# Augmenting Passage Representations with Query Generation for Enhanced Cross-Lingual Dense Retrieval

Shengyao Zhuang
The University of Queensland
Brisbane, QLD, Australia
s.zhuang@uq.edu.au

Linjun Shou
Microsoft STCA
Beijing, China
lisho@microsoft.com

Guido Zuccon
The University of Queensland
Brisbane, QLD, Australia
g.zuccon@uq.edu.au

## ABSTRACT

Effective cross-lingual dense retrieval methods that rely on multi-lingual pre-trained language models (PLMs) need to be trained to encompass both the relevance matching task and the cross-language alignment task. However, cross-lingual data for training is often scarcely available. In this paper, rather than using more cross-lingual data for training, we propose to use cross-lingual query generation to augment passage representations with queries in languages other than the original passage language. These augmented representations are used at inference time so that the representation can encode more information across the different target languages. Training of a cross-lingual query generator does not require additional training data to that used for the dense retriever. The query generator training is also effective because the pre-training task for the generator (T5 text-to-text training) is very similar to the fine-tuning task (generation of a query). The use of the generator does not increase query latency at inference and can be combined with any cross-lingual dense retrieval method. Results from experiments on a benchmark cross-lingual information retrieval dataset show that our approach can improve the effectiveness of existing cross-lingual dense retrieval methods. Implementation of our methods, along with all generated query files are made publicly available at https://github.com/ielab/xQG4xDR.

## CCS CONCEPTS

• **Information systems** → **Query representation**; **Language models**.

## KEYWORDS

Cross-lingual query generation; Cross-lingual retrieval; Dense retriever

## 1 INTRODUCTION

Pre-trained language model-based (PLM) dense retrievers (DRs) have achieved remarkable success in the task of English-only passage retrieval [12–14, 20, 21, 26, 28, 29, 40, 43, 44]. These models use a dual-encoder architecture that encodes both queries and passages with a PLM encoder into dense embeddings. They then perform approximate nearest neighbor (ANN) searching in the embedding space. Compared to traditional bag-of-words approaches, DRs benefit from semantic soft matching, which helps overcome the problem of word mismatch in passage retrieval [33, 45].

To leverage the semantic modelling power of DRs, recent research has extended English-only DRs to support cross-lingual settings [1, 2, 19, 22, 27, 31], i.e, where queries and passages are in different languages. This is achieved using multi-lingual PLMs, such as multilingual BERT [6], in place of the English-only PLMs. This approach is particularly important in this setting where traditional bag-of-words methods are ineffective due to the limited number of matching terms across languages. In contrast, cross-lingual DRs (xDRs) are able to encode queries and passages in different languages into a shared embedding space, enabling efficient ANN search across languages. However, such multi-lingual PLM-based xDRs usually are less effective on the cross-lingual passage retrieval task than DRs in the English-only setting [1]. The hypothesis to explain this result is that, in the English-only setting, a DR only needs to learn relevance matching between queries and passages. In contrast, a xDR not only needs to learn the relevance matching task, but also needs to learn how to align the embeddings of texts with similar semantic meaning but in different language [24, 41]. It is this language gap that makes cross lingual retrieval a relatively harder task for xDRs.

Based on this hypothesis, this paper proposes the use of cross-lingual query generation (xQG) to bridge the language gap for xDRs. Our approach is illustrated in Figure 1. Specifically, we fine-tune a multilingual T5 (mT5) model using language-specific prompts to generate several queries per passage for each target language. In the indexing phase, we use a given xDR model to encode passages and their associated generated queries in different languages into embeddings. Finally, we augment the original passage embeddings with the embeddings of the generated queries before adding them to the ANN index. By doing so, we move the passage embeddings to a space that is closer to the target user query language. Our approach does not add extra query latency and can be applied to any existing xDRs.

## 2 RELATED WORKS

***Cross-lingual dense retrievers.*** The development of cross-lingual pre-trained language models has significantly contributed to the
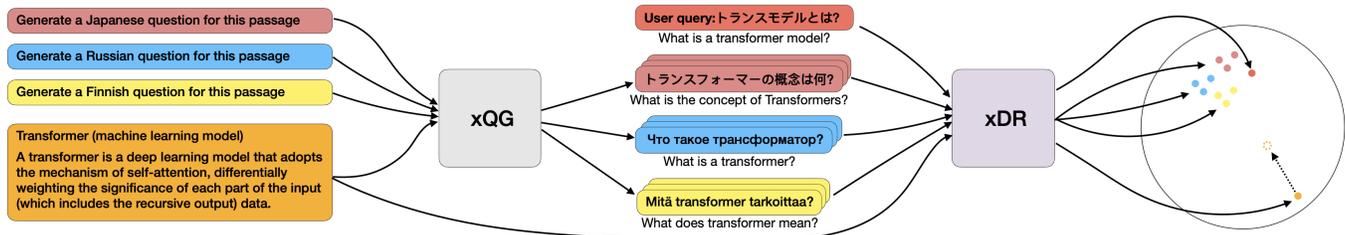
**Figure 1: Augmenting passage representations with cross-lingual generated query embeddings. The query examples shown in this figure were generated using our trained xQG model. For each query, we report the corresponding translation obtained using Google's translation service.**

progress of cross-lingual dense retrieval (xDR) [4, 5, 8, 39]. Notable examples of recent xDR methods include CORA [2], which employs a generator to facilitate retrieval training data mining, Sentri [31], which proposes a single encoder and self-training, and DR.DECR [19], which utilizes parallel queries and sentences for cross-lingual knowledge distillation. Among these, our work is most closely related to QuiCK [27], which also utilizes a cross-lingual query generator for xDR. However, unlike QuiCK, which uses the xQG as a teacher model in knowledge distillation for xDR in the training phase, our method directly augments passage embeddings with xQG queries without any xDR training involved.

***Query generation for information retrieval.*** Query generation is a well-established technique that has been widely used to improve retrieval performance in various retrieval models [11, 25]. In addition, it has been shown to be effective for domain adaptation in dense passage retrieval tasks [23, 32, 34], as well as for enhancing the effectiveness of other PLM-based rankers [37, 48]. In our approach, we also rely on a query generation model to generate high-quality queries for downstream passage embedding augmentation tasks.

***Augmenting embeddings for dense retrievers.*** Our method relies on effectively augmenting representations encoded by DR encoders – a direction recently explored also by other works. For instance, Li et al. [16, 17] use the top retrieved passages as pseudo-relevant feedback to augment query embeddings using the Rocchio aggregation function. Similarly, Zhuang et al. [47] extend this idea by using embeddings of clicked passages to augment query embeddings. Other PRF methods have also been extensively researched to enhance both English-only dense retrieval [18, 35, 36, 42, 46] and cross-lingual dense retrieval [3]. On the other hand, the HyDE method uses large pre-trained language models to generate hypothetical passages for a given query and directly uses the embedding of the hypothetical passage to perform the search [9]. While all these works focus on augmenting query embeddings for DRs at query time, our work focuses on augmenting passage embeddings *at indexing time*, thereby avoiding the extra overhead in terms of query latency.

## 3 METHOD

Our approach consists of two components: (1) a cross-lingual query generation model that, for each passage, generates high-quality queries in different languages, and (2) an embedding aggregation function that augments the original passage embedding with the embeddings of the generated queries.

To obtain a xQG model, we fine-tune a mT5 model with labeled relevant $(q_t, p)$ pairs, where $t$ is the target query language. We use a seq2seq objective akin to docTquery-T5's objective [25]; in our

case the input is the passage text with language-specific prompts:

$$prompt\,(t, p) = Generate\ a\ [t]\ question\ for\ this\ passage:\ [p], \quad (1)$$

where $[t]$ and $[p]$ are prompt placeholders for the target language and passage text, respectively. Once we have trained a xQG model, we can generate a set of queries for each target language and passage by replacing the placeholders accordingly:

$$Q_p = \bigcup_{t \in T} Q_p^t;\ Q_p^t = xQG\,(prompt\,(t, p)) \times n, \quad (2)$$

where $T$ is the target language set and $Q_p$ is the set of all generated queries for passage $p$ which includes $n$ generated queries for each target language. In our experiments, we use a top-$k$ sampling scheme [7] to generate queries, and set $k = 10$. We find that these simple language-specific prompts can effectively lead the T5 model to generate the desired output in that language.

For the embedding aggregation function, we use a Rocchio-like aggregation function that is similar to previous works that aggregated dense representations [16, 47]. We use this aggregation function to obtain the embeddings of all passages in the corpus:

$$emb(p, Q_p, \theta, \alpha) = (1 - \alpha) \cdot \theta(p) + \alpha \sum_{\hat{q}_t \in Q_p} \theta(\hat{q}_t), \quad (3)$$

where $\hat{q}_t$ is a generated query for target language $t$, $\theta$ is the xDR encoder model and $\theta(.)$ is the text embedding given by the xDR encoder. The hyper-parameter $\alpha$ is the augmentation ratio, $\alpha \in [0, 1]$. This ratio is used to control the weights assigned to the original passage embedding and the generated query embeddings.

## 4 EXPERIMENTAL SETTINGS

We design our experiments to answer the research questions:

**RQ1:** How does the number of generated queries affect xDR effectiveness?

**RQ2:** How does the augmentation ratio $\alpha$ impact xDR effectiveness?

**RQ3:** How do the queries generated for each language impact xDR effectiveness?

***Datasets and Evaluation.*** We train and evaluate our approach on XOR-TyDi [1], a cross-lingual open retrieval question answering benchmark dataset. The dataset contains approximately 15k annotated relevant passage-query pairs in the training set and 2k passage-answer pairs in the dev set. Queries in both the train and dev sets are in seven typologically diverse languages (Ar, Bn, Fi, Ja, Ko, Ru, and Te), while passages are in English. There are about 18M passages in the corpus. We use the training set to train both the

**Table 1: xDR models trained with different backbone PLMs. Zero-shot means trained with only the English subset of the NQ training data (thus, zero-shot for the cross-lingual task). Statistical differences against the base model (without xQG) are labelled with $\star$ ($p < 0.05$).**

(a) R@2kt

| Model | Ar | Bn | Fi | Ja | Ko | Ru | Te | Average |
|---|---|---|---|---|---|---|---|---|
| XLM-R | 28.5 | 26.3 | 29.3 | 22.4 | 34.0 | 17.3 | 34.9 | 27.5 |
| XLM-R + xQG | **30.1** | **29.6** | **31.2** | **23.7** | **36.1** | **19.4** | **38.2** | **29.8**$^\star$ |
| mBERT | 41.1 | 49.0 | 52.2 | **37.3**$^\star$ | 48.1 | 33.3 | 47.9 | 44.1 |
| mBERT + xQG | **42.4** | **54.9**$^\star$ | **54.1** | 33.6 | **52.3**$^\star$ | **33.8** | **52.5**$^\star$ | **46.2**$^\star$ |
| mBERT (zero-shot) | 32.6 | 25.0 | 38.2 | 29.5 | 38.9 | 27.0 | 39.9 | 33.0 |
| mBERT (zero-shot) + XQG | **33.9** | **28.9**$^\star$ | **43.0**$^\star$ | **32.8**$^\star$ | **41.4** | **29.5** | **42.0** | **36.0**$^\star$ |

(b) R@5kt

| Model | Ar | Bn | Fi | Ja | Ko | Ru | Te | Average |
|---|---|---|---|---|---|---|---|---|
| XLM-R | 38.5 | 33.6 | 37.9 | **32.8** | 42.8 | 28.3 | 47.5 | 37.3 |
| XLM-R + xQG | **38.8** | **41.1**$^\star$ | **39.8** | **32.8** | **43.2** | **30.8** | **49.6** | **39.4**$^\star$ |
| mBERT | 49.2 | 57.6 | **58.6** | 42.7 | 57.5 | **41.4** | 55.9 | 51.8 |
| mBERT + xQG | **51.5** | **60.9**$^\star$ | 58.3 | **43.6** | **58.6** | **41.4** | **60.1**$^\star$ | **53.5**$^\star$ |
| mBERT (zero-shot) | 38.5 | 36.5 | 47.5 | 38.2 | 48.1 | 35.0 | 48.7 | 41.8 |
| mBERT (zero-shot) + xQG | **43.7**$^\star$ | **40.8**$^\star$ | **50.0**$^\star$ | **40.2** | **49.8** | **39.7**$^\star$ | **51.7** | **45.1**$^\star$ |

xQG model and the xDRs. We evaluate the effectiveness of our approach and baselines using recall at $m$ kilo-tokens (R@$m$kt), which is the dataset's official evaluation metric. This metric computes the fraction of queries for which the minimal answer is contained in the top $m$ tokens of the retrieved passages. We consider $m = 2k, 5k$ (R@2kt and R@5kt), as per common practice for this dataset. Statistical significant differences are reported with respect to a two-tailed t-test with Bonferroni correction.

**Baselines.** Following common practice on XOR-TyDi [1], we adapt DPR [14] to the cross-lingual setting by initializing it with different multilingual pre-trained language models (PLMs). Specifically, in our experiments we use the multilingual variants of BERT [6] (mBERT) and XLM-RoBERTa [5] (XLM-R) for this purpose. In addition to the standard supervised baselines, we also explore how our xQG embedding augmentation approach can improve a zero-shot xDR model, where the xDR is initialized with mBERT but it is trained with English only passage-query pairs and is directly applied to the XOR-TyDi cross-lingual retrieval task.

**Implementation details.** We initialize our xQG model with the multilingual T5-base checkpoint provided by Google[1] and available in the Huggingface library [38]. We fine-tune this checkpoint with the passage-query pairs in the XOR-TyDi training set. We train the xQG model for 1600 steps using a batch size of 128 and a learning rate of 1e-4. After training, we generate 5 queries per passage for each target language, resulting in about 7 * 5 * 18M = 630M generated queries in total. We use four A100 GPUs to generate all queries; this process took about 70 hours to complete. We release our generated queries on the Huggingface hub [2] to allow the community to reuse this resource [30]. For training the xDRs, we use the Tevatron dense retriever training toolkit [10], which uses with BM25 hard negative passages. We train the xDRs with a batch size of 128, initializing them with mBERT base[3] or XLM-R base[4] checkpoints and training them on the XOR-TyDi training set for 40 epochs. For each training sample, we set the number of hard negative passages in the contrastive loss to 7 and applied in-batch negatives training. We use a learning rate of 1e-5

---

[1]https://huggingface.co/google/mt5-base
[2]https://huggingface.co/datasets/ielab/xor-tydi-xqg-augmented
[3]https://huggingface.co/bert-base-multilingual-cased
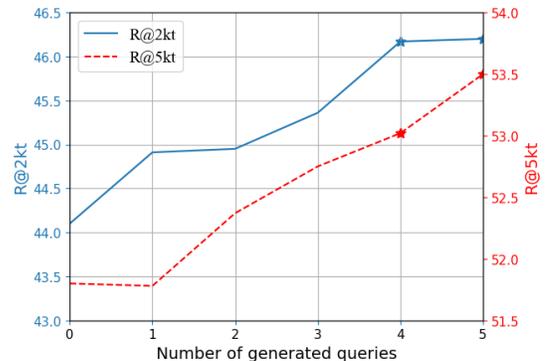[4]https://huggingface.co/xlm-roberta-base



**Figure 2: Impact of the amount of generated queries for target language. Scores are averaged across all languages. Statistical significant improvements over no augmentation (number of queries $n = 0$) are labelled with stars ($p < 0.05$).**

for mBERT-based xDRs and of 2e-5 for XLM-R-based xDRs. For the zero-shot xDR, we use the same training configurations as for the mBERT-based xDR trained on XOR-TyDi but using the Natural Questions (NQ) training data [15] which contains English-only query-passage training samples.

## 5  RESULTS

### 5.1  Main results

Table 1 presents the effectiveness of xDR models initialized with the XLM-R and mBERT backbone PLMs and trained on the XOR-TyDi dataset. Zero-shot denotes the models trained only on the English subset of the NQ dataset. In these experiments, we use all the queries generated by our xQG and set the augmentation ratio to $\alpha = 0.01$ for augmenting the passage embeddings of the xDRs.

For the R@2kt metric, the xDR initialized with mBERT outperforms the xDR initialized with XLM-R, achieving an average R@2kt score of 44.1, while XLM-R achieves an average score of 27.5. Our xQG passage embedding augmentation approach improves the XLM-R xDR, achieving an average score of 29.8, which is a statistically significant improvement compared to its baseline effectiveness ($p < 0.05$). Similarly, mBERT's effectiveness improves with xQG, achieving an average score of 46.2, which is also a statistically
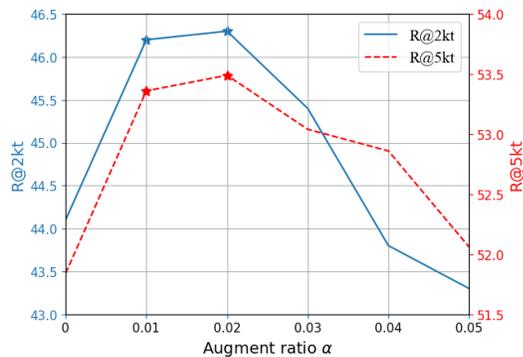
**Figure 3: Impact of the augmentation ratio $\alpha$. Scores are averaged across all languages. Statistical significant improvements over no augmentation ($\alpha = 0$) are labelled with stars ($p < 0.05$).**

significant improvement compared to its corresponding baseline ($p < 0.05$). The zero-shot mBERT model achieves an average R@2kt of 33.0; this also improves when combined with xQG, achieving an average score of 36.0. This improvement is statistically significant ($p < 0.05$). Similar trends are found for R@5kt. Overall, we find that our xQG can significantly improve all investigated xDR models. In terms of per language effectiveness, xQG improves almost all models across all languages with the exceptions of mBERT's R@2kt for Japanese (Ja) and mBERT's R@5kt for Finnish (Fi).

In summary, mBERT performs better than XLM-R for both R@2kt and R@5kt. The use of xQG embedding augmentation statistically significantly improves the effectiveness of both backbones.

## 5.2 RQ1: Impact of number of generated queries

Figure 2 reports the impact of using different amounts of generated queries to augment passage embeddings when using mBERT xDR. The results suggest that using more generated queries is beneficial for both R@2tk and R@5tk. The improvements become statistically significant when 4 or more generated queries are used for each target language. While the curves do not plateau, indicating that using even more generated queries could further improve the effectiveness, our experiments were limited to up to 5 generated queries per target language due to computational constraints.

## 5.3 RQ2: Impact of augmentation ratio

We report the impact of the augmentation ratio $\alpha$ on the effectiveness of xDR in Figure 3. Higher values of $\alpha$ correspond to assigning more weight to the generated query embeddings during embedding aggregation, and $\alpha = 0$ corresponds to no augmentation. As shown in the figure, even a small value of $\alpha$ (0.01) leads to a significant improvement in both R@2kt and R@5kt. The best effectiveness is achieved when $\alpha = 0.02$. However, using higher values of $\alpha$ does not result in further improvements – rather, it can even hurt the effectiveness when $\alpha > 0.05$. Based on these results, we conclude that the augmentation ratio in our embedding aggregation function has a significant impact on the effectiveness of xDR, and using a small value of $\alpha$ can be beneficial for improving effectiveness.
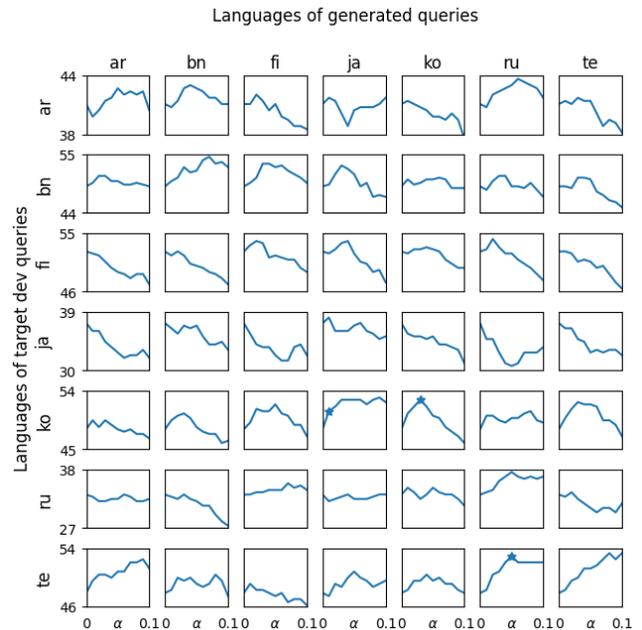


**Figure 4: R@2kt for target languages (rows) augmented by source generated queries (columns). We plot $\alpha$ from 0 to 0.1 with step size of 0.01 (x axis). Statistically significant better results with respect to no augmentation ($\alpha = 0$) are labelled with stars ($p < 0.05$).**

## 5.4 RQ3: Impact of each languages

In the previous experiments we used all the queries generated for a passage to augment the original passage embedding, irrespective of language of the generated query. Next, we investigate the impact on effectiveness of using generated queries from each of the languages separately. We analyze this in Figure 4, where we plot R@2kt for each target language (rows) against different values of $\alpha$ (x-axis).

The plots in the diagonal show that xDR effectiveness improves when passage embeddings are augmented with generated queries for the same language. Notably, we observe that the best value of $\alpha$ varies across different languages, and the improvements in effectiveness are not always statistically significant. We also observe that, for some languages, queries generated for other languages can improve the effectiveness of target queries in a different language. For instance, using queries generated for Japanese (Ja) can improve the effectiveness of target queries in Korean (ko), while using Russian (Ru) generated queries can help target queries in Telugu (Te). These results suggest that the embeddings of the queries generated for any single language potentially can also provide useful information for target queries in other languages.

## 6 CONCLUSION AND FUTURE WORK

In this paper we propose a passage embedding augmentation approach to enhance the effectiveness of cross-lingual DRs. Our method can be applied to any cross-lingual DR and it requires no further DR training nor changes to the retrieval pipeline. We empirically showed the method is effective for the cross-lingual DPR method across different backbones. However, a limitation of our empirical investigation is that we did not evaluate the method across other dense retriever architectures. We leave the extension of this

investigation to future work. We also note that our approach relies on a xQG model that can generate high quality queries. However, a recent work has shown that T5-based query generation is prone to hallucination and that along with highly effective queries, the generator also produces poor queries that negatively impact retrieval effectiveness [11]. They then propose the use of a cross-encoder ranker to filter out some ineffective generated queries; this practice can further improve effectiveness. We leave the adaptation of this approach in our xQG setting to future work.

## REFERENCES

[1] Akari Asai, Jungo Kasai, Jonathan H. Clark, Kenton Lee, Eunsol Choi, and Hannaneh Hajishirzi. 2021. XOR QA: Cross-lingual Open-Retrieval Question Answering. In *NAACL-HLT*.

[2] Akari Asai, Xinyan Yu, Jungo Kasai, and Hanna Hajishirzi. 2021. One question answering model for many languages with cross-lingual dense passage retrieval. *Advances in Neural Information Processing Systems* 34 (2021), 7547–7560.

[3] Ramraj Chandradevan, Eugene Yang, Mahsa Yarmohammadi, and Eugene Agichtein. 2022. Learning to Enrich Query Representation with Pseudo-Relevance Feedback for Cross-lingual Retrieval. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1790–1795.

[4] Zewen Chi, Li Dong, Furu Wei, Nan Yang, Saksham Singhal, Wenhui Wang, Xia Song, Xian-Ling Mao, He-Yan Huang, and Ming Zhou. 2021. InfoXLM: An Information-Theoretic Framework for Cross-Lingual Language Model Pre-Training. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 3576–3588.

[5] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Édouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised Cross-lingual Representation Learning at Scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 8440–8451.

[6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, 4171–4186. https://doi.org/10.18653/v1/n19-1423

[7] Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical Neural Story Generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 889–898.

[8] Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT Sentence Embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 878–891.

[9] Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2022. Precise Zero-Shot Dense Retrieval without Relevance Labels. *arXiv preprint arXiv:2212.10496* (2022).

[10] Luyu Gao, Xueguang Ma, Jimmy J. Lin, and Jamie Callan. 2022. Tevatron: An Efficient and Flexible Toolkit for Dense Retrieval. *ArXiv* abs/2203.05765 (2022).

[11] Mitko Gospodinov, Sean MacAvaney, and Craig Macdonald. 2023. Doc2Query--: When Less is More. *arXiv preprint arXiv:2301.03266* (2023).

[12] Sebastian Hofstätter, Sophia Althammer, Michael Schröder, Mete Sertkan, and Allan Hanbury. 2020. Improving efficient neural ranking models with cross-architecture knowledge distillation. *arXiv preprint arXiv:2010.02666* (2020).

[13] Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021. Efficiently teaching an effective dense retriever with balanced topic aware sampling. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 113–122.

[14] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 6769–6781. https://doi.org/10.18653/v1/2020.emnlp-main.550

[15] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural Questions: a Benchmark for Question Answering Research. *Transactions of the Association of Computational Linguistics* (2019).

[16] Hang Li, Ahmed Mourad, Shengyao Zhuang, Bevan Koopman, and Guido Zuccon. 2023. Pseudo relevance feedback with deep language models and dense retrievers: Successes and pitfalls. *ACM Transactions on Information Systems* (2023).

[17] Hang Li, Shengyao Zhuang, Xueguang Ma, Jimmy Lin, and Guido Zuccon. 2023. Pseudo-Relevance Feedback with Dense Retrievers in Pyserini. In *Proceedings of the 26th Australasian Document Computing Symposium (ADCS '22)*.

[18] Hang Li, Shengyao Zhuang, Ahmed Mourad, Xueguang Ma, Jimmy Lin, and Guido Zuccon. 2022. Improving query representations for dense retrieval with pseudo relevance feedback: A reproducibility study. In *Advances in Information Retrieval: 44th European Conference on IR Research, ECIR 2022, Stavanger, Norway, April 10–14, 2022, Proceedings, Part I*. Springer, 599–612.

[19] Yulong Li, Martin Franz, Md Arafat Sultan, Bhavani Iyer, Young-Suk Lee, and Avirup Sil. 2022. Learning Cross-Lingual IR from an English Retriever. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 4428–4436.

[20] Sheng-Chieh Lin, Jheng-Hong Yang, and Jimmy Lin. 2020. Distilling dense representations for ranking using tightly-coupled teachers. *arXiv preprint arXiv:2010.11386* (2020).

[21] Sheng-Chieh Lin, Jheng-Hong Yang, and Jimmy Lin. 2021. In-batch negatives for knowledge distillation with tightly-coupled teachers for dense retrieval. In *Proceedings of the 6th Workshop on Representation Learning for NLP (RepL4NLP-2021)*. 163–173.

[22] Shayne Longpre, Yi Lu, and Joachim Daiber. 2021. MKQA: A linguistically diverse benchmark for multilingual open domain question answering. *Transactions of the Association for Computational Linguistics* 9 (2021), 1389–1406.

[23] Ji Ma, Ivan Korotkov, Yinfei Yang, Keith Hall, and Ryan McDonald. 2021. Zero-shot Neural Passage Retrieval via Domain-targeted Synthetic Question Generation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. 1075–1088.

[24] Suraj Nair, Eugene Yang, Dawn Lawrie, Kevin Duh, Paul McNamee, Kenton Murray, James Mayfield, and Douglas W Oard. 2022. Transfer learning approaches for building cross-language dense retrieval models. In *Advances in Information Retrieval: 44th European Conference on IR Research, ECIR 2022, Stavanger, Norway, April 10–14, 2022, Proceedings, Part I*. Springer, 382–396.

[25] Rodrigo Nogueira and Jimmy Lin. 2019. From doc2query to docTTTTTquery.

[26] Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng Wang. 2021. RocketQA: An optimized training approach to dense passage retrieval for open-domain question answering. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 5835–5847.

[27] Houxing Ren, Linjun Shou, Ning Wu, Ming Gong, and Daxin Jiang. 2022. Empowering Dual-Encoder with Query Generator for Cross-Lingual Dense Retrieval. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. 3107–3121.

[28] Ruiyang Ren, Shangwen Lv, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021. PAIR: Leveraging Passage-Centric Similarity Relation for Improving Dense Passage Retrieval. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. 2173–2183.

[29] Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021. RocketQAv2: A Joint Training Method for Dense Passage Retrieval and Passage Re-ranking. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. 2825–2835.

[30] Harrisen Scells, Shengyao Zhuang, and Guido Zuccon. 2022. Reduce, Reuse, Recycle: Green Information Retrieval Research. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2825–2837.

[31] Nikita Sorokin, Dmitry Abulkhanov, Irina Piontkovskaya, and Valentin Malykh. 2022. Ask me anything in your native language. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 395–406.

[32] Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. BEIR: A Heterogeneous Benchmark for Zero-shot Evaluation of Information Retrieval Models. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*. https://openreview.net/forum?id=wCu6T5xFjeJ

[33] Nicola Tonellotto. 2022. Lecture Notes on Neural Information Retrieval. *arXiv preprint arXiv:2207.13443* (2022).

[34] Kexin Wang, Nandan Thakur, Nils Reimers, and Iryna Gurevych. 2022. GPL: Generative Pseudo Labeling for Unsupervised Domain Adaptation of Dense Retrieval. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2345–2360.

[35] Xiao Wang, Craig Macdonald, Nicola Tonellotto, and Iadh Ounis. 2021. Pseudo-relevance feedback for multiple representation dense retrieval. In *Proceedings of the 2021 ACM SIGIR International Conference on Theory of Information Retrieval*. 297–306.

[36] Xiao Wang, Craig Macdonald, Nicola Tonellotto, and Iadh Ounis. 2023. ColBERT-PRF: Semantic pseudo-relevance feedback for dense passage and document retrieval. *ACM Transactions on the Web* 17, 1 (2023), 1–39.

[37] Yujing Wang, Yingyan Hou, Haonan Wang, Ziming Miao, Shibin Wu, Qi Chen, Yuqing Xia, Chengmin Chi, Guoshuai Zhao, Zheng Liu, et al. [n.d.]. A Neural Corpus Indexer for Document Retrieval. In *Advances in Neural Information Processing Systems*.

[38] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, 38–45.

[39] Ning Wu, Yaobo Liang, Houxing Ren, Linjun Shou, Nan Duan, Ming Gong, and Daxin Jiang. 2022. Unsupervised context aware sentence representation pretraining for multi-lingual dense retrieval. *arXiv preprint arXiv:2206.03281* (2022).

[40] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N Bennett, Junaid Ahmed, and Arnold Overwijk. 2020. Approximate Nearest Neighbor Negative Contrastive Learning for Dense Text Retrieval. In *International Conference on Learning Representations*.

[41] Eugene Yang, Suraj Nair, Ramraj Chandradevan, Rebecca Iglesias-Flores, and Douglas W Oard. 2022. C3: Continued pretraining with contrastive weak supervision for cross language ad-hoc retrieval. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2507–2512.

[42] HongChien Yu, Chenyan Xiong, and Jamie Callan. 2021. Improving Query Representations for Dense Retrieval with Pseudo Relevance Feedback. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 3592–3596.

[43] Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2021. Optimizing Dense Retrieval Model Training with Hard Negatives. *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval* (2021).

[44] Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Min Zhang, and Shaoping Ma. 2020. RepBERT: Contextualized text embeddings for first-stage retrieval. *arXiv preprint arXiv:2006.15498* (2020).

[45] Wayne Xin Zhao, Jing Liu, Ruiyang Ren, and Ji-Rong Wen. 2022. Dense text retrieval based on pretrained language models: A survey. *arXiv preprint arXiv:2211.14876* (2022).

[46] Yunchang Zhu, Liang Pang, Yanyan Lan, Huawei Shen, and Xueqi Cheng. 2022. LoL: A Comparative Regularization Loss over Query Reformulation Losses for Pseudo-Relevance Feedback. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 825–836.

[47] Shengyao Zhuang, Hang Li, and Guido Zuccon. 2022. Implicit feedback for dense passage retrieval: A counterfactual approach. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 18–28.

[48] Shengyao Zhuang, Houxing Ren, Linjun Shou, Jian Pei, Ming Gong, Guido Zuccon, and Daxin Jiang. 2022. Bridging the gap between indexing and retrieval for differentiable search index with query generation. *arXiv preprint arXiv:2206.10128* (2022).