

Learning Discrete Representations via Constrained Clustering for Effective and Efficient Dense Retrieval

Jingtao Zhan¹, Jiaxin Mao², Yiqun Liu^{1*}, Jiafeng Guo³, Min Zhang¹, Shaoping Ma¹

¹ Department of Computer Science and Technology, Institute for Artificial Intelligence, Beijing National Research Center for Information Science and Technology, Tsinghua University, Beijing 100084, China

² Beijing Key Laboratory of Big Data Management and Analysis Methods, Gaoling School of Artificial Intelligence, Renmin University of China, Beijing 100872, China

³ CAS Key Lab of Network Data Science and Technology, Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China

jingtaozhan@gmail.com, maojiaxin@gmail.com, yiqunliu@tsinghua.edu.cn, guojiafeng@ict.ac.cn

ABSTRACT

Dense Retrieval (DR) has achieved state-of-the-art first-stage ranking effectiveness. However, the efficiency of most existing DR models is limited by the large memory cost of storing dense vectors and the time-consuming nearest neighbor search (NNS) in vector space. Therefore, we present RepCONC, a novel retrieval model that learns discrete Representations via CONstrained Clustering. RepCONC jointly trains dual-encoders and the Product Quantization (PQ) method to learn discrete document representations and enables fast approximate NNS with compact indexes. It models quantization as a constrained clustering process, which requires the document embeddings to be uniformly clustered around the quantization centroids and supports end-to-end optimization of the quantization method and dual-encoders. We theoretically demonstrate the importance of the uniform clustering constraint in RepCONC and derive an efficient approximate solution for constrained clustering by reducing it to an instance of the optimal transport problem. Besides constrained clustering, RepCONC further adopts a vector-based inverted file system (IVF) to support highly efficient vector search on CPUs. Extensive experiments on two popular ad-hoc retrieval benchmarks show that RepCONC achieves better ranking effectiveness than competitive vector quantization baselines under different compression ratio settings. It also substantially outperforms a wide range of existing retrieval models in terms of retrieval effectiveness, memory efficiency, and time efficiency.

CCS CONCEPTS

• Information systems → Search index compression; Retrieval models and ranking; Information retrieval.

KEYWORDS

index compression, dense retrieval, neural ranking

ACM Reference Format:

Jingtao Zhan¹, Jiaxin Mao², Yiqun Liu^{1*}, Jiafeng Guo³, Min Zhang¹, Shaoping Ma¹. 2021. Learning Discrete Representations via Constrained Clustering for Effective and Efficient Dense Retrieval. In *Proceedings of The 15th ACM International Conference on Web Search and Data Mining (WSDM'22)*. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

Dense Retrieval (DR) has become a popular paradigm for first-stage retrieval in ad-hoc retrieval tasks. Through embedding queries and documents in a latent vector space with dual-encoders and using nearest neighbor search to retrieve relevant documents, the DR paradigm avoids the vocabulary mismatch problem, which has been a great challenge for traditional bag-of-words (BoW) models [30]. With end-to-end supervised training, recent works have achieved state-of-the-art ranking performance and significantly outperforms BoW models [22, 29, 33, 37].

Despite the success in improving ranking performance, most existing DR models [19, 29, 33, 36] are inefficient in memory usage and computational time. For memory inefficiency, the size of the embedding index is usually an order of magnitude larger than that of BoW index [35]. At runtime, the vectors must be loaded to system memory or even GPU memory, which is both costly and highly limited in size. As for time inefficiency, many existing DR models [19, 29, 33, 36] do not use approximate vector search [17, 25]. They have to conduct exhaustive search, i.e., computing relevance scores between the submitted query and all documents, which is less time-efficient than BoW models with inverted indexes. As a result, these DR models cannot use CPUs for retrieval due to high latency and have to use much more expensive GPUs to accelerate the search [33, 35, 36].

A key solution for the efficiency issue of DR models is to learn discrete representations for document embeddings, which can be encoded into compact indexes and enable efficient vector search. Popular methods for learning this kind of discrete representation include Product Quantization (PQ) [11, 17] and Locality Sensitive Hashing (LSH) [16]. However, these methods usually learn discrete representations in an unsupervised way and cannot benefit from

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

WSDM'22, February 21-25, 2022, Phoenix, Arizona

© 2021 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

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

*Corresponding author

supervised signals. Directly adopting these techniques usually hurts ranking effectiveness [35, 38].

Therefore, jointly optimizing dual-encoders and the quantization methods with supervised labels is regarded as a promising direction in improving retrieval effectiveness. However, it is inherently challenging because the quantization operation is non-differentiable and the model cannot be trained in an end-to-end fashion. There exist a number of recent works (e.g., JPQ [35]) trying to solve this problem, but they usually suffer from significant performance loss while improving efficiency. Therefore, we believe it is still an unsolved but essential problem.

To tackle this problem, we present RepCONC, which stands for learning discrete **R**epresentations via **C**ONstrained Clustering¹. It jointly trains the dual-encoders and PQ by modeling quantization as a *constrained clustering* process. Specifically, *constrained clustering* involves a clustering loss and a uniform clustering constraint. The clustering loss is introduced to train the discrete codes with the requirement that document embeddings are clustered around the quantization centroids. We also employ a uniform clustering constraint, which requires the vectors to be equally assigned to all quantization centroids. We add the constraint because we find that unconstrained clustering tends to assign vectors to a few major clusters and makes the quantized vectors indistinguishable with each other. Since this constraint leads to a difficult combinatorial optimization problem, we derive an approximate solution by relaxing it to an instance of the optimal transport problem. Besides the two components of *constrained clustering*, RepCONC further employs vector-based inverted file system (IVF) [17], which enables efficient non-exhaustive vector search. With these designs, RepCONC can run on either GPU or CPU² and perform vector search in an efficient way.

We conduct experiments on two widely-adopted ad-hoc retrieval benchmarks [1, 5] and compare RepCONC with a wide range of baselines, including both vector compression methods and retrieval models. Experimental results show that: 1) RepCONC significantly outperforms competitive vector compression baselines with different compression ratio settings from tens of times to hundreds of times. 2) RepCONC substantially outperforms various retrieval baselines in terms of retrieval effectiveness, memory efficiency, and time efficiency. 3) The ablation study demonstrates that *constrained clustering* is the key to the effectiveness of RepCONC.

2 RELATED WORKS

DR represents queries and documents with embeddings and utilizes vector search to retrieve relevant documents. Most existing DR models [19, 29, 33, 36, 37] share the same BERT-base [9, 23] architecture and utilize brute-force vector search. They differ in training methods, which can be classified into two categories. One line of research is negative sampling [15, 19, 33, 36, 37]. According to Zhan et al. [36], utilizing hard negatives helps improve top ranking performance. The other line is knowledge distillation [14, 22, 29], which adopts a cross-encoder to generate pseudo labels. This paper uses negative sampling to train RepCONC and leaves training RepCONC with knowledge distillation to future work.

Since the models mentioned above utilize brute-force vector search, they incur very large embedding indexes and have to use costly GPUs to accelerate the search. How to address the efficiency issue has recently attracted researchers' attention. Several studies propose some workarounds [34, 35, 38]. BPR [34] binarizes dense vectors and is conducted on the OpenQA task. An obvious limitation is that the compression ratio is fixed to 32x. DPQ [4, 38] utilizes PQ [17] for compression and is designed for word embedding compression and recommendation systems. Most recently, Zhan et al. [35] propose JPQ for document ranking and achieve state-of-the-art results. JPQ utilizes fixed discrete codes (Index Assignments) generated by K-Means and only trains the query encoder and PQ Centroid Embeddings. RepCONC is different from JPQ in both joint learning framework and efficiency design. Firstly, with the help of constrained clustering, RepCONC is able to optimize discrete codes (Index Assignments) while JPQ cannot. Secondly, RepCONC additionally employs the inverted file system (IVF) [17] to accelerate search and thus, can efficiently retrieve documents on CPUs while JPQ has to rely on GPUs.

3 CONSTRAINED CLUSTERING MODEL

In this section, we propose RepCONC, which stands for learning discrete **R**epresentations via **C**ONstrained Clustering. We firstly introduce the preliminary of Production Quantization [17], a widely-used vector compression method for approximate nearest neighbor search (ANNS). Then we elaborate our model.

3.1 Revisiting Product Quantization

RepCONC is based on Product Quantization (PQ) [17]. For vectors of dimension D , PQ defines M sets of embeddings, each of which includes K embeddings of dimension D/M . They are called PQ Centroid Embeddings. Formally, let $\mathbf{c}_{i,j}$ be the j -th centroid embedding from the i -th set:

$$\mathbf{c}_{i,j} \in \mathbb{R}^{\frac{D}{M}} \quad (1 \leq i \leq M, 1 \leq j \leq K) \quad (1)$$

Given a document embedding $\mathbf{d} \in \mathbb{R}^D$, PQ firstly splits it into M sub-vectors.

$$\mathbf{d} = \mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_M \quad (2)$$

Then PQ independently quantizes each sub-vector to the nearest PQ Centroid Embedding. Formally, to quantize a sub-vector \mathbf{d}_i , PQ selects $\mathbf{c}_{i,\varphi_i(\mathbf{d})}$ which achieves the minimum quantization error:

$$\varphi_i(\mathbf{d}) = \arg \min_j \|\mathbf{c}_{i,j} - \mathbf{d}_i\|^2 \quad (3)$$

Let $\boldsymbol{\varphi}(\mathbf{d})$ be the concatenation of $\varphi_i(\mathbf{d})$:

$$\boldsymbol{\varphi}(\mathbf{d}) = \varphi_1(\mathbf{d}), \varphi_2(\mathbf{d}), \dots, \varphi_M(\mathbf{d}) \in \{1, 2, \dots, K\}^M \quad (4)$$

where comma denotes vector concatenation. $\boldsymbol{\varphi}(\mathbf{d})$ is called the Index Assignment of \mathbf{d} . Along with the PQ Centroid Embeddings, $\boldsymbol{\varphi}(\mathbf{d})$ can reconstruct the quantized document embedding $\hat{\mathbf{d}}$ as follows:

$$\hat{\mathbf{d}} = \mathbf{c}_{1,\varphi_1(\mathbf{d})}, \mathbf{c}_{2,\varphi_2(\mathbf{d})}, \dots, \mathbf{c}_{M,\varphi_M(\mathbf{d})} \in \mathbb{R}^D \quad (5)$$

PQ improves both memory efficiency and time efficiency. For memory efficiency, PQ does not explicitly store \mathbf{d} or $\hat{\mathbf{d}}$. Instead, it only stores the PQ Centroid Embeddings $\{\mathbf{c}_{i,j}\}$ and Index Assignments $\boldsymbol{\varphi}(\mathbf{d})$. Since K is usually less than or equal to 256, $\boldsymbol{\varphi}(\mathbf{d})$ can be encoded with M bytes. Therefore, the compression ratio is about

¹Code and models are available at <https://github.com/jingtaozhan/RepCONC>.

²Except that the user queries are still encoded on GPU.

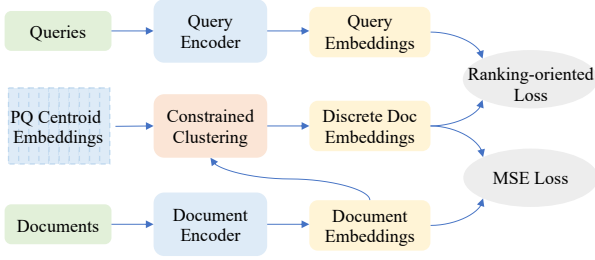


Figure 1: Training process of RepCONC.

$4D/M$. As for time efficiency, PQ enables efficient vector search. Given a query embedding, PQ splits it equally to M sub-vectors and pre-computes the similarities between the sub-vectors and PQ Centroid Embeddings. Then, PQ efficiently computes the similarities between the query embedding and each document embedding by aggregating the corresponding pre-computed similarities. Compared with directly computing vector similarity, the speedup ratio is about D/M .

3.2 Clustering and Representation Learning

Jointly optimizing dual-encoders and PQ parameters is challenging. RepCONC views it as a simultaneous clustering and representation learning problem. To solve it, RepCONC utilizes both the ranking-oriented loss [35] and a clustering loss for training. The ranking-oriented loss helps learn representations for ranking. The clustering loss, i.e., the mean square error (MSE) between the document embeddings and the quantization centroids, helps cluster document embeddings to centroid embeddings. We illustrate the training workflow in Figure 1.

The ranking-oriented loss [35] replaces the uncompressed document embeddings in the common DR ranking loss functions [29, 33, 36] with the quantized document embeddings. Therefore, it better evaluates the ranking performance with respect to the current compression parameters. Formally, ranking-oriented loss L_r is formulated as:

$$L_r = -\log \frac{e^{\langle q, \hat{d}^+ \rangle}}{e^{\langle q, \hat{d}^+ \rangle} + \sum_{d^-} e^{\langle q, \hat{d}^- \rangle}} \quad (6)$$

where d^+ and d^- are relevant and irrelevant documents, respectively. L_r facilitates effective representation learning by encouraging the relevant pairs to be scored higher than irrelevant pairs.

Although incorporating the ranking-oriented loss helps a recent joint learning work achieve state-of-the-art compression results [35], we argue that simply relying on this loss is problematic. Since the quantization error in Eq. (3) is not included in L_r , it may change arbitrarily and selecting Index Assignments based on Eq. (3) may lead to unexpected behaviors. Zhan et al. [35] avoids this pitfall by fixing the Index Assignments. However, fixed Index Assignments lead to sub-optimal ranking performance because they cannot benefit from supervised signals.

Different from Zhan et al. [35], RepCONC regards quantization as a clustering problem and introduces the MSE loss L_m :

$$L_m = \|d - \hat{d}\|^2 \quad (7)$$

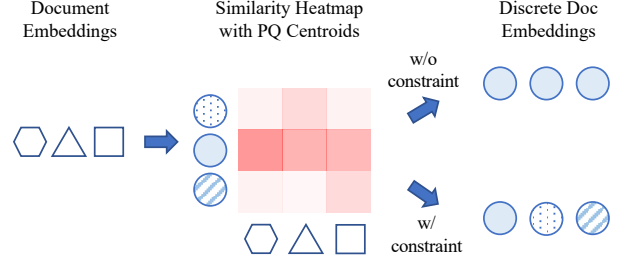


Figure 2: Illustration of Constrained Clustering. Darker colors in the heatmap indicate higher similarities (smaller distances). With the constraint, the discrete document embeddings are more diverse.

Minimizing L_m requires the vectors before and after quantization to be close to each other. In this way, the document embeddings are clustered around the centroid embeddings. Combining both L_r and L_m helps the model to cluster document embeddings based on ranking effectiveness. It is expected to produce better clustering compared with unsupervised training.

The final loss L is a weighted sum of ranking-oriented loss L_r and the MSE loss L_m .

$$L = L_r + \lambda L_m \quad (8)$$

λ is a hyper-parameter. If λ is too small, the documents are not clustered and the selected Index Assignments become arbitrary. If λ is too big, the MSE loss dominates the training process and harms ranking effectiveness. In practice, we find the model is relatively sensitive to λ , but becomes more robust and effective with the help of the uniform clustering constraint introduced in the next sections.

Since quantization involves some non-differentiable operations, we explicitly design the gradient back-propagation policy. The gradients of uncompressed document embeddings are defined as follows:

$$\frac{\partial L}{\partial d} := \frac{\partial L_r}{\partial d} + \lambda \frac{\partial L_m}{\partial d} \quad (9)$$

As the equation shows, we add the gradient of quantized document embeddings (the first term). The gradients are further back-propagated to dual-encoders. As for the PQ Centroid Embeddings, their gradients can be derived with chain rule. We formally show it as follows:

$$\begin{cases} \frac{\partial L}{\partial d} &= \frac{\partial L_r}{\partial d} + \lambda \frac{\partial L_m}{\partial d} \\ \frac{\partial L}{\partial c_{i,j}} &= 1_{\varphi_i(d)=j} \cdot \frac{\partial L}{\partial d} \end{cases} \quad (10)$$

In the following sections, we show how RepCONC selects Index Assignments (instead of Eq. (3)).

3.3 Importance of Uniform Clustering

It is non-trivial to simultaneously conduct clustering and representation learning because the two objects are conflicting to some extent. Although representation learning encourages vectors to be distinguishable, clustering encourages vectors to be identical. In practice, clustering tends to map vectors to several major clusters while some clusters are rarely used or even empty. The problem worsens with Eq. (10) where rarely used centroid embeddings are

less likely to be updated and may end up with arbitrary values. The unbalanced clustering distribution affects the distinguishability of the quantized vectors and compromises ranking effectiveness.

We tackle this challenge by imposing a uniform clustering constraint. It requires the document sub-vectors to be equally assigned to all PQ Centroid Embeddings. The learning object along with the constraint is formally expressed as:

$$\min L \quad \text{subject to } \forall i, j : P(\varphi_i(d) = j) = \frac{1}{K} \quad (11)$$

We illustrate constrained clustering in Figure 2. As the figure shows, the discrete document embeddings are selected by minimizing the quantization error (maximizing the similarity) given the uniform clustering constraint. Without the constraint, the discrete document embeddings become identical. In the following, we theoretically analyze the importance of uniform clustering.

We introduce several notations for our theoretical analysis. We define \mathcal{I} as all possible Index Assignments:

$$\mathcal{I} := \{\varphi(d)\} = \{1, 2, \dots, K\}^M \quad (12)$$

Let $\mathbf{I} \in \mathcal{I}$ be one Index Assignment and $I_i \in \{1, 2, \dots, K\}$ be the i_{th} value of \mathbf{I} .

Firstly, we show that maximizing the distinguishability of vectors is equivalent to forcing the vectors to be equally quantized to all possible Index Assignments. Let $\hat{\mathbf{d}}^s$ and $\hat{\mathbf{d}}^t$ be randomly sampled from all quantized document embeddings. We assume they are independent and identically distributed (i.i.d). The probability that they are equal satisfies:

$$\begin{aligned} P(\hat{\mathbf{d}}^s = \hat{\mathbf{d}}^t) &= P(\varphi(\mathbf{d}^s) = \varphi(\mathbf{d}^t)) \\ &= \sum_{\mathbf{I} \in \mathcal{I}} P(\varphi(\mathbf{d}^s) = \mathbf{I}) P(\varphi(\mathbf{d}^t) = \mathbf{I}) \\ &\stackrel{\text{i.i.d.}}{=} \sum_{\mathbf{I} \in \mathcal{I}} P^2(\varphi(d) = \mathbf{I}) \\ &\geq \frac{(\sum_{\mathbf{I} \in \mathcal{I}} P(\varphi(d) = \mathbf{I}))^2}{|\mathcal{I}|} \\ &= \frac{1}{|\mathcal{I}|} \end{aligned} \quad (13)$$

where AM-GM inequality is used. The equality is achieved if and only if

$$\forall \mathbf{I} : P(\varphi(d) = \mathbf{I}) = \frac{1}{|\mathcal{I}|} = \frac{1}{K^M} \quad (14)$$

That is to say, quantizing vectors equally to all possible Index Assignments helps representations to be distinguishable.

Next, we show uniformly clustering sub-vectors is the essential condition of Eq. (14). Given Eq. (14), the probability that a sub-vector is quantized to a centroid is a constant:

$$\forall i, j : P(\varphi_i(d) = j) = \sum_{\mathbf{I} : I_i = j} P(\varphi(d) = \mathbf{I}) = \frac{|\{\mathbf{I} : I_i = j\}|}{K^M} = \frac{1}{K} \quad (15)$$

We further show that if sub-vectors are independent, uniformly clustering sub-vectors (Eq. (15)) is the sufficient condition of Eq. (14):

$$P(\varphi(d) = \mathbf{I}) = \prod_{i=1}^M P(\varphi_i(d) = I_i) = \frac{1}{K^M} \quad (16)$$

Although independence among sub-vectors may not hold for practical dual-encoders, we believe constraining quantization with Eq. (15) is still helpful for distinguishing quantized vectors.

3.4 Constrained Clustering Optimization

This section shows how to incorporate the uniform clustering constraint during training. In previous works related to joint learning with PQ [4, 35, 38], the Index Assignments are selected based on Eq. (3). However, it cannot be applied to RepCONC because of the uniform clustering constraint. Next, we show how RepCONC incorporates the constraint to select Index Assignments during training.

We introduce a posterior distribution $q(j|\mathbf{d}_i)$, which is the probability that the sub-vector \mathbf{d}_i is quantized to the centroid $\mathbf{c}_{i,j}$. The Index Assignment, $\varphi_i(d)$, is the centroid with the maximum probability:

$$\varphi_i(d) = \arg \max_j q(j|\mathbf{d}_i) \quad (17)$$

For previous works [4, 35, 38] that use Eq. (3), $q(j|\mathbf{d}_i)$ can be regarded as being computed solely based on quantization error. Here for RepCONC, we compute $q(j|\mathbf{d}_i)$ by minimizing the quantization error given the uniform clustering constraint:

$$\begin{aligned} \forall i : \min_q \sum_{d \in \mathcal{D}} \sum_{j=1}^K q(j|\mathbf{d}_i) \|\mathbf{c}_{i,j} - \mathbf{d}_i\|^2 \quad \text{subject to} \\ \forall j, d : q(j|\mathbf{d}_i) \in \{0, 1\}, \sum_{j=1}^K q(j|\mathbf{d}_i) = 1, \text{ and } \sum_{d \in \mathcal{D}} q(j|\mathbf{d}_i) = \frac{|\mathcal{D}|}{K} \end{aligned} \quad (18)$$

where \mathcal{D} indicates the set of all documents. The first condition constrains $q(j|\mathbf{d}_i)$ to be binary, the second condition is a natural requirement for probability, and the third condition is exactly the uniform clustering constraint. Without the third condition, Eq. (17) and (18) degenerate to Eq. (3), i.e., selecting Index Assignments with minimum quantization error.

Solving Eq. (18) is particularly difficult because it is a combinatorial optimization problem with the scale of millions or even billions of documents. Therefore, we use an approximate solution by relaxing q to be continuous and focusing on uniformly clustering a mini-batch of documents \mathcal{B} :

$$\begin{aligned} \forall i : \min_q \sum_{d \in \mathcal{B}} \sum_{j=1}^K q(j|\mathbf{d}_i) \|\mathbf{c}_{i,j} - \mathbf{d}_i\|^2 \\ \text{subject to } \forall d : \sum_{j=1}^K q(j|\mathbf{d}_i) = 1 \text{ and } \forall j : \sum_{d \in \mathcal{B}} q(j|\mathbf{d}_i) = \frac{|\mathcal{B}|}{K} \end{aligned} \quad (19)$$

Since $\|\mathbf{c}_{i,j} - \mathbf{d}_i\|^2$ can be regarded as the cost of mapping \mathbf{d}_i to $\mathbf{c}_{i,j}$, this is an instance of the optimal transport problem and can be solved in polynomial time by linear program. In our implementation, we use Sinkhorn-Knopp algorithm [6] to efficiently solve Eq. (19).

3.5 Accelerating Search with IVF

Besides PQ, RepCONC employs the inverted file system (IVF) to accelerate vector search. After quantizing document embeddings, RepCONC uses k-means to generate n clusters. Each document

embedding belongs to the nearest cluster and is stored in the corresponding inverted list. Note that n is much smaller than the corpus size. Given a query embedding, RepCONC selects the nearest \tilde{n} clusters and only ranks the documents in them. The documents in other clusters are ignored. In this way, RepCONC approximately accelerates vector search by n/\tilde{n} .

Note that RepCONC does not include IVF in the joint learning framework and simply uses IVF after training. The clusters are generated in an unsupervised manner. In practice, we find this already yields satisfying results, and thus training IVF with supervised labels is not explored.

IVF only induces negligible memory overhead and does not harm memory efficiency. For example, on MS MARCO Passage Ranking dataset [1] which has 8 million passages, n is set to 5,000 and additional memory overhead is less than 3%.

3.6 Training/Inference Details

3.6.1 Warmup with OPQ. In order to accelerate convergence, we warmup the dual-encoders and PQ Centroid Embeddings as follows. We use the open-sourced STAR [36] model to initialize dual-encoders. STAR is trained without quantization. Given the document embeddings output by STAR, we use OPQ [11] to warmup PQ parameters, which is a popular unsupervised PQ variant.

3.6.2 Two-Stage Negative Sampling. Hard negative sampling is shown to be important for retrieval models [33, 36]. Following Zhan et al. [36], we train RepCONC in two stages. In the first stage, we retrieve static hard negatives using the initialized RepCONC. In the second stage, we use dynamic hard negatives, the top irrelevant documents retrieved at each training step. To enable end-to-end retrieval during training, we fix the Index Assignments and only train the query encoder and PQ Centroid Embeddings.

3.6.3 Efficient Encoding during Inference. During inference, we use Eq. (3) to quantize document embeddings instead of Eq. (18) or Eq. (19). In this way, we can quantize each document embedding online efficiently. Otherwise, computing with Eq. (18) is expensive and Eq. (19) introduces stochastic noise when batching documents.

4 EXPERIMENTAL SETUP

Here we present our experimental settings, including datasets, baselines, and implementation details.

4.1 Datasets and Metrics

We conduct experiments on two large-scale ad-hoc retrieval benchmarks from the TREC 2019 Deep Learning Track [1, 5], passage ranking and document ranking. They have been widely-adopted in previous works related to neural ranking. The passage ranking task has a corpus of 8.8M passages, 0.5M training queries, 7k development queries (henceforth, MARCO Passage), and 43 test queries (DL Passage). The document ranking task has a corpus of 3.2M documents, 0.4M training queries, 5k development queries (MARCO Doc), and 43 test queries (DL Doc). For both tasks, we report the official metrics and R@100 based on the full-corpus retrieval results.

4.2 Baselines

We exploit two types of baselines, vector compression methods and retrieval models.

4.2.1 Vector Compression Baselines.

Unsupervised methods include PQ [17], ScaNN [13], ITQ+LSH [12], OPQ [11], and OPQ+ScaNN. We use Faiss library [18] to implement those baselines except for ScaNN [13], which is implemented based on its open-sourced code.

Supervised methods include recently proposed DPQ [4, 38] and JPQ [35], both of which are also based on PQ. We re-implement DPQ since it is originally designed for word embedding compression [4] and recommendation systems [38]. We use the same warmup process as RepCONC. JPQ is lately proposed for document ranking and shares the same warmup process. Another compression method, BPR [34] binarizes dense vectors and thus is limited to a fixed compression ratio (32x). As RepCONC already achieves very small performance loss with a 64x compression ratio, we do not implement BPR for comparison.

4.2.2 Retrieval Models.

First-stage retrieval models involve BoW models and DR models. BoW models include BM25 [30] and its variants, such as DeepCT [7], HDCT [8], doc2query [28], and docT5query [27]. DR models include RepBERT [37], ANCE [33], STAR [36], and ADORE [36]. Their output embeddings are of dimension 768 and are not compressed. All of them utilize negative sampling methods for training as RepCONC. In our experiments related to time efficiency, we use IVF [17] with the same hyperparameters as RepCONC to accelerate ADORE [36], the most competitive uncompressed DR baseline.

Although several ranking models also conduct end-to-end retrieval, their latency is significantly higher than typical first-stage retrievers. Therefore, we classify them as complex end-to-end retrieval models. These models include ColBERT [20], COIL [10], uniCOIL [21], and DeepImpact [26]³. Note, for COIL [10], the authors uploaded a new model trained with hard negatives in the github repository, which is not included in its paper. We denote it as COIL-Hard.

4.3 Implementation Details

Here are our model settings. We build RepCONC based on huggingface transformers [32] and Faiss ANNS library [18]. The dual-encoders use RoBERTa-base [23] as the backbone, and the output embedding dimension is 768. Embedding similarity is computed with inner product. For PQ hyper-parameters, K is set to 256, and M is set to 4, 8, 12, 16, 24, 32, and 48 for different compression ratios. The compression ratio equals $4 \times 768/M$ since one vector is compressed to M bytes.

Training settings are as follows. Most training hyper-parameters are kept the same in both datasets except for batch size due to the limitation of GPU memory. Following ADORE [36], training is in two stages. In the first stage where static hard negatives are used, the optimizer is AdamW [24]; learning rates are 5×10^{-6} and 2×10^{-4} separately for encoders and centroid embeddings;

³Although uniCOIL [21] and DeepImpact [26] can leverage the inverted indexes like BM25 [30], they are much slower possibly due to much smaller vocabulary size (30k vs. 500k) and not removing stop words.

Table 1: Comparison with different compression methods on TREC 2019 Deep Learning Track. Compression ratio is set to 64x, i.e., 48 bytes per passage/document. */ denotes that RepCONC performs significantly better than baselines at $p < 0.05/0.01$ level using the two-tailed pairwise t-test. ‘Unsup. Compr.’ and ‘Sup. Compr.’ denote unsupervised compression methods and supervised compression methods, respectively. Best compression method in each column is marked bold.**

Model	Compr.	MARCO Passage		DL Passage		MARCO Doc		DL Doc	
	Ratio	MRR@10	R@100	NDCG@10	R@100	MRR@100	R@100	NDCG@10	R@100
Uncompressed									
ANCE [33]	1x	0.338	0.862	0.654	0.445	0.377**	0.894**	0.610	0.273
ADORE [36]	1x	0.347	0.876	0.683	0.473	0.405	0.919	0.628	0.317
Unsup. Compr.									
PQ [17]	64x	0.028**	0.193**	0.077**	0.067**	0.038**	0.205**	0.076**	0.043**
ScaNN [13]	64x	0.034**	0.402**	0.085**	0.121**	0.149**	0.563**	0.318**	0.137**
ITQ+LSH [12]	64x	0.271**	0.782**	0.501**	0.367**	0.322**	0.845**	0.543**	0.263**
OPQ [11]	64x	0.290**	0.830**	0.591**	0.417**	0.340**	0.880**	0.579	0.282**
OPQ+ScaNN	64x	0.310**	0.837**	0.586**	0.429**	0.357**	0.893**	0.564*	0.286
Sup. Compr.									
DPQ [4, 38]	64x	0.305**	0.840**	0.589**	0.440	0.353**	0.891**	0.576	0.302
JPQ [35]	64x	0.332**	0.863	0.644	0.447	0.384**	0.905*	0.608	0.302
RepCONC (Ours)	64x	0.340	0.864	0.668	0.492	0.399	0.911	0.600	0.305

λ in Eq. (8) is set to 0.05 for $M = 48/32/24$, 0.07 for $M = 16$, 0.1 for $M = 12$, 0.2 for $M = 8$, and 0.3 for $M = 4$; batch sizes are separately set to 1024 and 256 for passage and document ranking. In the second stage where dynamic hard negatives are used, the optimizer is AdamW [24]; learning rates are 2×10^{-6} and 2×10^{-5} separately for encoders and centroid embeddings; batch sizes are set to 128. For $M = 4$, we replace Eq. (6) with LambdaLoss [3] for better ranking performance. Training time is about 4 hours for passage ranking and 2 hours for document ranking.

Now we present our hardware settings and details about latency measurement. We use Xeon Gold 5218 CPUs and RTX 3090 GPUs. When training and measuring latency, we use one CPU thread and one GPU. Training time is about 9 hours for passage ranking and 2 days for document ranking on one RTX 3090 GPU. Additional notes about latency measurement are as follows. BoW search and vector search are both conducted on the CPU. For most neural retrieval models including RepCONC, query encoding is required and is performed on GPU. In our reranking experiments, the reranking models are also running on GPU.

5 EXPERIMENTS

We empirically evaluate RepCONC to address the following three research questions:

- **RQ1:** Can RepCONC substantially compress the index without significantly hurting retrieval effectiveness?
- **RQ2:** How does RepCONC perform compared with other retrieval models?
- **RQ3:** How does constrained clustering contribute to the effectiveness of RepCONC?

5.1 Comparison with Compression Methods

This section compares RepCONC with vector compression baselines to answer **RQ1**. We compare it in two ways, a fixed 64x compression ratio and different compression ratios ranging from 64x to 784x.

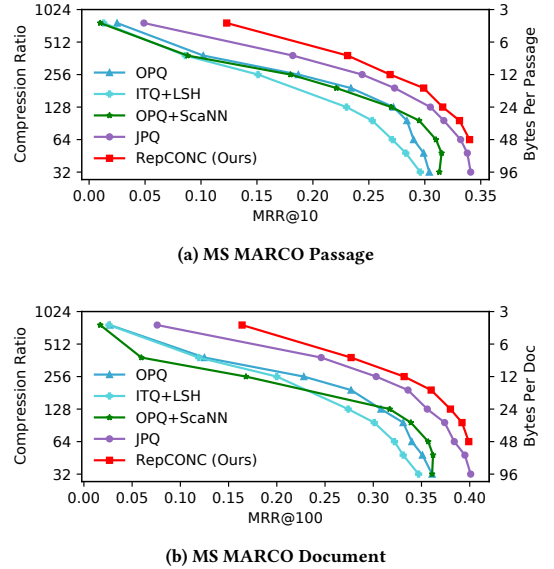
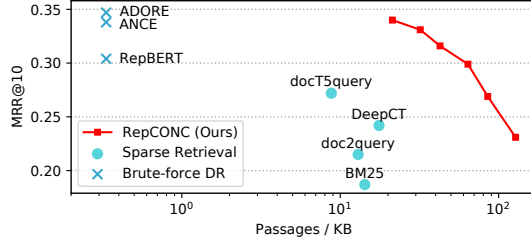
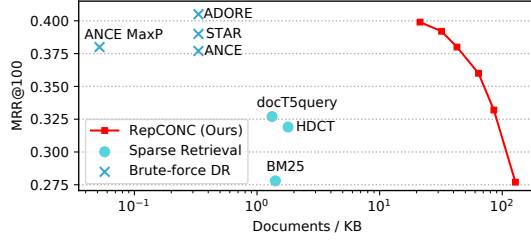


Figure 3: Comparison with compression methods. Up and right is better.

Ranking performances given a fixed 64x compression ratio are presented in Table 1. Even if the index is compressed by 64 times, RepCONC outperforms ANCE [33] and almost matches ADORE [36], the state-of-the-art DR model trained by negative sampling. Compared with unsupervised compression methods, RepCONC exhibits significant performance gains and demonstrates the importance of joint learning. As for supervised compression baselines, RepCONC significantly outperforms DPQ [4, 38] and especially outperforms the recently proposed state-of-the-art JPQ model [35] on most



(a) MS MARCO Passage



(b) MS MARCO Document

Figure 4: Comparison with first-stage retrieval models in terms of effectiveness-memory trade-off. Up and right is better. The x-axis indicates the average number of passages/documents stored in 1 kilobyte.

metrics. JPQ cannot train Index Assignments. It uses K-Means to generate them and fixes them during training. Our proposed RepCONC, on the contrary, is able to update Index Assignments during training and results show its effectiveness.

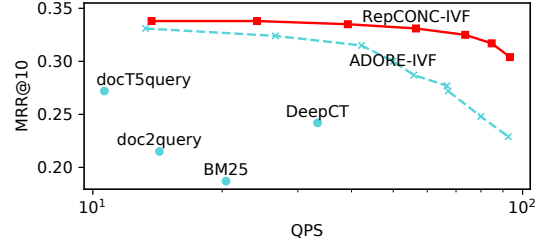
Ranking performances in terms of different compression ratios are plotted in Figure 3. The advantage of RepCONC is more significant when larger compression ratios are used. For example, its MRR score is more than twice the JPQ’s score when the compression ratio is 784x. We believe this is because RepCONC is able to generate high-quality Index Assignments specifically for ranking effectiveness, which becomes more important when fewer bytes are used. Instead, JPQ uses K-Means to produce task-blind Index Assignments and compromises ranking performance.

5.2 Comparison with Retrieval Models

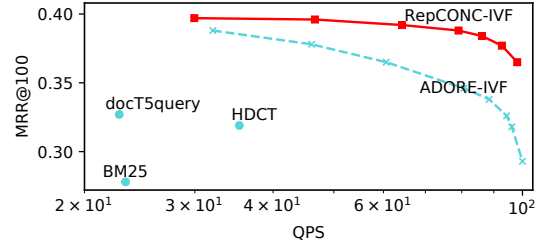
This section compares RepCONC with various retrieval models to address **RQ2**. We firstly compare it with first-stage retrievers, including BoW models and DR models. Then we compare it with complex (slow) end-to-end retrievers.

5.2.1 Comparison with First-Stage Retrievers.

Figure 4 summarizes the effectiveness-memory tradeoff. As the figure shows, although DR models are much more effective than BoW models, they incur severe memory inefficiency. By jointly training the dual-encoders and quantization methods, RepCONC substantially improves memory efficiency of DR while still being very effective in ranking. It outperforms RepBERT [37] and ANCE [33] in effectiveness, and is almost as effective as ADORE [36],



(a) MS MARCO Passage



(b) MS MARCO Document

Figure 5: Comparison with first-stage retrieval models in terms of effectiveness-latency trade-off. Up and right is better. The search is performed on CPU with one thread. QPS stands for ‘query per second’.

the state-of-the-art DR model trained by negative sampling. Compared with BoW models, it can build a much smaller index, especially on document dataset where text is much longer than that on passage dataset. For example, on document ranking task, it can build a 100x smaller index than BM25 while still being equally effective.

Figure 5 summarizes the effectiveness-latency tradeoff. To verify that RepCONC is more time-efficient than existing uncompressed DR models, we arm the state-of-the-art uncompressed DR model, ADORE, with the same IVF method [17] as RepCONC employs. As the figure shows, both RepCONC-IVF and ADORE-IVF substantially outperform BoW models with the help of IVF acceleration. Most importantly, RepCONC-IVF outperforms ADORE-IVF, especially at large QPS settings. This is because PQ already provides RepCONC with about 15x speedup compared with brute-force dense retrieval. Therefore, ADORE is more dependent on IVF than RepCONC and has to sacrifice more effectiveness for acceleration. The results demonstrate the time efficiency of RepCONC.

5.2.2 Comparison with Complex End-to-End Retrievers.

This section compares RepCONC with some complex (slow) end-to-end neural retrieval models. These models achieve better ranking performance with much higher query latency because of their complex model architecture. In consideration of fair comparison, we add a reranking stage to RepCONC and compare them in terms of effectiveness-latency tradeoff. The reranking models

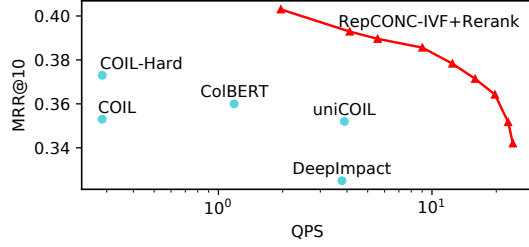


Figure 6: Comparison with complex (slow) end-to-end retrieval models in terms of effectiveness-latency tradeoff on MS MARCO Passage Ranking. The search is performed on CPU with one thread. Up and right is better. QPS stands for ‘query per second’.

Table 2: Ablation study on MSMARCO Passage Ranking dataset. BPP stands for ‘bytes per passage’

Models	MRR@10	
	BPP:16	BPP:48
Baselines		
DPQ [4, 38]	0.244	0.305
JPQ [35]	0.273	0.332
RepCONC		
OPQ [11]	0.234	0.290
+ Clustering	0.275	0.332
+ Constraint	0.284	0.337
+ Dynamic Neg	0.294	0.340

are MonoBERT and DuoT5 models open-sourced by the pygaggle library⁴. MonoBERT firstly reranks top passages retrieved by RepCONC-IVF. Then DuoT5 further reranks the top passages output by MonoBERT. We tune the IVF speedup ratio, the MonoBERT reranking depth, and the DuoT5 reranking depth to evaluate ranking performance at different query latency. Note, query encoding and reranking are performed on GPU while the search is performed on CPU with one thread.

Ranking performances are summarized in Figure 6. We can see that RepCONC-IVF+Rerank substantially outperforms all baselines in terms of both effectiveness and time efficiency. In fact, RepCONC is also more memory-efficient than these baselines. The index size of RepCONC is less than 0.5GB, while COIL [10], ColBERT [20], uniCOIL [21], and DeepImpact [26] separately consumes 60GB, 162GB, 1.3GB, and 1.5GB for storing indexes. Therefore, RepCONC’s efficient and effective retrieval is highly beneficial to second-stage reranking and helps the two-stage ranking achieve much better ranking performance than complex end-to-end retrieval models.

5.3 Ablation Study

This section conducts an ablation study to answer RQ3. We summarize the results in Table 2.

⁴<https://github.com/castorini/pygaggle>

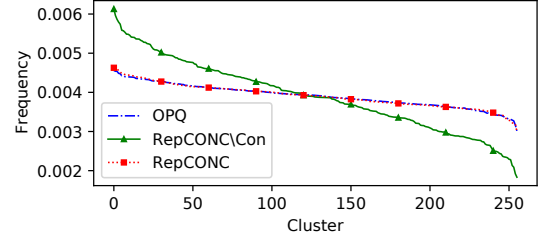


Figure 7: Cluster distribution on MS MARCO Passage Ranking. Clusters are sorted by the assigned frequency. Distributions across different sub-vector blocks are averaged. RepCONC\Con indicates RepCONC without constraint.

As the results show, clustering object helps RepCONC outperform DPQ [4, 38]. Although DPQ also utilizes a similar MSE loss, the gradients with respect to it only backpropagate to PQ centroids. Therefore, DPQ updates Index Assignments with a trick similar to Batch K-Means [2] instead of clustering. However, Batch K-Means is shown to converge slowly [31]. Besides, the target distribution of document embeddings is also changing during training, which makes convergence harder.

With the help of the uniform clustering constraint, RepCONC outperforms the state-of-the-art JPQ method [35], which uses fixed Index Assignments generated by OPQ [11]. It demonstrates that simply adding a clustering loss is risky to retrieval effectiveness and that the constraint helps to tackle this problem by distinguishing the quantized vectors. The ranking performance is further improved by employing dynamic hard negatives [36].

To further verify that the constraint helps produce balanced clustering results, we plot the frequencies of clusters being assigned in Figure 7. Without the constraint, RepCONC\Con generates unbalanced clustering distribution. With the help of the constraint, the distribution is more balanced. It is similar to that of OPQ [11], which uses K-Means for clustering. The distribution is not uniform because we do not use the constraint during inference as discussed in Section 3.6.3.

6 CONCLUSIONS

To solve the efficiency issue existing in brute-force DR models, we present RepCONC, which learns discrete representations by modeling quantization as constrained clustering in the joint learning process. The clustering object requires the document embeddings to be clustered around the quantization centroids and facilitates joint optimization of PQ parameters and dual-encoders. To tackle the risk that clustering assigns vectors to only a few major centroids and results in indistinguishable quantized vectors, we introduce a uniform clustering constraint that enforces the vectors to be equally quantized to all possible centroids during training. The constraint is approximately solved as an instance of the optimal transport problem. In addition to constrained clustering, RepCONC employs the inverted file system (IVF) to enable efficient vector search on CPUs. We conduct experiments on two widely-adopted ad-hoc retrieval benchmarks. Experimental results show that RepCONC significantly outperforms competitive quantization baselines and

substantially improves the memory efficiency and time efficiency of DR. It substantially outperforms various retrieval models in terms of retrieval effectiveness, memory efficiency, and time efficiency. The ablation study demonstrates that constrained clustering is the key to the effectiveness of RepCONC.

REFERENCES

- [1] Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamee, Bhaskar Mitra, Tri Nguyen, et al. 2016. Ms marco: A human generated machine reading comprehension dataset. *arXiv preprint arXiv:1611.09268* (2016).
- [2] Leon Bottou and Yoshua Bengio. 1995. Convergence properties of the k-means algorithms. In *Advances in neural information processing systems*. 585–592.
- [3] Christopher JC Burges. 2010. From ranknet to lambdarank to lambdamart: An overview. *Learning* 11, 23–581 (2010), 81.
- [4] Ting Chen, Lala Li, and Yizhou Sun. 2020. Differentiable product quantization for end-to-end embedding compression. In *International Conference on Machine Learning*. PMLR, 1617–1626.
- [5] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Ellen M. Voorhees. 2020. Overview of the TREC 2019 deep learning track. In *Text REtrieval Conference (TREC)*. TREC.
- [6] Marco Cuturi. 2013. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems* 26 (2013), 2292–2300.
- [7] Zhuyun Dai and Jamie Callan. 2019. Context-aware sentence/passage term importance estimation for first stage retrieval. *arXiv preprint arXiv:1910.10687* (2019).
- [8] Zhuyun Dai and J. Callan. 2020. Context-Aware Document Term Weighting for Ad-Hoc Search. *Proceedings of The Web Conference 2020* (2020).
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 4171–4186.
- [10] Luyu Gao, Zhuyun Dai, and Jamie Callan. 2021. COIL: Revisit Exact Lexical Match in Information Retrieval with Contextualized Inverted List. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Online, 3030–3042.
- [11] Tiezheng Ge, Kaiming He, Qifa Ke, and Jian Sun. 2013. Optimized product quantization. *IEEE transactions on pattern analysis and machine intelligence* 36, 4 (2013), 744–755.
- [12] Yunchao Gong, Svetlana Lazebnik, Albert Gordo, and Florent Perronnin. 2012. Iterative quantization: A procrustean approach to learning binary codes for large-scale image retrieval. *IEEE transactions on pattern analysis and machine intelligence* 35, 12 (2012), 2916–2929.
- [13] Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Kumar. 2020. Accelerating large-scale inference with anisotropic vector quantization. In *International Conference on Machine Learning*. PMLR, 3887–3896.
- [14] 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 (SIGIR '21)*. 113–122.
- [15] Jui-Ting Huang, Ashish Sharma, Shuying Sun, Li Xia, David Zhang, Philip Pronin, Janani Padmanabhan, Giuseppe Ottaviano, and Linjun Yang. 2020. Embedding-based Retrieval in Facebook Search. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2553–2561.
- [16] Piotr Indyk and Rajeev Motwani. 1998. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*. 604–613.
- [17] Herve Jegou, Matthijs Douze, and Cordelia Schmid. 2010. Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence* 33, 1 (2010), 117–128.
- [18] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data* (2019).
- [19] 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)*.
- [20] O. Khattab and M. Zaharia. 2020. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval* (2020).
- [21] Jimmy Lin and Xueguang Ma. 2021. A Few Brief Notes on DeepImpact, COIL, and a Conceptual Framework for Information Retrieval Techniques. *arXiv preprint arXiv:2106.14807* (2021).
- [22] 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).
- [23] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: a robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [24] Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).
- [25] Yu A Malkov and Dmitry A Yashunin. 2018. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and machine intelligence* 42, 4 (2018), 824–836.
- [26] Antonio Mallia, Omar Khattab, Torsten Suel, and Nicola Tonello. 2021. Learning Passage Impacts for Inverted Indexes. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21)*.
- [27] Rodrigo Nogueira, Jimmy Lin, and AI Epistemic. 2019. From doc2query to docTTTTTquery. *Online preprint* (2019).
- [28] Rodrigo Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. 2019. Document expansion by query prediction. *arXiv preprint arXiv:1904.08375* (2019).
- [29] 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.
- [30] Stephen E Robertson and Steve Walker. 1994. Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In *SIGIR '94*. Springer, 232–241.
- [31] David Sculley. 2010. Web-scale k-means clustering. In *Proceedings of the 19th international conference on World wide web*. 1177–1178.
- [32] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. HuggingFace's Transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771* (2019).
- [33] 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 *International Conference on Learning Representations*. <https://openreview.net/forum?id=zeFrfgYzIn>
- [34] Ikuya Yamada, Akari Asai, and Hannaneh Hajishirzi. 2021. Efficient Passage Retrieval with Hashing for Open-domain Question Answering. In *ACL*.
- [35] Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2021. Jointly Optimizing Query Encoder and Product Quantization to Improve Retrieval Performance. In *Proceedings of the 30th ACM International Conference on Information and Knowledge Management*.
- [36] Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2021. Optimizing Dense Retrieval Model Training with Hard Negatives. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21)*. 1503–1512.
- [37] 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).
- [38] Han Zhang, Hongwei Shen, Yiming Qiu, Yunjiang Jiang, Songlin Wang, Sulong Xu, Yun Xiao, Bo Long, and Wen-Yun Yang. 2021. Joint Learning of Deep Retrieval Model and Product Quantization Based Embedding Index. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21)*. 1718–1722.