

# A Dense Representation Framework for Lexical and Semantic Matching

SHENG-CHIEH LIN, University of Waterloo, Canada

JIMMY LIN, University of Waterloo, Canada

Lexical and semantic matching capture different successful approaches to text retrieval and the fusion of their results has proven to be more effective and robust than either alone. Prior work performs hybrid retrieval by conducting lexical and semantic matching using different systems (e.g., Lucene and Faiss, respectively) and then fusing their model outputs. In contrast, our work integrates lexical representations with dense semantic representations by densifying high-dimensional lexical representations into what we call low-dimensional dense lexical representations (DLRs). Our experiments show that DLRs can effectively approximate the original lexical representations, preserving effectiveness while improving query latency. Furthermore, we can combine dense lexical and semantic representations to generate dense hybrid representations (DHRs) that are more flexible and yield faster retrieval compared to existing hybrid techniques. In addition, we explore *jointly* training lexical and semantic representations in a single model and empirically show that the resulting DHRs are able to combine the advantages of the individual components. Our best DHR model is competitive with state-of-the-art single-vector and multi-vector dense retrievers in both in-domain and zero-shot evaluation settings. Furthermore, our model is both faster and requires smaller indexes, making our dense representation framework an attractive approach to text retrieval. Our code is available at <https://github.com/castorini/dhr>.

CCS Concepts: • **Information systems** → **Top-k retrieval in databases; Search engine indexing; Retrieval effectiveness; Retrieval efficiency.**

Additional Key Words and Phrases: Sparse Retrieval; Dense Retrieval; Hybrid Retrieval; Vector Compression

## 1 INTRODUCTION

Transformer-based bi-encoders have been widely used as first-stage retrievers for text retrieval. Compared to their multi-vector counterparts [12, 19, 27], single-vector representation learning approaches (with a few representative techniques listed in Table 1) are attractive due to their good balance between effectiveness and efficiency.

Semantic matching through dense representations [3, 26, 42] form one large successful class of models.<sup>1</sup> These dense semantic representations transform the relevance matching problem into nearest neighbor search in a semantic space, tackling vocabulary and semantic mismatches in ways that traditional lexical matching approaches (e.g., BM25) cannot. Subsequent work further improves these models through advanced training techniques [11, 18, 22, 32, 41, 50, 54]. However, as Sciavolino et al. [47] show, dense semantic representations still fail in some easy cases, and it remains challenging to interpret why they sometimes perform poorly. In addition, Thakur et al. [48] demonstrate that many existing dense retrievers still fall short in terms of generalization capability across different domains.

There is evidence [13, 28, 34] that lexical matching compensates for the weaknesses of semantic matching; these papers further propose hybrid retrieval techniques that fuse lexical and semantic representations. However, in practice, lexical and semantic matching are executed in very different ways: typically, lexical matching is conducted using inverted indexes, for example, in Lucene, while semantic matching is treated as nearest neighbor search using, for example, HNSW indexes in Faiss [24]. This means that a hybrid retrieval system requires two separate “software stacks”, running completely distinct retrieval operations in parallel before their outputs are post-processed to generate a final ranking (e.g., through linear combination of scores). Such a design makes

<sup>1</sup>In the literature, these are often just called dense retrieval models. However, we explicitly refer to these as dense *semantic* representations because we show that dense *lexical* representations exist as a separate class of models.

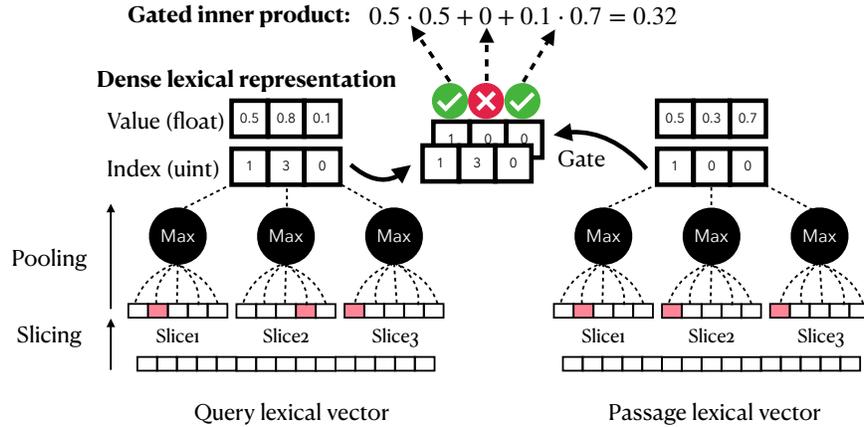


Fig. 1. Illustration of densified lexical representations (DLRs) and gated inner product (GIP). Both query and passage lexical representations are fixed-width vectors where the number of dimensions is equal to the vocabulary size. Our approach first groups these high-dimensional vectors into  $M$  slices, each with  $N$  dimensions (e.g.,  $M = 3$ ,  $N = 5$  here). For each slice, the maximum value is selected. These values from the original vector (Value) and their positions in each slice (Index) are recorded separately. When computing the gated inner product between two vectors, only the dimensions with the same index are considered.

real-world deployments more complicated than necessary, since, for example, the two separate indexes need to be maintained and be kept in sync.

Recently, another thread of work uses bi-encoders to learn sparse lexical (bag-of-words) representations for text retrieval. For example, Dai and Callan [7] demonstrate with DeepCT that replacing tf-idf with contextualized term weights from a transformer-based regression model significantly improves retrieval effectiveness. Subsequent work further combines term reweighting techniques with term expansion to address vocabulary and semantic mismatch issues with lexical representations. Some methods [28, 36, 56] leverage another model for expansion, which incurs additional costs in both training and inference. As an alternative, Formal et al. [9] exploit BERT’s masked language model (MLM) layer to train a single model for both expansion and term weighting.

Compared to dense semantic models, sparse lexical models appear to be more robust to domain shifts [48]. However, the optimization of sparse lexical retrievers must take into account the efficiency of query evaluation using inverted indexes. For example, to achieve better effectiveness, term expansion techniques tend to make sparse lexical representations more dense, sometimes rendering retrieval with inverted indexes much slower [35]. Further performance penalties are incurred when integrating sparse lexical and dense semantic representations into hybrid retrieval systems. Since it is impractical to directly compute dot products between high-dimensional vectors at scale in latency-sensitive retrieval applications (i.e., via brute-force approaches), using lexical representations with inverted indexes remains presently the only sensible choice, but this comes with the aforementioned limitations.

Motivated by these tradeoffs, we explore an approach to *directly* computing dot products between lexical vectors in a scalable and low-latency manner for retrieval applications. This is accomplished by densifying lexical representations, which can be applied to any existing lexical model. The approach described in this paper, which builds on our previous work [31], is comprised of two simple steps: (1) representation slicing (2) sliced representation pooling. This method can be viewed as compressing high-dimensional lexical representations into low-dimensional ones, as depicted in Fig. 1. One key feature of our approach is that it does not involve either unsupervised [24] or supervised training [4, 53, 55]. Any collection of lexical representations (from a “base model”)

can be converted into dense lexical representations (DLRs), and the nearest neighbor search problem between query and passage vectors can be performed with an operation we call the gated inner product (GIP). Standard dense representations and inner products can be considered a special case of DLRs and GIP, respectively. Thus, our representation and scoring function comprise a unified framework for *both* lexical and semantic matching.

An advantage of our framework is that end-to-end lexical and hybrid retrieval can be performed *directly* on GPUs with a single index structure in a unified execution environment. In fact, our implementation uses standard vector operations directly in the popular PyTorch open-source neural modeling toolkit. Retrieval efficiency, unlike with inverted indexes, is not sensitive to the sparsity of representations; thus, we can optimize retrieval effectiveness without any constraints. Building on these features, we propose a new model called Dense Lexical AnD Expansion (DeLADE) and explore a new matching approach, which we refer to as dense lexical matching (shown in Table 1).

In addition, we can integrate DLRs with dense semantic representations into what we call dense hybrid representations (DHRs). This can be accomplished in two ways: First, we can independently combine DLRs with any “off-the-shelf” dense semantic representations such as ANCE [50], TAS-B [18], etc. Second, we can *jointly* train DHRs that combine lexical as well as semantic components. Both DLRs and DHRs exhibit advantageous features for real-world applications: (1) end-to-end retrieval can be accomplished on GPUs using a single index and software framework, instead of, for example, using Lucene for lexical matching and Faiss for semantic matching; (2) vector densification can be applied to meet storage constraints without model retraining.

With our proposed dense representation framework, we explore the following research questions:

**RQ1** How well do DLRs approximate the original high-dimensional lexical representations?

We densify lexical representations from two lexical models based on whole word matching (BM25 and DeepImpact [36]) and two lexical models based on wordpiece token matching (uniCOIL [28] and SPLADE [9]). Experimental results show that our approach effectively compresses high-dimensional lexical representations (30K and even 3.5M dimensions) into 768- and 128-dimensional vectors with less than 1% and 5% retrieval effectiveness loss, respectively. In addition, this approach can be applied to DeLADE, our proposed dense lexical representation model, which is a SPLADE variant. More importantly, this compression enables us to perform lexical matching directly on GPUs, which substantially speeds up retrieval compared to using an inverted index and does not depend on the sparsity of the lexical representations.

**RQ2** How well do DHRs benefit from the independent fusion of DLRs and “off-the-shelf” dense semantic representations?

Within our framework, we propose hybrid dense representations (DHRs) by combining DLRs and other standard “off-the-shelf” dense semantic representations for hybrid fusion retrieval. We demonstrate that our method achieves comparable retrieval effectiveness to other existing hybrid retrieval methods but with lower query latency.

Next, we propose to *jointly* train lexical and semantic components within DHRs. Specifically, we combine lexical representations and the [CLS] embeddings from a transformer to capture lexical and semantic matching, and then conduct retrieval based on a fusion of these separate components in a single unified framework.

**RQ3** Can DHRs benefit from *joint* training of lexical and semantic components in a single model?

Experiments show that our approach to jointly training DHRs (i.e., single model fusion) outperforms both dense semantic and sparse lexical retrieval models and is competitive with state-of-the-art multi-vector approaches such as ColBERT [27] and COIL [12]. We achieve effectiveness on par with the state of the art on an in-domain benchmark (MS MARCO [2]) and obtain better generalization on an out-of-domain benchmark (BEIR [48]). In addition, the index size of our model can be tuned by controlling the vector dimensionality of the lexical component *without retraining*, making our approach attractive for real-world applications.

Table 1. Comparison of Single-Vector Representation Learning Approaches

matching type		Sparse	Dense
Lexical	whole words	DeepCT [7], DeepImpact [36]	<b>DLR</b> (our approach)
	wordpiece tokens	uniCOIL [28], SPLADE [10]	
Semantic		UHD [23]	DPR [26], ANCE [50] TASB [18], RocketQA [41]
Hybrid		-	<b>DHR</b> (our approach)

Further analysis identifies that GIP requires more operations compared to the standard inner product for dense vectors and this is one potential drawback of end-to-end retrieval with DLRs. To address this issue, we propose a two-stage retrieval approach: in the first stage, retrieval is conducted by *approximate* GIP, which is faster but less accurate; then, the retrieved sets are reranked by computing the exact GIP. Thus, we explore the following research question:

**RQ4** How effective is our proposed two-stage retrieval approach?

Our experiments show that approximate GIP is capable of retrieving sufficient relevant passages (i.e., achieves good recall) for the subsequent exact GIP reranking to achieve high end-to-end effectiveness. Specifically, we demonstrate that with approximate GIP, our proposed two-stage retrieval approach substantially reduces retrieval latency without sacrificing any effectiveness compared to end-to-end retrieval using more expensive but exact GIP.

Our contributions are summarized as follows:

- We propose a simple yet effective approach to densifying high-dimensional lexical representations for text retrieval, creating what we call dense lexical representations (DLRs).
- Building on DLRs, we introduce dense hybrid representations (DHRs) that combine lexical and semantic representations.
- While DHRs can combine arbitrary off-the-shelf lexical representations (e.g., BM25 and uniCOIL) and semantic representations (e.g., ANCE and TAS-B) *independently*, we demonstrate how to *jointly* train effective DHRs with complementary lexical and semantic components.
- We show how to efficiently conduct two-stage retrieval in our dense representation framework, with fast approximate GIP followed by exact GIP reranking.

Code to reproduce all experiments in this paper is available at <https://github.com/castorini/dhr>.

## 2 BACKGROUND AND RELATED WORK

Following Lin et al. [30], let us formulate the task of text (or *ad hoc*) retrieval as follows: Given a query  $q$ , the goal is to retrieve a ranked list of documents  $\{d_1, d_2, \dots, d_k\} \in C$  to maximize some ranking metric, such as nDCG, where  $C$  is the collection of documents.

Specifically, given a (query, passage) pair, we aim to maximize the following:

$$\text{sim}(q, d) \triangleq \phi(\eta_q(q), \eta_d(d)) = \langle \mathbf{q}, \mathbf{d} \rangle, \quad (1)$$

where  $\eta_q(\cdot)$  and  $\eta_d(\cdot) \in \mathbb{R}^h$  denote functions mapping the query and the passage into  $h$ -dimensional vector representations,  $\mathbf{q}$  and  $\mathbf{d}$ , respectively. The scoring function that quantifies the degree of relevance between the representations  $\mathbf{q}$  and  $\mathbf{d}$  is denoted  $\phi(\cdot, \cdot)$ , which can be a simple inner product or a more complex operation [12, 20, 27, 38]. We focus on single-vector representation learning approaches that apply the inner product as the

scoring function. We categorize single-vector representation learning approaches through “matching type”, as shown in Table 1. In the literature, there are two main lines of research: dense representations for semantic matching and sparse representations for lexical matching.

*Dense representations for semantic matching.* Pretrained transformers [8, 33] are able to encode sentences or passages into dense semantic representations, which have been shown to be effective for downstream tasks [3, 42]. In recent years, transformer-based bi-encoders have been widely applied to the task of passage retrieval [26] and further improved by advanced training techniques such as hard negative mining [50, 54], knowledge distillation [18, 32], pretraining [11, 22], or their combination [41]. These approaches encode queries and passages into dense vectors, using the inner product to capture the degree of relevance:

$$\text{sim}_{\text{semantic}}(q, d) \triangleq \langle \mathbf{q}_{[\text{CLS}]}, \mathbf{d}_{[\text{CLS}]} \rangle, \quad (2)$$

where  $\mathbf{q}_{[\text{CLS}]}$  and  $\mathbf{d}_{[\text{CLS}]}$  are typically 768-dimensional vectors taken from the [CLS] token in the final layer of a transformer model (or alternatively, pooling over the contextualized representations of the tokens).

*Sparse representations for lexical matching.* To our knowledge, Zamani et al. [52] was the first to demonstrate that neural networks can learn lexical representations for text retrieval. Recently, transformer-based bi-encoders have also been applied to lexical representation learning by replacing heuristic term weighting functions (e.g., BM25) with learned term weights. In the literature, these solutions can be classified into two broad classes: (1) linear layer and (2) MLM projection.

One early technique, DeepCT [7] uses a linear layer to project BERT token embeddings into contextualized term weights for input tokens. Subsequent work [28, 36, 56] further combines DeepCT and other document expansion methods [39, 57] to address vocabulary mismatch and missing terms in DeepCT. For example, DeepImpact [36] uses a trained sequence-to-sequence model [39] to expand the original passages in the corpus and learn contextualized term weights. Other models [28, 56] follow a similar approach to DeepImpact but term matching is performed in the space of BERT wordpiece tokens rather than whole words. They generate term weights only for tokens appearing in each (possibly expanded) query or passage and thus the lexical representations are sparse by design. However, to achieve competitive effectiveness, these techniques require additional models for term expansion (see Zhuang and Zuccon [56] for more discussion). In contrast, some researchers [1, 10] use the masked language model (MLM) layer in transformers such as BERT to learn term weighting and expansions at the same time. For these approaches, sparsity regularization must be applied during model training to ensure that the generated lexical representations are amenable to retrieval using inverted indexes.

Generally, these neural approaches to lexical term matching can be viewed as projecting queries and passages into  $|V_{\text{BERT}}|$ -dimensional vectors, where  $|V_{\text{BERT}}| = 30522$  is the vocabulary size of BERT wordpiece tokens:

$$\text{sim}_{\text{lexical}}(q, d) \triangleq \langle \mathbf{q}_{\text{BoW}}, \mathbf{d}_{\text{BoW}} \rangle, \quad (3)$$

where  $\mathbf{q}_{\text{BoW}}$  and  $\mathbf{d}_{\text{BoW}} \in \mathbb{R}^{30522}$ . The value in each dimension is the token’s term weight. As with dense retrieval, the relevance score between a query and a passage is computed by the inner product of their vector representations.

*Bridging the gap between the two worlds.* Although the inner product is a common operation for capturing relevance in the aforementioned two approaches, there are still two major differences between them:

- (1) Unlike semantic representations, lexical representations can be considered bags of words (or subwords) and thus are more interpretable, since dimensions of the representation vectors directly correspond to vocabulary items.
- (2) Text retrieval using lexical representations is usually performed using standard inverted indexes due to their high dimensionality, while text retrieval using semantic representations is usually performed using completely different infrastructure, e.g., HNSW indexes.

Table 2. Comparison of Different Vector Compression Approaches

Approach	unsupervised training	supervised training
PQ [25], OPQ [15], LSH [21]	✓	✗
JPQ [53], RepCONC [55]	✓	✓
SPAR [4]	✗	✓
Our work	✗	✗

The advantage of our approach is that it does not require *any* training.

Previous work [13, 28, 34] has demonstrated that semantic and lexical matching can compensate for each other; these papers typically implement dense–sparse hybrid retrieval by performing retrieval independently using different systems and then merging their results (e.g., by interpolating scores). To make such an approach “production-ready” for deployment in real-world applications, non-trivial software engineering effort is required to coordinate dense and sparse retrieval in parallel (on inverted and HNSW indexes) and the final fusion. The need for two entirely separate “software stacks” and the associated operational maintenance costs (e.g., of keeping indexes in sync) increase the complexity of hybrid retrieval systems.

To bridge the gap between semantic and lexical representations for text retrieval, some researchers extend their focus beyond the above two research threads. For example, Jang et al. [23] show that semantic matching can be executed using an inverted index by projecting semantic representations from BERT to sparse representations in an ultra-high-dimensional space. However, this projection operation requires additional training and the retrieval effectiveness of this approach still lags behind baseline dense retrieval approaches.

In contrast, another approach to bridging dense and sparse representations is to compress high-dimensional sparse lexical representations into low-dimensional dense ones. We can accomplish this using existing unsupervised approaches such as product quantization (PQ) [25], optimized product quantization (OPQ) [15], and locality-sensitive hashing (LSH) [21]. However, such approaches require substantial computational resources since the entire corpus (or large portions thereof) need to be loaded into CPU/GPU memory to perform the unsupervised training. These approaches are not practical to compress ultra-high-dimensional lexical vectors, especially when the corpus is large.<sup>2</sup> Supervised compression techniques [53, 55] that are built upon these unsupervised approaches similarly suffer from high resource requirements.

Closest to our own work, Chen et al. [4] distill lexical matching signals from existing sparse retrieval models such as BM25 and uniCOIL [28] into low-dimensional dense representations in an approach called SPAR. Nevertheless, SPAR requires massive amounts of training data and computational resources (e.g., 64 V100 GPUs for three days). In contrast, our approach simply performs max pooling over each ultra-high-dimensional lexical vector. Thus, we do not require any additional computational resources for unsupervised or supervised training. We summarize the comparison of different compression techniques in Table 2.

Beyond not needing any training (either supervised or unsupervised), our approach to lexical representation compression confers the additional advantage of interpretability. Specifically, we densify lexical representations in a reversible manner (except for some information loss from the max pooling operation). Thus, our dense lexical vectors still retain characteristics of the original lexical representations with respect to matching terms, as we will demonstrate.

<sup>2</sup>As an example, BM25 with a standard English tokenizer in the Lucene search library generates representations with 2.6M dimensions for the MS MARCO passage corpus (containing 8.8M passages).

### 3 METHODOLOGY

In this section, we first describe our approach to densifying lexical representations by slicing and pooling. We then introduce our dense representation framework, under which the densified lexical representations are captured in compact pairs of vectors. We then propose a new scoring function called gated inner product (GIP) for computing query–passage similarity. Next, we describe how to combine lexical and semantic representations using our framework and propose a two-stage approach for end-to-end retrieval. Finally, we introduce our dense lexical model, DeLADE.

#### 3.1 Dense Lexical Representations

Lexical representations can be considered vectors with  $|V|$  dimensions, i.e.,  $\mathbf{q}_{\text{BoW}} = (q_0, \dots, q_{|V|-1})$  and  $\mathbf{d}_{\text{BoW}} = (d_0, \dots, d_{|V|-1})$ . These representations come from an underlying “base model” such as uniCOIL or SPLADE, and our focus here is to “densify” such representations.

We first divide each vector into  $M$  slices, each of which is a smaller vector with  $N$  dimensions (i.e.,  $|V| = M \cdot N$ ). In terms of the standard “slice” notation used by Python:

$$\begin{aligned} S_m^q &= \mathbf{q}_{\text{BoW}}[mN : mN + N] \in \mathbb{R}^N; \\ S_m^d &= \mathbf{d}_{\text{BoW}}[mN : mN + N] \in \mathbb{R}^N, \end{aligned} \quad (4)$$

where  $m \in \{0, 1, \dots, M-1\}$ . Note that the slicing can be performed in different ways; for example, slicing randomly or with a fixed stride:  $[m : M(N-1) + m : N]$ . For simplicity, we use contiguous slicing as shown in Eq. (4) in our presentation. Thus, the inner product between  $\mathbf{q}_{\text{BoW}}$  and  $\mathbf{d}_{\text{BoW}}$  can be rewritten as the summation of all the dot products of their slices:

$$\langle \mathbf{q}_{\text{BoW}}, \mathbf{d}_{\text{BoW}} \rangle = \sum_{m=0}^{M-1} \langle S_m^q, S_m^d \rangle. \quad (5)$$

Intuitively, if a lexical representation is sparse enough, we can assume that for each slice, there is only one non-zero entry. Thus, we can approximate  $S_m^q$  ( $S_m^d$ ) by keeping only the entry with the maximum value in each slice:

$$\begin{aligned} S_m^q &\approx \max(S_m^q) \cdot \hat{\mathbf{u}}(e_m^q); \\ S_m^d &\approx \max(S_m^d) \cdot \hat{\mathbf{u}}(e_m^d), \end{aligned} \quad (6)$$

where  $\hat{\mathbf{u}}(e_m^q)$  is a unit vector with the only non-zero entry at the entry  $e_m^q = \text{argmax}(S_m^q)$ . Thus, the inner product of  $\mathbf{q}_{\text{BoW}}$  and  $\mathbf{d}_{\text{BoW}}$  lexical vectors in Eq. (5) can be approximated as follows:

$$\begin{aligned} \langle \mathbf{q}_{\text{BoW}}, \mathbf{d}_{\text{BoW}} \rangle &\approx \sum_{m=0}^{M-1} \max(S_m^q) \cdot \max(S_m^d) \cdot \langle \hat{\mathbf{u}}(e_m^q), \hat{\mathbf{u}}(e_m^d) \rangle \\ &= \sum_{m=0}^{M-1} \max(S_m^q) \cdot \max(S_m^d) \mathbb{1}_{\{e_m^q=e_m^d\}} \end{aligned} \quad (7)$$

Observing Eq. (7), in order to compute the approximate inner product of lexical vectors, each query (passage) can be alternatively represented as two  $M$ -dimension dense vectors:

$$\begin{aligned}\mathbf{q}_{\text{DLR}}^{\text{val}} &= (\max(S_0^q), \dots, \max(S_{M-1}^q)) \in \mathbb{R}^M \\ \mathbf{q}_{\text{DLR}}^{\text{idx}} &= (e_0^q, \dots, e_{M-1}^q) \in \mathbb{N}^M\end{aligned}\quad (8)$$

$$\begin{aligned}\mathbf{d}_{\text{DLR}}^{\text{val}} &= (\max(S_0^d), \dots, \max(S_{M-1}^d)) \in \mathbb{R}^M \\ \mathbf{d}_{\text{DLR}}^{\text{idx}} &= (e_0^d, \dots, e_{M-1}^d) \in \mathbb{N}^M,\end{aligned}\quad (9)$$

where  $\mathbf{q}_{\text{DLR}}^{\text{val}}$  ( $\mathbf{d}_{\text{DLR}}^{\text{val}}$ ) is the dense vector storing the  $M$  maximum values from the query (passage) slices, and  $\mathbf{q}_{\text{DLR}}^{\text{idx}}$  ( $\mathbf{d}_{\text{DLR}}^{\text{idx}}$ ) is the integer dense vector storing the entries with the maximum value in the corresponding slices.

Formally, a dense lexical representation (DLR) is a pair of vectors comprising a “value” vector and an “index” vector, as indicated in Figure 1. For a query, the DLR is  $(\mathbf{q}_{\text{DLR}}^{\text{val}}, \mathbf{q}_{\text{DLR}}^{\text{idx}})$  and for each passage,  $(\mathbf{d}_{\text{DLR}}^{\text{val}}, \mathbf{d}_{\text{DLR}}^{\text{idx}})$ . Intuitively, the “value” vector stores the most important term weight in each slice while the “index” vector stores the position of the corresponding terms in each slice. Using DLRs, Eq. (7) can be rewritten as:

$$\langle \mathbf{q}_{\text{BoW}}, \mathbf{d}_{\text{BoW}} \rangle \approx \sum_{m=0}^{M-1} \mathbf{q}_{\text{DLR}}^{\text{val}}[m] \cdot \mathbf{d}_{\text{DLR}}^{\text{val}}[m] \cdot \mathbb{1}_{\{\mathbf{q}_{\text{DLR}}^{\text{idx}}[m]=\mathbf{d}_{\text{DLR}}^{\text{idx}}[m]\}}, \quad (10)$$

where  $\mathbf{q}_{\text{DLR}}^{\text{val}}[m]$  ( $\mathbf{d}_{\text{DLR}}^{\text{val}}[m]$ ) and  $\mathbf{q}_{\text{DLR}}^{\text{idx}}[m]$  ( $\mathbf{d}_{\text{DLR}}^{\text{idx}}[m]$ ) is the  $m$ -th entry of the query (passage) DLR. Note that the query (passage) DLR represents approximations of the original lexical query (passage) vector,  $\mathbf{q}_{\text{BoW}}$  ( $\mathbf{d}_{\text{BoW}}$ ). Thus, the computation in Eq. (10) using DLRs is an approximation of the original inner product between the lexical vectors.

To simplify Eq. (10), we define a new operation, *gated inner product* (GIP) as follows:

$$\text{GIP}(\mathbf{q}, \mathbf{d}, \mathbf{g}) \triangleq \sum_{m=0}^{M-1} \mathbf{q}[m] \cdot \mathbf{d}[m] \cdot \mathbf{g}[m], \quad (11)$$

where  $\mathbf{g}$  is the gate vector with entries either equal to 0 or 1 (i.e.,  $\mathbf{g}[m] \in \{0, 1\}$ ). Thus, we can consider Eq. (10) as a GIP operation between query  $(\mathbf{q}_{\text{DLR}}^{\text{val}}, \mathbf{q}_{\text{DLR}}^{\text{idx}})$  and passage  $(\mathbf{d}_{\text{DLR}}^{\text{val}}, \mathbf{d}_{\text{DLR}}^{\text{idx}})$  DLRs:

$$\langle \mathbf{q}_{\text{BoW}}, \mathbf{d}_{\text{BoW}} \rangle \approx \text{GIP}(\mathbf{q}_{\text{DLR}}^{\text{val}}, \mathbf{d}_{\text{DLR}}^{\text{val}}, \mathbf{g}_{\text{DLR}}), \quad (12)$$

where the gate vector  $\mathbf{g}_{\text{DLR}}[m] = \mathbb{1}_{\{\mathbf{q}_{\text{DLR}}^{\text{idx}}[m]=\mathbf{d}_{\text{DLR}}^{\text{idx}}[m]\}}$  can be interpreted as the result of lexical matching between two DLRs. This is why lexical matching can still be performed in low-dimensional DLRs and is the key difference between GIP and standard inner product. In summary, DLRs and GIP (representing the original vectors and the inner product, respectively) capture the representation and scoring function in our framework.

### 3.2 Independent Model Fusion

In order to perform retrieval based on the fusion of lexical and semantic representations, we can compute their fusion scores as follows:

$$\text{sim}_{\text{hybrid}}(q, d) \triangleq \langle \mathbf{q}_{\text{BoW}}, \mathbf{d}_{\text{BoW}} \rangle + \lambda \cdot \langle \mathbf{q}_{[\text{CLS}]}, \mathbf{d}_{[\text{CLS}]} \rangle, \quad (13)$$

where  $\lambda$  is a hyperparameter. Here,  $\mathbf{q}_{[\text{CLS}]}$  and  $\mathbf{d}_{[\text{CLS}]}$  can refer to any dense semantic representation, including “off-the-shelf” ones such as ANCE, DPR, TAS-B, etc. In standard implementations, these inner products are computed using completely different systems, for example, using Lucene and Faiss for the first and second terms, respectively. Our work aims to avoid this via a unified framework.

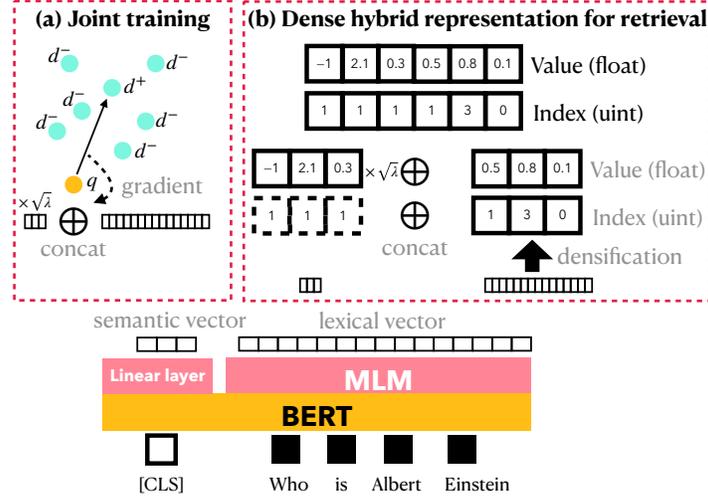


Fig. 2. Illustration of single model fusion. (a) During training, we directly concatenate semantic and lexical vectors to compute the relevance score between a query and a passage. (b) During retrieval, we concatenate the semantic and densified lexical vectors to form query and passage DHRs for end-to-end retrieval.

We approximate the first term in Eq. (13) using Eq. (12) and rewrite the second term as a special case of GIP when all the entries in the gate vector are equal to one:

$$\begin{aligned} \text{sim}_{\text{hybrid}}(q, d) &\approx \text{GIP}(\mathbf{q}_{\text{DLR}}^{\text{val}}, \mathbf{d}_{\text{DLR}}^{\text{val}}, \mathbf{g}_{\text{DLR}}) + \lambda \cdot \text{GIP}(\mathbf{q}_{[\text{CLS}]}, \mathbf{d}_{[\text{CLS}]}, \mathbf{1}) \\ &= \text{GIP}(\underbrace{\mathbf{q}_{\text{DLR}}^{\text{val}} \oplus \sqrt{\lambda} \cdot \mathbf{q}_{[\text{CLS}]}}_{\mathbf{q}_{\text{DHR}}}, \underbrace{\mathbf{d}_{\text{DLR}}^{\text{val}} \oplus \sqrt{\lambda} \cdot \mathbf{d}_{[\text{CLS}]}}_{\mathbf{d}_{\text{DHR}}}, \underbrace{\mathbf{g}_{\text{DLR}} \oplus \mathbf{1}}_{\mathbf{g}_{\text{DHR}}}), \end{aligned} \quad (14)$$

where  $\mathbf{1}$  is a vector of all ones with the same dimension as  $[\text{CLS}]$  and  $\oplus$  is vector concatenation. Observing Eq. (14), the fusion score can be considered a GIP operation between  $(\mathbf{q}_{\text{DHR}}^{\text{val}}, \mathbf{q}_{\text{DHR}}^{\text{idx}})$  and  $(\mathbf{d}_{\text{DHR}}^{\text{val}}, \mathbf{d}_{\text{DHR}}^{\text{idx}})$ :

$$\mathbf{q}_{\text{DHR}}^{\text{val}} = \mathbf{q}_{\text{DLR}}^{\text{val}} \oplus \sqrt{\lambda} \cdot \mathbf{q}_{[\text{CLS}]}; \mathbf{q}_{\text{DHR}}^{\text{idx}} = \mathbf{q}_{\text{DLR}}^{\text{idx}} \oplus \mathbf{1}, \quad (15)$$

$$\mathbf{d}_{\text{DHR}}^{\text{val}} = \mathbf{d}_{\text{DLR}}^{\text{val}} \oplus \sqrt{\lambda} \cdot \mathbf{d}_{[\text{CLS}]}; \mathbf{d}_{\text{DHR}}^{\text{idx}} = \mathbf{d}_{\text{DLR}}^{\text{idx}} \oplus \mathbf{1}. \quad (16)$$

We call the pair of vectors,  $\mathbf{q}_{\text{DHR}}^{\text{val}}$  and  $\mathbf{q}_{\text{DHR}}^{\text{idx}}$  ( $\mathbf{d}_{\text{DHR}}^{\text{val}}$  and  $\mathbf{d}_{\text{DHR}}^{\text{idx}}$ ), query (passage) dense hybrid representations (DHRs). Note that DHRs represent a special case of DLRs since they have the same form of representation and scoring function.

### 3.3 Single Model Fusion

Note that Eq. (13) is completely general and can be used to combine any arbitrary “off-the-shelf” lexical representation (e.g., uniCOIL, SPLADE, etc.) and semantic representation (e.g., ANCE, TAS-B, etc.) for hybrid retrieval. This approach can be described as fusion of *independent* models. We further study whether our framework can benefit from a single model; thus, we propose to *jointly* train lexical and semantic representations within a single model. Inspired by Gao et al. [13], our intuition is that such learned representations can better complement each other to perform both lexical and semantic matching.

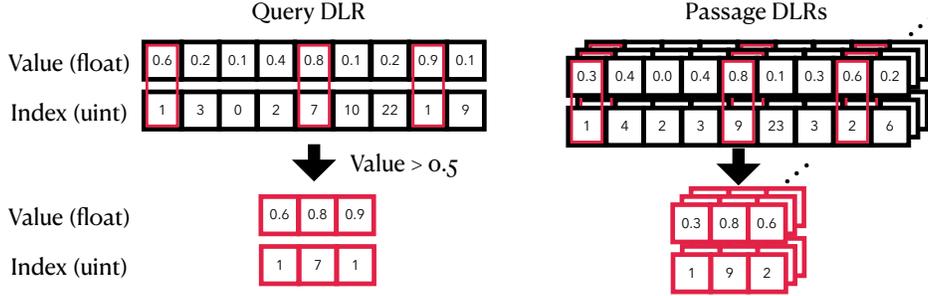


Fig. 3. Approximate GIP. DLR retrieval can be approximated by only computing the gated inner product from a few dimensions where the query value vector has values above a threshold (e.g., greater than 0.5 here). In two-stage retrieval, we can then precisely rerank the top- $K$  passages using all the dimensions.

Specifically, given a query  $q$ , its relevant passage  $d^+$  and a set of negative passages  $\{d_1^-, d_2^-, \dots, d_l^-\}$ , we train our model by minimizing negative log likelihood of the positive  $\{q, d^+\}$  pair over all the passages:

$$\text{NLL} = -\log \frac{e^{\text{sim}_{\text{hybrid}}(q, d^+)}}{e^{\text{sim}_{\text{hybrid}}(q, d^+)} + \sum_{j=1}^l e^{\text{sim}_{\text{hybrid}}(q, d_j^-)}}. \quad (17)$$

Following Karpukhin et al. [26], we also include both negative and positive passages from the other queries in the same batch as the negatives. Note that, as depicted in Figure 2, we use the exact fusion score  $\text{sim}_{\text{hybrid}}(\cdot, \cdot)$  in Eq. (13) for training, where the original 30522-dimensional lexical representations are used. However, when performing retrieval, we use the approximate fusion score in Eq. (14).

### 3.4 End-to-End Retrieval with DLRs

While it is possible to perform end-to-end retrieval for each query DLR through GIP computations against all passage DLRs in the corpus, we can identify one weakness. Unlike standard dense representations, GIP requires  $4 \cdot M$  operations, which is more than the standard inner product of  $M$ -dimensional dense vectors, which only requires  $2 \cdot M$  operations. When conducting brute-force search over corpus  $C$ , the difference becomes:  $4 \cdot M \cdot |C| > 2 \cdot M \cdot |C|$ .

To address this issue, we propose a two-stage retrieval approach inspired by previous work [27, 51]. We first retrieve the top- $K$  candidates (where  $K \ll |C|$ ) using approximate score computations and then rerank the  $K$  candidates based on the more accurate GIP computations. In this work, we propose *approximate* GIP for first-stage approximate retrieval. User queries usually contain only a few key terms, which means that when searching the entire corpus, we can perform GIP based on only a few dimensions of DLRs:

$$\text{GIP}(\mathbf{q}_{\text{DLR}}^{\text{val}}, \mathbf{d}_{\text{DLR}}^{\text{val}}, \mathbf{g}_{\text{DLR}}) \approx \sum_{m \in \mathcal{M}} \mathbf{q}_{\text{DLR}}^{\text{val}}[m] \cdot \mathbf{d}_{\text{DLR}}^{\text{val}}[m] \mathbf{g}_{\text{DLR}}[m], \quad (18)$$

where  $\mathcal{M} = \{m | \mathbf{q}_{\text{DLR}}^{\text{val}}[m] > \theta\}$  is the set of indices for the GIP computation, and  $\theta$  is a hyperparameter. This first-stage retrieval relies on the dimensions where  $\mathbf{q}_{\text{DLR}}^{\text{val}}[m]$  is above a threshold, as depicted in Figure 3.

Note that this approach can also be applied to DHRs, which inherit all their properties from DLRs (and thus can be applied to representations that combine lexical and semantic components). In the subsequent main experiments, we use this retrieve-and-rerank approach with approximate GIP as our first-stage retriever. See Section 5.4 for further analyses.

### 3.5 Choice of Lexical Representation Models

Although our approach can be applied to any off-the-shelf model for lexical matching, as described in Section 5.1, in practice, many of the models discussed in Section 2 still have limitations. For example, uniCOIL [28] and DeepImpact [36] require another model for passage expansion, which incurs additional costs for training and inference. Thus, in this paper, we choose SPLADE [10] as the basis of our lexical representation model, which addresses the above issue by directly learning term expansions and term weights together. However, Formal et al. [10] demonstrate that additional steps are required for tuning a good efficiency–effectiveness tradeoff using a sparsity regularization term [40] to enable retrieval with inverted indexes. To address this shortcoming, we propose a variant: playing off the name of SPLADE, which stands for **SP**arse **L**exical **AN**D **E**xpansion, we call this variant the **D**ense **L**exical **AN**D **E**xpansion (DeLADE) model. Instead of regularizing vector sparsity, we increase vector density—the exact opposite. There are two advantages of this design: (1) hyperparameter tuning is not required for model training; (2) the dense vectors are more robust to our densification approach, as we show in Section 5.1.

*SPLADE.* Following SparTerm [1], SPLADE generates lexical vectors based on the logits of the BERT pretrained masked language model (MLM). Consider an input (query or passage) wordpiece sequence  $S = (s_1, s_2, \dots, s_n)$  and its corresponding contextualized token embeddings  $H = (\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n) \in \mathbb{R}^{n \times 768}$ . The logit of each token  $s_i$  is computed as follows:

$$\mathbf{logit}_i = \text{transform}(\mathbf{h}_i)^T \cdot E_{\text{MLM}} + b_{\text{MLM}} \in \mathbb{R}^{|\mathcal{V}_{\text{BERT}}|}, \quad (19)$$

where  $\text{transform}(\cdot)$  is a linear layer with ReLU activation and LayerNorm,  $E_{\text{MLM}} \in \mathbb{R}^{|\mathcal{V}_{\text{BERT}}| \times 768}$  is the BERT embedding table, and  $b_{\text{MLM}}$  is the bias.  $\mathcal{V}_{\text{BERT}}$  is the vocabulary of BERT wordpiece tokens with size  $|\mathcal{V}_{\text{BERT}}| = 30522$ . SPLADE generates a single-vector embedding  $\mathbf{sp}_{\text{BoW}} \in \mathbb{R}^{|\mathcal{V}_{\text{BERT}}|}$  by max pooling over sequence token logits.

$$\mathbf{sp}_{\text{BoW}}[v] = \max_{i=1,2,\dots,n} \log(1 + \text{ReLU}(\mathbf{logit}_i[v])), \quad (20)$$

where  $\mathbf{sp}_{\text{BoW}}[v]$  is the vector value of the index (or vocabulary ID)  $v \in [0, |\mathcal{V}_{\text{BERT}}|)$  and  $\text{ReLU}(\cdot)$  is the activation function. This design naturally promotes vector sparsity since  $\log(1 + \text{ReLU}(\cdot))$  becomes zero for all negative input. Together with FLOP regularization loss within a mini-batch  $\mathcal{B}$ :

$$\mathcal{L}_{\text{FLOP}} = \sum_{v=1}^{|\mathcal{V}_{\text{BERT}}|} \left( \frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \mathbf{sp}_{\text{BoW}}[v] \right)^2, \quad (21)$$

the sparsity of  $\mathbf{sp}_{\text{BoW}}$  can be further increased. The overall training loss for SPLADE becomes:

$$\mathcal{L} = \text{NLL} + \lambda_q \cdot \mathcal{L}_{\text{FLOP}}^q + \lambda_d \cdot \mathcal{L}_{\text{FLOP}}^d. \quad (22)$$

Thus, SPLADE training requires tuning two hyperparameters,  $\lambda_q$  and  $\lambda_d$ , in order to obtain a good effectiveness–efficiency tradeoff. We refer readers to previous work [10, 40] for more details.

*DeLADE.* We make a slight revision to Eq. (20) by replacing the activation function  $\text{ReLU}(\cdot)$  with  $\text{softmax}(\cdot)$  to promote vector density.

$$\mathbf{ds}_{\text{BoW}}[v] = \max_{i=1,2,\dots,n} w_i \cdot \text{softmax}(\mathbf{logit}_i)[v], \quad (23)$$

where  $w_i = \mathbf{h}_i^T W + b$  is the linear transformation of  $\mathbf{h}_i$  as the term weight (i.e., capturing term importance) of token  $s_i$ ;  $\text{softmax}(\mathbf{logit}_i)$  can be interpreted as the contextualized representation of token  $s_i$ , as a probability distribution over  $\mathcal{V}_{\text{BERT}}$ . The intuition behind this design is that a sequence is represented as the max pooling of all contextualized representations, while giving more weight to important tokens. Furthermore, in contrast to ReLU, softmax is the activation function used in BERT MLM pretraining, which also ensures that the output

vector is dense. Since vector sparsity is not a concern for retrieval latency in our framework, DeLADE can focus on optimizing negative log likelihood loss without tuning additional hyperparameters to balance efficiency considerations.

## 4 EXPERIMENTAL SETUP

### 4.1 Dataset Descriptions

*In-domain IR datasets.* We use the MS MARCO passage ranking dataset introduced by Bajaj et al. [2], comprising a corpus of 8.8M web passages with the following public query sets for evaluation: (a) MS MARCO dev: 6980 queries comprise the development set for the MS MARCO passage leaderboard, with on average one relevant passage per query. Following established procedure, we report MRR@10 and R@1000 as the evaluation metrics. (b) TREC DL [5, 6]: the organizers of the 2019 (2020) Deep Learning Track at the Text REtrieval Conference (TREC) released 43 (53) queries with graded relevance labels, where (query, passage) pairs were annotated by NIST assessors. We report nDCG@10 for these two evaluation sets.

*Out-of-domain IR datasets.* We use BEIR, recently introduced by Thakur et al. [48], which contains 18 distinct IR datasets spanning diverse domains and tasks, including retrieval, question answering, fact checking, question paraphrasing, and citation prediction. Each individual dataset comprises its own corpus, queries, and relevance judgements. Following previous work [9, 46], we conduct zero-shot retrieval on 13 of the 18 datasets. Model retrieval effectiveness is evaluated in terms of nDCG@10 and R@100, except for TREC-COVID, where we follow Thakur et al. [48] and use “capped” Recall@100 instead of “regular” R@100.

### 4.2 Models

*Lexical retrieval models.* To evaluate our approach to densifying lexical representations, we conduct experiments on four lexical matching models: two models based on whole word matching, with large vocabulary sizes (millions of distinct terms), and three models that operate on the wordpiece vocabulary of BERT (30522 distinct tokens), including our proposed DeLADE model.

The whole word matching approaches in more detail:

- (1) BM25: We can characterize this classic retrieval model as generating heuristically assigned term weights.
- (2) DeepImpact [36]: Term expansion is first applied to each passage in the collection using doc2query-T5 [39]. Then, the encoder (a two-layer MLP with ReLU activation) projects token embeddings from the final layer of BERT into contextualized term weights for each input token in the query and expanded passage.

For BM25 and DeepImpact, we output term weights for each query and passage using Pyserini [29] and then randomly assign each unique term to a unique vocabulary ID to form a high-dimensional sparse vector.

We describe the wordpiece token matching approaches in more detail:

- (1) uniCOIL [28]: This model is similar to DeepImpact, but one main difference is that uniCOIL performs lexical matching on BERT wordpiece tokens during both training and retrieval.
- (2) SPLADE [9]: To be precise, we refer to the SPLADE-max model, which uses BERT pretrained MLM to project query (or passage) tokens into a  $|V_{\text{BERT}}|$ -dimensional sparse lexical representation, where  $|V_{\text{BERT}}| = 30522$  is the vocabulary size of BERT wordpiece tokens.
- (3) DeLADE: Our proposed variant of SPLADE encodes queries and passages into  $|V_{\text{BERT}}|$ -dimensional dense lexical representations.

For uniCOIL<sup>3</sup> and SPLADE,<sup>4</sup> we use checkpoints provided by the authors to generate vector representations for each query and passage. We refer to the DLRs from a base model as  $\text{model}_{\text{DLR}}(\text{dim})$ , where  $\text{dim}$  denotes the

<sup>3</sup><https://huggingface.co/castorini/unicoil-msmarco-passage>

<sup>4</sup>[https://github.com/naver/splade/tree/main/weights/splade\\_max](https://github.com/naver/splade/tree/main/weights/splade_max)

Table 3. Comparison of Vector Densification Settings for the MS MARCO Passage Corpus

Model	vocabulary size	discarded vocabulary	value vector type	index vector type
BM25	2,660,824	472	float16	uint16
DeepImpact	3,514,102	502	float16	uint16
uniCOIL/SPLADE/DeLADE	30,522	570	float16	uint8
Dense [CLS]	-	-	float16	*

\* Note that we do not have to store the index vector for [CLS] since it is a vector of all ones.

dimensionality of the DLR. For example,  $\text{uniCOIL}_{\text{DLR}}(768)$  refers to DLRs of 768 dimensions using uniCOIL as the base model.

*Densification.* We densify high-dimensional representations by first removing some dimensions to make the dimensionality a multiple of 768; then, we further divide each vector into 768 slices. For example, the 30522-dimensional representations are densified into 768-dimensional vectors by (1) discarding the first 570 unused tokens in the BERT vocabulary; (2) dividing the remaining 29952 tokens into 768 slices. Each of the slices contains 39 distinct vocabulary items; that is,  $M = 768$  and  $N = 39$ . In our main experiments, we use the stride slicing strategy, although in Section 5.1 we examine the effects of alternatives. In addition, we conduct experiments to densify into 256 and 128 dimensions, where there are 117 and 234 tokens in each slice, respectively. The “value” and “index” dense vectors (see Figure 1) are stored as `float16` and `uint8`, respectively. Note that since we randomly assign a vocabulary ID to each word for BM25 and DeepImpact, the discarded words are also randomly chosen. In addition, `uint8` is sufficient to represent the index vectors for the lexical models that rely on wordpiece token matching, while `uint16` is required for BM25 and DeepImpact. The detailed settings for densifying the vectors from different models are shown in Table 3. In addition, we also list the standard dense [CLS] vector storage setting under our framework for comparison.

*Single model fusion.* In practice, we can jointly train the [CLS] representation with any of the above lexical retrieval models using a single BERT model. However, as discussed in Section 3.5, DeepImpact and uniCOIL require an additional model for passage expansion, and thus they are not ideal base models. For SPLADE, FLOP regularization and hyperparameter tuning are required to achieve a good effectiveness–efficiency balance, also making it not an ideal base model. Instead, we choose our proposed DeLADE model as the base lexical model and refer to the single fusion model as DeLADE+[CLS]. The corresponding DHR is called  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}(\text{dim})$ , where  $\text{dim}$  denotes the dimensionality of the lexical component, DLR. We project the [CLS] vectors into 128 dimensions with a linear layer.

*Training and inference details.* Our proposed models, DeLADE and  $(\text{DeLADE}+[\text{CLS}])$ , are implemented using Tevatron [14] and trained on a single Tesla V100 GPU with 32 GB memory. We train our models using `distilbert-base-uncased` [44] for 6 epochs (around 100k steps) with learning rate  $7e-6$ . Each batch includes 24 samples, and for each query, we randomly sample one positive and seven negative passages. All the negatives are sampled from the MS MARCO “small” triples training set, which is created using BM25. We set the maximum input length for the query and the passage to 32 and 150 (including the special tokens [CLS] and [SEP]), respectively, at both training and inference stages, except for the BEIR dataset, where we set the maximum input length to 512 for both the query and the passage at inference time. For  $(\text{DeLADE}+[\text{CLS}])$ , we set  $\lambda$  to one during training and inference.

*Advanced training techniques.* To further compare with other state-of-the-art retrievers, we also train our model with knowledge distillation (KD) [16] and hard negative mining (HNM), denoted  $(\text{DeLADE}+[\text{CLS}])^+$ . Specifically, we use  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}(128)$  to retrieve the top-200 passages using the 8M queries in the training set and

then retrain (DeLADE+[CLS]) using ColBERT as the teacher model. Note that our ColBERT model (initialized from `distilbert-base-uncased`) is trained with the soft labels provided by Hofstätter et al. [17].<sup>5</sup> Following Lin et al. [32], we use KL divergence as the listwise KD loss, which considers all the in-batch negatives. For each batch, we include 288 triples (i.e., a query with positive and negative examples) by randomly sampling negatives from the top-200 hard negatives. We train (DeLADE+[CLS])<sup>+</sup> on three Tesla V100 GPUs.

### 4.3 Retrieval Implementation and Settings

Operations in our dense representation framework can be implemented by existing packages that support common array operations, which makes our approach easy to implement and to deploy in real-world settings. Specifically, our DLR and DHR retrieval experiments are performed using a custom PyTorch implementation, which means that training and retrieval experiments can be conducted within the same execution environment (on GPUs). Our experiments can be performed without an additional toolkit for nearest neighbor search such as Faiss [24].

For experimental results reported in Section 5.1, we use two-stage retrieval with the following settings,  $\theta = \{1, 1, 1, 1, 0.1\}$ , for the corresponding lexical models, {BM25, DeepImpact, uniCOIL, SPLADE, DeLADE}. For experimental results reported in Section 5.2 and Section 5.3, we use two-stage retrieval with  $\theta = 0.3$  for the independent fusion models, i.e., (BM25+ANCE)<sub>DHR</sub> and (uniCOIL+ANCE)<sub>DHR</sub>, and single model fusion, i.e., (DeLADE+[CLS])<sub>DHR</sub> and (DeLADE+[CLS])<sub>DHR</sub><sup>+</sup>. We set  $K = 10000$  in all experiments.

As already noted, our DLR and DHR retrieval experiments are conducted in PyTorch directly. For the other dense retrieval models, we use Faiss FlatIP GPU indexes as points of comparison. We perform all retrieval experiments using a single NVIDIA RTX A6000 with batch size one. For retrieval using inverted indexes, we use Pyserini (which is built on Lucene) and measure retrieval latency using a single thread on a Linux machine with two 2.1 GHz Intel Xeon Platinum 8160 CPUs and 944G of RAM. Note that in our main experiments, we primarily compare retrieval latency between different types of text representations (i.e., sparse lexical, dense semantic, and DLR/DHR). Thus, we exclude the encoding time of the query text from the neural models for simplicity. Query encoding latency on a GPU is around 10–20 ms per query for the backbones (`bert-base-uncased` or `distilbert-base-uncased`) used by the compared neural models (except for GTR<sub>xxl</sub> [37]); thus, this does not have much impact on the latency comparison between neural models. We refer readers to Table 12 for detailed online retrieval latency measurements of our models on the CPU and GPU.

## 5 RESULTS

In this section, we present our experimental results and discuss each research question in turn.

### 5.1 Quality of DLR Approximations

To begin, we densify lexical representations into DLRs of different dimensions to investigate our first research question:

**RQ1** How well do DLRs approximate the original high-dimensional lexical representations?

Table 4 shows the results of densifying different lexical representations. In these experiments, we use stride slicing, but examine the impact of different slicing strategies below. The first row in each block reports the retrieval effectiveness of the original lexical representations using inverted indexes, which can be considered an upper bound. We also report the effectiveness difference (shown as a percentage) between each method and the upper bound using inverted indexes. Note that since DeLADE is trained without any sparsity constraints, the output query and passage representations are quite dense and hence impractical for retrieval using inverted indexes. For

<sup>5</sup><https://github.com/sebastian-hofstaetter/neural-ranking-kd>

Table 4. Effectiveness/Efficiency Comparisons of DLRs with Different Base Models on MS MARCO (Dev)

Base model	Method	Dim	# tokens/doc	Quality			Storage	Latency
			number (diff.)	MRR@10 (diff.)	R@1K (diff.)	GB	ms/q	
whole word	BM25	Inv. index	2.6M	30.11	0.188	0.858	0.7	40
		SPAR [4]	768	-	0.173 (-8.0%)	0.831 (-3.1%)	26.0	64
		DLR	768	29.18 (-3.1%)	0.180 (-4.3%)	0.845 (-1.5%)	26.0	26
			256	28.41 (-5.7%)	0.177 (-5.9%)	0.834 (-2.8%)	8.6	23
			128	26.62 (-11.6%)	0.169 (-10.1%)	0.816 (-4.9%)	4.3	22
	DeepImpact	Inv. index	3.5M	71.61	0.327	0.948	1.4	285
		DLR	768	65.28 (-8.9%)	0.324 (-0.9%)	0.942 (-0.6%)	26.0	26
			256	59.90 (-16.4%)	0.316 (-3.4%)	0.933 (-1.6%)	8.6	24
			128	52.74 (-26.4%)	0.304 (-7.0%)	0.923 (-2.6%)	4.3	22
			Inv. index	30K	67.96	0.351	0.958	1.3
wordpiece token	uniCOIL	SPAR [4]	768	-	0.341 (-2.8%)	0.970 (+1.2%)	26.0	64
		DLR	768	64.15 (-5.6%)	0.349 (-0.6%)	0.957 (-0.1%)	20.0	25
			256	58.62 (-13.7%)	0.344 (-2.0%)	0.952 (-0.6%)	6.4	22
			128	52.48 (-22.8%)	0.335 (-4.6%)	0.944 (-1.5%)	3.3	22
			Inv. index	30K	91.50	0.340	0.965	2.6
	SPLADE	768	86.33 (-5.7%)	0.336 (-1.2%)	0.963 (-0.2%)	20.0	28	
		256	76.45 (-16.5%)	0.326 (-4.2%)	0.959 (-0.6%)	6.4	25	
		128	64.35 (-29.7%)	0.318 (-6.5%)	0.951 (-1.5%)	3.3	24	
		FlatIP*	30K	30522	0.347	0.957	1033.3	-
	DeLADE	768	768 (-97.5%)	0.345 (-0.6%)	0.953 (-0.4%)	20.0	25	
256		256 (-99.2%)	0.341 (-1.7%)	0.951 (-0.6%)	6.4	22		
128		128 (-99.6%)	0.335 (-3.5%)	0.945 (-1.3%)	3.3	21		

\* Index size with faiss.FlatIP is provided only as a reference; query latency is not comparable to retrieval with inverted indexes and hence omitted.

Table 5. Effectiveness Comparisons of DLRs with Different Unsupervised Compression Techniques on FiQA-2018 (Test).

Base model	Method	Dim	Quality		Storage
			nDCG@10 (diff.)	R@100 (diff.)	GB
DeLADE	FlatIP	30K	0.301	0.576	6.54
	LSH	-	0.210 (-30.2%)	0.545 (-5.4%)	0.75
	PQ768	-	0.281 (-6.6%)	0.557 (-3.3%)	0.07
	PQ256	-	0.253 (-15.9%)	0.535 (-7.1%)	0.04
	PQ128	-	0.221 (-26.6%)	0.496 (-13.9%)	0.03
	DLR	768	0.294 (-2.3%)	0.567 (-1.6%)	0.13
		256	0.287 (-4.7%)	0.558 (-3.1%)	0.04
		128	0.274 (-9.0%)	0.552 (-4.2%)	0.02

this condition, we use faiss.FlatIP brute-force search instead. Index size is provided only as a reference, and we omit query latency since it is not comparable to retrieval with inverted indexes. We also report the number of tokens per passage (i.e., vector indices with non-zero weights) for each vector densification condition.

Our method is able to densify high-dimensional lexical vectors into 768-dimensional vectors with only a small retrieval effectiveness drop. As the vectors are further densified into smaller dimensions, retrieval effectiveness drops more, which can be explained by information loss since the number of tokens are reduced through max pooling in each slice; i.e., collision between tokens appearing in the same slice. Compared to models that use whole word matching, BM25 and DeepImpact, wordpiece matching models appear to be more robust to vector

densification. For example, with 128-dimensional DLRs, uniCOIL sees only 4.6% and 1.5% degradation in MRR@10 and R@1K, respectively, while BM25 sees 10% and 4.9% degradation in MRR@10 and R@1K, respectively. In contrast to whole word matching models, wordpiece matching models represent many terms with multiple wordpiece tokens; thus, effectiveness is less sensitive to vector densification since there is more redundancy in the representation. In addition, BM25 is the only model that does not benefit from passage expansion; thus, reducing the token space causes greater effectiveness loss.

Among wordpiece matching models, uniCOIL sees only modest retrieval effectiveness degradation for 256 and 128 dimensions while SPLADE sees larger effectiveness drops. This is likely because SPLADE represents each passage with more wordpiece tokens than uniCOIL does; thus, there are more collisions as the vectors are densified into smaller dimensions. In contrast to SPLADE, although there are more collisions from vector densification with DeLADE<sub>DLR</sub>, our model sees less retrieval effectiveness degradation. We attribute this robustness to collisions to DeLADE’s full expansion over the entire wordpiece vocabulary space without any sparsity regularization; that is, there appears to be more “redundancy” in the representations. This result indicates that DeLADE is a better alternative to SPLADE for dense lexical matching.

Finally, a comparison of DLR and SPAR [4] shows the advantages of our approach. Without any training, our 256-dimensional DLRs applied to BM25 and uniCOIL are able to compete with SPAR. Furthermore, our approach consumes less space and exhibits lower retrieval latency. Note that the index vectors for the whole word matching models (BM25 and DeepImpact) are stored in `uint16` due to larger vocabulary sizes; thus, they consume more storage compared to the wordpiece matching models.

In addition, we observe that retrieval latency using inverted indexes is sensitive to the average number of tokens per passage (i.e., vector sparsity). For example, a 35% increase from 68 to 92 (uniCOIL vs. SPLADE) leads to more than 60% latency increase. In contrast, DLRs see lower retrieval latency, which appears to be insensitive to the average number of tokens per passage; for example, uniCOIL<sub>DLR</sub> and DeLADE<sub>DLR</sub> exhibit large differences in vector sparsity but have comparable latency when using vectors with the same dimensions. Even with different dimensions, our approach does not exhibit much variability in retrieval latency. We attribute this advantage to our two-stage retrieval approach, where computationally expensive end-to-end retrieval only relies on a few dimensions (see Section 5.4 for more details).

To compare our approach to the unsupervised vector compression techniques shown in Table 2, we use the FiQA-2018 test collection in BEIR [48], which is based on a medium-size corpus with 57K passages. Compared to MS MARCO (8.8M passages), it is less computationally demanding to perform unsupervised training for the various vector compression techniques. We use DeLADE as the base model and report results (nDCG@10 and R@100 as quality metrics and storage as the efficiency metric) in Table 5. Our experiments follow Faiss instructions<sup>6</sup> for applying the unsupervised vector compression techniques. The FlatIP index (without any compression) provides the performance upper bound. For locality-sensitivity hashing (LSH), we use `faiss.IndexLSH` with `nbits = 8 · 768`; for product quantization (PQ), we use `faiss.IndexPQ` with `nbits = 8` and `M = 768, 256, 128`.<sup>7</sup> We observe that PQ768 can retain most of the information from the original high-dimensional vectors (less than 10% retrieval effectiveness drop) while LSH cannot effectively compress even a 30K-dimensional vector into binary codes. When we further compress the original vectors with PQ256 and PQ128, retrieval effectiveness drops more than 10% (and even more for PQ128). In contrast, DLR shows retrieval effectiveness drops less than 10% and performs consistently better than product quantization with the same storage size. Furthermore, our approach does not require any supervised or unsupervised training.

<sup>6</sup><https://github.com/facebookresearch/faiss/wiki/Faiss-indexes>

<sup>7</sup>We also tried OPQ [15] training before performing PQ; however, the effectiveness is much worse.

Table 6. Comparison of Different DLR Slicing Strategies on MS MARCO (Dev).

Strategy	uniCOIL <sub>D<sub>DLR</sub></sub> (768)			SPLADE <sub>D<sub>DLR</sub></sub> (768)			DeLAD <sub>D<sub>DLR</sub></sub> (768)		
	# tokens/doc	MRR@10	R@1K	# tokens/doc	MRR@10	R@1K	# tokens/doc	MRR@10	R@1K
Contiguous	39.27	0.334	0.947	68.49	0.326	0.952	768	0.332	0.949
Stride	64.15	0.349	0.957	86.33	0.336	0.963	768	0.345	0.953
Random	64.32	0.349	0.957	87.91	0.336	0.963	768	0.345	0.953

Finally, we examine three different slicing strategies on the wordpiece matching models.<sup>8</sup> We densify lexical representations with 768 dimensions (i.e.,  $M = 768$ ). Table 6 reports results for the three lexical retrieval models, along with the number of tokens per passage for each condition. We notice that stride slicing has the same effectiveness as randomized slicing (only minor differences observed beyond four digits). However, surprisingly, the contiguous slicing strategy shows degradation in ranking effectiveness and the number of tokens per passage for this condition is smaller than the other slicing strategies. This indicates that BERT wordpiece tokens with adjacent token IDs may co-occur with higher probability than any two randomly chosen tokens. Thus, max pooling over contiguous slices of BERT wordpiece token IDs leads to more collisions compared to the other strategies. Based on this analysis, we use stride slicing as our default setting in the rest of our experiments.

## 5.2 Evaluation of Independent Model Fusion

In this section, we describe experiments on fusing different “off-the-shelf” lexical and semantic retrieval models and compare their effectiveness and efficiency to other hybrid retrieval methods to answer the following research question:

**RQ2** How well do DHRs benefit from the independent fusion of DLRs and “off-the-shelf” dense semantic representations?

Following Chen et al. [4], we conduct experiments on BM25–ANCE and uniCOIL–ANCE fusion on the MS MARCO dev set, reported in rows (1)–(3) and rows (4)–(6), respectively, in Table 7. The first entry in each main block, rows (1) and (4), represents the linear combination approach used in previous work [26, 28, 32], which requires two separate indexes and additional post-processing of the ranked lists. For these experiments, we measure latency in Pyserini and report a theoretically optimized system as the maximum latency between Lucene and Faiss retrieval plus 3 ms of post-processing time. We leave aside the engineering challenge of synchronizing CPU and GPU search necessary to achieve this performance under real-world conditions. For SPAR [4], rows (2) and (5), we directly report the ranking effectiveness from the paper and measure the retrieval latency of the Faiss FlatIP index in our environment. Note that their approach distills uniCOIL’s lexical representations into semantic representations and then concatenates them to ANCE for dense retrieval. Thus, dimensionality of the representation vectors is reported as  $2 \times 768$ . For DHRs, rows (3) and (6), we tune  $\lambda$  on MRR@10 using a subset of 100 queries in the training set and set  $\theta = 0.3$  for two-stage retrieval.

Overall, all three systems yield similar retrieval effectiveness but our DHRs achieve lower retrieval latency. Note that our model variant with 768-dimensional lexical and semantic vectors is faster than SPAR using Faiss GPU. We attribute this improvement to our two-stage retrieval approach (see Section 5.4 for more details). It is worth noting that with a negligible effectiveness drop we can further compress the lexical representations to 128 dimensions. Furthermore, we can convert lexical representations into DHRs of any width according to user design requirements *without any model retraining*. This flexibility is one major advantage of our approach.

<sup>8</sup>Recall that the vocabulary IDs for the whole word matching models are randomly assigned; thus, the slicing strategy is considered random.

Table 7. Effectiveness/Efficiency Comparisons of Independent Fusion of DLRs with “Off-the-Shelf” Dense Semantic Representations on MS MARCO (Dev)

Approach	Lexical		Semantic		quality		storage	latency
	Software	Dim	Software	Dim	MRR@10 (diff.)	R@1K (diff.)	(GB)	(ms/q)
(1) Linear combination (BM25, ANCE)	Lucene	30K	Faiss FlatIP	768	0.347	0.969	26	64
(2) SPAR [4] (lexical conversion to semantic)	n/a	n/a	Faiss FlatIP	2× 768	0.344 (−0.9%)	0.971 (+0.2%)	52	81
(3) (BM25+ANCE) <sub>DHR</sub>	PyTorch	768	PyTorch	768	0.349 (−0.6%)	0.967 (−0.2%)	39	56
	PyTorch	256	PyTorch	768	0.348 (+0.3%)	0.967 (−0.2%)	21	56
	PyTorch	128	PyTorch	768	0.347 (−0.0%)	0.967 (−0.2%)	17	53
(4) Linear combination (uniCOIL, ANCE)	Lucene	30K	Faiss FlatIP	768	0.375	0.976	27	291
(5) SPAR [4] (lexical conversion to semantic)	n/a	n/a	Faiss FlatIP	2× 768	0.369 (−1.6%)	0.981 (+0.5%)	52	81
(6) (uniCOIL+ANCE) <sub>DHR</sub>	PyTorch	768	PyTorch	768	0.378 (+0.8%)	0.975 (−0.1%)	32	60
	PyTorch	256	PyTorch	768	0.375 (−0.0%)	0.973 (−0.3%)	19	58
	PyTorch	128	PyTorch	768	0.369 (−1.6%)	0.971 (−0.5%)	16	57

Retrieval with DHRs uses our custom PyTorch implementation running on GPUs.

### 5.3 Evaluation of Single Model Fusion

With our framework and proposed DeLADE model, fusing lexical and semantic representations becomes easier. This motivates us to investigate:

**RQ3** Can DHRs benefit from *joint* training of lexical and semantic components in a single model?

Table 8 compares model performance in terms of retrieval effectiveness and efficiency. For efficiency, index size and retrieval latency are measured on the MS MARCO dev set as the point of reference. The comparison models across the columns are categorized as: (1) sparse lexical retrievers, including BM25, docT5q [39], and SPLADE [9]; (2) dense semantic retrievers, including our trained baseline dense retriever (denoted Dense with 768-dimensional [CLS] vectors) and ANCE [50]; (3) multi-vector retrievers, including ColBERT [27] and COIL [12].

For the MS MARCO datasets, we conduct experiments using Pyserini [29] for all models except for ColBERT<sup>9</sup> and COIL.<sup>10</sup> For the BEIR datasets, we directly copy numbers from Izacard et al. [22], except for SPLADE<sup>11</sup> and COIL.<sup>10</sup> To determine the statistical significance of our results, we perform paired *t*-tests ( $p < 0.05$ ) comparing all models except for ColBERT on the MS MARCO datasets. For a fair comparison, Table 8 only includes models that use the same baseline training strategy as ours. Thus, we exclude approaches that depend on other models for expansion [28, 36, 56], costly training techniques such as knowledge distillation [9, 18, 19, 41, 46, 49], or special pretraining [11, 22, 37], although see Table 11 for additional comparisons.

In terms of our proposed models, we report results on three (DeLADE+[CLS])<sub>DHR</sub> variants by densifying the lexical components into 128, 256, and 768 dimensions. The (DeLADE+[CLS])<sub>DHR</sub> 128- and 256-dimensional variants can be considered “small vectors”, to compare with single-vector (sparse lexical and dense semantic) retrievers, while the (DeLADE+[CLS])<sub>DHR</sub> 768-dimensional variant can be considered “large vectors”, to compare with multi-vector retrievers. Note that all three variants are derived from the same model. Finally, we report the performance of DeLADE<sub>DLR</sub>(768), i.e., without the incorporation of the [CLS] vector.

*DHRs vs DLRs.* We first compare (DeLADE+[CLS])<sub>DHR</sub>(768) with DeLADE<sub>DLR</sub>(768), columns (8) and (b), to examine the effectiveness of single model fusion. From the results, we see that (DeLADE+[CLS])<sub>DHR</sub>(768), which

<sup>9</sup>We copy numbers from Mallia et al. [36].

<sup>10</sup>We run COIL using the inference code from the authors’ repo at <https://github.com/luyug/COIL>.

<sup>11</sup>We run SPLADE-max using the inference code from the authors’ repo at <https://github.com/naver/splade>.

Table 8. Effectiveness/Efficiency Comparisons of DHRs Using Single Model Fusion

	sparse lexical			dense semantic		multi-vector		DeLADE <sub>DLR</sub>	(DeLADE+[CLS]) <sub>DHR</sub>			
	(1) BM25	(2) docT5q	(3) SPLADE	(4) Dense	(5) ANCE	(6) COIL	(7) ColBERT	(8) 768 dim	(9) 128 dim	(a) 256 dim	(b) 768 dim	
<b>Efficiency*</b>												
storage (GB)	0.67	0.98	2.6	26	26	60	154	20	5.4	8.6	22	
latency (ms/q)	40	64	475	64	64	40*	458*	30	28	31	33	
<b>MS MARCO</b>												
Dev	MRR@10	0.188	0.277 <sup>1</sup>	0.340 <sup>1245</sup>	0.307 <sup>12</sup>	0.330 <sup>124</sup>	0.354 <sup>1-58</sup>	<b>0.360</b> <sup>△</sup>	0.345 <sup>1245</sup>	0.351 <sup>1-58</sup>	0.355 <sup>1-589</sup>	0.357 <sup>1-58-a</sup>
	R@1K	0.858	0.947 <sup>1</sup>	0.965 <sup>12458</sup>	0.944 <sup>1</sup>	0.959 <sup>124</sup>	0.964 <sup>12458</sup>	<b>0.968</b> <sup>△</sup>	0.953 <sup>124</sup>	0.962 <sup>1248</sup>	0.965 <sup>124589</sup>	0.967 <sup>12458-a</sup>
DL 19	nDCG@10	0.506	0.642 <sup>1</sup>	0.683 <sup>14</sup>	0.631 <sup>1</sup>	0.646 <sup>1</sup>	<b>0.714</b> <sup>1245</sup>	0.694 <sup>△</sup>	0.691 <sup>14</sup>	0.691 <sup>145</sup>	0.696 <sup>145</sup>	0.693 <sup>145</sup>
DL 20		0.475	0.619 <sup>1</sup>	0.671 <sup>1</sup>	0.648 <sup>1</sup>	0.646 <sup>1</sup>	<b>0.688</b> <sup>125</sup>	0.676 <sup>△</sup>	0.668 <sup>12</sup>	0.683 <sup>12</sup>	0.686 <sup>12</sup>	<b>0.688</b> <sup>12</sup>
<b>BEIR</b>												
nDCG@10												
TREC-COVID	0.656	0.713	0.661	0.604	0.654	0.668	0.677	0.681	0.695	0.701	<b>0.727</b>	
NFCorpus	0.325	0.328	0.322	0.244	0.237	<b>0.331</b>	0.305	<u>0.331</u>	0.319	0.324	0.329	
NQ	0.329	0.399	0.469	0.410	0.446	0.519	<b>0.524</b>	0.471	0.476	0.487	0.497	
HotpotQA	0.603	0.580	0.640	0.441	0.456	<b>0.713</b>	0.593	0.666	0.664	0.679	<u>0.689</u>	
FiQA-2018	0.236	0.291	0.289	0.224	0.295	<u>0.313</u>	<b>0.317</b>	0.294	0.292	0.304	0.312	
ArguAna	0.315	0.349	<b>0.445</b>	0.323	0.415	0.295	0.233	0.360	<u>0.423</u>	0.405	0.360	
Touche-2020 (v2)	<b>0.367</b>	<u>0.347</u>	0.201	0.185	0.240	0.281	0.202	0.266	0.248	0.259	0.280	
Quora	0.789	0.802	0.834	0.750	<u>0.852</u>	0.838	<b>0.854</b>	0.755	0.829	0.830	0.831	
DBPedia	0.313	0.331	0.370	0.295	0.281	0.398	0.392	0.376	0.397	<u>0.402</u>	<b>0.404</b>	
Scidocs	0.158	<u>0.162</u>	0.149	0.103	0.122	0.155	<b>0.165</b>	0.148	0.146	0.148	0.150	
Fever	0.753	0.714	0.740	0.651	0.669	<b>0.840</b>	0.771	0.787	0.750	0.781	<u>0.810</u>	
Climate-Fever	0.213	0.201	0.187	0.167	0.198	0.216	0.184	0.202	0.215	<u>0.220</u>	<b>0.229</b>	
SciFact	0.665	0.675	0.633	0.479	0.507	<b>0.707</b>	0.671	0.674	0.651	0.670	<u>0.685</u>	
Avg. nDCG@10	0.440	0.453	0.458	0.375	0.413	<u>0.483</u>	0.453	0.462	0.470	0.478	<b>0.485</b>	
Avg. rank	7.00	6.07	7.00	10.46	8.15	<u>3.15</u>	5.54	5.39	5.92	4.46	<b>2.85</b>	
<b>BEIR</b>												
R@100												
TREC-COVID	0.498	<b>0.541</b>	0.502	0.386	0.457	0.531	0.464	0.502	0.489	<u>0.512</u>	<b>0.541</b>	
NFCorpus	0.250	0.253	0.266	0.228	0.232	<b>0.277</b>	0.254	0.260	0.268	0.272	<u>0.273</u>	
NQ	0.760	0.832	0.883	0.817	0.836	<b>0.917</b>	0.912	0.876	0.894	0.901	<u>0.909</u>	
HotpotQA	0.740	0.709	0.793	0.586	0.578	<b>0.837</b>	0.748	0.801	0.797	0.808	<u>0.818</u>	
FiQA-2018	0.539	0.598	0.576	0.500	0.581	<u>0.596</u>	<b>0.603</b>	0.567	0.584	0.592	0.591	
ArguAna	0.942	<b>0.972</b>	0.954	0.873	0.937	0.777	0.914	0.893	<u>0.957</u>	0.943	0.900	
Touché-2020 (v2)	<u>0.538</u>	<b>0.557</b>	0.450	0.409	0.458	0.512	0.439	0.499	0.462	0.480	0.500	
Quora	0.973	0.982	0.984	0.963	0.987	<b>0.998</b>	<u>0.989</u>	0.972	0.984	0.984	0.985	
DBPedia	0.398	0.365	0.493	0.353	0.319	<b>0.522</b>	0.461	0.495	0.491	0.506	<u>0.510</u>	
SCIDOCs	<u>0.356</u>	<b>0.360</b>	0.351	0.247	0.269	0.354	0.344	0.343	0.338	0.342	0.346	
FEVER	0.931	0.916	0.934	0.886	0.900	<b>0.959</b>	0.934	0.948	0.948	0.953	<u>0.956</u>	
Climate-FEVER	0.436	0.427	0.452	0.401	0.445	<u>0.508</u>	0.444	0.464	0.490	0.506	<b>0.525</b>	
SciFact	0.908	0.914	0.899	0.829	0.816	<b>0.931</b>	0.878	<u>0.920</u>	0.898	0.904	0.910	
Avg. R@100	0.636	0.648	0.657	0.575	0.601	<u>0.671</u>	0.645	0.657	0.661	0.669	<b>0.674</b>	
Avg. rank	7.15	5.54	6.08	10.62	8.62	<b>2.46</b>	6.23	5.92	5.77	4.31	<u>3.31</u>	

\* We report efficiency figures on the MS MARCO passage corpus for comparison.

\* These numbers are copied from the original papers, which are measured on multi-GPU systems; thus, they cannot be reproduced in our setup.

△ We do not have these run files; thus, no significance testing is performed against ColBERT.

This table compares against a selection of retrieval models trained with comparable baseline strategies. Bold (underline) denotes the best (second best) effectiveness for each row. For the MS MARCO datasets, superscripts denote significant improvements over the labeled model with paired  $t$ -test ( $p < 0.05$ ).

incorporates an additional 128-dimensional [CLS] vector, demonstrates significantly better retrieval effectiveness than DeLADE<sub>DLR</sub> (768) for both in-domain and out-of-domain datasets. Furthermore, (DeLADE+[CLS])<sub>DHR</sub> (256)

Table 9. Component Retrieval Effectiveness of (DeLADE+[CLS])<sub>DHR</sub>(768)

Component*	MS MARCO dev		TREC-COVID		FiQA-2018		SciFact	
	MRR@10	R@1K	NDCG@10	CapR@100	NDCG@10	R@100	NDCG@10	R@100
(1) DeLADE + [CLS]	0.357	0.967	0.727	0.541	0.312	0.591	0.685	0.910
(2) DeLADE	0.294	0.930	0.658	0.501	0.193	0.441	0.684	0.920
(3) [CLS]	0.045	0.494	0.125	0.091	0.037	0.185	0.118	0.454

\* Row (1) corresponds to (DeLADE+[CLS])<sub>DHR</sub>(768) in Table 8. Note that rows (2) and (3) are different from the Dense and (DeLADE+[CLS])<sub>DHR</sub>(768) conditions in Table 8, where the separate models are trained independently.

Table 10. Ablation of (DeLADE+[CLS])<sub>DHR</sub>(768) Varying the Semantic Component

[CLS] dimension	MS MARCO dev		TREC-COVID		FiQA-2018		SciFact	
	MRR@10	R@1K	NDCG@10	CapR@100	NDCG@10	R@100	NDCG@10	R@100
(1) 0*	0.345	0.953	0.681	0.502	0.294	0.567	0.674	<b>0.920</b>
(2) 128*	0.357 <sup>1</sup>	0.967 <sup>1</sup>	<b>0.727<sup>1</sup></b>	<b>0.541<sup>1</sup></b>	0.312 <sup>1</sup>	0.591 <sup>1</sup>	<b>0.685</b>	0.910
(3) 256	0.358 <sup>1</sup>	0.969 <sup>1</sup>	0.722 <sup>1</sup>	0.530 <sup>1</sup>	0.313 <sup>1</sup>	0.590 <sup>1</sup>	0.683	0.917
(4) 768	0.358 <sup>1</sup>	0.969 <sup>1</sup>	0.717	0.540 <sup>1</sup>	<b>0.318<sup>1</sup></b>	<b>0.601<sup>1</sup></b>	0.680	0.918

\* The 0 and 128 variants correspond to DeLADE<sub>DLR</sub>(768) and (DeLADE+[CLS])<sub>DHR</sub>(768) in Table 8, respectively.

Bold denotes the best effectiveness for each column. Superscript denotes significant improvement over the labeled model based on paired  $t$ -tests ( $p < 0.05$ ).

outperforms DeLADE<sub>DLR</sub>(768) for all metrics, which suggests that incorporating the [CLS] vector can mitigate information loss from our densification approach.

To further understand how joint training works for (DeLADE+[CLS])<sub>DHR</sub>(768), we conduct retrieval using the densified 768-dimensional lexical vectors from DeLADE and 128-dimensional [CLS] semantic vectors separately. Retrieval effectiveness on the MS MARCO dev set and three BEIR datasets (TREC-COVID, FiQA-2018, and SciFact) is reported in Table 9. We see that each component is far less effective individually than their hybrid fusion. Specifically, the DeLADE and [CLS] components in (DeLADE+[CLS])<sub>DHR</sub>(768) are less effective than the independently trained DeLADE<sub>DLR</sub>(768) and Dense models, respectively; see columns (4) and (8) in Table 8. This result indicates that joint training yields components that are highly complementary, as designed.

To further explore this complementarity of representation, we vary the dimensionality of the semantic component of (DeLADE+[CLS])<sub>DHR</sub>(768). We jointly train DeLADE<sub>DLR</sub>(768) with 0, 128, 256, and 768-dimensional [CLS] vectors separately. Retrieval effectiveness on the MS MARCO dev set and three BEIR datasets (TREC-COVID, FiQA-2018, and SciFact) is reported in Table 10. We see that fusing a small [CLS] vector (e.g., 128 dimensions) in training improves retrieval effectiveness, indicating that joint training yields complementary lexical and semantic representations within a single model. This is consistent with the component analysis above. However, further increasing the dimensionality of the [CLS] vector does not appear to yield obvious advantages. This result indicates that a relatively small [CLS] vector is sufficient to complement the lexical component.

*DHRs vs single-vector models.* First, a comparison of single-vector retrievers shows that sparse lexical representations have better storage efficiency than dense semantic representations. For example, docT5q, column (2), with an index less than 1 GB, exhibits better out-of-domain retrieval effectiveness than dense semantic retrievers, columns (4) and (5). Similarly, SPLADE, column (3), outperforms the dense retrievers for both in-domain and out-of-domain conditions with an index of only 2.6 GB. However, SPLADE requires over 7 times the retrieval latency of docT5q and other dense retrievers. As lexical-semantic hybrid representations, DHRs exhibit better

retrieval effectiveness with modest retrieval latency and index storage consumption. We see that the 128- and 256-dimensional variants of  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}$ , columns (9) and (a), outperform the single-vector retrievers, columns (1) to (5), for both in-domain and out-of-domain datasets in terms of retrieval effectiveness, and furthermore achieves lower query latency. In addition, the two  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}$  variants have modest index sizes; for example, the Dense model, column (4), requires 26 GB to store the MS MARCO passage corpus (using a Faiss FlatIP index) while  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}(128)$  only consumes 5.4 GB.

*DHRs vs multi-vector models.* We compare retrieval models with large vectors,  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}(768)$  and multi-vector retrieval models, shown in columns (b), (6), and (7). Although their in-domain retrieval effectiveness does not appear to be very different, ColBERT falls behind COIL and  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}(768)$  in out-of-domain evaluation. This result suggests that lexical matching remains a key component for generalization. Compared to COIL,  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}(768)$  exhibits comparable retrieval effectiveness but has a much smaller index (60 GB vs 22 GB).

It is worth mentioning that COIL has other variants with smaller indexes; for example, the configuration with 8-dimensional tokens plus 128-dimensional [CLS] embeddings consumes 14 GB and yields 0.347 (0.956)  $\text{MRR}@10$  ( $\text{R}@1\text{K}$ ) on the MS MARCO dev queries.<sup>12</sup> This variant is still slightly less effective than our models,  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}(256)$  and  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}(128)$ , with index sizes of 8.6 GB and 5.4 GB, respectively. Furthermore, all our variants are derived from the same model without retraining. This comparison demonstrates the advantages of single model fusion under our proposed dense representation framework.

*Summary.* We observe that dense retrievers generally perform well in domain, while sparse retrievers appear to yield stronger generalization capabilities. Our results demonstrate that DHRs inherit advantages of both lexical and semantic matching. Specifically, our  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}$  model achieves competitive retrieval effectiveness in both evaluation settings with low retrieval latency and modest index storage consumption. This advantageous effectiveness–efficiency tradeoff makes the design of single model fusion with DHRs attractive. It is also worth noting that  $\text{DeLADE}_{\text{DLR}}(768)$  outperforms SPLADE slightly,<sup>13</sup> especially in term of nDCG. This result shows that DeLADE is a good alternative under our framework compared to SPLADE using inverted indexes in the scenario where lower retrieval latency is more important than index size.

*DHRs vs more advanced retrieval models.* Finally, we compare  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}^+$ , which uses advanced training strategies, to the effectiveness and efficiency of existing state-of-the-art retrievers in Table 11. For sparse lexical retrieval models, we report results from SPLADEv2 [9], which uses two rounds (i.e., training with BM25 and hard negatives) of knowledge distillation from a cross-encoder teacher.

For dense semantic retrieval models, we include three representative models for comparison: (1) TAS-B [18] distills knowledge from multiple teachers using a topic-aware negative sampling strategy. (2) RocketQAv2 [43] leverages a student–teacher joint training approach to make the dense retriever better mimic a cross-encoder teacher. (3) Contriever [22] leverages pretraining by combining advanced contrastive learning techniques with an Inverse Cloze Task (ICT) variant. We also include dense retrievers trained with much larger backbone models (i.e., more model parameters): (1) GPL [49] trains an expert model for each target dataset in BEIR. (2) GTR [37] trains even larger encoder models; the authors’ T5-3B and T5-11B models correspond to  $\text{GTR}_{\text{XL}}$  and  $\text{GTR}_{\text{XXL}}$ , which also combine pretraining, knowledge distillation (KD), and hard negative mining (HNM).<sup>14</sup> Finally, for multi-vector retrieval models, we report ColBERTv2 [46], which also combines KD and HNM, and at the same time compresses the multi-vectors into a smaller index.

<sup>12</sup>We refer readers to Gao et al. [12] for more details.

<sup>13</sup>Our reproduced numbers are slightly better than Formal et al. [9] except for Tóuche-2020, where we use v2 instead of v1.

<sup>14</sup>GTR is fine-tuned on MS MARCO training queries with hard negatives denoised by a cross-encoder.

Table 11. Effectiveness/Efficiency Comparisons with Existing State-of-the-Art Models Using Advanced Training Techniques

	sparse lexical	dense semantic						multi-vector	(DeLADE+[CLS]) <sub>DHR</sub> <sup>+</sup>			
	(c) SPLADEv2	(d) TAS-B	(e) Contriever	(f) RocketQAv2	(g) GPL	(h) GTR <sub>XL</sub>	(i) GTR <sub>XXL</sub>	(j) ColBERTv2	(k) 128 dim	(l) 256 dim	(m) 768 dim	
model size	66M	66M	110M	110M	66M×13*	1.24B	4.8B	110M	66M			
IR pretrain	✗	✗	✓	✗	✗	✓	✓	✗	✗			
KD	✓	✓	✗	✓	✓	✓	✓	✓	✓			
HNM	✓	✗	✓	✓	✓	✓	✓	✓	✓			
batch size >1K	✗	✗	✓	✓	✗	✓	✓	✗	✗			
<b>Efficiency*</b>												
storage (GBs)	5.0	26	26	26	26	26	26	29	5.4	8.6	22	
latency (ms/q)	2864	64	64	64	64	64	64	260	28	31	33	
<b>MS MARCO</b>												
Dev	MRR@10	0.368	0.347	0.341	0.381	-	0.385	<u>0.388</u>	<b>0.397</b>	0.366	0.370	0.371
	R@1K	<u>0.979</u>	0.978	0.979	0.981	-	<u>0.989</u>	<b>0.990</b>	0.984	0.973	0.975	0.977
DL 19	nDCG@10	<b>0.729</b>	<u>0.717</u>	0.678	-	-	-	-	0.703	0.711	0.708	
DL 20	nDCG@10	<b>0.710</b>	0.686	0.661	-	-	-	-	0.684	0.696	<u>0.700</u>	
<b>BEIR</b>												
		nDCG@10										
TREC-COVID	0.710	0.481	0.596	0.675	0.700	0.584	0.501	<b>0.738</b>	0.686	0.702	<u>0.735</u>	
NFCorpus	0.334	0.319	0.328	0.293	<b>0.345</b>	<u>0.343</u>	0.342	0.338	0.327	0.332	0.337	
NQ	0.521	0.463	0.498	0.505	0.483	0.559	<b>0.568</b>	<u>0.562</u>	0.497	0.512	0.523	
HotpotQA	<b>0.684</b>	0.584	0.638	0.533	0.582	0.591	0.599	0.667	0.659	<u>0.673</u>	<b>0.684</b>	
FiQA-2018	0.336	0.300	0.329	0.302	0.344	<u>0.444</u>	<b>0.467</b>	0.356	0.320	0.326	0.335	
ArguAna	0.479	0.429	0.446	0.451	<b>0.557</b>	0.531	<u>0.540</u>	0.463	0.475	0.458	0.436	
Touche-2020 (v2)	<b>0.272</b>	0.162	0.230	0.247	0.255	0.230	0.256	<u>0.263</u>	0.226	0.237	0.254	
Quora	0.838	0.835	0.865	0.749	0.836	<u>0.890</u>	<b>0.892</b>	0.852	0.846	0.848	0.849	
DBPedia	<u>0.435</u>	0.384	0.413	0.356	0.384	0.396	0.408	<b>0.446</b>	0.402	0.409	0.413	
Scidocs	0.158	0.149	<u>0.165</u>	0.131	<b>0.169</b>	0.159	0.161	0.154	0.156	0.158	0.158	
Fever	0.786	0.700	0.758	0.676	0.759	0.717	0.740	0.785	0.764	<u>0.794</u>	<b>0.815</b>	
Climate-Fever	0.235	0.228	0.237	0.180	0.235	<b>0.270</b>	<u>0.267</u>	0.176	0.222	<u>0.232</u>	0.239	
SciFact	0.693	0.643	0.677	0.568	0.674	0.635	0.662	0.693	0.674	<u>0.694</u>	<b>0.699</b>	
Avg.nDCG@10	<b>0.499</b>	0.437	0.475	0.436	0.486	0.488	0.493	<b>0.499</b>	0.481	0.490	<u>0.498</u>	
Avg. rank	<u>4.15</u>	10.00	6.00	9.62	5.46	5.31	4.23	<u>4.15</u>	7.31	5.38	<b>4.00</b>	

\* Wang et al. [49] trained one expert model for each BEIR dataset using pseudo-relevant labels from a cross-encoder model.

\* We report efficiency figures on the MS MARCO passage corpus for comparison. As detailed in Mackenzie et al. [35], SPLADEv2 is much slower than the other models with Lucene; even with PISA, a much faster query evaluation implementation, latency is still 220 ms/q. As for ColBERTv2, the number is copied from Vanilla ColBERTv2 (p=4, c=2<sup>16</sup>) measured by Santhanam et al. [45].

Bold (underline) denotes the best (second best) effectiveness for each row.

We point out that it is not easy to fairly compare models with more advanced and costly training techniques since, as shown in Table 11, there are many substantive differences that cannot be captured by general descriptive labels. Even considering a general method such as knowledge distillation, there are many different implementations. For example, we use a lightweight ColBERT teacher, while GPL, ColBERTv2, and SPLADEv2 use a more expensive cross-encoder teacher. While end-to-end results are comparable since the evaluations use the same test collections, it is difficult to attribute effectiveness differences to specific components.

Nevertheless, it is possible to draw some conclusions from these experiments. We first observe that dense semantic vectors from smaller models, in columns (d)–(f), do not appear to perform well in both in-domain and out-of-domain evaluations. For example, Contriever performs well on BEIR but lags behind most models on in-domain evaluation and the reverse trend can be observed for RocketQAv2. To improve the generalization capability of dense retrieval models, existing work either leverages multiple expert models for multi-domain datasets as GPL, column (g), or larger pre-trained models as GTR, columns (h)–(i).

On the other hand, SPLADEv2 and ColBERTv2, columns (c) and (j), leverage representations with more expressive capacity than dense semantic vectors and appear to show equally good generalization capability without adding model parameters. However, these models sacrifice retrieval efficiency to gain this generalization capability. For example, both SPLADEv2 and ColBERTv2 are slower than dense semantic models. GTR also sacrifices query encoding latency, which is excluded in our latency measurement.<sup>15</sup>

In contrast, (DeLADe+[CLS])<sub>DHR</sub><sup>+</sup>(768) not only yields competitive zero-shot retrieval effectiveness on BEIR compared to GTR<sub>XXL</sub>, ColBERTv2, and SPLADEv2, but also maintains low retrieval latency. In addition, our (DeLADe+[CLS])<sub>DHR</sub><sup>+</sup>(256), column (l), still yields better overall performance (i.e., effectiveness and efficiency) compared to the single-vector dense retrievers, TAS-B, GPL, and Contriever. We notice that GTR performs particularly well in some QA datasets (e.g., NQ and FiQA-2018), likely because GTR includes more training data such as QA pairs mined from the web and the NaturalQuestions dataset. We also note that (DeLADe+[CLS])<sub>DHR</sub> cannot compete with RocketQAv2, GTR, and ColBERTv2 on in-domain evaluation since the first two exploit additional QA training pairs and all of them use more expensive cross-encoder teachers, which is orthogonal to our work. It is likely that we can further boost (DeLADe+[CLS])<sub>DHR</sub> effectiveness by leveraging more data and more expensive teachers; however, we leave these explorations for future work.

#### 5.4 Performance of Two-Stage Retrieval

In this section, we further study our approach to end-to-end retrieval with DLRs/DHRs proposed in Section 3.4 to investigate our final research question:

**RQ4** How effective is our proposed two-stage retrieval approach?

To illustrate how well our proposed approximate GIP operation compares to the more expensive exact GIP operation between DLRs (or DHRs), we use DeLADe<sub>DLR</sub>(768) and (DeLADe+[CLS])<sub>DHR</sub>(768) to conduct end-to-end retrieval experiments on the MS MARCO dev queries.

Figure 4 illustrates recall at different cutoffs using approximate GIP with various settings of the parameter  $\theta$ . Exact GIP, which is equivalent to  $\theta = 0$ , is shown as the black dashed line and represents the upper bound; as  $\theta$  increases, recall drops, as expected. However, the results show that the top-10000 candidates retrieved using approximate GIP include more relevant passages (i.e., has higher recall) than GIP at cutoff 1000. That is, the effectiveness drop from approximate (first-stage) retrieval can be remedied by reranking the top-10000 candidates with GIP. In addition, we note that (DeLADe+[CLS])<sub>DHR</sub>(768) suffers a smaller recall drop at the approximate retrieval stage. This suggests that [CLS] vectors can help capture relevance and compensate for lower recall from the lexical component with larger values of  $\theta$ . To be clear,  $\theta$  is also applied to the [CLS] vector to guide the approximate computations in (DeLADe+[CLS])<sub>DHR</sub>. In addition, we also perform retrieval by computing the standard inner product (IP) between two DLRs' or DHRs' value vectors (without considering their index vectors) for comparison. We observe that IP, surprisingly, also retrieves more relevant passages at cutoff 10000 than GIP at cutoff 1000, especially for DHRs. This indicates that standard inner product between two DLRs' (or DHRs') value vectors can be an alternative to approximate GIP. We show below that using IP as first-stage retrieval can benefit DLRs (or DHRs) in end-to-end retrieval on the CPU.

We further examine effectiveness–efficiency tradeoffs after reranking the top-10000 candidates with exact GIP. The red lines in Figures 5(a) and 5(b) plot the DeLADe<sub>DLR</sub>(768) performance tradeoff curve (ranking effectiveness vs retrieval latency) of approximate GIP and reranking with different values of  $\theta$ , while the blue dashed line depicts the performance curve of approximate GIP without reranking for comparison. Retrieval effectiveness with exact GIP (the black triangle) can be considered the upper bound. From the blue dashed lines, we observe that approximate GIP substantially reduces retrieval latency but sacrifices effectiveness. However, reranking the

<sup>15</sup>GTR<sub>XXL</sub> and GTR<sub>XXXL</sub> query encoding consume 96 and 349 ms, respectively, as opposed to 10 ms for TAS-B, reported by Ni et al. [37]

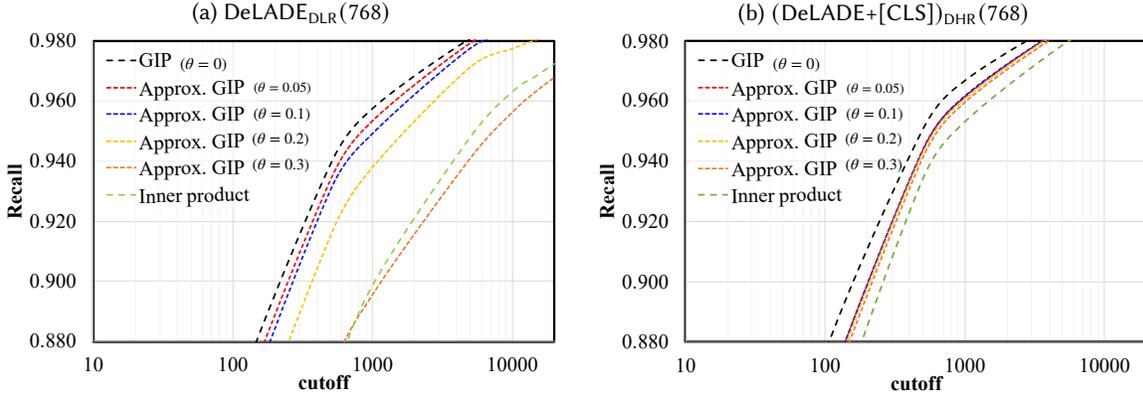


Fig. 4. Recall comparison of approximate retrieval approaches at different cutoffs. Approx. GIP refers to approximate gated inner product with threshold  $\theta$ . Inner product refers to standard inner product between query and passage value vectors without involving the index vectors.

Table 12. End-to-End Retrieval Latency for  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}(768)$  on MS MARCO (Dev)

Two-stage retrieval		CPU latency (ms/q)				GPU latency (ms/q)				Quality		
1st	2nd	Q Enc.	1st	2nd	Total	Q Enc.	1st	2nd	Total	MRR@10	R@1K	
(1)	GIP	-	164	72560	0	72724	12	670	0	682	0.357	0.967
(2)	Approx. GIP	GIP	164	9013	102	9279	12	23	10	45	0.357	0.967
(3)	Inner Product	GIP	164	3428	102	3694	12	37	10	59	0.357	0.967
(4)	Inner Product (w/ PQ128)	GIP	164	562	102	828	12	-*	10	-*	0.357	0.965
(5)	Inner Product (w/ PQ64)	GIP	164	284	102	650	12	-*	10	-*	0.357	0.961

\* Faiss. IndexPQ does not support GPU search.

Retrieval latency is measured in terms of ms/q with a single thread and batch size 1. With the exception of row (1), for each query, we retrieve the top-10000 passages at the first stage and then rerank the candidates using the second stage. Approx. GIP refers to approximate gated inner product with threshold  $\theta = 0.3$ . Inner product refers to standard inner product between query and passage value vectors without involving the index vectors.

top-10000 candidates with GIP mostly recovers the effectiveness loss (except for  $\theta = 0.3$ ) and requires only an additional 5–10 ms.

By incorporating the [CLS] vector, the performance tradeoff curves appear to be even better, as shown in Figures 5(c) and 5(d), which are organized in the same manner as Figures 5(a) and 5(b). For  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}(768)$ , we see no obvious retrieval effectiveness drop, even at  $\theta = 0.3$ . This result is consistent with our observation in Figure 4 that  $\text{DeLADE}_{\text{DLR}}(768)$  and  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}(768)$  retrieve enough relevant passages in the top-10000 candidates using approximate GIP with  $\theta \leq 0.3$ . Finally, we observe that IP only shows a minor R@1K degradation after GIP reranking, indicating that the standard inner product can provide an alternative to approximate GIP. This means that end-to-end retrieval with DLRs or DHRs can also be implemented with first-stage IP followed by GIP reranking.

Next, we measure the end-to-end retrieval latency of  $(\text{DeLADE}+[\text{CLS}])_{\text{DHR}}$ , including both query encoding (including vector densification) and retrieval components, using the same GPU and CPU environments described in Section 4.3. We report the query latency averaged over the 6980 MS MARCO dev queries in Table 12. With

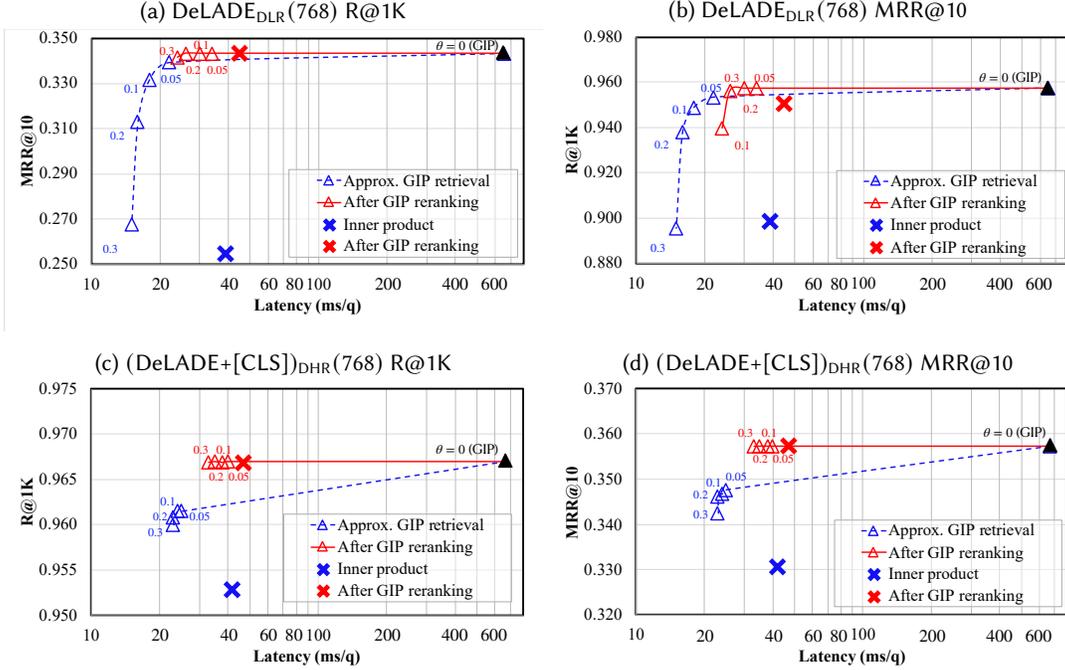


Fig. 5. Two-stage retrieval performance comparisons. We compare retrieval latency and effectiveness between approximate gated inner product (Approx. GIP) and standard inner product retrieval (between query and passage value vectors without involving the index vectors) and their performance after GIP reranking the top-10000 retrieved candidates. The labeled numbers denote the threshold  $\theta$  for Approx. GIP retrieval. Note that approx. GIP with  $\theta = 0$  is equal to GIP retrieval without any approximation.

the exception of row (1), for each query, we retrieve the top-10000 passages at the first stage and rerank the candidates using the second stage.

On the GPU, query encoding consumes 12 ms per query. With approximate GIP retrieval and GIP reranking, row (2), our lowest end-to-end retrieval latency on the GPU is 45 ms per query. With inner product retrieval and GIP reranking, query latency rises slightly to 59 ms per query, row (3). End-to-end retrieval latency on the CPU can be reduced to less than one second per query by performing inner product first-stage retrieval on the quantized query and passage value vectors, shown in rows (4) and (5). Although approximate GIP can improve the latency of first-stage retrieval over 80× compared to GIP on the CPU, row (2) vs row (1), approximate GIP is still over 2× slower than inner product, row (2) vs row (3). This stands in contrast to the GPU, where approximate GIP is *faster* than inner product. We note that approximate GIP requires selection of certain dimensions, i.e.,  $\mathcal{M}$  in Eq. (18), from the passage value and index vectors across the entire corpus, and this operation appears to be the latency bottleneck in our implementation using PyTorch on the CPU with a single thread. Thus, on the CPU, using the inner product between query and passage vectors is a better choice for first-stage retrieval, which can be further accelerated with approximate nearest neighbor (ANN) search algorithms.

In our final analysis, we demonstrate the interpretability of DLRs, which contrasts with dense semantic representations, where it is often difficult to understand why certain query–passage pairs obtain high scores. Since the “index” vector stores the position of the most important term in each slice, a DLR can be “reverted”

Table 13. Effects of  $\theta$  on Approximate GIP with the Reconstructed Expanded Query from DeLADE<sub>DLR</sub>(768)

Original query: where was the bauhaus built		
$\theta$	rank	Reconstructed (expanded) query terms
0.00	2	location, built, ##uh, ##aus, ba, was, house, build, building, site, were, school, home, church, the, store, originally, place, construction, studio, founded, headquarters, structure, later, city, is, origin, be, theater, college, first, hotel, villa, manufacture ... (omit)
0.05	3	location, built, ##uh, ##aus, ba, was, house, build, building, site, were, school
0.10	6	location, built, ##uh, ##aus, ba, was, house, build, building, site
0.20	16	location, built, ##uh, ##aus, ba, was, house, build
0.30	19	location, built, ##uh, ##aus, ba, was, house
Original passage (relevant)		
So the built output of Bauhaus architecture in these years is the output of Gropius: the Sommerfeld house in Berlin, the Otte house in Berlin, the Auerbach house in Jena, and the competition design for the Chicago Tribune Tower, which brought the school much attention. taatliches Bauhaus, commonly known simply as Bauhaus, was an art school in Germany that combined crafts and the fine arts, and was famous for the approach to design that it publicised and taught. It operated from 1919 to 1933.		
Reconstructed relevant (expanded) passage terms		
ba, ##uh, ##aus, ta, output, house, was, berlin, year, 1933, school, help, ##liche, germany, tribune, architecture, jena, chicago, ##at, design, 1919, ##ius, tower, competition, ot, somme, known, ##op, info, art, ##bach, attention, an, ##s, ##rf, location, commonly, ##a, simply, is, for, au, combined, ##te, built, operated, ##er, as, famous, ##ised, brought, and, approach, period, public, build, in, fine, building, be, reason, crafts, taught, bring, greatest, much, were, that, combine, skyscraper, architectural, fact, most, production, largest, are, origin, definition, arts, term, so, during, german, name, institution, called, abbreviation, it, type, ##t, which, highest, widely, history, of, important, work, operate, war, the, include, state, difference, organization, its, designed, ##st, company, great, century, 1920, recent, purpose, sculpture, acronym, meaning, studio, produced, last, lot, notable, run, literally, concept, 1934, college, place, created, record, historical, 1918, clock, major, see, time, decorative, produce, say, being, biggest, villa, fifty, many, founded, structure, culture ... (omit)		

The matching terms between the queries and the original (expanded) passage are colored with blue (red). The “rank” column denotes the rank of the relevant passage under different  $\theta$  settings. The reconstructed query and passage terms are ordered by their term weights in descending order.

back into a bag of words with term weights. We showcase how  $\theta$  impacts approximate GIP in Table 13, where we reconstruct expanded query terms from DeLADE<sub>DLR</sub>(768), shown in the top portion of the table. In addition, we reconstruct the passage judged as relevant to the query in the bottom portion of the table. For simplicity, we do not show term weights.

We observe that when  $\theta = 0$ , there are more than 10 matching terms (colored terms) between the query and the passage. As  $\theta$  increases, terms with lower weights are filtered out; see Eq. (18). Thus, the relevance score between the query and the passage decreases. For example, from  $\theta = 0$  to 0.05, terms with weights lower than 0.05 are removed for retrieval and the rank of the passage degrades slightly. At  $\theta = 0.3$ , four important terms (e.g., ‘build’, ‘building’, ‘were’, ‘school’) are removed, and the passage is retrieved much lower in the ranked list (at rank to 19).<sup>16</sup> This example illustrates that tuning  $\theta$  determines how many query terms (dimensions) are used for approximate GIP.

<sup>16</sup>We notice some “wacky” terms from the reconstructed passage with high term weights, originally observed by Mackenzie et al. [35].

## 6 CONCLUSIONS AND FUTURE WORK

We present a simple yet effective approach to densifying lexical representations for passage retrieval. This work introduces a dense representation framework and proposes a new scoring function to compute relevance scores between dense lexical representations (DLRs) derived from queries and passages. Using our framework, we can combine lexical and semantic representations into dense hybrid representations (DHRs) for hybrid retrieval. Our experiments show that DLRs can accurately approximate any “off-the-shelf” lexical model. Furthermore, when combined with other semantic representations (as DHRs), the resulting models can achieve comparable effectiveness to existing state-of-the-art hybrid retrieval methods.

The main advantage of our framework is that we can execute end-to-end retrieval using DLRs/DHRs on GPUs with a single index structure in a uniform execution environment. In our implementation, retrieval latency is insensitive to the sparsity of the vectors, unlike lexical representations that use inverted indexes. This feature makes our approach both fast and easy to deploy, especially for modern lexical retrieval models, which generate representations that are quite dense and require additional tuning (e.g., by introducing a sparsity constraint) to enable practical deployment using inverted indexes. Furthermore, we propose to better model content using a single model by jointly training semantic and lexical representations, which are then combined into hybrid representations for retrieval. We demonstrate that this combination, (DeLADE+[CLS])<sub>DHR</sub>, outperforms most models (e.g., ANCE, ColBERT, COIL, etc.) in both in-domain and out-of-domain evaluations while requiring smaller indexes and achieving lower query latency. Finally, we examine our two-stage retrieval approach, uncovering how the proposed design achieves low query latency without sacrificing accuracy.

One future research direction is to combine various supervised techniques [51, 53, 55] for further vector compression. In addition, since our proposed single model fusion approach can combine both lexical and semantic matching capabilities, it would be interesting to further explore methods that can integrate the strengths of each in a complementary manner. Finally, we believe that joint *self* training of semantic and lexical components is a promising future direction, especially in domain transfer scenarios.

## ACKNOWLEDGEMENTS

This research was supported in part by the Canada First Research Excellence Fund and the Natural Sciences and Engineering Research Council (NSERC) of Canada. We acknowledge Cloud TPU support from Google’s TPU Research Cloud (TRC). We thank the anonymous referees who provided useful feedback to improve this work.

## REFERENCES

- [1] Yang Bai, Xiaoguang Li, Gang Wang, Chaoliang Zhang, Lifeng Shang, Jun Xu, Zhaowei Wang, Fangshan Wang, and Qun Liu. 2020. SparTerm: Learning Term-based Sparse Representation for Fast Text Retrieval. *arXiv:2010.00768* (2020).
- [2] Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, et al. 2016. MS MARCO: A human generated MACHine Reading COMprehension dataset. *arXiv:1611.09268* (2016).
- [3] Wei-Cheng Chang, Felix X. Yu, Yin-Wen Chang, Yiming Yang, and Sanjiv Kumar. 2020. Pre-training Tasks for Embedding-based Large-scale Retrieval. In *Proc. ICLR*.
- [4] Xilun Chen, Kushal Lakhotia, Barlas Oğuz, Anchit Gupta, Patrick Lewis, Stan Peshterliev, Yashar Mehdad, Sonal Gupta, and Wen-tau Yih. 2021. Salient Phrase Aware Dense Retrieval: Can a Dense Retriever Imitate a Sparse One? *arXiv:2110.06918* (2021).
- [5] Nick Craswell, Bhaskar Mitra, and Daniel Campos. 2019. Overview of the TREC 2019 Deep Learning Track. In *Proc. TREC*.
- [6] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, and Daniel Campos. 2020. Overview of the TREC 2020 Deep Learning Track. In *Proc. TREC*.
- [7] Zhuyun Dai and Jamie Callan. 2020. Context-Aware Term Weighting For First Stage Passage Retrieval. In *Proc. SIGIR*. 1533–1536.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proc. NAACL*. 4171–4186.
- [9] Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2021. SPLADE v2: Sparse Lexical and Expansion Model for Information Retrieval. *arXiv:2109.10086* (2021).

- [10] Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2021. SPLADE: Sparse Lexical and Expansion Model for First Stage Ranking. In *Proc. SIGIR*. 2288–2292.
- [11] Luyu Gao and Jamie Callan. 2021. Condenser: a Pre-training Architecture for Dense Retrieval. In *Proc. EMNLP*. 981–993.
- [12] Luyu Gao, Zhuyun Dai, and Jamie Callan. 2021. COIL: Revisit Exact Lexical Match in Information Retrieval with Contextualized Inverted List. In *Proc. NAACL*. 3030–3042.
- [13] Luyu Gao, Zhuyun Dai, Tongfei Chen, Zhen Fan, Benjamin Van Durme, and Jamie Callan. 2021. Complement Lexical Retrieval Model with Semantic Residual Embeddings. In *Proc. ECIR*. 146–160.
- [14] Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2022. Tevatron: An Efficient and Flexible Toolkit for Dense Retrieval. *arXiv:2203.05765*.
- [15] Tiezheng Ge, Kaiming He, Qifa Ke, and Jian Sun. 2014. Optimized Product Quantization. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2014), 744–755.
- [16] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the Knowledge in a Neural Network. In *Proc. NeurIPS: Deep Learning and Representation Learning Workshop*.
- [17] Sebastian Hofstätter, Sophia Althammer, Michael Schröder, Mete Sertkan, and Allan Hanbury. 2020. Improving Efficient Neural Ranking Models with Cross-Architecture Knowledge Distillation. *arXiv:2010.02666* (2020).
- [18] 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 *Proc. SIGIR*. 113–122.
- [19] Sebastian Hofstätter, Omar Khattab, Sophia Althammer, Mete Sertkan, and Allan Hanbury. 2022. Introducing Neural Bag of Whole-Words with ColBERT: Contextualized Late Interactions using Enhanced Reduction. *arXiv:2203.13088* (2022).
- [20] Samuel Humeau, Kurt Shuster, Marie-Anne Lachaux, and Jason Weston. 2020. Poly-encoders: Architectures and Pre-training Strategies for Fast and Accurate Multi-sentence Scoring. In *Proc. ICLR*.
- [21] Piotr Indyk and Rameez Motwani. 2000. Approximate Nearest Neighbors: Towards Removing the Curse of Dimensionality. *Proc. STOC*.
- [22] Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Towards Unsupervised Dense Information Retrieval with Contrastive Learning. *arXiv:2112.09118* (2021).
- [23] Kyoung-Rok Jang, Junmo Kang, Giwon Hong, Sung-Hyon Myaeng, Joohee Park, Taewon Yoon, and Heecheol Seo. 2021. Ultra-High Dimensional Sparse Representations with Binarization for Efficient Text Retrieval. In *Proc. EMNLP*. 1016–1029.
- [24] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2021. Billion-Scale Similarity Search with GPUs. *IEEE Transactions on Big Data* (2021), 535–547.
- [25] Herve Jégou, Matthijs Douze, and Cordelia Schmid. 2011. Product Quantization for Nearest Neighbor Search. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2011), 117–128.
- [26] 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 *Proc. EMNLP*. 6769–6781.
- [27] Omar Khattab and Matei Zaharia. 2020. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. In *Proc. SIGIR*. 39–48.
- [28] Jimmy Lin and Xueguang Ma. 2021. A Few Brief Notes on DeepImpact, COIL, and a Conceptual Framework for Information Retrieval Techniques. *arXiv:2106.14807* (2021).
- [29] Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo Nogueira. 2021. Pyserini: A Python Toolkit for Reproducible Information Retrieval Research with Sparse and Dense Representations. In *Proc. SIGIR*. 2356–2362.
- [30] Jimmy Lin, Rodrigo Nogueira, and Andrew Yates. 2021. *Pretrained Transformers for Text Ranking: BERT and Beyond*. Morgan & Claypool.
- [31] Sheng-Chieh Lin and Jimmy Lin. 2021. Densifying Sparse Representations for Passage Retrieval by Representational Slicing. *arXiv:2112.04666* (2021).
- [32] Sheng-Chieh Lin, Jheng-Hong Yang, and Jimmy Lin. 2021. In-Batch Negatives for Knowledge Distillation with Tightly-Coupled Teachers for Dense Retrieval. In *Proc. Repl4NLP*. 163–173.
- [33] 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:1907.11692* (2019).
- [34] Yi Luan, Jacob Eisenstein, Kristina Toutanova, and Michael Collins. 2021. Sparse, Dense, and Attentional Representations for Text Retrieval. *Trans. Assoc. Comput. Linguistics* (2021), 329–345.
- [35] Joel Mackenzie, Andrew Trotman, and Jimmy Lin. 2021. Wacky Weights in Learned Sparse Representations and the Revenge of Score-at-a-Time Query Evaluation. *arXiv:2110.11540* (2021).
- [36] Antonio Mallia, Omar Khattab, Torsten Suel, and Nicola Tonellotto. 2021. Learning Passage Impacts for Inverted Indexes. In *Proc. SIGIR*. 1723–1727.
- [37] Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernández Ábrego, Ji Ma, Vincent Y. Zhao, Yi Luan, Keith B. Hall, Ming-Wei Chang, and Yinfei Yang. 2021. Large Dual Encoders Are Generalizable Retrievers. *arXiv:2112.07899* (2021).
- [38] Rodrigo Nogueira and Kyunghyun Cho. 2019. Passage Re-ranking with BERT. *arXiv:1901.04085* (2019).
- [39] Rodrigo Nogueira and Jimmy Lin. 2019. From doc2query to docTTTTTquery.

- [40] Biswajit Paria, Chih-Kuan Yeh, Ian E.H. Yen, Ning Xu, Pradeep Ravikumar, and Barnabás Póczos. 2020. Minimizing FLOPs to Learn Efficient Sparse Representations. In *Proc. ICLR*.
- [41] 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 *Proc. NAACL*. 5835–5847.
- [42] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proc. EMNLP*. 3982–3992.
- [43] 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 *Proc. EMNLP*. 2825–2835.
- [44] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv:1910.01108* (2019).
- [45] Keshav Santhanam, Omar Khattab, Christopher Potts, and Matei Zaharia. 2022. PLAID: An Efficient Engine for Late Interaction Retrieval. *arXiv:2205.09707* (2022).
- [46] Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. 2021. ColBERTv2: Effective and Efficient Retrieval via Lightweight Late Interaction. *arXiv:2112.01488* (2021).
- [47] Christopher Scovel, Zexuan Zhong, Jinhyuk Lee, and Danqi Chen. 2021. Simple Entity-Centric Questions Challenge Dense Retrievers. In *Proc. EMNLP*. 6138–6148.
- [48] 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 *Proc. NIPS*.
- [49] Kexin Wang, Nandan Thakur, Nils Reimers, and Iryna Gurevych. 2021. GPL: Generative Pseudo Labeling for Unsupervised Domain Adaptation of Dense Retrieval. *arXiv:2112.07577* (2021).
- [50] 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 *Proc. ICLR*.
- [51] Ikuya Yamada, Akari Asai, and Hannaneh Hajishirzi. 2021. Efficient Passage Retrieval with Hashing for Open-domain Question Answering. In *Proc. ACL*. 979–986.
- [52] Hamed Zamani, Mostafa Dehghani, W. Bruce Croft, Erik Learned-Miller, and Jaap Kamps. 2018. From Neural Re-Ranking to Neural Ranking: Learning a Sparse Representation for Inverted Indexing. In *Proc. CIKM*. 497–506.
- [53] 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 *Proc. CIKM*. 2487–2496.
- [54] Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2021. Optimizing Dense Retrieval Model Training with Hard Negatives. In *Proc. SIGIR*. 1503–1512.
- [55] Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2022. Learning Discrete Representations via Constrained Clustering for Effective and Efficient Dense Retrieval. In *Proc. WSDM*. 1328–1336.
- [56] Shengyao Zhuang and Guido Zuccon. 2021. Fast Passage Re-ranking with Contextualized Exact Term Matching and Efficient Passage Expansion. *arXiv:2108.08513* (2021).
- [57] Shengyao Zhuang and Guido Zuccon. 2021. TILDE: Term Independent Likelihood MoDEL for Passage Re-Ranking. In *Proc. SIGIR*. 1483–1492.