

Enhancing Sparse Retrieval via Unsupervised Learning

Anonymous Author(s)

ABSTRACT

Recent work has shown that neural retrieval models excel in text ranking tasks in a supervised setting when given large amounts of manually labeled training data. However, it remains an open question how one might train effective *unsupervised* retrieval models that are unequivocally better than lexical-matching baselines such as BM25. While some progress has been made in unsupervised *dense* retrieval models within a bi-encoder architecture, unsupervised *sparse* retrieval models remain unexplored. In this work, we propose BM26, to our knowledge the first such model. BM26 is trained in an unsupervised manner without the need for any human relevance judgments but keeps the same retrieval paradigm as BM25. Evaluations across multiple modern test collections, including MS MARCO, NaturalQuestions, TriviaQA, and BEIR, show that BM26 alone outperforms Contriever, the current state-of-the-art unsupervised dense retriever, in general, and performs on par with BM25. We further demonstrate two promising avenues to enhance existing lexical retrieval systems using such unsupervised lexical retrieval models: 1) Enhancing the BM25 system by hybrid with BM26 based on simple vector concatenation (that we dub “BM51”), demonstrating a significant and consistent improvement over BM25 alone, while remains unsupervised system. 2) Enhancing the supervised lexical retrieval model with improved initialization using BM26, as seen with the SPLADE model, generating significant improvements in both in-domain and zero-shot retrieval effectiveness.

1 INTRODUCTION

Traditional text retrieval methods such as TF-IDF and BM25 treat documents as “bags of words” and assign term weights using a heuristic function [28]. In general, these methods can be characterized as *unsupervised* lexical-matching models.¹ Although such methods date back many decades, they remain strong baselines [19] in various ranking tasks [1, 2], even in the age of deep neural networks [20]. Deployed in popular search platforms such as Elasticsearch, traditional lexical retrieval methods such as BM25 are widely used in industry for real-world search applications due to their robustness to domain and query variations.

There has been much recent progress in neural retrieval models that adopt a bi-encoder architecture to encode queries and documents independently into a representation space [15] using pre-trained language models such as BERT [5]. These representations can comprise either dense low-dimensional vectors [7, 15, 21, 30] or sparse high-dimensional vectors [6, 22, 24, 32]. Various models have been shown to be more effective than BM25 under the in-domain supervised learning setting for various document retrieval tasks. That is, given a sufficient amount of labeled training data comprising query-document pairs that have been (manually) judged for relevance, there is no doubt that we can train highly effective models. However, whether we can obtain an effective neural retrieval

model that performs better than BM25 in an *unsupervised* setting remains an open question [12, 29].

Existing work on unsupervised neural retrieval models have focused on dense retrievers such as ICT [17], Contriever [12], cpt-text [25]. At a high level, these models demonstrate how to craft pseudo relevant query-document pairs and how to obtain a large negative candidate pool, two important factors in the effectiveness of unsupervised dense retrievers. To date, though, we are not aware of any unsupervised model demonstrating effectiveness that is unequivocally better than BM25.

We noticed that in the evolution of unsupervised retrieval models, unsupervised dense retrievers introduce two main innovations compared to a traditional heuristic model such as BM25, which is used as a point of reference: (1) they change heuristic weighting functions to deep neural networks. (2) they change sparse lexical representations to dense semantic representations. Although changing the representation space gives such models more freedom to fit target labels, they lose the ability to perform exact lexical matches, which are more robust to noisy data and domain shifts. Moreover, sparse retrieval models are amenable to efficient retrieval using standard inverted indexes, a well-established technology with decades of research that has produced sophisticated query evaluation techniques. Thus, we hypothesize that unsupervised retrieval which learn sparse lexical representations has advantage on robustness than those that learn dense semantic representations. Our chain of reasoning is easy to see if we understand BM25 as a bi-encoder model with an unsupervised (i.e., heuristic) encoder [18]. Given this starting point, we propose a method to train a sparse retrieval model in an unsupervised manner. We call our model BM26 (i.e., what comes after BM25).

Our experiments show that BM26 alone is more effective than the existing state-of-the-art unsupervised dense retriever Contriever [12] in general and performs on par with BM25 in terms of effectiveness across a broad range of modern test collections. Furthermore, a retrieval model based on the simple concatenation of BM25 and BM26 representations significantly and consistently outperforms BM25 alone. We call this model BM51 (since $25 + 26 = 51$). A key feature is that it remains a lexical retrieval model and thus is compatible with infrastructure for “bag of words” retrieval based on inverted indexes. This provides a major advantage over dense-sparse hybrid systems, which require (separate) approximate nearest neighbor search libraries for efficient top- k retrieval. We also explored unsupervised lexical neural retrieval using “token expansion” to reduce token mismatch issues for lexical retrieval. These variants, BM26e and BM51e, show potential for further improvements in effectiveness. Moreover, we demonstrate the unsupervised lexical neural retrieval model can act as enhanced initialization for supervised lexical models such as SPLADE.

Contributions. We view this work as having three main contributions.

- Firstly, we introduce an unsupervised sparse retrieval model that we call BM26. This, to our knowledge, the first attempt

¹While it is possible to tune parameters (in the case of BM25, b_1 and k) with training data, in this work we simply adopt default parameters in all experiments.

to learn neural lexical retriever in unsupervised manner. The BM26 performs more effectively than existing dense retrieval models in the unsupervised setting.

- Secondly, we suggest how to combine BM25 and BM26 via simple vector concatenation to create BM51, a novel lexical retrieval model that is also unsupervised and fully compatible with existing inverted indexing infrastructure such as the Lucene search library. The method has the potential to improve “cold start” system where no labeled data are available with an inverted index only.
- Thirdly, we show supervised lexical models can be enhanced by using the unsupervised lexical neural retrieval models as initialization.

2 BACKGROUND AND RELATED WORK

2.1 Dense Retrieval Models

Throughout this paper we assume the standard definition of the *ad hoc* retrieval problem, where given a corpus $C = \{D_1, D_2, \dots, D_n\}$ comprised of an arbitrarily large collection of documents and a query Q , the system’s task is to return a top- k ranking of documents that maximizes some metric of quality such as nDCG, MRR, etc. “Documents” here is used generically to refer to the unit of retrieval, even though in actuality the system may be retrieving passages or other units of content (e.g., images).

Dense Passage Retriever (DPR) [15], which we take as an exemplar of a popular and large class of models known as dense retrieval models, adopts a bi-encoder structure to encode queries and documents separately into low-dimensional (e.g., 768 dimensions) dense vector representation as follows:

$$\mathbf{E}_Q = \text{Encoder}_Q(Q), \mathbf{E}_D = \text{Encoder}_D(D)$$

where the encoders are initialized with a pretrained language model such as BERT [5]. The query representation \mathbf{E}_Q and the document representation \mathbf{E}_D are taken from the last layer output of the [CLS] token of the corresponding encoder. The relevance between a query and a document is measured by the dot product of their representations, $\text{Sim}(Q, D) = \langle \mathbf{E}_Q, \mathbf{E}_D \rangle$.

Dense retrieval models are typically trained (more precisely, their underlying transformers are fine-tuned) using large amounts of supervised data comprising human-labeled query–document pairs. For DPR, during training, given a query Q , a labeled relevant document D^+ , and n non-relevant documents $D_1^-, D_2^-, \dots, D_n^-$, the model is optimized by contrastive learning using infoNCE loss:

$$\begin{aligned} \mathcal{L}(Q, D^+, D_1^-, D_2^-, \dots, D_n^-) \\ &= -\log p(D = D^+ | Q) \\ &= -\log \frac{\exp(\text{Sim}(Q, D^+))}{\exp(\text{Sim}(Q, D^+)) + \sum_{i=1}^n \exp(\text{Sim}(Q, D_i^-))}. \end{aligned}$$

Once the model has been trained, the document encoder can be applied to generate document representations for every document in the corpus (as a preprocessing step). At retrieval time, inference is applied to the query to generate the query representation, and the top- k most similar documents (in terms of dot products) are

retrieved. This is operationalized as a nearest neighbor search problem in dense vector space, which can be accomplished by existing libraries such as Faiss [13].

2.2 Unsupervised Dense Retrieval Models

Training dense retrieval models requires large amounts of labeled data. Recently, however, researchers have begun to explore training dense retrieval models in an *unsupervised* manner. The main challenge is how to automatically generate positive *pseudo* query–document pairs on which a model can be trained.

Two existing methods of creating such pseudo pairs are the Inverse Cloze Task (ICT) [17] and Independent Cropping (IC) [12]. Given a text span S composed of a sequence of tokens $\{t_1, t_2, \dots, t_n\}$, ICT randomly samples a sub-span M from S as the pseudo query and uses the rest of the sequence $S \setminus M$ as the pseudo document to form a positive query–document pair. The Independent Cropping method introduced by the recent Contriever model samples two random sub-spans S_1, S_2 (with replacement) to be the pseudo query and document respectively, which has been shown to be another simple yet promising way to create positive query–document pairs. The Contriever experiments found that IC achieves better effectiveness than ICT for unsupervised dense retriever training [12].

Another challenge of training unsupervised dense retrievers is how to find hard negative documents for each query. Training with hard negatives can improve the discriminability of the model [30]. In a supervised learning setting, the hard negatives for a query, in the simplest case, can be sampled from negative documents in the top results of an existing ranker (e.g., BM25), although the community has since developed much more sophisticated techniques. However, in the unsupervised setting, the most common approach to obtain harder negative documents is to increase the negative candidates pool for a higher chance of the model seeing hard negatives. Two existing methods to achieve this are increasing the batch size with in-batch negative training [25] or caching a large negative representation queue with momentum update [11, 12].

3 METHODS

3.1 BM25 as a Bi-Encoder

The starting point of our approach is the representation learning framework for IR articulated by Lin [18]. Recall that given a query Q and a document D , the BM25 ranking model computes the following similarity score:

$$\text{Sim}(Q, D) = \sum_{i=1}^n \text{IDF}(q_i) \cdot \frac{\text{TF}(q_i, D) \cdot (k_1 + 1)}{\text{TF}(q_i, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{\text{avgdl}})}$$

where n is the number of tokens in the query and q_i is the i -th token. This formula can be abstracted as:

$$\text{Sim}(Q, D) = \sum_{i=1}^n f(q_i, D)$$

where f is a function of token statistics of q_i in D and other global statistics. Instead of summing across the tokens in the query, we can rewrite the formula as a summation across all terms in the entire vocabulary space, where $\mathbb{1}$ is an indicator function for a term

appearing in the query:

$$\text{Sim}(Q, D) = \sum_{j=1}^{|V|} \mathbb{1}(v_j, Q) \cdot f(v_j, D)$$

This is in fact equivalent to the dot product of the query representation and the document representation in a vector space:

$$\text{Sim}(Q, D) = \langle \mathbf{E}_Q, \mathbf{E}_D \rangle$$

This formulation is exactly equivalent to a bi-encoder architecture that underlies recent transformer-based dense retrieval models such as DPR [15]. This generic design, however, admits a number of parametric differences, which can be characterized as design choices:

- Basis of the representation vectors: They can be dense (i.e., latent semantic dimensions learned by transformers) or sparse (i.e., the document vocabulary space).
- How weights are assigned: The weight of each dimension can be assigned by a heuristic (i.e., unsupervised) function or a function that has been learned (for example, a transformer).

In other words, BM25 adopts a bi-encoder structure where the query and document representations are generated independently by heuristic function “encoders” that operate on the document vocabulary space.

Similarly, dense retrieval models—exemplified by DPR [15]—fall into the category of generating dense semantic representations using encoders that are learned in a supervised setting. Many existing models have demonstrated the effectiveness of a *supervised* approach to learning encoders that operate on the document vocabulary space, e.g., DeepCT [4], DeepImpact [24], uniCOIL [22], and SPLADE [6]. Among these existing models, we adopt the architecture proposed in uniCOIL [22] for training our *unsupervised* lexical retrieval model as its similarity function has the same parametric form as BM25. However, we change the heuristic functions of both query and document “encoders” into neural models. This design also provides a natural point of comparison: Can we build an unsupervised lexical retrieval model that outperforms BM25 utilizing pretrained transformers such as BERT?

3.2 BM26: An Unsupervised Sparse Bi-Encoder

Given the general framework introduced above, BM26 can be characterized as an unsupervised (sparse) lexical retrieval model. The basis of the vector representation is the BERT token space, and the vector weights are assigned by a transformer that has been fine-tuned *without any human-labeled relevance judgments*.

In our model, the similarity between a query Q and a document D is computed by:

$$\text{Sim}(Q, D) = \sum_{j=1}^{|V|} w_{v_j, Q} \cdot w_{v_j, D}$$

$$w_{v_j, Q} = \text{ReLU}(P \cdot \text{BERT}(Q)_{v_j}) \text{ if } v_j \in Q \text{ else } 0$$

$$w_{v_j, D} = \text{ReLU}(P \cdot \text{BERT}(D)_{v_j}) \text{ if } v_j \in D \text{ else } 0$$

where $\text{BERT}(\cdot)_{v_j}$ is the contextual token representation from the last layer of BERT for token v_j and P is a linear projection that maps the token representation into a scalar weight. Note that if v_j occurs in a query or a document multiple times, we select the

maximum scalar weight for v_j as the weight of the token in the text. Exactly as above, this formula is equivalent to the dot product of the query and the document representation in a vector space with a basis defined by the vocabulary: $\text{Sim}(Q, D) = \langle \mathbf{E}_Q, \mathbf{E}_D \rangle$.

Following Contriever, we use independent cropping to create pseudo-positive pairs for contrastive learning. However, we choose not to use the MoCo [11] method to create a large negative candidate pool, as the representations in the pool are not generated by the latest model parameters. As an alternative, we increase the number of negative examples by increasing the batch size to a large number (16384 in our experiments). To achieve such a large batch size with limited GPU memory, we adopt the Gradient Caching [10] method. The training objective for our unsupervised model is the same as training a dense retriever, i.e., using the infoNCE loss.

3.3 BM25 \oplus BM26 = BM51

Since BM25 and BM26 are both unsupervised lexical retrieval models, we can combine their results to create a hybrid retrieval model that itself remains in the space of unsupervised methods. While linear combination of retrieval scores or reciprocal rank fusion are popular methods for fusion, we create a novel “sparse–sparse” hybrid by vector concatenation in this work. We name this model “BM51” as it is the fusion of BM25 and BM26 (25+26=51).

Specifically, the query representation and document representation of the BM51 retrieval model are computed by:

$$\mathbf{E}_{Q_{\text{BM51}}} = \mathbf{E}_{Q_{\text{BM25}}} \oplus \alpha \cdot \mathbf{E}_{Q_{\text{BM26}}}$$

$$\mathbf{E}_{D_{\text{BM51}}} = \mathbf{E}_{D_{\text{BM25}}} \oplus \beta \cdot \mathbf{E}_{D_{\text{BM26}}}$$

where \oplus represents the vector concatenation operation and α, β adjust the relative contributions of each representation to the final model. Since the token weights from BM25 and BM26 are on different scales, we first normalize the token weights of the two systems into integers in the range(0, 256); i.e., these become, in essence, impact scores. For simplicity and to retain the unsupervised setup, we keep $\alpha = \beta = 1$ after normalizing the token weights, which means that BM25 and BM26 are weighted equally. The BM51 query–document similarity is also computed as a dot product.

A detail worth noting: the vector bases of the individual representation vectors from BM25 and BM26 are different, since they are determined by the tokenizer used to process the text. In our implementation of BM25, the vocabulary space is defined by a Lucene analyzer (which, in the case of the MS MARCO passage collection, contains 2.7M unique tokens). In contrast, BM26 operates in the vocabulary space of BERT subwords, which contains 30,522 unique tokens. Thus, in our implementation, the vocabulary size of BM51 is only marginally bigger than that of BM25.

A key feature of our fusion-via-vector-concatenation approach of BM51 is that it remains yet another lexical retrieval model. This means that top- k retrieval can take advantage of all existing software infrastructure that has been developed around inverted indexes. In addition, vector concatenation presents an advantage over other popular fusion techniques such as the linear combination of scores or reciprocal rank fusion, both of which involve performing retrieval twice and post-processing two ranked lists. In other words, BM51 can serve as a drop-in replacement for BM25. That is,

Method	Rep. S/D	MS MARCO		NQ		TriviaQA	
		MRR@10	R@1k	Top20	Top100	Top20	Top100
(1) BM25	S	18.4	85.3	62.9	78.3	76.4	83.2
(2) BM25 _{wp}	S	17.5	82.6	63.2	77.8	73.8	81.4
(3) BM25' = BM25 \oplus BM25 _{wp}	S	18.9	86.9	64.4	79.4	76.5	83.5
(4) BM26	S	19.0 [▲]	88.6	63.9 [▲]	78.8	74.3 [▼]	82.5
(5) BM26e	S	18.1	92.1	70.0 [▲]	83.3	78.7 [▲]	85.2
(6) BM51 = BM25 \oplus BM26	S	22.2[▲]	92.0	68.1 [▲]	81.6	77.6 [▲]	83.9
(7) BM51e = BM25 \oplus BM26e	S	20.2 [▲]	93.5	71.8[▲]	84.2	79.3[▲]	85.5
(a) Contriever	D	16.0	88.1	67.8	82.1	74.2	83.2
(b) cpt-text-S	D	19.9	-	65.5	77.2	75.1	81.7

Table 1: The retrieval effectiveness of BM26 and BM51 compared to baselines and other unsupervised retrieval models on MS MARCO, NQ and TriviaQA. We performed significance tests comparing BM26 and BM51 to BM25 indicated by [▲]/_▼.

it can serve as a first-stage ranker to provide candidates for further reranking, it can be further combined with supervised models, etc.

3.4 BM26e & BM51e: Improving Lexical Match By Token Expansion

BM26 adopts the same parametric form as BM25, which means the lexical representation is restricted to the tokens that appeared in the original text. However, vocabulary mismatch issue is common to see in the lexical retrieval setting. The SPLADE work [6] utilizes the MLM layer of BERT to conduct expansion for the tokens in the original text. The purpose is to assign weights to the tokens that have close meaning to the existing tokens.

SPLADE was originally focused on supervised learning settings. We adopt its architecture here to achieve the token expansion for unsupervised lexical representation learning. Following the same name convention, we denote this unsupervised neural lexical retriever with expansion as BM26e. Specifically, the similarity between the pseudo query-document pairs are computed by:

$$\text{Sim}(Q, D) = \sum_{j=1}^{|V|} w_{v_j, Q} \cdot w_{v_j, D}$$

$$w_{v_j, Q} = \max_{i \in |Q|} \log(1 + \text{ReLU}(\text{MLM}(\text{BERT}(Q))_{v_{i,j}}))$$

$$w_{v_j, D} = \max_{i \in |D|} \log(1 + \text{ReLU}(\text{MLM}(\text{BERT}(D))_{v_{i,j}}))$$

where the term impact of a token v_j from the vocabulary space V is determined by the maximum-pooled MLM layer output, corresponding to the token, across all input tokens in a query Q or document D .

During training, additional FLOPS [26] loss was added to the infoNCE loss to control the sparsity of query and document expansions:

$$\ell_{\text{FLOPS}} = \sum_{j \in V} \bar{w}_j^2$$

which is calculated as the sum of the squares of the average impacts of each token within a batch. The FLOPS loss for query and document are weighted by λ_q and λ_d . We set them as 0.001 in all our experiments based on BM26e.

Similar to the way we build BM51 = BM25 \oplus BM26, we craft the “with expansion” version BM51e = BM25 \oplus BM26e. We follow the same vector concatenation and token weights processing as described in Sec. 3.3.

4 EXPERIMENT: ENHANCING UNSUPERVISED SPARSE RETRIEVAL

In this section, we show how BM26/BM26e act as an effective unsupervised neural lexical retriever and how it enhances the existing BM25 retrieval system.

Data. Following the data crafting method in the Contriever work, we use Independent Cropping on the CCNet corpus [3] to create the unsupervised training data. Specifically, we concatenate all the natural documents in CCNet together and treat each 500-token split as a document. Each training pair is obtained by independent cropping from a random span of 256 tokens in a document. We also apply 10% random token deletion on both the pseudo query and the pseudo document as data augmentation.

Models. We train BM26 and BM26e following the method proposed in Section 3.2 and 4. We use bert-base-uncased as initialization and train 10k steps with batch size 16384. Our model is trained on a single AWS EC2 instance with 8 \times A100 40 GB GPUs. The training code is modified based on the open-source toolkit Tevatron [9].

Our retrieval experiments are performed using the Pyserini IR toolkit [19] using the “fake words” trick, which is a common technique that allows sparse retrieval models to transparently reuse inverted indexes and existing query evaluation machinery [23].

BM51 is implemented exactly in the manner described in Section 3.3: We materialize the document vectors from BM25 and BM26 and then feed the concatenated vectors back into Pyserini for indexing and retrieval.

Finally, we evaluate two variants that allow us to precisely attribute differences in model effectiveness:

- BM25_{wp}: This represents the BM25 weighting function applied directly to wordpiece tokens. That is, BM25_{wp} and BM26 share exactly the vector basis, except BM26 weights are computed by a transformer model.

	AA	CF	DB	FE	FQ	HQ	NF	NQ	QU	SD	SF	TC	T2	Avg.
(1) BM25	39.7	16.5	31.8	65.1	23.6	63.3	32.2	30.6	78.9	14.9	67.9	59.5	44.2	43.7
(2) BM25 _{wp}	36.4	15.8	28.4	65.8	21.8	59.3	31.4	30.5	73.0	13.8	67.2	56.5	46.6	42.0
(3) BM25'	39.0	17.1	31.1	67.3	24.2	63.4	32.8	31.8	78.0	15.0	69.4	60.0	50.2	44.7
(4) BM26	38.1 [▼]	15.8	28.6 [▼]	74.4 [▲]	25.8 [▲]	63.6	32.4	27.9 [▼]	77.5 [▼]	15.0	66.7	54.4	26.1 [▼]	42.0
(5) BM26e	42.2 [▲]	18.9 [▲]	33.0	75.1 [▲]	28.9 [▲]	64.0 [▲]	33.3	29.9	82.4 [▲]	16.0 [▲]	67.9	34.9 [▼]	23.0 [▼]	42.3
(6) BM51	41.7 [▲]	18.9 [▲]	34.0 [▲]	75.2 [▲]	28.3 [▲]	68.2[▲]	34.2 [▲]	33.9 [▲]	82.3 [▲]	15.9 [▲]	69.8 [▲]	65.6	34.7 [▼]	46.4
(7) BM51e	43.1[▲]	20.0[▲]	36.5[▲]	77.9[▲]	30.1 [▲]	67.9 [▲]	34.7[▲]	33.4[▲]	83.5[▲]	16.8[▲]	70.9[▲]	53.9	28.6 [▼]	46.0
(a) Contriever	37.9	15.5	29.3	68.2	24.5	48.1	31.7	25.4	83.5	14.9	64.9	27.4	19.3	37.7
(b) cpt-text-S	38.7	15.8	27.2	57.1	34.1	51.5	32.0	-	68.1	-	65.4	52.9	21.0	-

Table 2: The nDCG@10 retrieval effectiveness of BM26 and BM51 compared to baselines and other unsupervised retrieval models on BEIR. We performed significance tests comparing BM26 and BM51 to BM25 indicated by [▲]/[▼]. Dataset Legend: AA=ArguAna, CF=Climate-FEVER, DB=DBpedia, FE=FEVER, FQ=FiQA, HQ=HotpotQA, NF=NFCorpus, NQ=NaturalQuestions, QU=Quora, SD=SCIDOCS, SF=SciFact, TC=TREC-COVID, T2=Touché-2020 (v2).

- BM25': This represents BM25 \oplus BM25_{wp}, or the concatenation of "standard" BM25 and BM25_{wp} vectors. That is, BM25' and BM51 share exactly the same vector basis, differing only in how the BM26 portion of the weights are assigned.

Evaluation. We evaluate our models on the following datasets:

- MS MARCO Passage Ranking [1]: a web search dataset that contains 8.8 million passages and 6980 queries in the development set for evaluation. Effectiveness is measured by MRR@10 and Recall@1k.
- NaturalQuestions [16] and TriviaQA [14]: open-domain question answering datasets that use English Wikipedia as the corpus. Effectiveness is measured by top-20/100 retrieval accuracy.
- BEIR [29]: a collection of datasets for zero-shot evaluation of search and related tasks. We evaluate our models on the 13 publicly available datasets; effectiveness is measured in terms of nDCG@10.

4.1 Results

Table 1 presents results on MS MARCO, NQ and TriviaQA, and Table 2 shows results on the 13 BEIR datasets. The results of the BM25 baselines are obtained by following the reproduction guides in the open-source Anserini IR toolkit [31]. The results of Contriever and cpt-text-S are copied from the original papers.

Evaluation of BM26. First, we compare BM26 with BM25. The high level conclusion is that they are roughly on par, which can be seen by comparing row (4) with row (1) in Table 1 on MS MARCO, NQ and TriviaQA, and by comparing column (4) with column (1) in Table 2 for the BEIR datasets. In both tables, we performed significance testing using paired t -tests ($p < 0.05$); the symbols [▲]/[▼] indicate significant differences (better/worse). Despite the statistical significance of some differences, overall the effect sizes are relatively small: less than a point in most cases.

For BM25, the default Lucene analyzer used by Anserini generates 2.7M unique tokens for the MS MARCO passage corpus, which is orders of magnitude larger than the vocabulary space of the bert-base-uncased tokenizer (30,522 unique subwords). To isolate the effects of this difference, we turn to BM25_{wp}, which applies BM25 term weighting on wordpiece tokens. Thus, BM25_{wp}

shares exactly the same representation space as BM26, providing a fair comparison between neural and heuristic encoders. Here, we see that BM26 outperforms BM25_{wp}.

Turning our attention to effectiveness on the BEIR datasets in Table 2, column (4) vs. column (1), we see that BM26 is significantly better than BM25 for two datasets and worse for five datasets. Overall, BM26 is 1.7 points worse than BM25, but BM26 and BM25_{wp} achieve the same level of effectiveness.

Comparing BM26 to Contriever and cpt-text-S on the datasets in Table 1, shown in row (a) and row (b), respectively, we see some cases where our method performs better and other cases where our method performs worse. At a high level, it would be fair to characterize effectiveness as "comparable". On BEIR, however, BM26 scores four points higher than Contriever. Although cpt-text-S results are not available for all BEIR datasets, it does appear that BM26 obtains higher scores in most cases. Note that significance tests are not possible here because we do not have access to the raw run files from these models.

Evaluation of BM51. Secondly, we investigate how BM26 helps improve baseline BM25 via our vector concatenation technique. The main comparison condition is between BM51 (= BM25 \oplus BM26) and BM25. We performed significance testing and denote better/worse results exactly in the same manner as above.

From Table 1, we see significant improvements across all three datasets, comparing row (6) to row (1). In the case of MS MARCO, we observe a nearly four point gain, which translates into a 21% relative improvement. From Table 2, comparing column (6) with column (1), we see that BM51 is significantly more effective than BM25 for all BEIR datasets but Touché-2020. We emphasize that we obtain these gains without the use of manual labels.

On BEIR, the only case where BM51 falls short is on the Touché-2020 dataset. This corpus contains many long documents, which might be the reason why our unsupervised sparse neural retriever performs worse than any BM25 variant. However, Touché-2020 appears to be an outlier, as other existing models [6, 12, 25] also perform poorly.

To untangle the effectiveness improvements due to BM26 from improvements from simply having a hybrid vocabulary base, we can

	AA	CF	DB	FE	FQ	HQ	NF	NQ	QU	SD	SF	TC	T2	Avg.
(1) BM25	39.7	16.5	31.8	65.1	23.6	63.3	32.2	30.6	78.9	14.9	67.9	59.5	44.2	43.7
(2) BM51e	43.1	20.0	36.5	77.9	30.1	67.9	34.7	33.4	83.5	16.8	70.9	53.9	28.6	46.0
(3) uniCOIL	39.6	18.2	33.8	81.2	28.9	66.7	33.3	42.5	66.3	14.4	68.6	64.0	29.8	45.3
(4) COIL-full	29.5	21.6	39.8	84.0	31.3	71.3	33.1	51.9	83.8	15.5	70.7	66.8	28.1	48.3
(5) SPLADE	43.9	19.9	36.6	73.0	28.7	63.6	31.3	46.9	83.5	14.5	62.8	67.3	31.6	46.4
(6) SPLADE _{repro}	47.4	21.2	38.8	78.1	30.0	63.7	33.9	48.6	75.5	14.5	65.8	72.7	28.0	47.6
(7) SPLADE _{BM26e}	53.5	21.9	41.8	81.4	34.2	70.7	34.6	52.2	84.4	15.6	70.8	70.1	29.8	50.8

Table 3: The zero-shot retrieval effectiveness of the SPLADE model initialized by BM26e compared with baseline sparse retrieval methods on BEIR. Supervised Models in the table are trained with only BM25 hard negatives.

	MRR@10	Recall@1k
(1) BM25	18.4	85.3
(2) BM51e	20.2	93.5
(3) COIL-full	35.3	96.7
(4) SPLADE	34.0	96.5
(5) SPLADE _{repro}	33.5	97.0
(6) SPLADE _{BM26e}	35.2	98.0

Table 4: The in-domain retrieval effectiveness of SPLADE model initialized by BM26e compared with baseline sparse retrieval methods on MS MARCO passage ranking task.

compare BM51 to $BM25' = BM25 \oplus BM25_{wp}$. Recall that the latter represents the concatenation of BM25 and $BM25_{wp}$: the BM25 components are identical, and as explained above, $BM25_{wp}$ and BM26 share exactly the same representation space, differing only in how term weights are computed. We see that BM51 consistently outperforms $BM25'$, which demonstrates that BM26 provides additional relevance signals beyond a simple vocabulary space hybrid.

We see that BM51 is more effective than Contriever and cpt-text-S across most datasets. On MS MARCO, BM51 outperforms Contriever by over six points, and on BEIR, nine points.

Evaluation of BM26e and BM51e. Finally, we investigate the performance of BM26e and BM51e, which incorporate token expansion to address the vocabulary mismatch issue. We compare BM26e and BM51e with their non-expanded counterparts, BM26 and BM51, respectively.

For example in Table 1, comparing row (5) to row (4), we observe that BM26e performs better than BM26 on NQ, and TriviaQA. This indicates that token expansion is beneficial for improving retrieval effectiveness. In Table 2, comparing column (7) with column (6), we find that BM51e is more effective than BM51 on 9 out of the 13 BEIR datasets. These results highlight the value of token expansion in addressing the vocabulary mismatch issue and improving the retrieval performance of unsupervised sparse neural retrievers.

5 EXPERIMENT: ENHANCING SUPERVISED SPARSE RETRIEVAL

In this section, we show how the unsupervised neural lexical retriever can enhance the supervised lexical retriever as better backbone initialization.

Data We use the train set of MS MARCO passage ranking task [1] to train the model. During training, the positive label is from human judgment and the in-batch negative document is from BM25 negatives. The in-domain retrieval effectiveness is evaluated on the dev set, measured by MRR@10 and Recall@1k. The zero-shot retrieval effectiveness was evaluated on the aforementioned BEIR datasets, measured by nDCG@10.

Model As BM26e has shown higher unsupervised retrieval effectiveness than BM26 in Sec.4.1. In this experiment, we study the effectiveness of using BM26e as initialization to train SPLADE [6], denoted as SPLADE_{BM26e}. Specifically, we replace. We finetune the model on MS MARCO train set for 3 epochs with batch size 256. Each training example contains 1 positive passage and 7 sampled BM25 hard negative passages.

Baselines We aim to compare our proposed SPLADE_{BM26e} with other existing representative lexical neural retrieval models like COIL [8], uniCOIL [22], and SPLADE [6]. We also have our own implementation of SPLADE denoted as SPLADE_{repro} in comparison, which has the same hyperparameter setting with SPLADE_{BM26e} but initialized with bert-base-uncased. While recent studies suggest that the effectiveness of the retrieval can be further enhanced through hard negative mining and distillation from a cross-encoder reranker [6, 21, 27, 30], we have only compared models trained on BM25 hard negatives, excluding distillation and multi-round hard negative mining, to maintain a fair comparison.

5.1 Results

Evaluation of in-domain retrieval. In Table 4, we present the in-domain retrieval results on MS MARCO passage ranking task. The SPLADE_{BM26e} model outperformed other baseline models in terms of both MRR@10 and Recall@1k, thus illustrating the benefits of our proposed approach for model initialization. The SPLADE_{BM26e} model achieved an MRR@10 of 35.2 and Recall@1k 0.98, which is a significant improvement over the SPLADE. These results demonstrate that using unsupervised neural lexical retrieval models, specifically the BM26e in this case, to initialize the weights of the supervised models can improve the performance of supervised lexical retrieval.

Evaluation of zero-shot retrieval. In Table 3, we provide the results for zero-shot retrieval effectiveness on the BEIR dataset.

The SPLADE_{BM26e} model consistently outperformed other baseline models in most retrieval tasks, thus demonstrating its superior generalizability in a zero-shot retrieval scenario. On average, SPLADE_{BM26e} achieved an overall score of 50.8, which is 3 points higher than the SPLADE model initialized by bert-base-uncased. This further solidifies the effectiveness of our proposed approach in using unsupervised lexical retrieval model as the initialization of the supervised model for enhancing retrieval performance in a zero-shot setting.

6 CONCLUSION

In this work, we propose an unsupervised sparse retrieval model called BM26 that adapts techniques originally designed for dense retrieval to a sparse lexical representation space. We found that our unsupervised lexical approach is on par with BM25 but has great potential to enhance the existing lexical retrieval systems. Our unsupervised lexical retriever exhibits its utility when there is a scarcity of relevance judgement. It successfully augments heuristic retrieval methods such as BM25 through the creation of a hybrid vocabulary vector space. This allows for a "cold start" retrieval system with increased effectiveness. Furthermore, given a sufficient availability of relevance judgement, the unsupervised lexical retriever (for instance, BM26e) can act as better initialization for supervised finetuning.

REFERENCES

- [1] Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamea, Bhaskar Mitra, Tri Nguyen, Mir Rosenberg, Xia Song, Alina Stoica, Saurabh Tiwary, and Tong Wang. 2018. MS MARCO: A Human Generated Machine Reading Comprehension Dataset. *arXiv:1611.09268v3* (2018).
- [2] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading Wikipedia to Answer Open-Domain Questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vancouver, Canada, 1870–1879.
- [3] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised Cross-lingual Representation Learning at Scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 8440–8451.
- [4] Zhuyun Dai and Jamie Callan. 2020. Context-Aware Term Weighting For First Stage Passage Retrieval. In *Proceedings of the 43rd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2020)*. 1533–1536.
- [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota, 4171–4186.
- [6] Thibault Formal, C. Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2021. SPLADE v2: Sparse Lexical and Expansion Model for Information Retrieval. *arXiv:2109.10086* (2021).
- [7] Luyu Gao and Jamie Callan. 2021. Unsupervised Corpus Aware Language Model Pre-training for Dense Passage Retrieval. *arXiv:2108.05540* (2021).
- [8] Luyu Gao, Zhuyun Dai, and Jamie Callan. 2021. COIL: Revisit Exact Lexical Match in Information Retrieval with Contextualized Inverted List. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Online, 3030–3042. <https://doi.org/10.18653/v1/2021.naacl-main.241>
- [9] Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2022. Tevatron: An Efficient and Flexible Toolkit for Dense Retrieval. *arXiv:2203.05765* (2022).
- [10] Luyu Gao, Yunyi Zhang, Jiawei Han, and Jamie Callan. 2021. Scaling Deep Contrastive Learning Batch Size under Memory Limited Setup. In *Proceedings of the 6th Workshop on Representation Learning for NLP (RePL4NLP-2021)*. Online, 316–321.
- [11] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. 2020. Momentum Contrast for Unsupervised Visual Representation Learning. In *Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 9726–9735.
- [12] Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised Dense Information Retrieval with Contrastive Learning. *arXiv:2112.09118* (2021).
- [13] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data* 7, 3 (2019), 535–547.
- [14] Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vancouver, Canada, 1601–1611.
- [15] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Online, 6769–6781.
- [16] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural Questions: A Benchmark for Question Answering Research. *Transactions of the Association for Computational Linguistics* 7 (2019), 452–466.
- [17] Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. Latent Retrieval for Weakly Supervised Open Domain Question Answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy, 6086–6096.
- [18] Jimmy Lin. 2021. A Proposed Conceptual Framework for a Representational Approach to Information Retrieval. *arXiv:2110.01529* (2021).
- [19] 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 *Proceedings of the 44th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2021)*. 2356–2362.
- [20] Jimmy Lin, Rodrigo Nogueira, and Andrew Yates. 2021. *Pretrained Transformers for Text Ranking: BERT and Beyond*. Morgan & Claypool Publishers.
- [21] Sheng-Chieh Lin, Jheng-Hong Yang, and Jimmy Lin. 2021. In-Batch Negatives for Knowledge Distillation with Tightly-Coupled Teachers for Dense Retrieval. In *Proceedings of the 6th Workshop on Representation Learning for NLP (RePL4NLP-2021)*. Online, 163–173.
- [22] Xueguang Ma, Ronak Pradeep, Rodrigo Nogueira, and Jimmy Lin. 2022. Document Expansions and Learned Sparse Lexical Representations for MS MARCO V1 and V2. In *Proceedings of the 45th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2022)*. Madrid, Spain, 3187–3197.
- [23] 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).
- [24] Antonio Mallia, Omar Khattab, Torsten Suel, and Nicola Tonellotto. 2021. Learning Passage Impacts for Inverted Indexes. In *Proceedings of the 44th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2021)*. 1723–1727.
- [25] Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, Johannes Heidecke, Pranav Shyam, Boris Power, Tyna Eloundou Nekoul, Girish Sastry, Gretchen Krueger, David Schnurr, Felipe Petroski Such, Kenny Hsu, Madeleine Thompson, Tabarak Khan, Toki Sherbakov, Joanne Jang, Peter Welinder, and Lilian Weng. 2022. Text and Code Embeddings by Contrastive Pre-Training. *arXiv:2201.10005* (2022).
- [26] Biswajit Paria, Chih-Kuan Yeh, Ning Xu, Barnabás Póczos, Pradeep Ravikumar, and Ian En-Hsu Yen. 2020. Minimizing FLOPs to Learn Efficient Sparse Representations. *ArXiv abs/2004.05665* (2020).
- [27] Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng Wang. 2021. RocketQA: An Optimized Training Approach to Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Online, 5835–5847. <https://doi.org/10.18653/v1/2021.naacl-main.466>
- [28] Stephen E. Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Foundations and Trends in Information Retrieval* 3, 4 (2009), 333–389.
- [29] Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. BEIR: A Heterogeneous Benchmark for Zero-shot Evaluation of Information Retrieval Models. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- [30] 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 *Proceedings of the 9th International Conference on Learning Representations (ICLR 2021)*.

- [31] Peilin Yang, Hui Fang, and Jimmy Lin. 2017. Anserini: Enabling the Use of Lucene for Information Retrieval Research. In *Proceedings of the 40th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2017)*. Tokyo, Japan, 1253–1256.

- [32] Shengyao Zhuang and Guido Zuccon. 2021. Fast Passage Re-ranking with Contextualized Exact Term Matching and Efficient Passage Expansion. *arXiv:2108.08513* (2021).

SIGIR RR Explanation

SIGIR Short Paper Submission ID: 2399

In the review of our SIGIR short paper submission the reviewers all mentioned that our work is well motivated and interesting. The comments from reviewers about the weakness of our SIGIR submission can be summarized as three key points:

1. too strong claim about drop-in replacement of BM25

We initially claimed that our BM51/BM51e setting was a drop-in replacement. However, since BM25 does not require any fine-tuning or parameter tuning, we now suggest using BM51 to enhance an existing BM25 system for better retrieval effectiveness in the "cold start" setting for our current submission.

2. lack of technical detail for expansion model BM26e

We have included additional technical and implementation details about the model with expansion, BM26e. Furthermore, we conducted a series of new experiments that utilized BM26e as initialization to train a supervised SPLADE model. Our results demonstrate that an unsupervised lexical neural retrieval model serves as a good initialization for its supervised counterpart.

3. baselines should including other learned sparse retrieval model such as SPLADE and COIL.

SPLADE and COIL are supervised sparse retrieval models. In this submission, we compare SPLADE, COIL, and uniCOIL while evaluating the SPLADE model initialized by BM26e. The SPLADE model initialized using BM26e achieves the best zero-shot performance among the models.

Overall, we have made significant revisions based on valuable reviewer comments from SIGIR. We have reframed the main contributions as enhancing both unsupervised and supervised retrieval systems with our proposed method. Additionally, we have added a series of evaluations and experiments.

BM51: Improving BM25 via Unsupervised Learning

Anonymous Author(s)

ABSTRACT

Recent work has shown that neural retrieval models excel in text ranking tasks in a supervised setting when given large amounts of manually labeled training data. However, it remains an open question how one might train effective *unsupervised* retrieval models that are unequivocally better than lexical-matching baselines such as BM25. While some progress has been made in unsupervised *dense* retrieval models within a bi-encoder architecture, unsupervised *sparse* retrieval models remain unexplored. In this work, we propose BM26, to our knowledge the first such model. BM26 is trained in an unsupervised manner without the need for any human relevance judgments. Evaluations across multiple modern test collections, including MS MARCO, NaturalQuestions, TriviaQA, and BEIR, show that BM26 alone outperforms Contriever, the current state-of-the-art unsupervised dense retriever, in general, and performs on par with BM25. We further demonstrate that a hybrid of BM25 and BM26 based on simple vector concatenation (that we dub “BM51”), significantly and consistently outperforms BM25 alone. As a sparse model, BM51 can serve as a drop-in replacement for BM25 in scenarios that already use inverted indexes, and we encourage the adoption of our model in this capacity.

1 INTRODUCTION

Traditional text retrieval methods such as TF-IDF and BM25 treat documents as “bags of words” and assign term weights using a heuristic function [23]. In general, these methods can be characterized as *unsupervised* lexical-matching models.¹ Although such methods date back many decades, they remain strong baselines [17] in various ranking tasks [1, 2], even in the age of deep neural networks [18]. Deployed in popular search platforms such as Elasticsearch, traditional lexical retrieval methods such as BM25 are widely used in industry for real-world search applications due to their robustness to domain and query variations.

There has been much recent progress in neural retrieval models that adopt a bi-encoder architecture to encode queries and documents independently into a representation space [13] using pre-trained language models such as BERT [5]. These representations can comprise either dense low-dimensional vectors [7, 13, 19, 25] or sparse high-dimensional vectors [6, 20, 21, 27]. Various models have been shown to be more effective than BM25 under the in-domain supervised learning setting for various document retrieval tasks. That is, given a sufficient amount of labeled training data comprising query–document pairs that have been (manually) judged for relevance, there is no doubt that we can train highly effective models. However, whether we can obtain an effective neural retrieval model that performs better than BM25 in an *unsupervised* setting remains an open question [11, 24].

Existing work on unsupervised neural retrieval models have focused on dense retrievers such as ICT [15], Contriever [11], cpt-text [22]. At a high level, these models demonstrate how to craft

pseudo relevant query–document pairs and how to obtain a large negative candidate pool, two important factors in the effectiveness of unsupervised dense retrievers. To date, though, we are not aware of any unsupervised model demonstrating effectiveness that is unequivocally better than BM25.

We noticed that in the evolution of unsupervised retrieval models, unsupervised dense retrievers introduce two main innovations compared to a traditional heuristic model such as BM25, which is used as a point of reference: (1) they change heuristic weighting functions to deep neural networks. (2) they change sparse lexical representations to dense semantic representations. Although changing the representation space gives such models more freedom to fit target labels, they lose the ability to perform exact lexical matches, which are more robust to noisy data and domain shifts. Moreover, sparse retrieval models are amenable to efficient retrieval using standard inverted indexes, a well-established technology with decades of research that has produced sophisticated query evaluation techniques. Thus, we hypothesize that unsupervised retrieval which learn sparse lexical representations will be more effective than those that learn dense semantic representations. Our chain of reasoning is easy to see if we understand BM25 as a bi-encoder model with an unsupervised (i.e., heuristic) encoder [16]. Given this starting point, we propose a method to train a sparse retrieval model in an unsupervised manner. We call our model BM26 (i.e., what comes after BM25).

Our experiments show that BM26 alone is more effective than the existing state-of-the-art unsupervised dense retriever Contriever [11] in general and performs on par with BM25 in terms of effectiveness across a broad range of modern test collections. Furthermore, a retrieval model based on the simple concatenation of BM25 and BM26 representations significantly and consistently outperforms BM25 alone. We call this model BM51 (since $25 + 26 = 51$). A key feature is that it remains a lexical retrieval model and thus is compatible with infrastructure for “bag of words” retrieval based on inverted indexes. This provides a major advantage over dense models, which require (separate) approximate nearest neighbor search libraries for efficient top- k retrieval. We also explored unsupervised lexical neural retrieval using “token expansion” to reduce token mismatch issues for lexical retrieval. These variants, BM26e and BM51e, show potential for further improvements in effectiveness.

Contributions. We view this work as having two main contributions. Firstly, we introduce an unsupervised sparse retrieval model that we call BM26. This, to our knowledge, is the first model of its type. The BM26 model learns sparse lexical representations and performs more effectively than existing dense retrieval models in the unsupervised setting. Secondly, we suggest how to combine BM25 and BM26 via simple vector concatenation to create BM51, a novel lexical retrieval model that is also unsupervised and fully compatible with existing inverted indexing infrastructure such as the Lucene search library. Our methods can consistently improve “cold start” system where no labeled data are available.

¹While it is possible to tune parameters (in the case of BM25, b_1 and k) with training data, in this work we simply adopt default parameters in all experiments.

2 BACKGROUND AND RELATED WORK

Training neural retrieval models via contrastive learning requires large amounts of labeled data [1, 13, 19, 25]. Recently, however, researchers have begun to explore training dense retrieval models in an *unsupervised* manner [11, 22]. The main challenge is how to automatically generate positive *pseudo* query–document pairs on which a model can be trained.

Two existing methods of creating such pseudo pairs are the Inverse Cloze Task (ICT) [15] and Independent Cropping (IC) [11]. Given a text span S composed of a sequence of tokens $\{t_1, t_2, \dots, t_n\}$, ICT randomly samples a sub-span M from S as the pseudo query and uses the rest of the sequence $S \setminus M$ as the pseudo document to form a positive query–document pair. The Independent Cropping method introduced by the recent Contriever model samples two random sub-spans S_1, S_2 (with replacement) to be the pseudo query and document respectively, which has been shown to be another simple yet promising way to create positive query–document pairs. The Contriever experiments found that IC achieves better effectiveness than ICT for unsupervised dense retriever training [11].

Another challenge of training unsupervised dense retrievers is how to find hard negative documents for each query. Training with hard negatives can improve the discriminability of the model [25]. In a supervised learning setting, the hard negatives for a query, in the simplest case, can be sampled from negative documents in the top results of an existing ranker (e.g., BM25), although the community has since developed much more sophisticated techniques. However, in the unsupervised setting, the most common approach to obtain harder negative documents is to increase the negative candidates pool for a higher chance of the model seeing hard negatives. Two existing methods to achieve this are increasing the batch size with in-batch negative training [22] or caching a large negative representation queue with momentum update [10, 11].

3 METHODS

3.1 BM25 as a Bi-Encoder

The starting point of our approach is the representation learning framework for IR articulated by Lin [16]. Recall that given a query Q and a document D , the BM25 ranking model computes the following similarity score:

$$\text{Sim}(Q, D) = \sum_{i=1}^n \text{IDF}(q_i) \cdot \frac{\text{TF}(q_i, D) \cdot (k_1 + 1)}{\text{TF}(q_i, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{\text{avgdl}})}$$

where n is the number of tokens in the query and q_i is the i -th token. This formula can be abstracted as:

$$\text{Sim}(Q, D) = \sum_{i=1}^n f(q_i, D)$$

where f is a function of token statistics of q_i in D and other global statistics. Instead of summing across the tokens in the query, we can rewrite the formula as a summation across all terms in the entire vocabulary space, where $\mathbb{1}$ is an indicator function for a term appearing in the query:

$$\text{Sim}(Q, D) = \sum_{j=1}^{|V|} \mathbb{1}(v_j, Q) \cdot f(v_j, D)$$

This is in fact equivalent to the dot product of the query representation and the document representation in a vector space:

$$\text{Sim}(Q, D) = \langle \mathbf{E}_Q, \mathbf{E}_D \rangle$$

This formulation is exactly equivalent to a bi-encoder architecture that underlies recent transformer-based dense retrieval models such as DPR [13]. This generic design, however, admits a number of parametric differences, which can be characterized as design choices:

- Basis of the representation vectors: They can be dense (i.e., latent semantic dimensions learned by transformers) or sparse (i.e., the document vocabulary space).
- How weights are assigned: The weight of each dimension can be assigned by a heuristic (i.e., unsupervised) function or a function that has been learned (for example, a transformer).

In other words, BM25 adopts a bi-encoder structure where the query and document representations are generated independently by heuristic function “encoders” that operate on the document vocabulary space.

Many existing models have demonstrated the effectiveness of a *supervised* approach to learning encoders that operate on the document vocabulary space, e.g., DeepCT [4], DeepImpact [21], uniCOIL [20], and SPLADE [6]. Among these existing models, we adopt the architecture proposed in uniCOIL [20] for training our *unsupervised* lexical retrieval model as its similarity function has the same parametric form as BM25. However, we change the heuristic functions of both query and document “encoders” into neural models. This design also provides a natural point of comparison: Can we build an unsupervised lexical retrieval model that outperforms BM25 utilizing pretrained transformers such as BERT?

3.2 BM26: An Unsupervised Sparse Bi-Encoder

Given the general framework introduced above, BM26 can be characterized as an unsupervised (sparse) lexical retrieval model. The basis of the vector representation is the BERT token space, and the vector weights are assigned by a transformer that has been fine-tuned *without any human-labeled relevance judgments*.

In our model, the similarity between a query Q and a document D is computed by:

$$\text{Sim}(Q, D) = \sum_{j=1}^{|V|} w_{v_j, Q} \cdot w_{v_j, D}$$

$$w_{v_j, Q} = \text{ReLU}(P \cdot \text{BERT}(Q)_{v_j}) \text{ if } v_j \in Q \text{ else } 0$$

$$w_{v_j, D} = \text{ReLU}(P \cdot \text{BERT}(D)_{v_j}) \text{ if } v_j \in D \text{ else } 0$$

where $\text{BERT}(\cdot)_{v_j}$ is the contextual token representation from the last layer of BERT for token v_j and P is a linear projection that maps the token representation into a scalar weight. Note that if v_j occurs in a query or a document multiple times, we select the maximum scalar weight for v_j as the weight of the token in the text. Exactly as above, this formula is equivalent to the dot product of the query and the document representation in a vector space with a basis defined by the vocabulary: $\text{Sim}(Q, D) = \langle \mathbf{E}_Q, \mathbf{E}_D \rangle$.

Following Contriever, we use independent cropping to create pseudo-positive pairs for contrastive learning. However, we choose not to use the MoCo [10] method to create a large negative candidate

pool, as the representations in the pool are not generated by the latest model parameters. As an alternative, we increase the number of negative examples by increasing the batch size to a large number (16384 in our experiments). To achieve such a large batch size with limited GPU memory, we adopt the Gradient Caching [9] method. The training objective for our unsupervised model is the same as training a dense retriever, i.e., using the infoNCE loss.

3.3 BM25 \oplus BM26 = BM51

Since BM25 and BM26 are both unsupervised lexical retrieval models, we can combine their results to create a hybrid retrieval model that itself remains in the space of unsupervised methods. While linear combination of retrieval scores or reciprocal rank fusion are popular methods for fusion, we create a novel “sparse-sparse” hybrid by vector concatenation in this work. We name this model “BM51” as it is the fusion of BM25 and BM26 (25+26=51).

Specifically, the query representation and document representation of the BM51 retrieval model are computed by:

$$\mathbf{E}_{Q_{BM51}} = \mathbf{E}_{Q_{BM25}} \oplus \alpha \cdot \mathbf{E}_{Q_{BM26}}$$

$$\mathbf{E}_{D_{BM51}} = \mathbf{E}_{D_{BM25}} \oplus \beta \cdot \mathbf{E}_{D_{BM26}}$$

where \oplus represents the vector concatenation operation and α, β adjust the relative contributions of each representation to the final model. We keep $\alpha = \beta = 1$ after normalizing the token weights, which means that BM25 and BM26 are weighted equally. The BM51 query-document similarity is also computed as a dot product.

A key feature of our fusion-via-vector-concatenation approach of BM51 is that it remains yet another lexical retrieval model. This means that top- k retrieval can take advantage of all existing software infrastructure that has been developed around inverted indexes. In addition, vector concatenation presents an advantage over other popular fusion techniques such as the linear combination of scores or reciprocal rank fusion, both of which involve performing retrieval twice and post-processing two ranked lists. In other words, BM51 can serve as a drop-in replacement for BM25. That is, it can serve as a first-stage ranker to provide candidates for further reranking, it can be further combined with learned models, etc.

3.4 BM26e & BM51e: Improving Lexical Match By Token Expansion

BM26 adopts the same parametric form as BM25, which means the lexical representation is restricted to the tokens that appeared in the original text. However, vocabulary mismatch issue is common to see in the lexical retrieval setting. The SPLADE work [6] utilizes the MLM layer of BERT to conduct expansion for the tokens in the original text. The purpose is to assign weights to the tokens that have close meaning to the existing tokens. SPLADE was originally focused on supervised learning settings. We adopt its architecture here to achieve the token expansion for unsupervised lexical representation learning. Following the same name convention, we denote this unsupervised neural lexical retriever with expansion as BM26e. Similar to the way we build BM51 = BM25 \oplus BM26, we craft the “with expansion” version BM51e = BM25 \oplus BM26e. We follow the same vector concatenation and token weights processing as described in Sec. 3.3.

4 EXPERIMENTAL SETUP

Data. Following the data crafting method in the Contriever work, we use Independent Cropping on the CCNet corpus [3] to create the unsupervised training data. Specifically, we concatenate all the natural documents in CCNet together and treat each 500-token split as a document. Each training pair is obtained by independent cropping from a random span of 256 tokens in a document. We also apply 10% random token deletion on both the pseudo query and the pseudo document as data augmentation.

Models. We train BM26 and BM26e following the method proposed in Section 3.2 and 3.3. We use bert-base-uncased as initialization and train 10k steps with batch size 16384. Our model is trained on a single AWS EC2 instance with 8 \times A100 40 GB GPUs. The training code is modified based on the open-source toolkit Tevatron [8].

BM51 is implemented exactly in the manner described in Section 3.3: We materialize the document vectors from BM25 and BM26 and then feed the concatenated vectors back into Pyserini for indexing and retrieval.

As additional baselines, we evaluate two variants that allow us to precisely attribute differences in model effectiveness:

- BM25_{wp}: This represents the BM25 weighting function applied directly to wordpiece tokens. That is, BM25_{wp} and BM26 share exactly the vector basis, except BM26 weights are computed by a transformer model.
- BM25': This represents BM25 \oplus BM25_{wp}, or the concatenation of “standard” BM25 and BM25_{wp} vectors. That is, BM25' and BM51 share exactly the same vector basis, differing only in how the BM26 portion of the weights are assigned.

Evaluation. We evaluate our models on the following datasets:

- MS MARCO Passage Ranking [1]: a web search dataset that contains 8.8 million passages and 6980 queries in the development set for evaluation. Effectiveness is measured by MRR@10 and Recall@1k.
- NaturalQuestions [14] and TriviaQA [12]: open-domain question answering datasets that use English Wikipedia as the corpus. Effectiveness is measured by top-20/100 retrieval accuracy.
- BEIR [24]: a collection of datasets for zero-shot evaluation of search and related tasks. We evaluate our models on the 13 publicly available datasets; effectiveness is measured in terms of nDCG@10.

5 RESULTS

Table 1 presents results on MS MARCO, NQ and TriviaQA, and Table 2 shows results on the 13 BEIR datasets. The results of the BM25 baselines are obtained by following the reproduction guides in the open-source Anserini IR toolkit [26]. The results of Contriever and cpt-text-S are copied from the original papers.

Evaluation of BM26. First, we compare BM26 with BM25. The high level conclusion is that they are roughly on par, which can be seen by comparing row (4) with row (1) in Table 1 and by comparing column (4) with column (1) in Table 2. In both tables, we performed significance testing using paired t -tests ($p < 0.05$); the symbols $\blacktriangle/\blacktriangledown$ indicate significant differences (better/worse). Despite the statistical significance of some differences, overall the effect sizes are relatively small: less than a point in most cases.

Method	Rep. S/D	MS MARCO		NQ		TriviaQA	
		MRR@10	R@1k	Top20	Top100	Top20	Top100
(1) BM25	S	18.4	85.3	62.9	78.3	76.4	83.2
(2) BM25 _{wp}	S	17.5	82.6	63.2	77.8	73.8	81.4
(3) BM25' = BM25 \oplus BM25 _{wp}	S	18.9	86.9	64.4	79.4	76.5	83.5
(4) BM26	S	19.0 [*]	88.6	63.9 [*]	78.8	74.3 [▼]	82.5
(5) BM26e	S	18.1	92.1	70.0 [*]	83.3	78.7 [*]	85.2
(6) BM51 = BM25 \oplus BM26	S	22.2[*]	92.0	68.1 [*]	81.6	77.6 [*]	83.9
(7) BM51e = BM25 \oplus BM26e	S	20.2 [*]	93.5	71.8[*]	84.2	79.3[*]	85.5
(a) Contriever	D	16.0	88.1	67.8	82.1	74.2	83.2
(b) cpt-text-S	D	19.9	-	65.5	77.2	75.1	81.7

Table 1: The retrieval effectiveness of BM26 and BM51 compared to baselines and other unsupervised retrieval models on MS MARCO, NQ and TriviaQA. We performed significance tests comparing BM26 and BM51 to BM25 indicated by ^{*}/[▼].

	(1) BM25 sparse	(2) BM25 _{wp} sparse	(3) BM25' sparse	(4) BM26 sparse	(5) BM26e sparse	(6) BM51 sparse	(7) BM51e sparse	(a) Contriever dense	(b) cpt-text-S dense
NFCorpus	32.2	31.4	32.8	32.4	33.3	34.2 [*]	34.7[*]	31.7	32.0
NQ	30.6	30.5	31.8	27.9 [*]	29.9	33.9 [*]	33.4[*]	25.4	-
HotpotQA	63.3	59.3	63.4	63.6	64.0 [*]	68.2[*]	67.9 [*]	48.1	51.5
FiQA	23.6	21.8	24.2	25.8 [*]	28.9 [*]	28.3 [*]	30.1 [*]	24.5	34.1
ArguAna	39.7	36.4	39.0	38.1 [*]	42.2 [*]	41.7 [*]	43.1[*]	37.9	38.7
Quora	78.9	73.0	78.0	77.5 [*]	82.4 [*]	82.3 [*]	83.5[*]	68.1	68.1
DBPedia	31.8	28.4	31.1	28.6 [*]	33.0	34.0 [*]	36.5[*]	29.3	27.2
SCIDOCS	14.9	13.8	15.0	15.0	16.0 [*]	15.9 [*]	16.8[*]	14.9	-
SciFact	67.9	67.2	69.4	66.7	67.9	69.8 [*]	70.9[*]	64.9	65.4
FEVER	65.1	65.8	67.3	74.4 [*]	75.1 [*]	75.2 [*]	77.9[*]	68.2	57.1
Climate-FEVER	16.5	15.8	17.1	15.8	18.9 [*]	18.9 [*]	20.0[*]	15.5	15.8
TREC-COVID	59.5	56.5	60.0	54.4	34.9 [▼]	65.6	53.9	27.4	52.9
Touche-2020	44.2	46.6	50.2	26.1 [*]	23.0 [▼]	34.7 [▼]	28.6 [▼]	19.3	21.0
average	43.7	42.0	44.7	42.0	42.3	46.4	46.0	37.7	-

Table 2: The nDCG@10 retrieval effectiveness of BM26 and BM51 compared to baselines and other unsupervised retrieval models on BEIR. We performed significance tests comparing BM26 and BM51 to BM25 indicated by ^{*}/[▼].

For BM25, the default Lucene analyzer used by Anserini generates 2.7M unique tokens for the MS MARCO passage corpus, which is orders of magnitude larger than the vocabulary space of the bert-base-uncased tokenizer (30,522 unique subwords). To isolate the effects of this difference, we turn to BM25_{wp}, which applies BM25 term weighting on wordpiece tokens. Thus, BM25_{wp} shares exactly the same representation space as BM26, providing a fair comparison between neural and heuristic encoders. Here, we see that BM26 outperforms BM25_{wp}.

Comparing BM26 to Contriever and cpt-text-S, We see BM26 scores four points higher than Contriever on BEIR dataset. Although cpt-text-S results are not available for all BEIR datasets, it does appear that BM26 obtains higher scores in most cases.

Evaluation of BM51. Secondly, we investigate how BM26 helps improve baseline BM25 via our vector concatenation technique. The main comparison condition is between BM51 (= BM25 \oplus BM26) and BM25. From Table 1, we see significant improvements across all three datasets, comparing row (6) to row (1). In the case of MS MARCO, we observe a nearly four point gain, which translates into a 21% relative improvement. From Table 2, comparing column (6) with column (1), we see that BM51 is significantly more effective than BM25 for all BEIR datasets but Touche-2020. However, Touche-2020 appears to be an outlier, as other existing models [6, 11, 22] also perform poorly. We emphasize that we obtain these gains without the use of manual labels.

To untangle the effectiveness improvements due to BM26 from improvements from simply having a hybrid vocabulary base, we can compare BM51 to BM25' = BM25 \oplus BM25_{wp}. Recall that the latter represents the concatenation of BM25 and BM25_{wp}: the BM25 components are identical, and as explained above, BM25_{wp} and BM26 share exactly the same representation space, differing only in how term weights are computed. We see that BM51 consistently outperforms BM25', which demonstrates that BM26 provides additional relevance signals beyond a simple vocabulary space hybrid.

We see that BM51 is more effective than Contriever and cpt-text-S across most datasets. On MS MARCO, BM51 outperforms Contriever by over six points, and on BEIR, nine points.

Evaluation of BM26e and BM51e. Finally, we see that allowing terms expansion gives a further improvement in most cases. For example, BM26e improves the top-20 accuracy of NQ and TriviaQA by 6.1 and 4.3 points over BM26. And on BEIR datasets, the BM51e variant achieves the best effectiveness across the unsupervised retrieval models on 9 out of 13 datasets.

6 CONCLUSION

In this work, we propose an unsupervised sparse retrieval model called BM26 that adapts techniques originally designed for dense retrieval to a sparse lexical representation space. We show that our unsupervised lexical approach has the potential to help build a "better BM25" system in the age of neural models.

REFERENCES

- [1] Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, Mir Rosenberg, Xia Song, Alina Stoica, Saurabh Tiwary, and Tong Wang. 2018. MS MARCO: A Human Generated MACHine Reading COMprehension Dataset. *arXiv:1611.09268v3* (2018).
- [2] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading Wikipedia to Answer Open-Domain Questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vancouver, Canada, 1870–1879.
- [3] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised Cross-lingual Representation Learning at Scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 8440–8451.
- [4] Zhuyun Dai and Jamie Callan. 2020. Context-Aware Term Weighting For First Stage Passage Retrieval. In *Proceedings of the 43rd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2020)*. 1533–1536.
- [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota, 4171–4186.
- [6] Thibault Formal, C. Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2021. SPLADE v2: Sparse Lexical and Expansion Model for Information Retrieval. *arXiv:2109.10086* (2021).
- [7] Luyu Gao and Jamie Callan. 2021. Unsupervised Corpus Aware Language Model Pre-training for Dense Passage Retrieval. *arXiv:2108.05540* (2021).
- [8] Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2022. Tevatron: An Efficient and Flexible Toolkit for Dense Retrieval. *arXiv:2203.05765* (2022).
- [9] Luyu Gao, Yunyi Zhang, Jiawei Han, and Jamie Callan. 2021. Scaling Deep Contrastive Learning Batch Size under Memory Limited Setup. In *Proceedings of the 6th Workshop on Representation Learning for NLP (Repl4NLP-2021)*. Online, 316–321.
- [10] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. 2020. Momentum Contrast for Unsupervised Visual Representation Learning. In *Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 9726–9735.
- [11] Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised Dense Information Retrieval with Contrastive Learning. *arXiv:2112.09118* (2021).
- [12] Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vancouver, Canada, 1601–1611.
- [13] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Online, 6769–6781.
- [14] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural Questions: A Benchmark for Question Answering Research. *Transactions of the Association for Computational Linguistics* 7 (2019), 452–466.
- [15] Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. Latent Retrieval for Weakly Supervised Open Domain Question Answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy, 6086–6096.
- [16] Jimmy Lin. 2021. A Proposed Conceptual Framework for a Representational Approach to Information Retrieval. *arXiv:2110.01529* (2021).
- [17] 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 *Proceedings of the 44th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2021)*. 2356–2362.
- [18] Jimmy Lin, Rodrigo Nogueira, and Andrew Yates. 2021. *Pretrained Transformers for Text Ranking: BERT and Beyond*. Morgan & Claypool Publishers.
- [19] Sheng-Chieh Lin, Jheng-Hong Yang, and Jimmy Lin. 2021. In-Batch Negatives for Knowledge Distillation with Tightly-Coupled Teachers for Dense Retrieval. In *Proceedings of the 6th Workshop on Representation Learning for NLP (Repl4NLP-2021)*. Online, 163–173.
- [20] Xueguang Ma, Ronak Pradeep, Rodrigo Nogueira, and Jimmy Lin. 2022. Document Expansions and Learned Sparse Lexical Representations for MS MARCO V1 and V2. In *Proceedings of the 45th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2022)*. Madrid, Spain, 3187–3197.
- [21] Antonio Mallia, Omar Khattab, Torsten Suel, and Nicola Tonello. 2021. Learning Passage Impacts for Inverted Indexes. In *Proceedings of the 44th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2021)*. 1723–1727.
- [22] Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, Johannes Heidecke, Pranav Shyam, Boris Power, Tyna Eloundou Nekoul, Girish Sastry, Gretchen Krueger, David Schnurr, Felipe Petroski Such, Kenny Hsu, Madeleine Thompson, Tabarak Khan, Toki Sherbakov, Joanne Jang, Peter Welinder, and Lilian Weng. 2022. Text and Code Embeddings by Contrastive Pre-Training. *arXiv:2201.10005* (2022).
- [23] Stephen E. Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Foundations and Trends in Information Retrieval* 3, 4 (2009), 333–389.
- [24] Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. BEIR: A Heterogeneous Benchmark for Zero-shot Evaluation of Information Retrieval Models. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- [25] 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 *Proceedings of the 9th International Conference on Learning Representations (ICLR 2021)*.
- [26] Peilin Yang, Hui Fang, and Jimmy Lin. 2017. Anserini: Enabling the Use of Lucene for Information Retrieval Research. In *Proceedings of the 40th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2017)*. Tokyo, Japan, 1253–1256.
- [27] Shengyao Zhuang and Guido Zuccon. 2021. Fast Passage Re-ranking with Contextualized Exact Term Matching and Efficient Passage Expansion. *arXiv:2108.08513* (2021).