# NV-Retriever: Improving text embedding models with effective hard-negative mining

Gabriel de Souza P. Moreira*
NVIDIA
São Paulo, Brazil
gmoreira@nvidia.com

Radek Osmulski*
NVIDIA
Brisbane, Australia
rosmulski@nvidia.com

Mengyao Xu*
NVIDIA
Santa Clara, USA
mengyaox@nvidia.com

Ronay Ak*
NVIDIA
Sarasota, USA
ronaya@nvidia.com

Benedikt Schifferer*
NVIDIA
Berlin, Germany
bschifferer@nvidia.com

Even Oldridge*
NVIDIA
Vancouver, Canada
eoldridge@nvidia.com

## Abstract

Text embedding models have been popular for information retrieval applications such as semantic search and Question-Answering systems based on Retrieval-Augmented Generation (RAG). Those models are typically Transformer models that are fine-tuned with contrastive learning objectives. One of the challenging aspects of fine-tuning embedding models is the selection of high quality hard-negative passages for contrastive learning. In this paper we introduce a family of positive-aware mining methods that use the positive relevance score as an anchor for effective false negative removal, leading to faster training and more accurate retrieval models. We provide an ablation study on hard-negative mining methods over their configurations, exploring different teacher and base models. We further demonstrate the efficacy of our proposed mining methods at scale with the NV-Retriever-v1 model, which scores 60.9 on MTEB Retrieval (BEIR) benchmark and placed 1st when it was published to the MTEB Retrieval on July, 2024.

## Keywords

Text retrieval, embedding models, hard-negative mining, contrastive learning, transformers.

## 1 Introduction

Text retrieval is critical for a range of information retrieval applications such as search, question answering, semantic textual similarity, and item recommendation. It is also vital to the field of Retrieval Augmented Generation (RAG)[16, 27], which enables Large Language Models (LLM) to access external context without modifying model parameters.

Dense embedding models are a key component of text retrieval, semantically representing queries and passages (pieces of content) with low token overlap and generalize for out-of-domain corpus. Retrieval systems that index passages into embeddings can efficiently retrieve relevant passages for a query using Maximum Inner Product Search (MIPS) [16].

There has been a growing interest in text embedding models from both academia and industry leading to the recent release of many models including E5[33], GTE[17], and Jina[9]. To consistently compare the accuracy of available text embedding models the MS-MARCO[1], BEIR[32] and MTEB[22] benchmarks on public

(HuggingFace[36]) leaderboards provide an important focal point of comparison for different embedding based tasks.

Embedding models are trained with Contrastive Learning (CL)[4] to maximize the similarity between the embeddings of a query and the relevant passages (positives) while minimizing the similarity with embeddings of passages not relevant to the query (negatives)[11].

When preparing data for training embedding models, *hard-negative mining* is typically used to select negative passages for queries. It leverages a teacher retrieval model to find passages somewhat relevant to the query, making it harder for the contrastive loss to differentiate positives from negatives, resulting in more efficient and effective fine-tuning of the embedding model.

Despite its importance for embedding model fine-tuning, the *hard-negative mining* methods remain underexplored or poorly detailed, particularly in papers [12, 17, 17, 23, 34] introducing models topping the MTEB leaderboard, as these studies primarily focus on model architectures, fine-tuning approaches and data blends for training.

The main contributions of this research are threefold:

- **Positive-aware hard-negative mining methods**. We introduce a family of mining methods that are capable of improving contrastive learning, removing potential false negatives and increasing the accuracy of text embedding models, as described in Section 3.1;
- **Investigation on hard-negative mining best practices**. We present in Section 3.2 our methodology and research questions, and in Section 4.1 compare different hard-negative mining methods over their configurations, with different teacher and base models, demonstrating how sensitive retrieval models are to hard-negative mining choices;
- **Scaling the mining methods to achieve state-of-the-art text retrieval: *NV-Retriever-v1***. We demonstrate in Section 4.4 the effectiveness of our proposed mining methods at scale in the training of *NV-Retriever-v1*, a state-of-the-art embedding model which achieved 1st place on MTEB Retrieval leaderboard[1] when published on MTEB on July

---

[1]https://huggingface.co/spaces/mteb/leaderboard

11th, 2024. We describe *NV-Retriever-v1* train blends, hyper-parameters, architecture, and particularly how our positive-aware hard-negative mining methods were crucial to its leading retrieval accuracy.

## 2 Background

In this section we discuss related work on text embedding models and hard-negative mining.

### 2.1 Text embedding models

Text embedding models represent variable-length text as a fixed dimension vector that can be used for downstream tasks.

A seminal work on embedding sentences was Sentence-BERT [28], which proposed a modification to BERT network to represent pairs of related short texts in the same embedding space by using siamese networks (query,positive passages) or triplet networks (query,positive,negative passages). They also explored different options of objective functions and embedding pooling.

Contrastive learning was popularized by SimCLR[4] superior performance than classification-based losses[28] for embeddings. The (DPR)[11] proposed a bi-encoder architecture where separate BERT encoders (with no shared weights) represented query and passage, whose output embeddings are used for CL.

The E5 [33] family of models leveraged two-stage training of embedding models, pre-training with unsupervised pairs of text (e.g. neighboring text spans, title-abstract) and fine-tuning with supervised data (e.g. question-answer, text-summary, search-relevant passages). E5 models are available in different sizes, depending on the base model: MiniLM[35] and BERT[7]. E5 was the first dense retrieval model to beat BM25[30] sparse model baseline on BEIR[32] without using labeled data (with its unsupervised version). On second stage, they fine-tune E5 unsupervised model with labeled data from MS-MARCO dataset [1], Natural Questions (NQ) and NLI for superior performance.

E5-Mistral [34] embedding model proposed using a decoder model instead of encoder model (like BERT) as a base model. They choose Mistral-7B[10], which is already extensively pre-trained on web-scale data. It was fine-tuned in a single round with synthetic data generated by LLMs for different tasks and languages and with a small amount of labeled data.

The BEIR benchmark[32] has become the standard evaluation for zero-shot text retrieval with 18 retrieval datasets. Later MTEB [22] was introduced as a more comprehensive text embeddings benchmark, with 56 datasets distributed over 7 tasks for the English subset: retrieval, reranking, classification, clustering, pair classification, summarization, and semantic textual similarity. The MTEB retrieval task is composed by 15 selected datasets from BEIR. We evaluate on MTEB Retrieval in Section 4.4.

### 2.2 Hard-negative mining for fine-tuning embedding models

Contrastive Learning (CL) requires triples of query, positive passage and negative passages. Negative passages can be either labeled by humans or, more commonly be sampled from the corpus.

A basic approach for selecting negative passages is using the positive passages from other examples (queries) in the batch (*in-batch negatives*) [4, 11]. This is computationally efficient because the embeddings for those passages were already generated by the model forward pass, although, the batch size limits the number of negatives. Some proposals to increase the number of negatives are keeping a memory bank with embeddings from past batches [33, 40] or combining batches from different GPUs [25] (*cross-batch*).

The *in-batch* negatives are random with respect to the query and easy to discriminate. Thus, they are very uninformative for CL, as their loss / gradients are low and contribute little to model convergence [39]. On the other hand, using challenging (*hard*) negatives can lift the upper bound of gradient norm, reduce the variance of stochastic gradient estimation, and lead to faster learning [39].

Hard-negatives can be mined from the corpus of passages by retrieving examples similar to the query that are not labeled positive[8, 11, 37, 39]. Sparse or dense (embedding) retrieval models can be used for mining. DPR [11] uses one or two *hard* negatives mined with BM25[30] in addition to the in-batch negatives that help them to finetune models with better accuracy.

Some methods have been proposed for incremental hard-negative mining during training. ANCE[39] asynchronously refreshes the ANN index with the updated passage embeddings and re-mines the hard-negatives for each question. NGAME[5] clusters queries and positive embeddings and uses those clusters to prepare *negative mining-aware mini-batches*. The batches contain data points close to each other in the embedding space, so that hard-negatives can be found among the in-batch samples, whose embeddings are efficiently computed. Incremental mining methods such as ANCE and NGAME are in general complex to implement and costly to compute for large corpus, which is why for most of the top models in MTEB [12, 17, 23, 34] the negatives are mined upfront (before training) using a pre-trained model.

*2.2.1 False negatives.* In [25] they found that naive mining of hard-negatives may select many false negatives. Their experiment on mining negatives from the MS-Marco dataset found that about 70% of passages most similar to the queries should be actually labeled as positive. They propose to *denoise* hard-negatives, i.e., ignore the potential false negatives by filtering out the ones with high relevance score to the query.

Some works try to denoise hard-negatives retrieved from embedding models with cross-encoder models, like in RocketQA [25, 29], or with more powerful decoder LLMs[15]. These approaches might be costly to run on large train sets depending on the model size, as they require inference for every (query, negative) pair.

While most top performing models on MTEB like *e5-mistral-7b-instruct* [34], *Linq-Embed-Mistral* [12], *NV-Embed-v1* [14], *gte-large-en-v1.5* [17], *nomic-embed-text-v1* [23] leverage hard-negative mining in their fine-tuning, they do not explore or describe in detail their methodology to decide which model and method to use for mining. In *snowflake-arctic-embed-l* [19], the authors include an ablation study of three maximum negative score thresholds (0.4, 0.5 and 0.8) for hard-negative mining. The *SFR-Embedding-Mistral* [18] blog provides an ablation on sampling hard-negatives from three different ranges of top-k candidates (0-100, 30-100, 50-100) and reports that the range 30-100 helps eliminate false negatives and

improve model performance. They also compare different teacher models for mining (BM25, BGE-base [38], E5-Mistral[34] and their own SFR-Embedding-Mistral) and show that more powerful models can yield more effective hard-negatives, which we also observe in our ablation presented in Section 4.1.

## 3 Methodology

In this section, we introduce our main contribution - the positive-aware hard-negative mining methods - the research questions, and experiments setup.

### 3.1 Positive-aware hard-negative mining methods

A popular contrastive learning loss for retrieval models is *InfoNCE* [24] shown in Equation 1, where $sim(\cdot)$ is a similarity function (e.g., cosine similarity or dot product), $d^+$ is a positive relevant passage, $d^-$ is one of the $N$ negative passages and $\tau$ is the temperature parameter. The objective is to maximize the embedding similarity between the query and positive passage, whilst minimizing the similarity between the query and negative passages.

$$\mathcal{L}(q, d^+, d_N) = -\log p(d = d^+ \mid q)$$

$$= -\log \frac{\exp(sim(q, d^+)/\tau)}{\sum_{d_i \in \{d^+\} \cup d_N} \exp(sim(q, d_i)/\tau)}, \qquad (1)$$

As discussed in Section 2.2, random negatives are uninformative for training and hard-negatives (with higher relevance score $sim(q, d_i)$) can provide faster convergence.

The basic method for mining hard-negatives is to select the top-k most similar candidates to the query (ignoring positive passages) which we name **Naive Top-K**.

However, the hard-negative mining process may introduce some false negatives, i.e., passages $d_i$ are relevant to the query but not annotated as positives $d^+$, which add noise to constrastive learning. False negatives can be present after negative mining when annotation of positive passages is not comprehensive enough or when mining is done on a large corpus. For example, in open-domain question answering (OpenQA) datasets such as MS MARCO [1] and Natural Questions [13], the question answer might be supported by many paragraphs from Wikipedia[2, 13] or the web [1].

Some methods for filtering out false negatives have been proposed in the literature:

- **Top-K shifted by N** - Selects the top-k negatives after a rank $N$, e.g., Top-10 shifted by 5 would ignore the first 5 negatives and consider negatives between rank 5 and 15 [18, 38];
- **Top-k with absolute threshold (*TopK-Abs*)** - Ignores negatives with relevance score higher than an absolute threshold [14, 19, 25];

These methods have some important limitations. *Top-K shifted by N* is a basic method that does not take into account the relevance score of the negative with respect to the query and might even throw away valuable hard-negatives or keep false negatives. *TopK-Abs* uses absolute thresholds as maximum negative scores with respect to the query, regardless the positive passage relevance.

Motivated by those limitations, we have designed a family of *positive-aware hard-negative mining methods*. Our methods are simple, generic and can be applied on top of any teacher model (e.g., embedding, reranker) that retrieves top-k candidates semantically relevant to the query. Our methods take advantage of the information from the positive relevance score to help identify and eliminate potential false negatives. The base method is described in Algorithm 1. It iterates over the retrieved top-k negatives and filters out potential false negatives using one of the following filter criteria:

- **Top-k with margin to positive threshold (*TopK-MarginPos*)** - The maximum threshold for negative scores is the positive score minus an absolute margin (Algorithm 2).
- **Top-k with percentage to positive threshold (*TopK-PercPos*)** - The maximum threshold for negative scores is a percentage of the positive score (Algorithm 3).

---

**Algorithm 1** Positive-aware hard-negative mining base method

1: **procedure** POSITIVEAWARENEGATIVEMINING($p$, $N$)
2:     valid_negatives = []
3:     **for** all $n$ in negatives $N$ **do**
4:         **if** *filter_fn(p,n)* **then**   ▷ Call one of the positive-aware mining filters: *TopKMarginPosFilter* or *TopKPercPosFilter*
5:             valid_negatives.append(n)
6:         **end if**
7:     **end for**
8:     **return** valid_negatives
9: **end procedure**

---

**Algorithm 2** TopK-MarginPos negatives filter

1: **procedure** TOPKMARGINPOSFILTER($p$, $n$)
2:     abs_margin ← <CONFIG_MARGIN>
3:     **return**← ($n$.rel_score < $p$.rel_score - abs_margin)
4: **end procedure**

---

**Algorithm 3** TopK-PercPos negatives filter

1: **procedure** TOPKPERCPOSFILTER($p$, $n$)
2:     perc_margin ← <CONFIG_MARGIN>
3:     **return** ← ($n$.rel_score < $p$.rel_score * perc_margin)
4: **end procedure**

---

In Section 4, we compare mining methods by fine-tuning different base models on negatives mined using distinct teacher models, and demonstrate the effectiveness of our proposed *positive-aware mining methods* at scale to achieve state-of-the-art retrieval accuracy. We also provide an ablation study on the margin configurations of our methods to serve as a reference.

### 3.2 Research Questions

In this paper, we investigate the following Research Questions with respect to hard-negative mining:

- *RQ1*. How much does mining hard-negatives with different teacher models affect the accuracy of the downstream fine-tuned embedding models?
- *RQ2*. Can ensembling hard-negatives mined from different teacher models improve results?
- *RQ3*. How does different hard-negative mining methods for fine-tuning compare on the evaluation accuracy?

The next sections present the experiments results for those research questions, comparing different teacher models for mining and the effect of ensembling hard-negatives from distinct models. We also provide a comprehensive experimentation on hard-negative mining methods over their thresholds for filtering out false negatives.

## 3.3 Experiments setup

*3.3.1 Training.* In our general setup, embedding models are fine-tuned[2] with constrastive learning from *e5-large-unsupervised* or *Mistral-7B-v0.1* base models, with hard-negatives mined using selected mining methods and teacher models of different sizes.

Train set is composed of Natural Questions (NQ) [13] [3], Stack Exchange (2023 dump) [4] and SQUAD [26] [5] datasets ( 287k examples).

*3.3.2 Evaluation.* We selected three Question-Answering datasets from MTEB Retrieval / BEIR benchmark - NQ, HotpotQA and FiQA-2018 [32] - as they are more relevant for Q&A RAG systems.

For *RQ3*, we also perform a scaled experiment setup using a larger train set (same from *NV-Retriever-v1* model) and evaluate on full MTEB Retrieval benchmark, as described in Section 4.4.

## 4 Experiment results and discussion

Here we investigate the research questions presented in Section 3.2.

## 4.1 RQ1. Using different teacher models for mining

We have selected a number of popular text embedding models as teacher models for mining hard-negatives. Those models represent different architectures, model sizes and retrieval accuracy.

- *e5-large-unsupervised* [6] (334M params) - The *E5* model pre-trained on unsupervised data with CL [33];
- *e5-large-v2* [7] (334M params) - An *E5* model fine-tuned on top of *e5-large-unsupervised* with supervised data[33]
- *snowflake-arctic-embed-l* [8] (334M params) - Member of the *artic-embed* models which is trained in two rounds (with supervised and unsupervised data) like *E5* model, with improvements on data and training that lead to higher retrieval accuracy [19]

- *e5-mistral-7b-instruct* [9] (7.1B params) - The decoder-only Mistral model fine-tuned with CL to create an embedding model [34];
- *NV-Embed-v1* [10] (7.8B params) - A Mistral-based embedding model with some modifications including bi-directional and latent attention[14].

We mine 4 hard-negatives[11] with those different teacher models for every question in the train set. This process results in one train set per teacher model, which is then used for fine-tuning the base model (*E5-large-unsupervised*).

We can see from the results on Table 1 that the worse retrieval accuracy was obtained by using negatives mined with the BM25 sparse retrieval model[30] followed by random negatives, which was surprising and opposite to what was found in [11]. The *E5-large-unsupervised* dense retrieval model, pre-trained only on unsupervised data from the web, provided better accuracy.

The next teacher models are trained on retrieval supervised data, particularly for Question-Answering RAG systems. The *e5-large-v2* and *snowflake-arctic-embed-l* use the *E5* architecture (334M params) and perform better than the baselines. The best teacher models were the larger *NV-Embed-v1* and *e5-mistral-7b-instruct*, both based on the Mistral 7B architecture. They provide better hard-negatives for CL, which result on higher accuracy for the fine-tuned models.

**Table 1: Evaluation (NDCG@10) of fine-tuned *e5-large-unsupervised* embedding models with hard-negatives mined using different teacher models**

| Teacher models | Avg. | NQ | HotpotQA | FiQA |
|---|---|---|---|---|
| BM25 | **0.5002** | 0.5307 | 0.5774 | 0.3923 |
| random | **0.5248** | 0.5123 | 0.6151 | 0.4471 |
| e5-large-unsupervised | **0.5494** | 0.5541 | 0.6247 | 0.4694 |
| e5-large-v2 | **0.5704** | 0.6058 | 0.6435 | 0.4618 |
| snowflake-arctic-embed-l | **0.5728** | 0.6118 | 0.6331 | 0.4735 |
| NV-Embed-v1 | **0.5744** | 0.6092 | 0.6355 | 0.4785 |
| e5-mistral-7b-instruct | **0.5810** | 0.6241 | 0.6434 | 0.4757 |

## 4.2 RQ2. Ensembling hard-negatives from different teacher models

Ensembling outputs from different models is a common practice in machine learning to improve predictions accuracy, for providing a more robust estimator [6, 20, 31].

In particular, we investigated the similarity of the top-4 hard-negatives mined by four teacher models and noticed a low level of agreement (*jaccard similarity* lower than 30%), as you can see on Appendix A. For this reason, we decided to explore ensembling to try improving the quality of the hard-negatives.

We explored two methods to combine hard-negatives mined from four different *E5* and *Mistral* based embedding models - *e5-large-v2*, *snowflake-arctic-embed-l*, *NV-Embed-v1*, and *e5-mistral-7b-instruct* -

---

[2]We use for fine-tuning https://github.com/microsoft/unilm/tree/master/simlm
[3]https://ai.google.com/research/NaturalQuestions
[4]https://archive.org/details/stack-exchange-data-dump-2023-09-12
[5]https://rajpurkar.github.io/SQuAD-explorer/
[6]https://huggingface.co/intfloat/e5-large-unsupervised
[7]https://huggingface.co/intfloat/e5-large-v2
[8]https://huggingface.co/Snowflake/snowflake-arctic-embed-l

[9]https://huggingface.co/intfloat/e5-mistral-7b-instruct
[10]https://huggingface.co/nvidia/NV-Embed-v1
[11]For the experiments comparing different mining models, in order to remove potential false negatives, we use the *TopK-PercPos* mining method configured to set the maximum threshold for hard-negative scores as 95% of the corresponding positive score. That mining method and configuration performed best in our ablation, as described in Section 4.3.

which are described next. Each ensembling method returns 4 hard-negatives for each example (query,positive).

- *Cross-sample ensembling* - For each example, samples a teacher model to obtain all the negatives.
- *Intra-sample ensembling* - For each example, selects the top-1 mined negative from each teacher model.

The evaluation results are shown in Table 2. In the first row, we have as baseline the model trained with hard-negatives from the best teacher model (from *RQ1*): *e5-mistral-7b-instruct*. We can see that the *cross-sample ensembling* method does not lead to better negatives than the best teacher model.

The *Intra-sample ensembling* method turned out to be more effective. It may lead to duplicate hard-negatives, as for some examples the teacher models might agree on the 1st hard-negative. Thus, we tried two variations: keeping duplicates (*no-dedup*) or removing duplicates (*dedup*) and replacing them by the next unique hard-negative from teacher models sorted by their accuracy. Surprisingly, we found that it was better to keep the duplicate hard-negatives for training. A possible explanation could be that if models agree on the 1st hard-negative, keeping it duplicate will increase its importance in the cross-entropy loss.

**Table 2: Evaluation (NDCG@10) of fine-tuned *e5-large-unsupervised* embedding model with ensembles of hard-negatives mined using 4 different models: *e5-large-v2, snowflake-arctic-embed-l, NV-Embed-v1, e5-mistral-7b-instruct***

| Teacher Model / Ensemble method | Avg. | NQ | HotpotQA | FiQA |
|---|---|---|---|---|
| e5-mistral-7b-instruct (baseline) | **0.5810** | 0.6241 | 0.6434 | 0.4757 |
| *Ensembled hard-negatives from 4 models* | | | | |
| Cross-sample ensembling | **0.5806** | 0.6279 | 0.6384 | 0.4611 |
| Intra-sample ensembling (dedup) | **0.5804** | 0.6324 | 0.6302 | 0.4716 |
| Intra-sample ensembling (no-dedup) | **0.5825** | 0.6357 | 0.6298 | 0.4820 |

## 4.3 RQ3.a Comparing methods for mining hard-negatives

To investigate RQ3 we performed a comprehensive number of experiments testing the different hard-negative mining methods described in Section 3.1.

*4.3.1 Ablation study on mining methods configurations.* We performed an ablation study where we fine-tune the *e5-large-unsupervised* model (*Base setup*) with hard-negatives mined by different mining methods, for a range of their configuration.

For *TopK-Abs*, *TopK-MarginPos* and *TopK-PercPos* methods, the interval for the threshold/margin values config was [0, 1], with increments of 0.05[12]. We use as teacher model the *e5-mistral-7b-instruct*, which performed best for mining hard-negatives (Section 4.1).

We present plots with the results for the different configuration choices of each mining method. The evaluation metric reported is the average of NDCG@10 for the three selected BEIR Q&A datasets (NQ, HotpotQA and FiQA).

---

[12]For *TopK-PercPos*, we included 0.98 and 1.05 thresholds for more fine-grained analysis around the best threshold.

We can see that the basic *Top-k shifted by N* (Figure 1.a) method provides its best accuracy when discarding the top-10 ranked hard-negatives from the training, as they may be too strong and have a higher chance of containing false negatives.
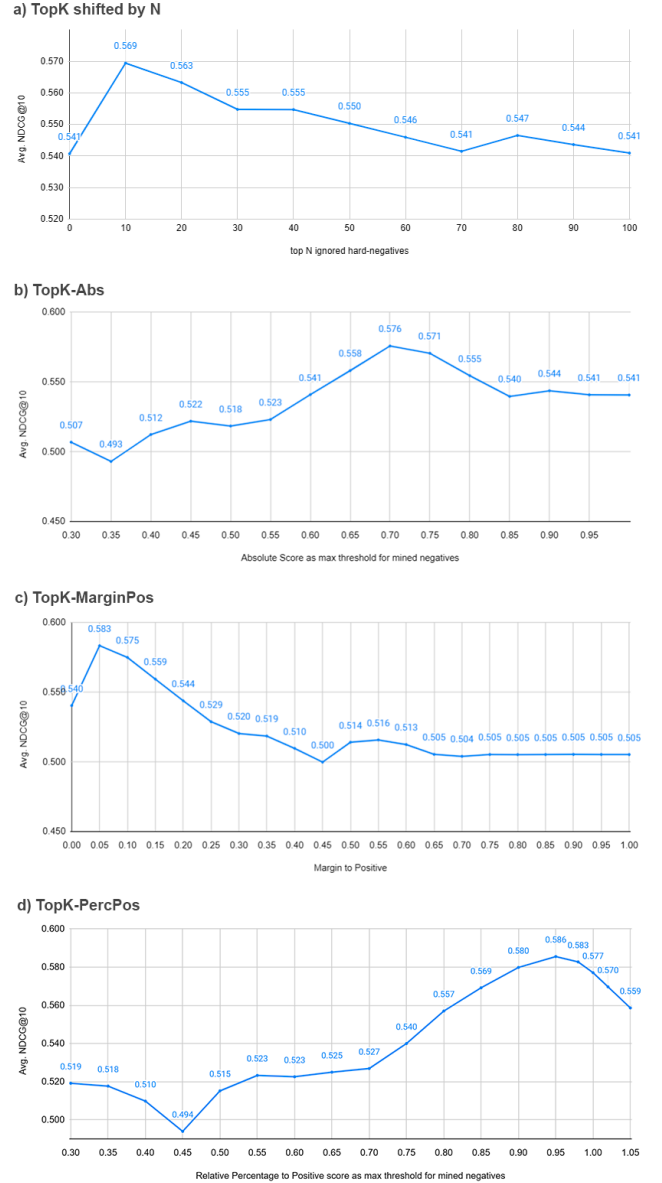


**Figure 1: Ablation study of the negative mining methods over their configuration.**

For *TopK-Abs* (Figure 1.b) – a traditional method to remove mined false negatives – the best configuration was to use the absolute score of 0.7 as the maximum threshold for negative scores. It seems that a higher threshold could include more false negatives, and a lower threshold could mine weaker negatives that would not be very informative for CL.

We propose in this paper *positive-aware mining methods*: *TopK-MarginPos* and *TopK-PercPos*. For *TopK-MarginPos* mining method (Figure 1.c), we observe that subtracting a small margin (0.05) from the positive score and setting it as the maximum threshold for the hard-negatives helped remove potential false negatives, and larger margins end up penalizing the model accuracy. For the *TopK-PercPos* mining method (Figure 1.d), we find that setting the maximum threshold for mined negatives as 95% of the positive scores was the optimum configuration. We can see that training on high-scoring negatives with respect to positives is detrimental to model accuracy, as it decreases after relative margin goes higher than 95% threshold.

*4.3.2 Sampling from mined negatives.* After a mining method returns an ordered list of negatives, the typical approach is selecting the top-k candidates, but some research works have proposed sampling among top-k to add some relevance diversity among the selected hard-negatives:

- **Sampled Top-k** - Samples *n* negatives from the top-k most relevant ones [3, 15, 25] or from a range of negatives based on its relevance rank, e.g. from 30-100 range like in [18];
- **Top-1+sampled top-k** - Selects the top-1 hard-negative (to secure a strong one) and samples $n-1$ negatives like described in the *Sampled Top-k* method.

We performed experiments of those two sampling approaches using *TopK-PercPos* (at 95% config)[13]. We experimented with *k* in the range of [10,100] with increments of 10, as shown in Figure 2.
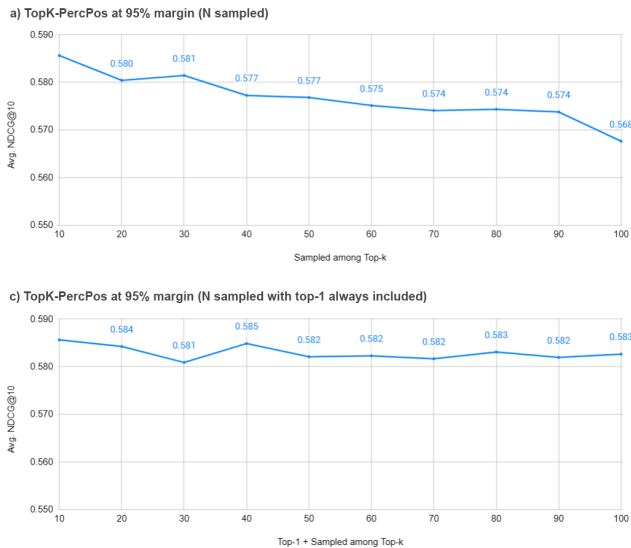


**Figure 2: Ablation study of *TopK-PercPos* negative mining method with 95% threshold + (*Sampled Top-k*) or (*top1 + Sampled*) sampling method. Four negatives are sampled among different top *k***

The best configuration was sampling the four negatives from top-10. Increasing the $k > 10$ for sampling dropped the retrieval

accuracy for *Sampled Top-k* (Figure 2.a), but not much for *Top-1+sampled top-k* sampling method (Figure 2.b), as the latter ensures the top-1 mined negative is always included.

*4.3.3 Comparing mining methods at their best configuration from the ablation.* We summarize in Table 3 the results of the best configurations of each mining method found in the ablation (Section 4.3).

**Table 3: NDCG@10 of fine-tuned *e5-large-unsupervised* model with top-4 hard-negatives from different mining methods, at their best configuration found in the ablation.**

| Mining Method | Config. | Avg. | NQ | HotpotQA | FIQA |
|---|---|---|---|---|---|
| *Naive Top-K* | - | 0.5407 | 0.5445 | 0.6120 | 0.4658 |
| *Top-K shifted by N* | N=10 | 0.5695 | 0.6007 | 0.6384 | 0.4693 |
| *TopK-Abs* | 0.7 | 0.5759 | 0.6133 | 0.6396 | 0.4748 |
| *TopK-MarginPos* | 0.05 | 0.5835 | 0.6338 | 0.6400 | 0.4766 |
| *TopK-PercPos* | 95% | **0.5856** | **0.6369** | **0.6414** | **0.4784** |
| *Sampling method on negatives from TopK-PercPos (at 95% threshold)* | | | | | |
| *TopK-PercPos (sampled)* | top-10 | 0.5856 | 0.6369 | 0.6414 | 0.4786 |
| *TopK-PercPos (top1+sampled)* | top-10 | 0.5857 | 0.6369 | 0.6414 | 0.4787 |

We can notice from the lower metrics of *Naive Top-K* compared to the other mining models how important it is to remove the potential false negatives (noise) for fine-tuning. The basic *Top-K shifted by N* method improves the obtained retrieval accuracy by ignoring the top (10) mined negatives, followed by *TopK-Abs*, that ignores any negative scoring higher than a max. threshold (0.7).

Our *positive-aware hard-negative mining methods* – *TopK-MarginPos* and *TopK-PercPos* perform the best. In particular, the *TopK-PercPos* was the most effective hard-negative mining method, with its threshold set to 95% of the positive score.

Sampling among top-k on top of *TopK-PercPos* mining method did not help to further improve results for *e5-large-unsupervised* base model, but it might help in some cases, e.g. for Mistral 7B v0.1 base model, as we show in Table 4 and discuss in the next section.

*4.3.4 Replicating mining methods experiments with Mistral-7B-v0.1.* In order to verify if our comparison of mining methods would generalize for base models other than *e5-large-unsupervised* (334M params), we ran experiments fine-tuning another base model: *Mistral-7B-v0.1* [14] (7.1B params). As Mistral is a much larger model, we had to use only 1 hard-negative per example to avoid out of memory errors, and re-used the best configuration for each of the mining methods from the ablation with *e5-large-unsupervised* (as running the full ablation with Mistral would require a lot of compute).

The mining methods comparison for *Mistral-7B-v0.1* base model can be seen in Table 4. It is noticeable the retrieval accuracy improvement when using Mistral larger base model compared to *e5-large-unsupervised* (Table 3) for fine-tuning. For *Mistral-7B-v0.1* base model, we can see that *Top-K shifted by N* also improves upon the *Naive Top-K*, while the traditional *TopK-Abs* mining method performs worse. On the other hand our *positive-aware mining methods* perform much better, with *TopK-PercPos* method also providing the best hard-negatives for finetuning.

---

[13]As we observed empirically that it is important to have strong negatives, in our sampling implementation we use *softmax* to create a probability distribution based on the negatives relevance score and sample based on that distribution.

[14]https://huggingface.co/mistralai/Mistral-7B-v0.1

**Table 4: NDCG@10 of *Mistral-7B-v0.1* base model with top-1 hard-negatives from different mining methods, at the best configuration found in the ablation**

| Mining Method | Config. | Avg. | NQ | HotpotQA | FIQA |
|---|---|---|---|---|---|
| *Naive Top-K* | - | 0.6214 | 0.6450 | 0.6733 | 0.5458 |
| *Top-K shifted by N* | N=10 | 0.6342 | 0.6247 | 0.7126 | 0.5653 |
| *TopK-Abs* | 0.7 | 0.6184 | 0.6416 | 0.6744 | 0.5391 |
| *TopK-MarginPos* | 0.05 margin | 0.6457 | 0.6694 | 0.7113 | 0.5565 |
| *TopK-PercPos* | 95% | **0.6479** | **0.6766** | **0.7053** | **0.5618** |
| *Sampling method on negatives from TopK-PercPos (at 95% threshold)* | | | | | |
| *TopK-PercPos (sampled)* | top-10 | 0.6499 | 0.6763 | 0.7063 | 0.5671 |

The last row in Table 4 applies the *Sampled Top-k*[15] sampling method on top of *TopK-PercPos* mining method. We can see that for Mistral base model it was slightly better to sample hard-negatives among top-10 instead of always using the top-1 negative.

## 4.4 RQ3.b Comparing mining methods at scale with *NV-Retriever-v1*

In this section, we present experiments for RQ3 performed at scale. We demonstrate that our proposed mining methods are the secret sauce to build the *NV-Retriever-v1* model, 1st place on the MTEB leaderboard when published. We leverage for these experiments the same setup we used for fine-tuning *NV-Retriever-v1*. It uses *Mistral-7B-v0.1* as base model and *E5-Mistral-7B* as teacher model for mining hard-negatives according to each mining method (at their best configuration from the ablation study). We trained models on the 15 retrieval datasets used to train the first stage of NV-Retriever-v1 (Appendix B.2) with a total of 728,160 examples. The model and training hyperparameters are the same from *NV-Retriever-v1*, as detailed in Appendix B.

This set of experiments demanded a lot of compute resources for requiring (1) mining negatives with a large teacher model (*E5-Mistral-7B*) for each mining method and train set, (2) fine-tuning large Mistral-7B models and (3) evaluating on full MTEB Retrieval.

We present in Table 5 the results for the models trained with the different mining methods, evaluated on the 15 datasets of the MTEB Retrieval benchmark. We can see in this scaled setup an even larger retrieval accuracy advantage of our proposed positive-aware mining methods compared to the baseline methods.

In particular, the model trained with *TopK-PercPos* mining method achieves the highest Avg. NDCG@10 (60.55), which would place 1st on the MTEB Retrieval / BEIR leaderboard leaderboard when *NV-Retriever-v1* model was published (see Appendix B.6). The *NV-Retriever-v1* model scored only slightly higher (60.9) on MTEB because after the first stage training on those retrieval datasets, it is further fine-tuned with some classification and clustering datasets[16]. Our *TopK-PercPos* method provided the best retrieval accuracy for 13 of the 15 BEIR datasets, for the other 2 datasets (*SCIDOCS* and *Touche-2020*) our *TopK-MarginPos* method was the best.

This set of large experiments makes it is evident how sensitive the accuracy is to the choice of a hard-negative mining method,

and how our positive-aware methods are effective in helping to achieve state-of-the-art results in MTEB Retrieval benchmark, as demonstrated by *NV-Retriever-v1*.

## 4.5 The effects of positive-aware neg. mining

In this section, we visualize the effects of positive-aware hard-negative mining on false negatives removal and on shifting scores and loss distributions for faster constrastive learning.

*4.5.1 False negatives removal.* To provide some insight on how effective our methods are in removing false negatives, we used LLM-as-a-judge (*Llama 3.1 70b instruct*[17]) with a prompt to classify, for a sample of the train set with mined negatives, whether a context is relevant to the question (false negative).

We present in Figure 3 the percentage of relevant (false) negatives, together with the percentage of relevant (true) positives for reference. We observed that false negatives rate is generally proportional[18] to the number of unique passages in the corpus, as larger number of candidates increases the chance of relevant passages to be mined. We can see that *Naive Top-k* has a high false negative rate for NQ (38.8%) and StackExchange (47%) train sets as it selects as negatives the most similar contexts to the question without any filter. Our methods successfully mined 57% and 50% less false negatives than *Naive Top-k* for NQ and StackExchange datasets, respectively.
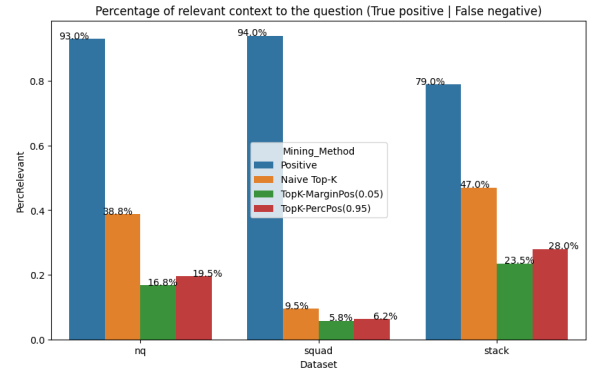


**Figure 3: Percentage of relevant context - (true) positives and (false) mined negatives - as classified by LLM-as-a-judge (*Llama 3.1 70b instruct*).**

*4.5.2 Visualizing the effect of positive-aware mining.* To provide some intuition on the effect of removing potential false negatives, in Figure 4 we provide histograms of the contexts scores and loss for *Naive Top-k* and *TopK-PercPos* mining methods on the train set. We can see in Figure 4.a that *TopK-PercPos* helps to separate the distributions of positive and negative scores, in Figure 4.b that it avoids negatives scoring higher (negative difference) than the corresponding positive score, and in Figure 4.c that it limits the max. cross-entropy loss (Equation 1), making training more stable.

---

[15]We don't include the results for *Top-1+sampled top-k* sampling method in Table 4 because they match the baseline *TopK-PercPos* method, as for the experiments with Mistral base model we could only use 1 hard-negative due to memory limitations

[16]The reason for not including those non-retrieval train sets in the experiments is because hard-negative mining is performed only for retrieval datasets.

[17]We obtained very similar context relevance classification results using the same prompt with Mixtral 8x 22b, which increases the confidence of our LLM-as-judge approach.

[18]Train set sizes: StackExchange (99,974), NQ (75,215), SQUAD (18,891)

**Table 5:** *Scaled setup*: Evaluation (NDCG@10) on full MTEB Retrieval of models trained with the same setup as *NV-Retriever-v1* using different hard-negative mining techniques

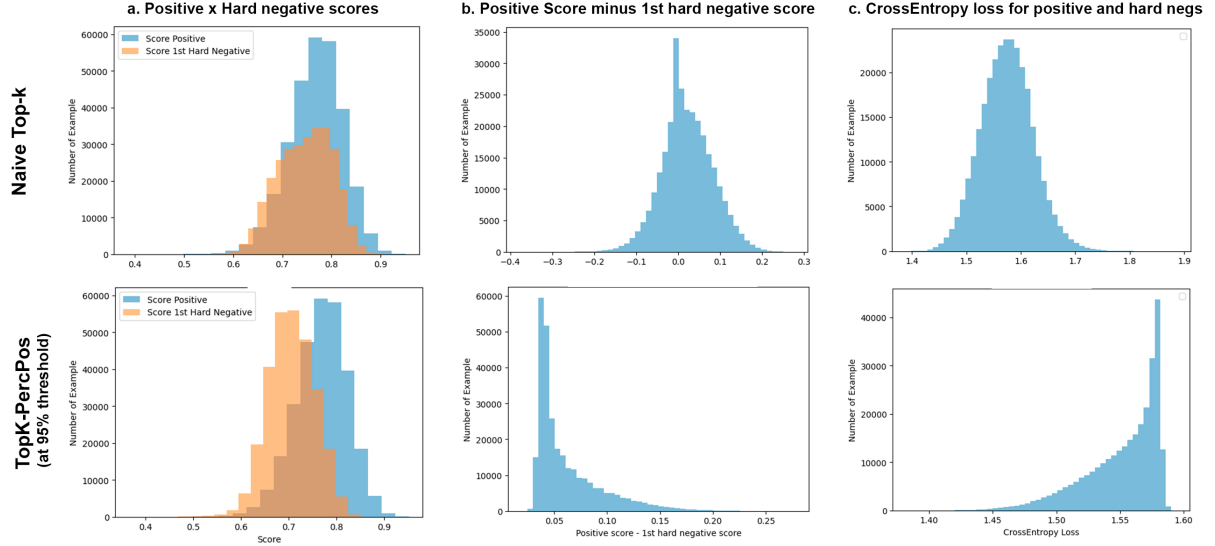| Mining method | Avg | ArguAna | Climate-FEVER | CQA-Dupstack-Retrieval | DBPedia | FEVER | FiQA2018 | Hotpot-QA | MS-MARCO | NF-Corpus | NQ | Quora-Retrieval | SCIDOCS | SciFact | Touche-2020 | TREC-COVID |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Naive Top-K | 51.44 | 63.8 | 38.0 | 41.6 | 42.6 | 89.2 | 48.5 | 73.9 | 38.2 | 13.1 | 69.8 | 84.8 | 16.5 | 77.6 | 19.8 | 54.2 |
| Top-K shifted by (10) | 54.66 | 57.3 | 36.3 | 45.8 | 48.7 | 87.6 | 54.1 | 75.2 | 40.7 | 31.1 | 64.7 | 86.4 | 20.5 | 73.0 | 20.4 | 78.0 |
| TopK-Abs (0.7) | 55.81 | 61.7 | 39.4 | 49.7 | 48.8 | 90.4 | 58.1 | 74.1 | 41.1 | 24.1 | 68.4 | 88.2 | 21.5 | 76.2 | 23.8 | 71.5 |
| TopK-MarginPos (0.05) | 59.77 | 61.7 | 39.9 | 48.7 | 50.5 | 92.6 | **61.5** | 77.7 | 44.4 | 44.6 | 71.3 | 88.4 | **22.5** | 79.2 | **28.6** | 84.8 |
| TopK-PercPos (95%) | **60.55** | **67.8** | **41.8** | 49.3 | 50.6 | 93.2 | 61.5 | 79.0 | 44.9 | 44.8 | 72.0 | 88.8 | 22.1 | **80.0** | 26.1 | **86.5** |



**Figure 4: Histograms comparing *Naive Top-k* and *TopK-PercPos* mining methods**

## 5 Conclusion

In this work, we introduce a novel family of *positive-aware hard-negative mining methods*, that leverage the positive relevance score for mining better negatives for embedding models finetuning.

We provide a comprehensive ablation study comparing hard-negative mining methods (under their many configurations), different teacher models and the ensembling of their hard-negatives. We also investigate the effect of our methods on false negative removal, loss stabilization and accuracy improvement. Finally, we demonstrate how effective our *positive-aware hard-negative mining methods* are at scale, leading to the state-of-the-art *NV-Retriever-v1* embedding model.

We recommend leveraging our mining methods for constrastive learning beyond text embedding models, as in some of our preliminary experiments on fine-tuning multi-modal (e.g., image,text) embedding models they also show to be effective in improving the retrieval accuracy. We suggest to future research works on text retrieval being aware of the sensitivity of models accuracy with respect to the negatives mined for contrastive learning. We also encourage them to disclose their methodology for hard-negative mining for better reproducibility and replicability.

## References

[1] 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 preprint arXiv:1611.09268* (2016).

[2] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading wikipedia to answer open-domain questions. *arXiv preprint arXiv:1704.00051* (2017).

[3] Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. *arXiv preprint arXiv:2402.03216* (2024).

[4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*. PMLR, 1597–1607.

[5] Kunal Dahiya, Nilesh Gupta, Deepak Saini, Akshay Soni, Yajun Wang, Kushal Dave, Jian Jiao, Gururaj K, Prasenjit Dey, Amit Singh, et al. 2023. Ngame: Negative mining-aware mini-batching for extreme classification. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 258–266.

[6] Chris Deotte, Bo Liu, Benedikt Schifferer, and Gilberto Titericz. 2021. GPU accelerated boosted trees and deep neural networks for better recommender systems. In *Proceedings of the Recommender Systems Challenge 2021*. 7–14.

[7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

[8] Daniel Gillick, Sayali Kulkarni, Larry Lansing, Alessandro Presta, Jason Baldridge, Eugene Ie, and Diego Garcia-Olano. 2019. Learning dense representations for entity retrieval. *arXiv preprint arXiv:1909.10506* (2019).

[9] Michael Günther, Jackmin Ong, Isabelle Mohr, Alaeddine Abdessalem, Tanguy Abel, Mohammad Kalim Akram, Susana Guzman, Georgios Mastrapas, Saba Sturua, Bo Wang, et al. 2023. Jina embeddings 2: 8192-token general-purpose

text embeddings for long documents. *arXiv preprint arXiv:2310.19923* (2023).

[10] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7B. *arXiv preprint arXiv:2310.06825* (2023).

[11] Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906* (2020).

[12] Junseong Kim, Seolhwa Lee, Jihoon Kwon, Sangmo Gu, Yejin Kim, Minkyung Cho, Jy yong Sohn, and Chanyeol Choi. 2024. Linq-Embed-Mistral:Elevating Text Retrieval with Improved GPT Data Through Task-Specific Control and Quality Refinement. Linq AI Research Blog. https://linqalpha.com/Blog/linq-embed-mistral

[13] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics* 7 (2019), 453–466.

[14] Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. 2024. NV-Embed: Improved Techniques for Training LLMs as Generalist Embedding Models. *arXiv preprint arXiv:2405.17428* (2024).

[15] Jinhyuk Lee, Zhuyun Dai, Xiaoqi Ren, Blair Chen, Daniel Cer, Jeremy R Cole, Kai Hui, Michael Boratko, Rajvi Kapadia, Wen Ding, et al. 2024. Gecko: Versatile text embeddings distilled from large language models. *arXiv preprint arXiv:2403.20327* (2024).

[16] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* 33 (2020), 9459–9474.

[17] Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. 2023. Towards general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281* (2023).

[18] Rui Meng, Ye Liu, Shafiq Rayhan Joty, Caiming Xiong, Yingbo Zhou, and Semih Yavuz. 2024. SFR-Embedding-Mistral:Enhance Text Retrieval with Transfer Learning. Salesforce AI Research Blog. https://blog.salesforceairesearch.com/sfr-embedded-mistral/

[19] Luke Merrick, Danmei Xu, Gaurav Nuti, and Daniel Campos. 2024. Arctic-Embed: Scalable, Efficient, and Accurate Text Embedding Models. *arXiv preprint arXiv:2405.05374* (2024).

[20] Gabriel de Souza P Moreira, Sara Rabhi, Ronay Ak, Md Yasin Kabir, and Even Oldridge. 2021. Transformers with multi-modal features and post-fusion context for e-commerce session-based recommendation. *arXiv preprint arXiv:2107.05124* (2021).

[21] Niklas Muennighoff, Hongjin Su, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, and Douwe Kiela. 2024. Generative Representational Instruction Tuning. arXiv:2402.09906 [cs.CL]

[22] Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. MTEB: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316* (2022).

[23] Zach Nussbaum, John X Morris, Brandon Duderstadt, and Andriy Mulyar. 2024. Nomic embed: Training a reproducible long context text embedder. *arXiv preprint arXiv:2402.01613* (2024).

[24] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018).

[25] Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daixiang Dong, Hua Wu, and Haifeng Wang. 2020. RocketQA: An optimized training approach to dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2010.08191* (2020).

[26] Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for SQuAD. *arXiv preprint arXiv:1806.03822* (2018).

[27] Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models. *Transactions of the Association for Computational Linguistics* 11 (2023), 1316–1331.

[28] Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084* (2019).

[29] Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021. Rocketqav2: A joint training method for dense passage retrieval and passage re-ranking. *arXiv preprint arXiv:2110.07367* (2021).

[30] Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends® in Information Retrieval* 3, 4 (2009), 333–389.

[31] Benedikt Schifferer, Chris Deotte, Jean-Francois Puget, Gabriel de Souza Pereira Moreira, Gilberto Titericz, Jiwei Liu, and Ronay Ak. 2021. Using Deep Learning to Win the Booking. com WSDM WebTour21 Challenge on Sequential Recommendations.. In *WebTour@ WSDM*. 22–28.

[32] Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. Beir: A heterogenous benchmark for zero-shot evaluation of

[33] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533* (2022).

[34] Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2023. Improving text embeddings with large language models. *arXiv preprint arXiv:2401.00368* (2023).

[35] Wenhui Wang, Hangbo Bao, Shaohan Huang, Li Dong, and Furu Wei. 2020. Minilmv2: Multi-head self-attention relation distillation for compressing pre-trained transformers. *arXiv preprint arXiv:2012.15828* (2020).

[36] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771* (2019).

[37] Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian Riedel, and Luke Zettlemoyer. 2019. Scalable zero-shot entity linking with dense entity retrieval. *arXiv preprint arXiv:1911.03814* (2019).

[38] Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. 2023. C-Pack: Packaged Resources To Advance General Chinese Embedding. arXiv:2309.07597 [cs.CL]

[39] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. 2020. Approximate nearest neighbor negative contrastive learning for dense text retrieval. *arXiv preprint arXiv:2007.00808* (2020).

[40] Wenhan Xiong, Xiang Lorraine Li, Srini Iyer, Jingfei Du, Patrick Lewis, William Yang Wang, Yashar Mehdad, Wen-tau Yih, Sebastian Riedel, Douwe Kiela, et al. 2020. Answering complex open-domain questions with multi-hop dense retrieval. *arXiv preprint arXiv:2009.12756* (2020).

information retrieval models. *arXiv preprint arXiv:2104.08663* (2021).

# A  Similarity of hard-negatives mined by different teacher models

In this appendix, we investigate the level of agreement of different teacher models on the hard-negatives.

In particular, we computed the Jaccard Similarity between the top-4 hard-negatives for pairs of different teacher models. Table 6 presents the similarity matrices for each train sets used in the ablation study: NQ, SQUAD and StackExchange. In general, it is possible to observe from the similarity matrices that the teacher models agree poorly on the top-4 mined hard-negatives, i.e., Jaccard similarity lower than 30%.

**Table 6: Jaccard similarity of hard-negatives mined with different teacher models for NQ | SQUAD | StackExchange datasets**

|  | e5-large-v2 | arctic-embed-l | NV-Embed-v1 | e5-mistral-7b |
|---|---|---|---|---|
| e5-large-v2 | - | 0.15 \| 0.14 \| 0.01 | 0.11 \| 0.13 \| 0.01 | 0.11 \| 0.14 \| 0.01 |
| arctic-embed-l | 0.15 \| 0.14 \| 0.01 | - | 0.28 \| 0.19 \| 0.05 | 0.18 \| 0.18 \| 0.04 |
| NV-Embed-v1 | 0.11 \| 0.13 \| 0.01 | 0.28 \| 0.19 \| 0.05 | - | 0.23 \| 0.30 \| 0.09 |
| e5-mistral-7b | 0.11 \| 0.14 \| 0.01 | 0.18 \| 0.18 \| 0.04 | 0.23 \| 0.30 \| 0.09 | - |

# B  NV-Retriever-v1

In this section, we describe the architecture, methods and techniques for training the state-of-the-art *NV-Retriever-v1* embedding model, which was placed first on the MTEB Retrieval leaderboard when it was published. The *positive-aware mining methods* we propose in this paper were crucial to obtaining high retrieval accuracy for *NV-Retriever-v1*, as we demonstrate in Section 4.4.

## B.1 Model architecture

The *NV-Retriever-v1* uses the Mistral 7B[10][19] as base model, as originally proposed by *e5-mistral-7b-instruct* [34] and followed by many of the top performing embedding models on MTEB such as *Linq-Embed-Mistral* [12], *GritLM* [21], *SFR-Embedding-Mistral* [18], and *NV-Embed-v1* [14]. We convert the decoder *Mistral-7b* model into an encoder model. We replace the causal attention of the base model with bi-directional attention, which empirically demonstrated higher accuracy in our experiment. Mean pooling is used to combine the outputs of the last Transformer layer, which corresponds to averaging the hidden states across the sequence length to form the sentence embedding. This approach was inspired by *GritLM* [21] and subsequent models using this approach like *NV-Embed-v1* [14] and *gte-Qwen2-7B-instruct.*

## B.2 Train sets

As MTEB contains different tasks like retrieval, reranking, classification, clustering among others, a diverse set of training datasets is required for a good overall performance. The train sets used for fine-tuning *NV-Retriever-v1*, listed in Table 7, are based on those used to train *E5-Mistral* [34] and *NV-Embed-v1* [14]. We use the same instruction prompts from [34] for each dataset.

## B.3 Hard-negative mining

We used the *E5-Mistral-7B* embedding model for hard-negative mining[20] with maximum sequence length of 4096.

To ignore potential false negatives, we leverage our proposed *TopK-PercPos* method and set the maximum threshold for the negative relevance score as 95% of the positive score. We explain this choice in the ablation study presented in Section 4.3.

To better leverage the base LLM model pre-training and adjust for datasets specific domain and task, [34] proposed and designed specific natural language instructions for each train set. The instructions prefixes are added to the query but not to the passages, so that for the latter no re-indexing is required for different instructions. We also adopt instruction prefixes with a small difference in the implementation: instead of the original template from [34] "*Instruct: {task_definition} \n Query: {query}*" , we use "*{task_definition}: {query}*". As in *NV-Embed-v1* [14] we mask out the instruction tokens for the average pooling during both training and evaluation, which can still affect the output due to self-attention.

## B.4 Training setup

We perform two stages of instruction tuning as in [14]. In the first stage, we only use supervised retrieval datasets with in-batch negatives in addition to the mined hard-negatives. In the second stage, we blend that retrieval data with other tasks' datasets (e.g. classification, regression, semantic sentence similarity).

**Table 7: Retrieval train datasets for *NV-Retriever-v1***

| Dataset | # of samples |
|---|---|
| *Retrieval datasets* | |
| ArguAna | 4065 |
| BioASQ | 2495 |
| FEVER | 50000 |
| FiQA2018 | 14166 |
| GOOAQ | 20000 |
| HotpotQA | 85000 |
| MS-MARCO | 200000 |
| NFCorpus | 3685 |
| Natural Language Inference | 20000 |
| Natural Questions | 100231 |
| PAQ | 20000 |
| SciFacts | 919 |
| SQUAD | 87599 |
| StackExchange | 100000 |
| TriviaQA | 20000 |
| *Non-retrieval datasets* | |
| Banking77Classification | 10000 |
| AmazonCounterfactualClassification | 4018 |
| AmazonReviewsClassification | 20000 |
| EmotionClassification | 16000 |
| ImdbClassification | 15000 |
| MTOPIntentClassification | 10000 |
| ToxicConversationsClassification | 40000 |
| TweetSentimentExtractionClassification | 27481 |
| STS12 | 1868 |
| STS22 | 416 |
| STSBenchmark | 2812 |
| ArxivClusteringP2P | 35000 |
| ArxivClusteringS2S | 35000 |
| BiorxivClusteringP2P | 4070 |
| BiorxivClusteringS2S | 4070 |
| MedrxivClusteringP2P | 1160 |
| MedrxivClusteringS2S | 1160 |

Our training implementation is based on the Hugging Face Transformers[21] and PEFT[22] libraries. The model and training hyperparameters are shown in Table 8 and Table 9.

**Table 8: Model hyperparameters for *NV-Retriever-v1***

| Hyperparameter | Value |
|---|---|
| Base model | mistralai/Mistral-7B-v0.1 |
| Layers | 32 |
| Attention | Bi-directional |
| Embedding dim | 4096 |
| Embedding pooling | Average at last layer |
| LoRA Rank | 16 |
| LoRA Alpha | 32 |
| Query max length | 192 |
| Passage max length | 512 |

## B.5 Computational cost

NV-Retriever-v1 is trained in two stages, with 12 epochs each. Full training takes about 90 hours on 8x A100 GPUs. The hard-negative mining process depends on the corpus size and teacher model complexity. As a reference, for the *e5-mistral-instruct* teacher model, it requires approximately 1.5 hours with 8x A100 GPUs to embed 1 million documents with max token length set to 4096.

Evaluating an embedding model on the 15 MTEB BEIR retrieval datasets takes approximately 120 hours of 8x A100 GPUs.

---

[19]https://huggingface.co/mistralai/Mistral-7B-v0.1 (we use the Mistral base version, without instruction fine-tuning)
[20]The only exception was ArguAna dataset, for which we found that using an internal Mistral embedding model for mining instead of *E5-Mistral-7B* trained with causal attention improved significantly the NDCG@10 for Arguana test set.

[21]https://github.com/huggingface/transformers
[22]https://github.com/huggingface/peft

**Table 9: Training hyperparameters for *NV-Retriever-v1***

| Hyperparameters | 1st stage | 2nd stage |
|---|---|---|
| Optimizer | AdamW | |
| Learning rate | 1.00E-05 | |
| Learning rate warm-up steps | 100 | |
| Negatives source | hard-negatives + in-batch negatives | hard-negatives |
| Number of hard-negatives | 1 | 5 |
| Batch-size | 32 | 8 |
| Gradient accumulation steps | 4 | 16 |
| Training Steps per epoch | 2844 | 829 |
| Epochs | 12 | 12 |

## B.6 *NV-Retriever-v1* results on MTEB Retrieval

*NV-Retriever-v1* model achieves at MTEB leaderboard an average NDCG@10 score of 60.9 and placed 1st place when published in July 11, 2024. Table 10 reports the top embedding models on the MTEB Retrieval leaderboard when it was published. Mistral 7B was the foundation model for 4 of the top 5 places, that used similar train sets. As discussed before, the main difference in *NV-Retriever-v1* training approach was the usage of the *positive-aware mining methods* introduced in this paper.

**Table 10: Top embedding models on MTEB Retrieval leaderboard as of 2024-07-11**

| MTEB Retrieval Dataset | NV-Retriever-v1 | gte-Qwen2-7B-instruct | Linq-Embed-Mistral | SFR-Embedding-2_R | NV-Embed-v1 | SFR-Embedding-Mistral |
|---|---|---|---|---|---|---|
| **Average** | **60.90** | **60.25** | **60.19** | **60.18** | **59.36** | **59.00** |
| ArguAna | 68.28 | 64.27 | 69.65 | 62.34 | 68.20 | 67.17 |
| ClimateFEVER | 43.47 | 45.88 | 39.11 | 34.43 | 34.72 | 36.41 |
| CQADupstackRetrieval | 49.36 | 46.43 | 47.27 | 46.11 | 50.51 | 46.49 |
| DBPedia | 50.82 | 52.42 | 51.32 | 51.21 | 48.29 | 49.06 |
| FEVER | 93.15 | 95.11 | 92.42 | 92.16 | 87.77 | 89.35 |
| FiQA2018 | 61.18 | 62.03 | 61.20 | 61.77 | 63.10 | 60.40 |
| HotpotQA | 79.12 | 73.08 | 76.24 | 81.36 | 79.92 | 77.02 |
| MSMARCO | 44.89 | 45.98 | 45.21 | 42.18 | 46.49 | 43.41 |
| NFCorpus | 45.06 | 40.60 | 41.62 | 41.34 | 38.04 | 41.88 |
| NQ | 72.44 | 67.00 | 70.63 | 73.96 | 71.22 | 69.92 |
| QuoraRetrieval | 88.78 | 90.09 | 90.27 | 89.58 | 89.21 | 89.78 |
| SCIDOCS | 22.55 | 28.91 | 21.93 | 24.87 | 20.19 | 19.91 |
| SciFact | 81.31 | 79.06 | 78.32 | 85.91 | 78.43 | 77.66 |
| Touche2020 | 26.60 | 30.57 | 30.61 | 28.18 | 28.38 | 29.00 |
| TRECCOVID | 86.44 | 82.26 | 87.10 | 87.27 | 85.88 | 87.60 |