

# colbert-kit Package Documentation

Thuong Dang, Qiqi Chen

April 2025

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>ColBERT and MaxSim score</b>	<b>2</b>
<b>3</b>	<b>Indexing</b>	<b>2</b>
<b>4</b>	<b>Retrieving</b>	<b>3</b>
<b>5</b>	<b>Re-ranking</b>	<b>4</b>
<b>6</b>	<b>Training</b>	<b>5</b>
<b>7</b>	<b>Evaluating</b>	<b>6</b>
<b>8</b>	<b>Choosing negatives</b>	<b>9</b>

## 1 Introduction

This is documentation for our work on ColBERT. In this document, we provide a detailed introduction to ColBERT and the functionalities of our ColBERT package. The documentation is structured as follows:

In the first section, we briefly review ColBERT and its MaxSim scoring function. In the second section, we discuss our indexing and retrieval methods for ColBERT, including the inputs and outputs used in our package. The third section covers the re-ranking functionality. In the fourth section, we explain how to train a ColBERT model and evaluate the training using standard information retrieval (IR) datasets. Our package includes training and evaluation scripts to compare the performance of ColBERT with BM25.

Finally, we address the general difficulty of selecting effective negatives for training or fine-tuning ranking models. To this end, we have developed several methods for negative sampling, which are also included in our package.

## 2 ColBERT and MaxSim score

ColBERT is a dense retrieval model that uses multiple embeddings per query and document to enable fine-grained similarity matching between their tokens.

Let  $q, d$  be a query, and document, respectively. Let  $[q_1 \cdots q_m]$  be the embeddings of tokens of  $q$ , and  $[d_1 \cdots d_n]$  be the embeddings of tokens of  $d$ . The *MaxSim score* is defined as follows

$$\text{MaxSim}(q, d) = \sum_{i=1}^m \max_j S(q_i, d_j),$$

where  $S(q_i, d_j)$  is the cosine similarity between the token  $q_i$  and  $d_j$ . As an intuition, this can be seen as a soft extension of keyword matching: if we define  $S(q_i, d_j) = 1$  if  $q_i = d_j$  and 0 otherwise, then the MaxSim score becomes a ranking function based on keyword matching. It means the more keyword matching, the higher ranking a document will have. In ColBERT, this hard match is replaced by semantic similarity, allowing tokens with similar meaning to contribute to the relevance score.

## 3 Indexing

To index with our ColBERT package, we need to define

```
1 # Place and name where we want to save the token to doc map
2 token_to_doc_map_name = 'token_to_doc_map.npy'
3 token_to_doc_map_path = './indices_test/'
4 # Place and name for saving the document dictionary
5 # This information is for computing ColBERT scores between a
   query and a document after we retrieve the nearest
   tokens
6 doc_dict_name = 'documents.json'
7 doc_dict_path = './indices_test/'
8 # We can load the model from huggingface or local
9 embedding_model_path = './model/colbert'
```

We then initialize the method for embedding and create index

```
1 from colbert_kit.embedding.colbert_embedding import
   colBERTEmbedding
2
3 colbertEmbed = colBERTEmbedding(
4     model_name_or_path=embedding_model_path,
5     token_to_doc_map_name=token_to_doc_map_name,
6     doc_dict_name=doc_dict_name,
7     token_to_doc_map_path=token_to_doc_map_path,
8     doc_dict_path=doc_dict_path,
9     device="cpu"
10 )
11
```

```

12 # texts is a list of string (or a list of our documents for
    retrieving)
13 embeddings = colbertEmbed.encode(texts) #output: numpy.
    ndarray
14
15 # We then to create FAISS index to search for nearest tokens
16 # Place and name to save the FAISS index
17 index_output_path = './indices_test/'
18 index_name = 'colbert_test_index.index'
19 dimension = embeddings.shape[1]
20
21 # Create an HNSW index
22 # There are three options: flat, IVF and HNSW
23 hnsw_index = HNSWfaissIndex(index_name=index_name, dimension
    =dimension, use_gpu=False, index_output_path=
    index_output_path)
24 hnsw_index.create_index(embeddings)

```

## 4 Retrieving

To retrieve, the main steps are

- Loading token to doc map, FAISS index, the document dictionary
- For an input query, use ColBERT model to have token embeddings of the query
- For each query's token embedding, retrieve top  $k'$  nearest tokens
- Map them back to original documents and compute ColBERT MaxSim scores
- Sort the results, and return the candidate documents

```

1 from colbert_kit.indexing.faiss.flat_faiss_index import
    FlatFaissIndex
2 from colbert_kit.indexing.faiss.HNSW_faiss_index import
    HNSWfaissIndex
3 from colbert_kit.indexing.faiss.IVF_faiss_index import
    IVFFaissIndex
4
5 colbertSearchFaiss = colBERTSearchFaiss(
6     model_name_or_path=embedding_model_path,
7     index_path=index_output_path,
8     index_name=index_name,
9     index_type='HNSW',
10    use_gpu=False,
11    token_to_doc_map_name=token_to_doc_map_name,
12    doc_dict_name=doc_dict_name,

```

```

13     token_to_doc_map_path=token_to_doc_map_path,
14     doc_dict_path=doc_dict_path,
15     device="cpu"
16 )
17
18 colbertSearchFaiss.load_index_map_and_dict()
19
20 # Search
21 scores, colbert_indices = colbertSearchFaiss.
    colbert_search_results(query, top_k_results=5,
    top_k_search_tokens=10)

```

## 5 Re-ranking

Re-ranking with ColBERT in our package are more straight-forward. The inputs for this is quite simple, we just need to define document candidates, document ids, batch size and number of top results.

```

1 from colbert_kit.reranking.colbert_batch_reranker import
    colBERTReRankerBatch
2
3 embedding_model_path = "./model/colbert"
4 batchreranker = colBERTReRankerBatch(model_name_or_path=
    embedding_model_path, device='cpu')

```

We can then define our query and candidate documents and candidate ids:

```

1     query = "Was ist die Hauptstadt von Frankreich?"
2
3 doc_candidates = [
4     "Paris ist die Hauptstadt von Frankreich.",
5     "Berlin ist die Hauptstadt von Deutschland.",
6     "Madrid ist die Hauptstadt von Spanien.",
7     "Rom ist die Hauptstadt von Italien.",
8     "Der Eiffelturm befindet sich in Paris."
9 ]
10
11 candidate_idx = list(range(len(doc_candidates)))

```

And then, finally, re-rank with ColBERT for top results

```

1 top_scores, top_indices = batchreranker.reranker(query,
    doc_candidates, candidate_idx, batch_size=8, top_n=10)
2 print("Top Scores:", top_scores)
3 print("Top Indices:", top_indices)
4 for id in top_indices:
5     print(doc_candidates[id])

```

## 6 Training

To train a ColBERT model, we follow the training loss proposed in the original paper. It is the softmax-cross-entropy loss on pairwise triplets. Given a triplet  $[q, d^+, d^-]$ , where  $q$  is a query,  $d^+$  is a positive document, and  $d^-$  is a negative document. The loss function is defined as follows

$$L = -\log \frac{\exp \text{MaxSim}(q, d^+)}{\exp \text{MaxSim}(q, d^+) + \exp \text{MaxSim}(q, d^-)}$$

The training script is simple, for really scaled dataset like MSMARCO, saving checkpoints and evaluating after each epoch is important. Therefore, we also included script to save and train from checkpoints.

```
1 import pandas as pd
2 import torch
3 from colbert_kit.training import ColBERTTrainer
4
5 # Create simple test dataset
6 test_data = [
7     ("A man is playing a guitar.", "A person is playing a
8     musical instrument.", "A bird is flying in the sky."),
9     ("A woman is cooking food.", "A person is preparing a
10    meal.", "A car is driving down the road."),
11    ("A child is reading a book.", "A young person is
12    engaged with a novel.", "The sun is setting behind
13    the mountains.")
14 ]
15 # Or we can train with triplet dataset with three columns:
16 # sentence, positive_sentence and negative_sentence
17 # test_df = pd.DataFrame(test_data, columns=["sentence", "
18 # positive_sentence", "negative_sentence"])
19 # test_df.to_csv("data_test/test_triplets.csv", sep='\t',
20 # index=False)
21
22 # Load BERT model for training
23 bert_model_name = "model/bert"
24 device = 'cuda' if torch.cuda.is_available() else 'cpu'
25
26 trainer = ColBERTTrainer(
27     model_name_or_path=bert_model_name,
28     device=device,
29     triplets=test_data,
30     # Or define triplet data from a csv file as above
31     # triplets_path="data_test/test_triplets.csv",
32     batch_size=2,
33     lr=2e-5,
34     num_epochs=4,
```

```

29     checkpoint_dir="./test_checkpoints/"
30 )
31 trainer.train()
32 trainer.save_model("./model/saved_model_test/")

```

We can also reload models from checkpoints and continue the training process.

```

1 trainer = ColBERTTrainer(
2     model_name_or_path=bert_model_name,
3     triplets=test_data,
4     # triplets_path="data_test/test_triplets.csv",
5     device = device,
6     batch_size=256,
7     lr=2e-5,
8     num_epochs=3,
9     checkpoint_dir="./checkpoints/",
10    checkpoint_path="./checkpoints/checkpoint_epoch_2.pth"
11 )
12 trainer.train()
13 trainer.save_model("./model/saved_model_test_3/")

```

## 7 Evaluating

To evaluate information retrieval (IR) metrics, we need standard IR datasets. They consist of three components

- queries component: a dataset with two columns query\_id and text
- documents component: a dataset with two columns doc\_id and text
- qrel component: a dataset with three columns query\_id, doc\_id and relevant\_score

To compute IR metric for our retrieving method, we need to

- create an index from documents component for retrieving
- for each query\_id in qrel component, map it back to the original query, and then retrieve from document index, collect all the results and compute metrics.

For more detail, the reader can take a look on our sample evaluation codes [...]. Here are parts of the script. For inputs of our IR evaluation, we need no more than a data frame of results and qrel component.

```

1 from colbert_kit.evaluating.ir_eval import IREvaluation
2 # Computing IR metrics for BM25
3
4 results = []

```

```

5 query_ids = df_qrels['query_id'].unique()
6 top_k = 50
7
8 for query_id in tqdm(query_ids, desc="Processing queries"):
9     query_row = df_queries.loc[df_queries['query_id'] ==
10                                query_id, text_column]
11     if query_row.empty:
12         raise ValueError(f"Query ID {query_id} not found in
13                             df_queries.")
14     query = query_row.values[0]
15
16     # Get top document indices from BM25
17     bm25_scores, bm25_top_idx = BM25search.search(query,
18                                                    top_k)
19
20     top_ranked_ids = []
21     for idx in bm25_top_idx:
22         idx = int(idx)
23         top_ranked_ids.append(df_docs.iloc[idx]['doc_id'])
24
25     for rank, doc_id in enumerate(top_ranked_ids):
26         results.append({'query_id': query_id, 'doc_id':
27                         doc_id, 'rank': rank})
28
29 df_results = pd.DataFrame(results)
30
31 evaluation = IREvaluation()
32
33 k_set = [1, 5, 10, 20, 50]
34 for k in k_set:
35     evaluation.compute_evaluation_metrics(df_qrels,
36                                         df_results, k)
37
38 # Computing metrics for ColBERT
39 results = []
40 query_ids = df_qrels['query_id'].unique()
41 top_k = 50
42
43 # ColBERT retrieving metrics
44
45 for query_id in tqdm(query_ids, desc="Processing queries"):
46     query_row = df_queries.loc[df_queries['query_id'] ==
47                                query_id, text_column]
48     if query_row.empty:
49         raise ValueError(f"Query ID {query_id} not found in
50                             df_queries.")
51     query = query_row.values[0]
52
53     # Semantic search with colBERT
54     colbert_top_idx, scores = colbertSearchFaiss.

```

```

colbert_search_results(query, top_k_results=top_k,
top_k_search_tokens=10)
48
top_ranked_ids = []
49
for idx in colbert_top_idx:
50
    idx = int(idx)
51
    top_ranked_ids.append(df_docs.iloc[idx]['doc_id'])
52
53
for rank, doc_id in enumerate(top_ranked_ids):
54
    results.append({'query_id': query_id, 'doc_id':
55
                    doc_id, 'rank': rank})
56
df_results = pd.DataFrame(results)
57
evaluation = IREvaluation()
58
60
k_set = [1, 5, 10, 20, 50]
61
for k in k_set:
62
    evaluation.compute_evaluation_metrics(df_qrels,
63
                                        df_results, k)

```

We can also compute metrics for ColBERT as a re-ranker. For more detail, the reader can look at [...].

```

1 from colbert_kit.evaluating.ir_eval import IREvaluation
2 from colbert_kit.reranking.colbert_batch_reranker import
  colBERTReRankerBatch
3 from colbert_kit.evaluating.ir_eval import IREvaluation
4
5 results = []
6 query_ids = df_qrels['query_id'].unique()
7 top_k = 100
8
9 # Load colBERT
10 model_name_or_path = "./model/colbert"
11
12 colbert_reranker_batch = colBERTReRankerBatch(
    model_name_or_path=model_name_or_path, device="cpu")
13
14 for query_id in tqdm(query_ids, desc="Processing queries"):
15     query_row = df_queries.loc[df_queries['query_id'] ==
16                               query_id, text_column]
17     if query_row.empty:
18         raise ValueError(f"Query ID {query_id} not found in
19                           df_queries.")
20     query = query_row.values[0]
21
22     # Get top document indices from BM25
23     bm25_scores, bm25_top_idx = BM25search.search(query,
24                                                    top_k)

```



```

22
23     doc_candidates = [texts[idx] for idx in bm25_top_idx]
24     candidate_idx = list(bm25_top_idx)
25
26     top_n_scores, top_n_idx = colbert_reranker_batch.
27         reranker(query, doc_candidates, candidate_idx,
28             batch_size = 1, top_n=50)
29
30     top_ranked_ids = []
31     for idx in top_n_idx:
32         idx = int(idx)
33         top_ranked_ids.append(df_docs.iloc[idx]['doc_id'])
34
35     for rank, doc_id in enumerate(top_ranked_ids):
36         results.append({'query_id': query_id, 'doc_id':
37             doc_id, 'rank': rank})
38
39 df_results = pd.DataFrame(results)
40
41 evaluation = IREvaluation()
42
43 k_set = [1, 5, 10, 20, 50]
44 for k in k_set:
45     evaluation.compute_evaluation_metrics(df_qrels,
46         df_results, k)

```

## 8 Choosing negatives

We employed two techniques for choosing negatives. In real applications, it is difficult in general to choose negatives for training or fine-tuning ranking models. The first technique is choosing random negative. The inputs include a data frame with two columns (EXACTLY) "sentence" and "positive\_sentence", number of negatives, and triplet output path. .

```

1 import pandas as pd
2 from colbert_kit.choosing_negatives.
3     choosing_random_negatives import build_random_triplets
4
5 # Step 1: Sample DataFrame
6 data = {
7     "sentence": [
8         "A man is playing guitar.",
9         "A woman is cooking.",
10        "Children are playing soccer.",
11        "A cat is sleeping on the couch.",
12        "He is reading a book."
13    ],
14     "positive_sentence": [

```

```

14         "Someone plays a musical instrument.",
15         "She is preparing food.",
16         "Kids are having fun outdoors.",
17         "An animal rests on furniture.",
18         "A person is reading."
19     ]
20 }
21
22 df = pd.DataFrame(data)
23
24 # Step 3: Run it and save to CSV
25 output_path = "triplets_example.csv"
26 build_random_triplets(
27     df=df,
28     n_negatives=2,
29     triplet_output_path=output_path
30 )
31
32 triplet_df = pd.read_csv("triplets_example.csv")
33
34 # Step 4: Load and show results
35 print("Sample triplets:")
36 print(triplet_df.head())

```

The second technique we used is to choose hard negatives. This is a little more complicated, as we use a sentence-transformer model to retrieve negatives that are very similar to the positive sentence. Our indexing method is HNSW. The inputs include

- a data frame with EXACTLY two columns named "sentence" and "positive\_sentence"
- a sentence-transformer model name or path
- a start index: an index from which to start selecting hard negatives among the top retrieved results
- number of negatives : number of hard negatives to retrieve per anchor sentence
- triplet\_output\_path: File path to save the generated triplets as a CSV
- index\_name and index\_output\_path: a place to save the index with sentence transformer
- use\_gpu: Whether to use GPU acceleration for sentence transformer encoding and FAISS. Default is False.

```

1 import pandas as pd
2 from colbert_kit.choosing_negatives.choosing_hard_negatives
   import build_hard_triplets

```

```

3
4 # Step 1: Sample DataFrame
5 data = {
6     "sentence": [
7         "A man is playing guitar.",
8         "A woman is cooking.",
9         "Children are playing soccer.",
10        "A cat is sleeping on the couch.",
11        "He is reading a book."
12    ],
13    "positive_sentence": [
14        "Someone plays a musical instrument.",
15        "She is preparing food.",
16        "Kids are having fun outdoors.",
17        "An animal rests on furniture.",
18        "A person is reading."
19    ]
20 }
21
22 df = pd.DataFrame(data)
23
24 # Step 3: Run it and save to CSV
25 build_hard_triplets(
26     df=df,
27     model_name_or_path='all-MiniLM-L6-v2',
28     start_index=2,
29     n_negatives=2,
30     triplet_output_path='hard_triplets.csv',
31     index_name='hard.index',
32     index_output_path='./indices_test',
33     use_gpu=False
34 )
35
36 triplet_df = pd.read_csv("hard_triplets.csv")
37
38 # Step 4: Load and show results
39 print("Sample triplets:")
40 print(triplet_df)

```