

Inverted Index Construction

The file *invidx.cons.cpp* implements the creation of the inverted index. I have implemented the **Single In Memory Pass Indexing** algorithm to create the inverted index. This algorithm can work for arbitrarily large collection sizes, if there is enough disk space despite having limited RAM.

Algorithm 1 SPIMI-INVERT Algorithm

```

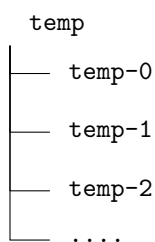
1: output_file ← new file()
2: dictionary ← new hash()
3: while free memory available do
4:   token ← next(token_stream)
5:   if term(token) ∉ dictionary then
6:     postings_list ← AddToDictionary(dictionary, term(token))
7:   else
8:     postings_list ← GetPostingsList(dictionary, term(token))
9:   end if
10:  AddToPostingsList(postings_list, docID(token))
11: end while
12: sorted_terms ← SortTerms(dictionary)
13: WriteBlockToDisk(sorted_terms, dictionary, output_file)
14: return output_file

```

I parse all the documents available in the collection, to get the token stream and till I have enough memory I keep using the algorithm mentioned above. As soon as memory is close to getting full, I write the postings list to disk and clear the memory.

From the next document onwards, I repeat the algorithm again till we are out of memory and so on.

This way at the end of this algorithm, we would have multiple temporary output files storing partial postings list. These files are stored in a folder named *temp*.



After this a merge algorithm is employed which uses a merge tree to merge files in pairs, till we have a single temporary file representing the postings list and vocabulary of the entire collection.

This single file is then processed, creating the dictionary and the index while employing the appropriate compression algorithm. We also store the mappings from the document IDs to the integer space and also store the norms of the document vectors which will be used during search.

Estimating Memory of Postings List

The postings list is implemented as a

```
1 std::map<std::string, std::vector<std::pair<int, int>>>
```

This way the tokens are always sorted, To estimate the size of this Data Structure, I use the following Heuristic

$$\text{size of postings list} = 8 * \text{Number of (docID, frequency pairs)} + 10 * \text{Number of Tokens}$$

Byte Pair Encoding

Since Training the Byte Pair Tokenizer, is an expensive and time consuming operation, I pretrained my own implementation of the BPE Tokenizer, and stored the merges learnt in order which can be loaded later for tokenization. The learned merges look as follows:

```
1 ...
2 yugoslav ia
3 at ta
4 dis count
5 at l
6 relig ious
7 j i
8 obser vers
9 grad u
10 bar ri
11 gree k
12 e a
13 compon ents
14 occas ion
15 new ly
16 ...
```

These are stored in the file named *bpe_merges* and are loaded automatically for tokenization if needed.

To speed up tokenization I remember, the splits of a particular token which prevents iterating all merges again and again.

Search

During Search, we first load the dictionary, document ID mappings and the norms of the document vectors in memory.

Then we parse the query file and load all the queries into memory.

We generate the tokens from the query title and content. We maintain a map to store the partial document scores, for every term in the query we load the postings list from the index file using the offsets specified by the dictionary file, we use the tf-idf formulation to update the partial scores of all the documents.

Once all Query Terms are processed we normalize the document scores using the precomputed document vector norms that we stored during index construction.

To get the top $K=100$ documents, we maintain a min heap which always maintains the best 100 documents, this way we can retrieve the top-k documents in $O(N \log k)$ time.

Heuristics and Optimizations

We treat the Query Title and Content differently, all terms in the title are considered important for scoring the documents, but for the content we try to pick only relevant terms, this speeds up query processing time significantly.

For all tokens in the query content, we sort them according to the increasing order of the count of documents they appear in.

We only consider the top 50% tokens in the query content according to this ordering.

Also since loading large postings list from the disk can be expensive, if a query term appears in more than 25% of the documents, we skip it during query processing to prevent loading large postings list from disk.

Results

Inverted Index Construction Time

Results are for all the entire training data, consisting of 2296 files containing 527939 documents

Compression	Tokenizer	Create Temporary Postings	Merge Partial Postings	Write Dictionary & Index File	Total Time
No Compression	Simple	784 s	120 s	897 s	1801 s
Variable Byte Encoding	Simple	717 s	69 s	229 s	1015 s
No Compression	BPE	5781 s	71 s	558 s	6410 s
Variable Byte Encoding	BPE	5823 s	36 s	281 s	6140 s

Index and Dictionary Sizes

Compression	Tokenization	Index File	Dictionary File
No Compression	Simple	811 MB	27 MB
Variable Byte Encoding	Simple	223 MB	26 MB
No Compression	BPE	860 MB	15 MB
Variable Byte Encoding	BPE	251 MB	15 MB

Average Query Processing Time

Compression	Tokenization	Average Time
No Compression	Simple	0.39 s
Variable Byte Encoding	Simple	0.69 s
No Compression	BPE	0.36 s
Variable Byte Encoding	BPE	0.65 s

Number of Correct Query Results

Simple Tokenization

On Average, 8 % results were relevant for the queries tested.

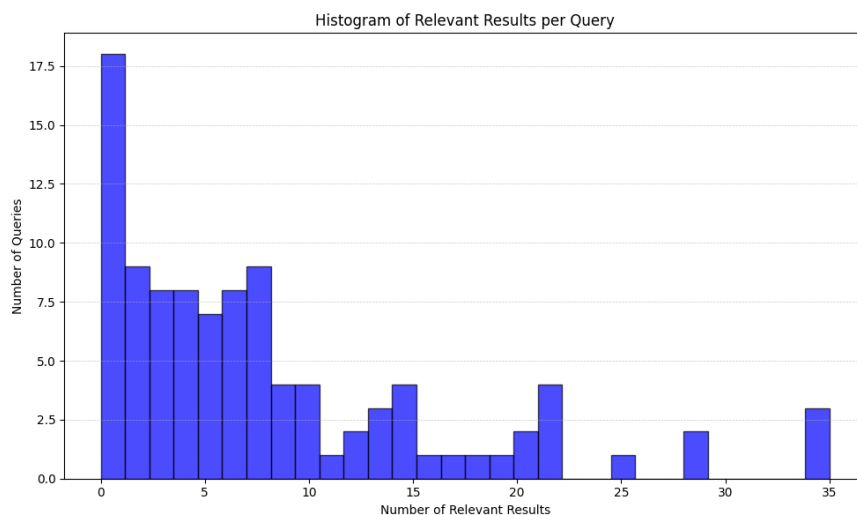


Figure 2: Number of Relevant Results in Top 100 Results

Byte Pair Encoding

Average Correct: 6

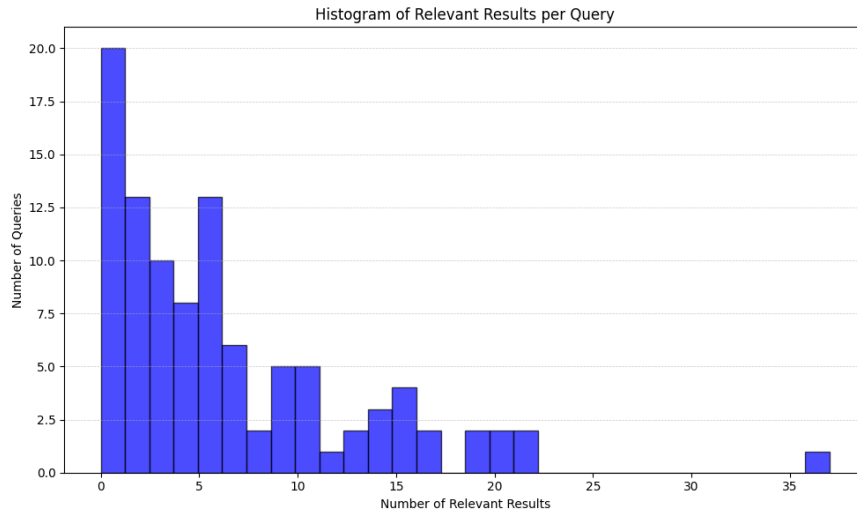


Figure 3: Number of Relevant Results in Top 100 Results

Number of Correct Query Results

Simple Tokenization

I computed Average F1@(100,50,20,10) for each query and then computed the average of these average F1 values Average F1: 0.0904361

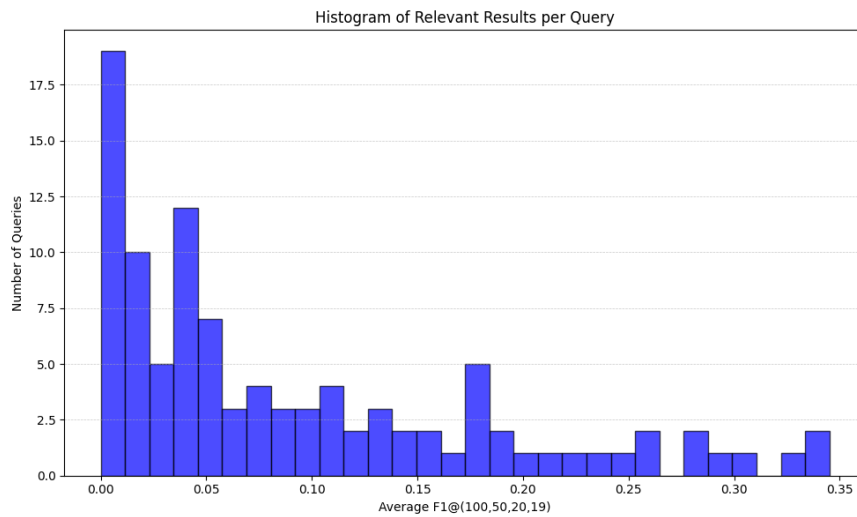


Figure 4: F1@(100,50,20,10)

Byte Pair Encoding

Average F1: 0.066461

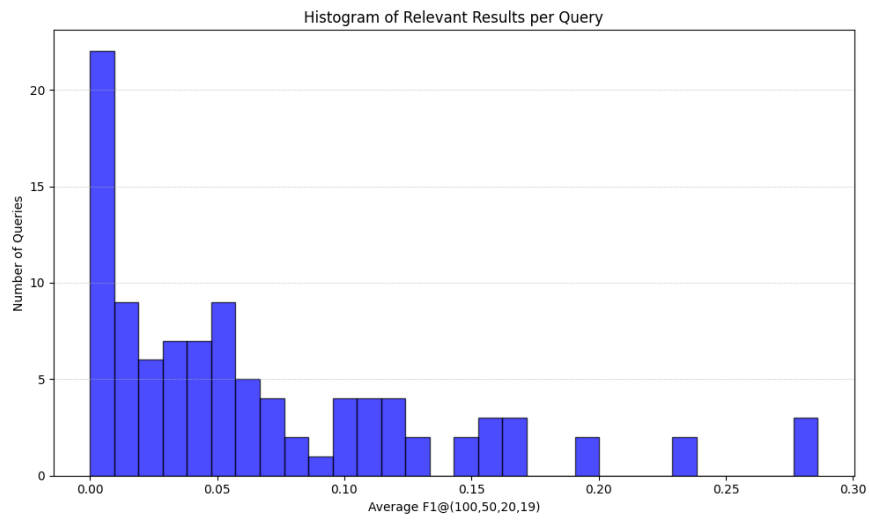


Figure 5: F1@(100,50,20,10)

Note: My Index and Dictionary files can be found at `/home/cse/btech/cs1200336/COL764/A1/submission_files`