

Information Retrieval - Assignment 2



DEPARTAMENTO DE ELETRÓNICA, TELECOMUNICAÇÕES E INFORMÁTICA

MESTRADO INTEGRADO EM ENG. DE COMPUTADORES E TELEMÁTICA

ANO 2018/2019

Beatriz Coronha 92210

Pedro Santos 76532

1. Introduction

The first part of this report will explain the development of the first two exercises of the assignment number two. In this assignment was extended the indexer to apply tf-idf weights to terms using the improved tokenizer and also it was created a class that implements a ranked retrieval method. Then, in the second part it will be explained the exercise number three that consists in implementation of a positional index and phrase search.

2. Class Diagram

The class diagram for the project is presented below (Figure 1) and it is in the folder of the assignment to be easier to analyse it.

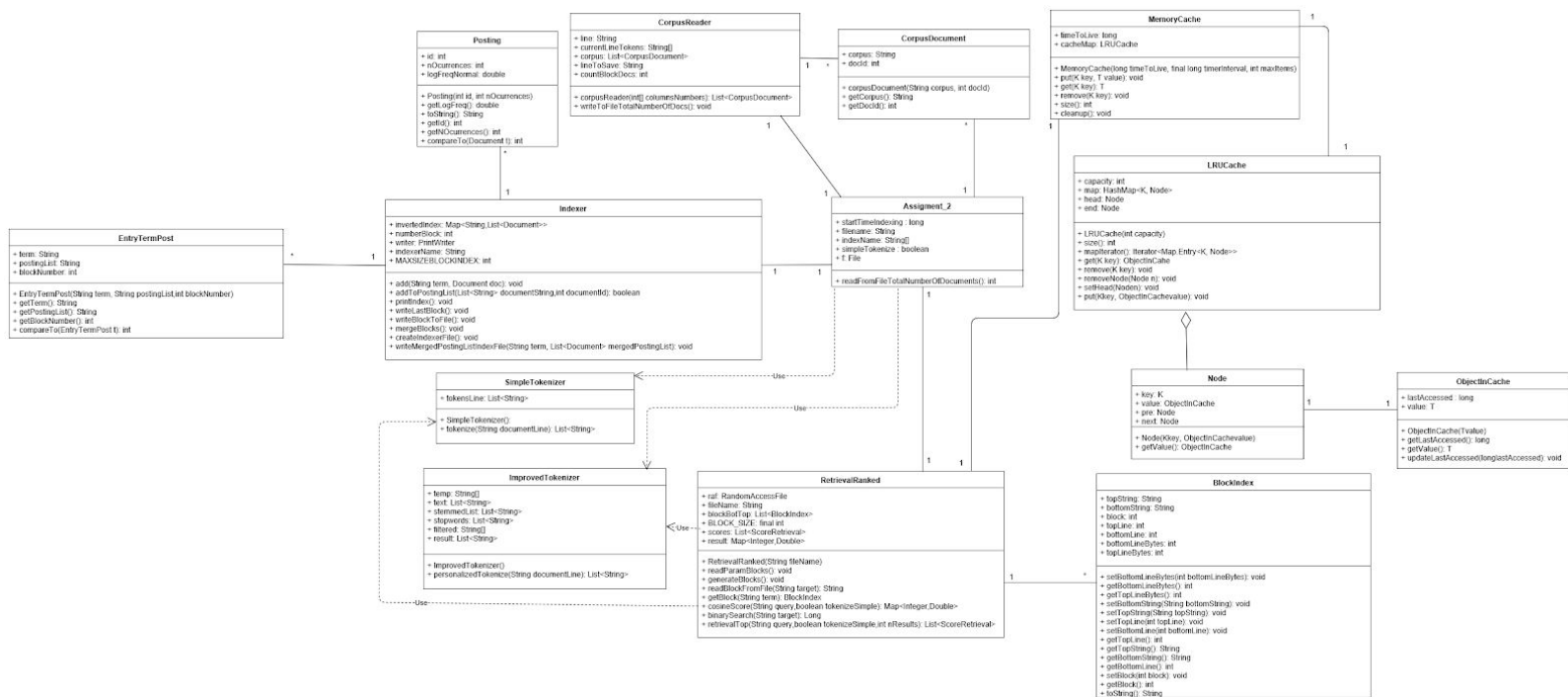


Figure 1: Class Diagram

3. Description of Class/Methods

This section will explain the goal of each class, the methods in them and how they're connected. The classes that were created for this assignment will be fully

explained and the changes made to the class from the assignment number one will be explained.

a. Assignment_2.java

This is the main class of the program. Like in the previous assignment there is a variable filename where it is put the path of the file that has the documents to index and a boolean variable to choose if it is to use the simple tokenizer or the improved one. It was added the code to verify if the indexer is created. If the indexer is not created it will be created like explained in the previous report, however now it will apply tf-idf weights to terms. If the indexer is already created it will skip that part and the next step will be the part of ranked retrieval method of the indexer. In this part the configurations that could be changed are the size of the cache that will be used, the time interval that the cache will verify for objects that are in the cache for more than a certain time and it is also possible to define that. If the objects are in the cache more than the define time, they will be erased. So, in this part it is instantiated the class RetrievalRanked and after the instantiation the function to generate the sub blocks of the indexer will be called, after that the retrievals can be made using the function *cosineScore* or the function *retrievalTop*. The function cosineScore will return all the results in an hashmap. The function retrievalTop will return the results in an array and it will be sorted by the high score, also this function receives an argument to define the maximum number of results.

b. CorpusReader.java

This class has only a difference from the last assignment. Now, it has a function that will be called after reading all the documents and this function will write in a file named *indexer_number_of_docs.txt* the total number of documents. This file is read in a function in the main class and the returned value will be used to instantiate the RetrivalRanked class.

c. Indexer.java

The goal of the first exercise on this assignment was to use tf-idf weights to terms. So, to accomplish this task in the function that adds to the indexer the term and the correspondent posting list it is necessary to do some extra operations. In the assignment number one, was only necessary to count the times that each term appears in the document and then add that information to the indexer. Now, it was necessary to count also the number of times that the term appears in the document. Then, it is used that result to calculate the Log-frequency weighting for each term that has a different count. Doing like this it is not necessary perform operations that we already know the result. So, basically it is saved in a hashmap the result where each count will have an result associated that corresponds to the $1 + \log_{10}(tf)$. The next step is to calculate the

value that will be used to normalize the results. To do that it is calculated the square root of the sum of the squares of the Log-frequency weight of each term in the document. The last step is to calculate the normalized values. Here it is done the same thing to not repeat operations and a lot of the values are the same, so this will optimize significantly the calculation. Therefore, basically for each different count will be calculated the division between the Log-frequency weight and the normalization value. After that the process is the same, however it is saved the tf-idf weights and not only the frequency of each term. So, for accomplish beyond this changes were made some other minor things, such as the Posting class was changed to have an attribute to save the weight.

d. Posting.java

To this class was added the attribute to save the normalized value of log frequency and the function to get the respective value. Also, it was added a different constructor to be able to instantiate an object using the document id and the weight. So, the changes made didn't compromise the use of this class in older versions of the project.

e. BlockIndex.java

This class is used for a more efficient search, the general idea is that the indexer is divided in blocks and this class is used to keep the number and the first and last word of it.

To do that, the indexer is read line by line and copied to another file that will be the block, this files are identified by a number, the function responsible for doing that is *generateBlocks*. After the blocks are generated, the objects of "blockIndex" are saved in an array, this objects will have the number of the block and its first and last words. That way, when a search is been made and the term is not in cache it's possible to know in which block the term is and only this block will be read to search for the query term.

f. RetrievalRanked.java

This is the class that will actually do the retrieval and the block operations mentioned previously. This class reads the indexer and then does the *generateBlocks* and the searches for the term and then calculates the scores.

This class has an constructor that accepts the name of the indexer, the total number of docs that were read to index the terms and the three settings for the cache that already were explained in the final of the section of the class `assignment_2.java`. So, the constructor will instantiate the cache. After that the main class calls the function *generateBlocks*. This function will check if the sub blocks of the indexer were created. If the blocks were not created the function will create the blocks,

like was described in the previous class `BlockIndex`. Then, after generating all the blocks it will serialize using the function *serializeBlocksArray* the array that has the collection of objects `BlockIndex`. By doing this is not necessary to read the blocks again to know the first and last word that each block contains. Therefore, if the blocks were already created this function will call *deserializeBlocksArray function* and it will deserialize the array that contains the objects `BlockIndex`. The main function of this class is *cosineScore*. This function has two arguments and it is the method to make the ranked retrieval, so it will return a hashmap where the key is the id of the document and the value the score. The first argument that the function has is the query and the second argument is the boolean value to such the tokenizer that it is supposed to use. The first step is to tokenize the query. Then, for each term of the query will be seen if it is in cache the value. If the value is in cache it is received immediately the posting list of that term. Otherwise, it will call the function *readBlockFromFile*. This function will call the function *getBlock* that will return the object `BlockIndex` that belongs to the respective term. Therefore, the function *readBlockFromFile* will read line by line until find the term and then it will return the line. The next step in the *cosineScore* function is to add the received line to the cache and then for it will calculate the score of each document. Then, it was create a extra function called *retrievalTop* that has the same arguments of the *cosineScore* and the maximum numbers of results. This function will call the *cosineScore* and it will order the results by the score and it will return a maximum of results.

g. ObjectInCache.java

The cache is generic, so it was not only developed to this project. This class is the object that will be each entry of the cache. This object has two attributes. One attribute is the value of the last time that was accessed and the other is the object to save, in our case will be a string.

h. LRUCache.java

A LRU cache is is a hash table of keys and double linked nodes. This class has the methods to put, remove and get elements in the cache. The cache will have a hashmap where the key will be associated to `Node`. Also, it will save the head and the end of the added nodes to the cache. This Class `Node` has three attributes, one is the object of the class `ObjectInCache` and the others two are the before and next node. In this way it is very fast to remove the older elements in the cache.

i. MemoryCache.java

This class is the “top level” of the cache and it is the class instantiated by the `RetrievalRanked` class. The constructor of this class receives three arguments. One of them is the size of the cache and this will be used to instantiate the `LRUCache`. The

other one is the timer interval to use to sleep the thread that will clean the memory by calling the function *cleanup*. Then, this function will check for all the objects in memory that are in the memory more time than supposed. This time is defined by the first argument in the constructor. Then, the class has the functions to remove, add and get from the LRU Cache. So, It was created a cache that has a limit of objects that can save, if it is full and it is necessary to put more objects will delete the older ones and it will preserve in the cache the most used values and it will delete the values that are not usually used.

j. ScoreRetrieval.java

This class was created to use to save the results of the scores of each document in a list and then order the results. So, the class implements the Comparable interface to be possible to order the elements easily and it has two attributes the id of the document and the score.

4. Data Flow

The data flow for this project is presented below (Figure 2). In it is possible to see the most important steps of the code, first there is the “Initialization”, where the “CorpusReader” and the “Indexer” are initialized. Then, the “Processing” happens, where the Indexer is created, as explained in the last assignment. After this part, it is necessary to merge all the different blocks, calculate the score and then return the rank. Firstly, it is verified if the indexer is created. If the indexer is not created the CorpusReader reads 10000 documents and then the tokenize will process that documents and the tokens will be consumed by the indexer. Then, another 10000 documents will be processed and this process happens until there are no more documents. After this part, it is necessary to merge all the different blocks and create the final indexer. If the indexer is already created the RetrievalRanked will be instantiated and the function generateBlocks will be called. If the blocks are already created it is only necessary to deserialize the array that contains the information, like was explained before. If the blocks are not created, the blocks will be created using the indexer and the array that saves the information will be serialized to not be necessary read all the times the blocks. Then, it is possible to do the retrievals. It is used the function cosineScore to receive all the results in a hashmap or the function retrievalTop to receive the results in an ordered list and it is possible to choose the maximum number of elements that will be returned. When the search for the term is being made before searching the terms in the block, it will be checked if the term is in the memory cache.

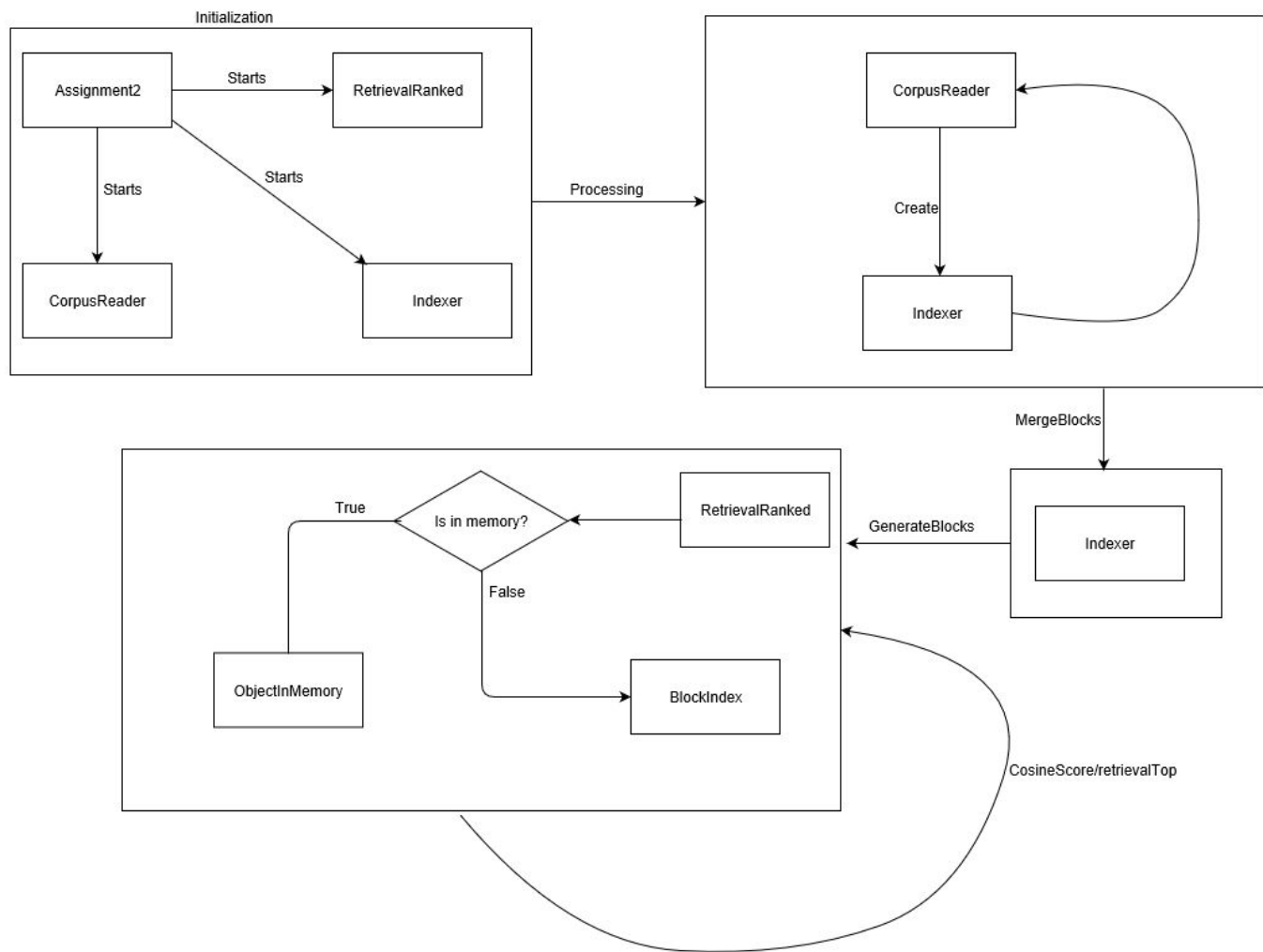


Figure 2: Data Flow

5. Results/Tests

a. Indexing Results

Total Indexing time = 33:37.771s

Maximum memory used in Indexing = 1249738832 bytes ~ 1.16 Gb

Total Indexer Size on disk = 3980859642 bytes ~ 4.00 Gb

b. Retrieval Results

To test the search of one term it was done 21 tests using 21 different random terms from the indexer, so the cache will not have any of the terms and it will be necessary to search on the index. The average time for each retrieval was 178 milliseconds. If all the terms were in cache the average time is 108 milliseconds. One thing to have in mind is if we take of the prints of the results the time without any term in the cache is 57 milliseconds and the average time when the terms are in cache is 14 milliseconds. It is possible to conclude that the retrieval is so fast that the prints increase a lot the time, also because we are printing all the results and some terms appear in a lot of documents. In all of the next tests the time of the prints will not be considered.

Then, it was done the test with query of 3 terms. It was tested with 8 different queries and all the terms were different, so the cache will be empty and all the terms belongs to indexers. The average time without prints was of 189 milliseconds. If all the terms are in cache the average time decreases to 45 milliseconds.

The last test was done with queries of 5 terms that are in the indexer and the total number of tests were 4. The average time with no terms in the cache was 315 milliseconds. If all the terms are in cache the average time was 78 milliseconds. In this test also all the terms were different and they are in the indexer.

6. How to Run

There aren't more dependencies or requirements in relation to the assignment number one. Like before, It is only necessary to install the porter stemmer and in the main class define the file that has the documents to be read.

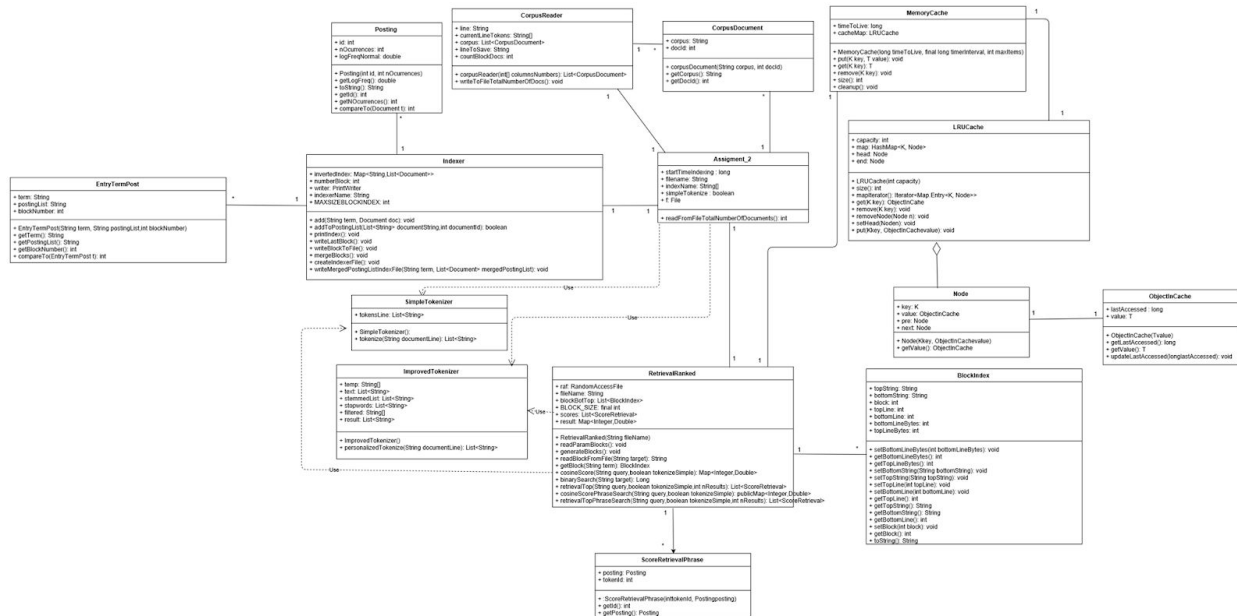
7. Extra Exercise

a. Problem Context

The exercise number three consists in transforming the index to be possible to make phrase search. In the next subsections will be explained the changes that were made to the previous versions of the index to be possible to do this type of search.

b. Class Diagram

The diagram is very similar to the previous one. This diagram has three more functions in the class `RetrievalRanked` and it was created a new class `ScoreRetrievalPhrase` to to handle the phrase search. The classes will be explain in the next subsection.



c. Description of Class/Methods

In this subsection will be explained each class and methods that suffer changes to be able to create the positional index and also the ones that were created.

i. Indexer.java

In this class several changes had to be made. When it is been add to the inverted index each document, it is necessary to save the position of each token in the document and if the token appears several times it is necessary to save all the different positions. The changes to handle this problem were made in the function `addToPostingList()`. To solve this problem is created a hashmap where the key will be the token and the value a list of integers and each element of this list will be the several positions of the token in the document. Then, it was only necessary to create a new attribute on `Posting` class to save the list. The function `writeBlockToFile()` suffer minor changes in how it separates the term from the posting list when it is writing the block of the index to a file to be easier to do the merge. The function `mergeBlocks()` suffer also

changes to be able to read each line of the blocks with the new format that supports positional lists. The two functions `writeMergedPostingListIndexFile()` and `writeIndexFile()` were changed to write the final index with the new format.

ii. Posting.java

Like was explained in the previous section it was added an attribute that it is a list of integers to save the different positions of a term in a document. Also, it was added two more attributes. One of them is a string to save the weight of the term in a doc. The other one is a string to save the list. These attributes were created to be optimized the merged process, because like this we don't have to convert the string to double and also the string read from the blocks with the list of number to a list to write again in the final index. So, it were added all the constructors and functions to be able to use this attributes. Also, the function `toString()` was changed to support the new index format that supports positional search.

iii. RetrievalRanked.java

In this Class were added three functions to handle the phrase search. These functions are the `cosineScorePhraseSearch()`, the `retrievalTopPhaseSearch()` and `checkIfTermsNextToOther()`. The `cosineScorePhraseSearch()` is the main function to do the phrase search. This function receives two arguments. One of the arguments is the query and the other one is the boolean value to check if it is to use the simple or the improved tokenized. The first thing that this function does is to tokenize the query. Then, for each term on the query it will get the posting list from the index or from the cache and it will for each term create an arraylist of the posting of each one and then I will save the list of postings of each term in a list. After this process, we have in memory everything necessary to calculate the scores and also to check the documents that have that specific phrase. There is an array named DIF. This array has the size of the number of terms on the query less one and it is initialized each position with one. This array can be used to make search with specific distances between terms. So, for example if the query has two terms, the array will be of size one and the position will say how many words can be between the words of the query. However, the array is there and can be changed in the code to test, it was not implemented to be manipulated by the user. In the next phase of the function each doc id of the first term will be compared with the doc id of the second term, (Also, to be faster we know that the posting list is ordered by the doc id, so if the doc id of first term is bigger of the second term it will stop.) Then, if one of the docs id of the first term of the query is equal to one of the second term it will compare the posting lists. Then, if any of the positions of the second term is one element higher it means that the terms are following each other

(here the DIF array also it is used to increment the difference that we want between each term) and they are added to a temporary list of objects ScoreRetrievalPhrase. This objects saves the postings and the tokenId. The tokenId corresponds to the index of the arraylist of postings of the term of the query. If the query has more than two terms it will enter in a while loop that will call the function checkIfTermsNextToOther(). This function returns a list of ScoreRetrievalPhrase. If this list is not empty it means that the next term of the query is also following that document that was being checked. So, if all the next terms of the query are in the same order has the document that is being checked this function will return always a arraylist of ScoreRetrievalPhrase and this result will be append to the arraylist with all the results. Then, it will exit from the while loop and if all the terms are in the document are in the same order, it will be calculated the score of each document and added to a hashmap that has keys the doc id and the value is the score. In the end after all the cycles it will be returned the hashmap with all the scores of each document. The function retrievalTopPhaseSearch() is equal to the same as the one created to do the normal search. The only difference is that it calls the function cosineScorePhraseSearch() to do the phrase search.

Some changes were done in the other functions to be able to read the index with the new format, however nothing special was changed or added.

iv. ScoreRetrievalPhrase.java

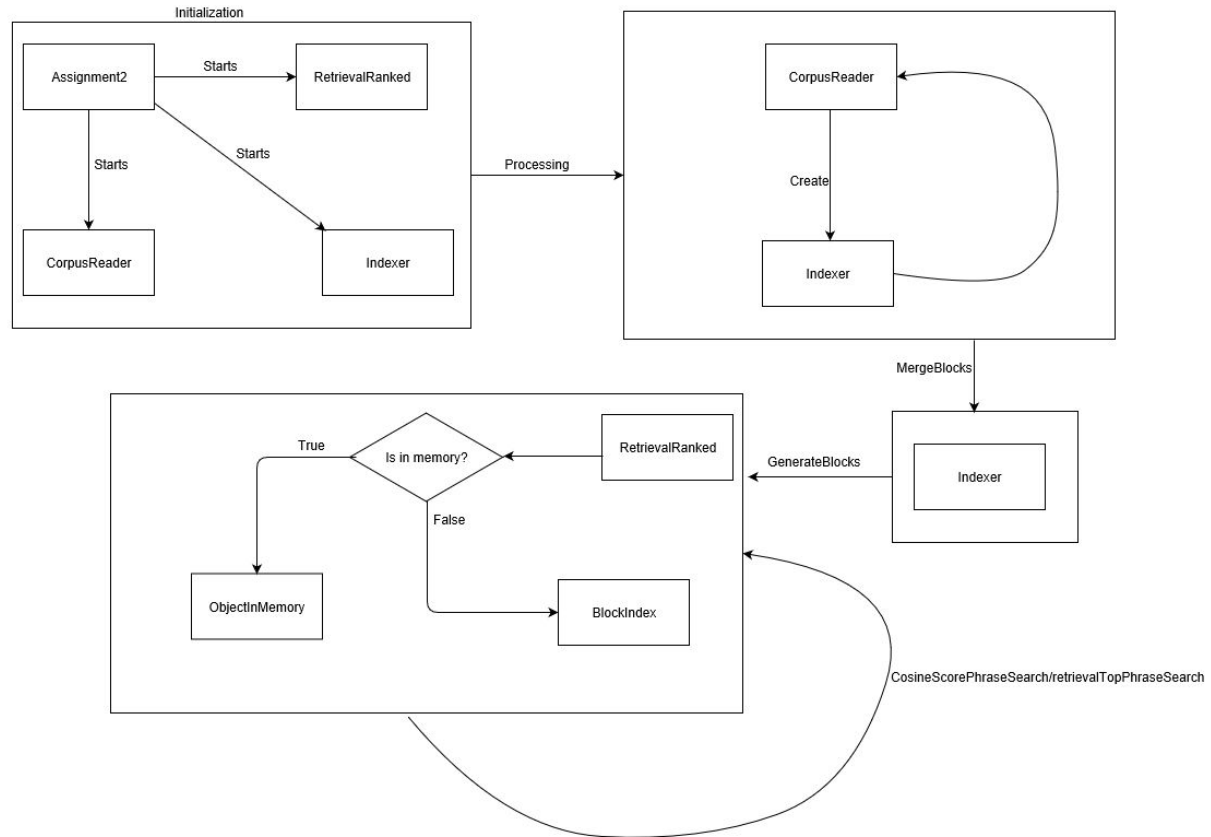
Class created that has as attributes the tokenId and the posting to be able to known of which term belongs the given posting to be able to calculate the score.

v. Assignment_2.java

Only added code to test and try the phrase search. Also, it was added a function that receives an arraylist of strings. This function can be used to do a lot of phrase searches. Each positions of the arraylist will be a query and it will be returned the ten more relevant documents by order of each search and also the time that took to calculate that.

d. Data Flow

The dataflow is equal to the one used in the exercise one and two. The only difference is in the called functions to do the search.



e. Results/tests

i. Indexing Results

Total Indexing time = 46:37.771s

Maximum memory used in Indexing = 1438536768 bytes ~ 1.33 Gb

Total Indexer Size on disk = 4722789710 bytes ~ 4.7 Gb

ii. Retrieval Results

In the phase test using the simple tokenizer and the 400 mb file was surprisingly fast, the exact time was not registered but the search made depending of the number of terms and the used terms would be between 300 ms and 3 seconds. And some of the results were verified by reading the specific line of file.

Using the dataset of 4 gb and using the improved tokenizer results:

Query: “eye-catching in reality”

Results:

docId: 198

time: 5939 milliseconds

```
root@rute-X550JK:~/Documents/cadeiras/Sano/ri$ sed -n 200p amazon_reviews_us_Wireless_v1_00.tsv
JS 20320537 RN6AFF0TC60NI B08PDSV6K2 414487636 HTC One M8 Case, StarCity HTC One M8 (Model 2014) Luxury Leather Back Hard Case Cover with Card Holder Snap on Cas
H 4 N Y Not as eye-catching in reality. Expected the case to be of slightly better quality and the graphics on the phone to be more vivid but when I got it, the cas
ap and not to my liking. However, you get what you paid for. 2015-08-31
```

It is possible to see that the document in the line 200 has that phrase. One of reasons of the offset between the line and the document id is first that is skipped. The other reason of the offset being bigger than one is unknown.

Query: “ GoPro is perfect for two cameras”

Results:

docId: 98, 650788

Time: 61659 ms

Query: “couple hours before”

Results:

docId: 2587801, 2010704 , 1050666, 3800406, 4206934, 1276741, 3302098, 214882, 3902913, 4107039

Time: 266790 ms

Query: “Switzerland or China”

Results:

docId : no results

Time: 253 ms

Query: “official Apple”

Results:

docId: 476273, 3550724, 545938, 1869739, 7165595, 5184351, 5663775, 4337571, 2745134, 6850335

Time: 145861 ms

Query: “cable from milklor”

Results:

docId: 921

Time: 950 ms

```
root@rute-X550JK:~/Documents/cadeiras/Sano/ri$ sed -n 923p amazon_reviews_us_Wireless_v1_00.tsv
JS 27289246 R50B06IQ8WB7 B00E4KLY34 153717096 BaoFeng UV-82 Dual Band Two-Way Radio 136-174MHz VHF & 400-520MHz UHF (Black) Wireless 5 0
ally nice equipment. Nice radio, perfect size, works well, programing requires some research if you haven't done it before, I spent 2 hrs to get my windows based computer syste
t 30 min trying to figure out how to make windows let me down load the required programing software, then I looked at videos on you tube and was programing these baby's likity spli
r cable from milklor, do not get the cheep cable seen on a-z, etc. the program you need is called chlrp and you will need the daily build. I'm using these for simplex in wildern
co fire and rescue, and of course noaa weather. These radios have the solid feel of professional equipment, some people complained about the chargers but mine worked perfectly, m
08-31
```

Query: "short distance"

Results:

docId: 4084082, 5224110, 2390832, 3894961, 4418378, 5033718, 7185907, 4482188, 4215662, 759207

Time: 11025 ms

The results were not the expected after testing in the small file and with the simple tokenizer, however the tests done in the simple tokenizer could be with terms with postings lists much more small. Nevertheless, by the tests runned we can see that is working and returning the right id of the documents. The id returned with the line of the file has an offset, but it was expectable an offset.

8. Conclusion

To have this performance and also to achieve this solution it were made a lot of other solutions, however not so good in performance. The first approach was to do binary search to the indexer, to take advantage of the fact that all the terms are sorted, but the indexer was too big to do that with good performance. So, the next test was to not do binary search in all the document, but to make a pre processing to know where are located the terms that begin with a given letter and then, it was taken the first letter of the term that it was to search in the indexer and it would be done the binary search only in that block. After this test, it was noticed that they are for example a lot more terms starting with an 'a' than with a 'z', so the time to retrieval was not consistent and to search a term starting with 'a' was much more slower than a letter by 'z'. So, the next approach was to divide the indexer in blocks and then do the binary search on that block. This solution was the best until now, but it was not good enough, because to do the binary search it is necessary to use the random access file and this is quite slow. Only after so many attempts it was found the solution that would have a very good performance with low memory use and with the implementation of the cache the performance increased even more.

The last exercise to implement the positional index and phrase research was also done with success. The performance could be more improved, however we were no expecting such big difference between using the small dataset with the simple tokenizer and the big file and the improved tokenized. Also, it would be possible to improve the phrase research by accepting in the queries tags to be possible defining a distance between the terms, even more because the algorithm is prepared to deal with this case, but it was not developed code to extract from the query that information.

9. Bibliography

The cache was developed with the support of this tutorial

<https://crunchify.com/how-to-create-a-simple-in-memory-cache-in-java-light-weight-cache/>

and this page

<https://www.programcreek.com/2013/03/leetcode-lru-cache-java/>