# Search Engine

## I.  Data structures & functions

### 1.  General

This search engine program is composed of 4 methods and a *main()* function. The inverted list is made is set up using an *unordered_map*, allowing us to be faster and more efficient because it uses a hash table. The unordered_map takes a string as the key and an array of integers as the value. The string represents the word we search, and the integers represent the index numbers of the lines in which they are.

### 2.  openFile()

The method *openFile()* allows us to initialize our unordered_map. We give a sting which is the path to a .txt file as the first parameter, and the empty unordered_map as the second parameter. Then, the function fills the unordered_map with the data in the file. At the end, it prints on the screen the time it took for the method to create the inverted list.

### 3.  printMap()

This method takes the unordered_map as parameter and print everything in the format: "word: list of lines the word is in".
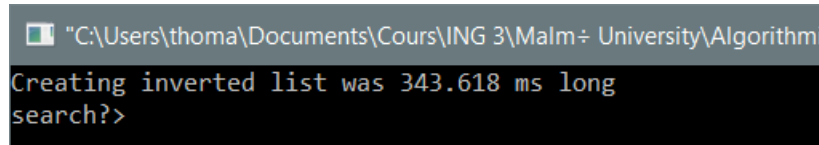
### 4.  query() and skipTo()

The method *query()* is the central method of the program. It is the one that allow us to search and find the lines in which the word(s) we want is/are situated. It takes the inverted list (unordered_map) as parameter. Then the user needs to input the word/bag of words he wants to find. Each word of the input will then be searched in the inverted list and their data copied in another inverted list (this could have been skipped, but it makes the program easier to understand and the speed results are still very good). If wanted, it is possible to show the results of every words by uncommenting blocs in the if/else. Then, the program search for the word that is in the smaller number of lines. The final results are stored in an array of integers, this array is filled by the content of the array of integers of the word precedingly found (with the smaller number of lines). Then, we run through the arrays of the other words entered by the user and we keep those that are in every word's list. For that we compare the *results* array with the array of each words (that are not the word with the smaller number of lines), and we put every word that is in both in a temporary array. When we have analyzed all the array, we copy the content of the temporary array in the results one and we go to the next word.

To do this efficiently we use the *skipTo()* function. It takes an int vector, a number to find (integer) and a start (unsigned int). Its purpose is to find the "number to find" in the vector. The start is a number, and we are sure that the number we are looking for is not before it. Finally, we browse the vector one by one in order to find or pass the number we want. When we find it, we return its position in the vector. If we didn't find or pass it, we return the size of the vector. We also could have used interesting methods we saw in class such as skipping 2 times more number each time (1->2->4->8…) using the "< 1". But as the results with this technique were very good, I decided that it would not be useful to implement it.

## II.    Observations

### 1.    Results and displayed content

When we launch the program, we have a first sentence displayed telling us how much time it took for the program to create the inverted list. It also writes "search?>" and give you the right to write something.
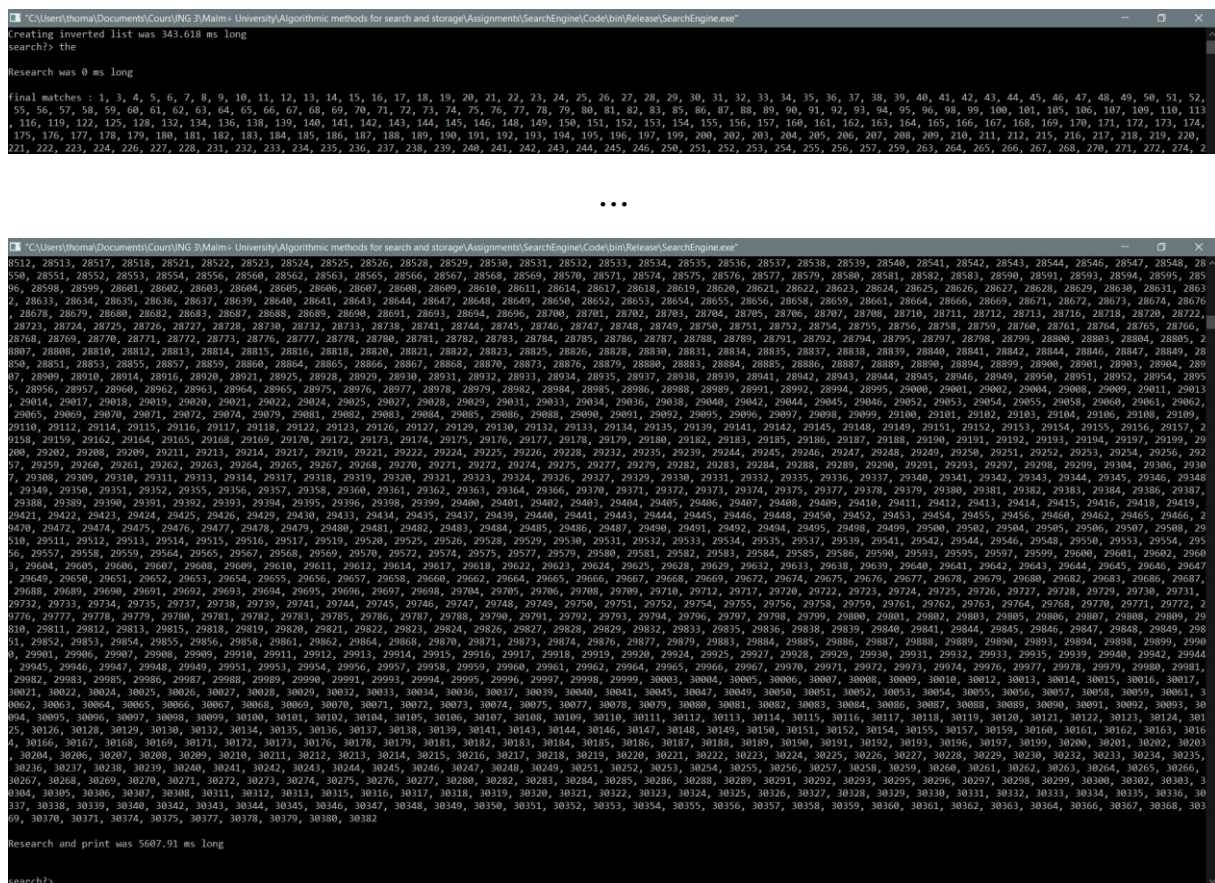


On this screen you can enter your query. 3 things will be displayed:

- The time it took to find the results
- The word and its results
- The time it took to find and print the results

Here are some examples of results:

- For "the"



…

- For "for those of an"

```
search?> for those of an

Research was 0 ms long

final matches : 6994, 17609, 18301, 22152, 22976, 29052, 29366

Research and print was 0 ms long
```

- For "in the beginning earth"

```
search?> in the beginning earth

Research was 0 ms long

final matches : 1, 29254

Research and print was 0 ms long
```

- For "foretteap the beginning god in an"

```
search?> foretteap the beginning god in an

Research was 0 ms long

No match found

Research and print was 0 ms long
```

- For "and god said let us make man in our image after our likeness and let them have dominion over the fish of the sea and over the fowl of the air and over the cattle and over all the earth and over every creeping thing that creepeth upon the earth"

```
search?> and god said let us make man in our image after our likeness and let them
 have dominion over the fish of the sea and over the fowl of the air and over the
cattle and over all the earth and over every creeping thing that creepeth upon the
 earth

Research was 0 ms long

final matches : 25

Research and print was 0 ms long
```

2. **Analyze of the results**

As we can see, all the tests we tried were very fast: less than 1ms to find the results. Even with long sentences the results were very fast. The print time was sometimes a little bit long, but the algorithm in its own is very efficient.

## III.    How can we be sure that there is no query that would make it perform much worse?

There are a lot of verifications that have been put in place in the query function that permits to make sure everything works well. First of all, if we find a word that has no match at all while browsing to find the word with the smaller number of occurrences, we leave the loop and we go straight to the printing part. Moreover, the loop that compares the results with each inverted list can be ended either when we have looked in every array OR when the results array is empty which allows us to end quickly if we don't have results while there is still 1000 word to look up. Also, the skipTo function starts its research at the index that was used in the calling function which allows us to save a lot of time browsing the array knowing nothing will be there.