

# Assignment Justification

## Data Structures

Since the data of this assignment has a fixed size in term of the number of words. Moreover, the Application should have to store all the data to perform search tasks, so that I decided to use ArrayList and HashTable as the core Data Structures to maximize the performance of the application by making use of the fact that they have the constant time complexity for most of the frequently used operations such as Access Item (into an ArrayList), Add Item (at the end of an ArrayList), Insert a Key (into a HashTable) and Search for a Key (in a HashTable) [1].

On the other hand, the main purpose of this assignment tasks are to find the position of words and phrase, so I decided to use a Standard Trie [1, pp. 586] to store all the position of every word. Moreover, the entire document consists of multiple sessions so I created a data structure to store all the content of each section so that I can easily search for words or phrases in a particular section using this data structure.

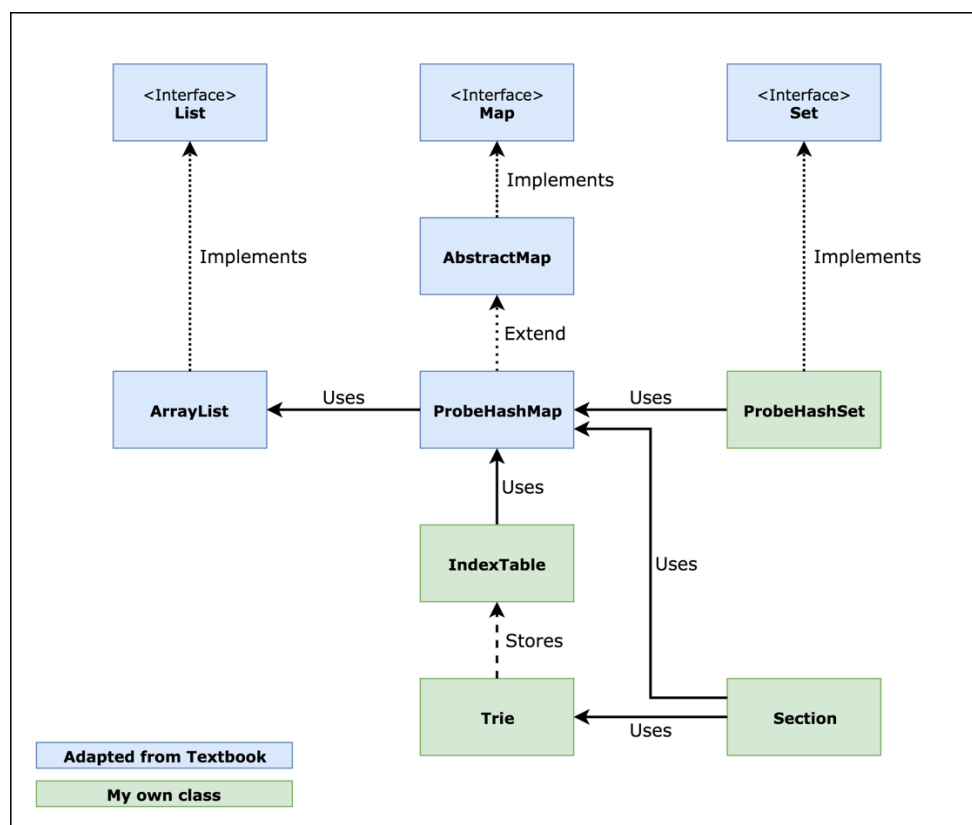


Figure 1. Data Structure Diagram

Each section contains a Map that stores all the lines, whose key is the line number and the value is the line text. Moreover, In order to enhance the search speed, I created an IndexTable for each word in a Trie which can be considered as an occurrence list of that word. The IndexTable is basically a MultiMap of line numbers and a List of columns that represents the word's position. This idea is based on the Inverted Index File concept which is mentioned in the textbook [1, pp. 594]

## Justifications

Firstly, I used ArrayList to store the section indexes because ArrayList reserves the ordering of the sections while reading from the file, so I don't have to sort it again. Then I added all the text of line into a Map which maps line numbers to the text of that line, I also add the position of each word into an IndexTable of the associated section. So that I can search on it later.

By using the Trie Data Structures, I can get the all the positions of a word with the time complexity of  $O(m)$  with  $m$  is the length of the word.

For the WordCount method, I only have to return the size of the index tables which is very efficient with the time complexity of  $O(1)$

For the PhraseOccurrence method, since I store all the text in my data structure, so I only have to search for the first word of the phrase and then loop through the phrase characters to compare texts. by using this method, it runs with the complexity of  $O(mn)$  with  $m$  is the number of characters in phrase and  $n$  is the number of first word occurrences.

For the PrefixOccurrence method, I use the recursive method to find all the descendants of a node in the Trie, so that I can get all the Index Table of the words that match the prefix with the time complexity of  $O(n+m)$  with  $n$  is the number of characters in the word and  $m$  is the number words that match.

For the Search Methods that include basic logic operations (AND, OR, NOT) I decided to use the Set data structure incorporated with set operations such as *addAll*, *removeAll* and *retainAll* [1, pp. 445] to manipulate the compound result of multiple words and logic operations. By using Sets, I can get the result of complex search with the time complexity of  $O(sn)$  with  $s$  is the number of sections that we are searching on, and  $n$  is the number of occurrences of all words in the query. Since the number of sections  $s$  can be considered as constants, so the time complexity is close to  $O(n)$ .

## Alternative Approaches

In order to enhance the performance, we could use the Compressed Trie instead of the Standard one because the compressed tried can reduce the redundant nodes in the Trie. One other potential improvement method is using the text compress technique [1] (e.g Huffman's Algorithms) to reduce the size of the text of line before storing them.

At first, I tried to store a Line as an object and has a LinkedList of words, and then an IndexTable can store the reference to the line object instead of line number as in the current implementation. However, I think it's a waste of computational resources, so I decided to simplify the data structure by just storing a sanitized string of line.

## References

- [1] M. T. Goodrich, R. Tamassia, and M. H. Goldwasser, Data structures and algorithms in Java. John Wiley & Sons, 2014