**CMT219 – Part 1 Report**

1. Functionality

A) Figure 1 is a screenshot of the result for running my code. It shows an ArrayList of strings as required, which contains the valid words in the given document based on the given vocabulary.
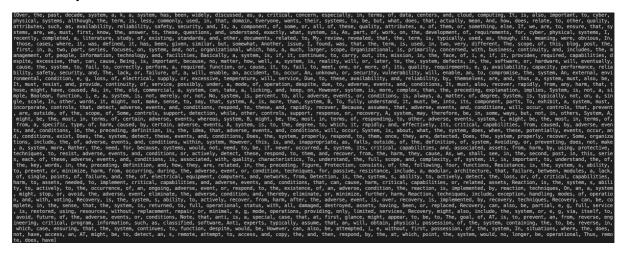


Figure 1. Printed ArrayList of valid words in the document based on the vocabulary.

B) Figure 2 shows the results of calling implemented merge sort algorithm to sort the first 100, 200, … words in the valid word ArrayList. For each sorting, the time used and count of moves/comparisons during sorting is printed.

Figure 3 show the sorted ArrayList after the last sorting, i.e., soring all the elements in the ArrayList in alphabetical order, confirming the correctness of the implemented merge sort algorithm. The results contain multiple same words since ArrayList can contain elements with the same values.



Figure 2. Calculated time used and counts of moves/comparisons during sorting.

Figure 3. Sorted ArrayList of valid words in alphabetical order.

## 2. Design

A) As mentioned in the requirement document, the "two text files are relatively big, you should consider how to make your program efficient". Therefore, I choose to use the **HashSet<String>** data structure to store the vocabulary words, as shown in Figure 4.

 Using HashSet could avoid redundant words in the set, and it could also provide a constant O(1) time complexity for testing whether a word in the document is contained in the HashSet or not. Since the document is relatively big (many words), using HashSet to improve the speed of checking each word is significant for the algorithm efficiency of the program.

```java
// vocabulary is a HashSet of strings, to store the key words in a given file.
// Using a HashSet will provide O(1) time for checking whether a single
// word is in the vocabulary

HashSet<String> vocabulary = new HashSet<String>();
```

Figure 4. Using HashSet data structure for the vocabulary.

B) For sorting, I implemented the merge sort algorithm as required, as shown in Figure 5. It partitioned the whole list into two half parts, and then recursively call merge sort algorithm to sort each of the two parts. Finally, a merge function is used to merge the sorted two half parts into one part.  The termination condition of recursion is start >= end, i.e., whenever there is at most one element, which means there is no need to further partition or sort.

The argument of int[] num_compartion in the definition of mergeSort function is used to maintain the count of moves/comparisons during sorting. Using int[] data type here is because of Java passes by values for int type, which cannot be obtained after calling the functions.

```
/*

The mergeSort function

Input: an ArrayList of strings, start and end position of mergesort, and a helper list for merging
Output: the ArrayList of strings will be sorted in ascending alphabetical order

*/

public static void mergeSort(ArrayList<String> words, int start, int end, String[] helper, int[] num_comparison) {

    // If start >= end, then there is no need to sort

    if (start < end) {

        // Set a mid position, and recursively call mergeSort twice
        // to sort the two half parts: start to mid, and mid+1 to end

        int mid = (start + end) / 2;

        mergeSort(words, start, mid, helper, num_comparison);
        mergeSort(words, mid+1, end, helper, num_comparison);

        // Merge the sorted two half parts after sorting them

        merge(words, start, end, helper, num_comparison);
    }
}
```

Figure 5. Implementation of merge sort.

To calculate the time used during sorting, I record the time before and after sorting, and calculate their difference, as shown in Figure 6.

```
// Record the time before sorting

start_time = System.nanoTime();

// Call the merge sort function to sort the first num_words_sorted words in valid_words

mergeSort(valid_words, start: 0, num_words_sorted-1, helper, num_comparison);

// Record the time after sorting

end_time = System.nanoTime();

// Calculate the time used for sorting

time_elapsed = end_time - start_time;
```

Figure 6. Calculate time used during sorting.

3. Ease of use and documentation

I use try and catch to make the code deal with invalid user input, such as file names which do not exist in the devices, as shown in Figure 7. In case of using non-existing file names, the program will not crash, and it will output a description of exception for fixing the issues, (e.g, ./google-10000-english-no-swears-wrong-filename.txt (No such file or directory)).

```java
// First, try to open the given vocabulary file,
// and add each word (line) in the vocabulary file into the HashSet of strings,
// using try and catch for situations like files not existing (e.g., wrong file names),

try {
    FileInputStream filestream = new FileInputStream(name: "./google-10000-english-no-swears.txt");
    BufferedReader br = new BufferedReader(new InputStreamReader(filestream));

    String strline;

    while ((strline = br.readLine()) != null) {
        vocabulary.add(strline.toLowerCase());
    }

    filestream.close();

} catch(Exception exception) {
    System.err.println(exception.getMessage());
}
```

Figure 7. Calculate time used during sorting.

As shown in above screenshots, the output is formatted in a easy to understand way, and the code is well commented with clear explanations.