# CMP-5014Y Coursework 2 - Word Auto Completion with Tries

Student number: 100225776. Blackboard ID: gny17hvu

## Contents

# 1 Part 1: Form a Dictionary and Word Frequency Count

Part one required me to implement the methods formDictionary and saveToFile. For the former, after trying a few different methods with different data structures I've settled on using a single hash table to store the dictionary, as this allows me to store both the word and it's count as a pair. A hash table data structure also allows for a very good average time complexity for searching and inserting of $\theta(1)$ (constant time) while maintaining a decent space complexity of O(n).

The design of saveToFile followed a similar path as formDictionary, whereby I experimented with designing and implementing many different solutions simultaneously with those of formDictionary before settling on what I feel is the most efficient implementation. This final version takes the key-value pairs one by one from the hash table previously mentioned, forms a string from them and adds them to a list data structure. This list is then sorted and sent to the saveCollectionToFile method given.

In addition to the above methods, DictionaryFinder also contains readWordsFromCSV, which produces a list from a given CSV format file.

---

**Algorithm 1** formDictionary(**A**,$l$) **return D**

---

**Require:** list **A** of length $l$ containing strings.
**Ensure:** **D**, the dictionary formed from **A**.
1:  **D** ← new map                                  ▷ *initialise map to store dictionary*
2:  **for** $i ← 1$ **to** $l$ **do**
3:     $found ←$ **false**
4:     **for each** $(k,v)$ pair $\in$ **D do**
5:        **if** $A_i := k$ **then**
6:           $v ← v + 1$                     ▷ *add one to the value on this entry*
7:           $found ←$ **true**
8:           break
9:     **if** $found :=$**false then**
10:       **D** $← (k ← A_i, v ← 1)$              ▷ *add a new (k,v) pair to the map*
    **return D**

---

In analysing formDictionary, we see that the fundamental operations are ($v ← v + 1$, $found ←$, break). The case we are interested in is the worst case - Order O. The Runtime Complexity function is as follows - the three fundamental operations only occur if there is a match therefore only a check of the map entry happens every loop, and this happens within a loop. Therefore for any given value of $l$ and for any size of map D (represented as $m$), we get

$$t_{formDictionary} = \sum_{i=1}^{l} \sum_{j=1}^{m} 1 \tag{1}$$

For solving the summations, we first remove the rightmost summation and apply the summation rule

$$\sum_{i=1}^{p} 1 = p \tag{2}$$

on the innermost loop to give $m$, which leaves

$$t_{formDictionary} = \sum_{i=1}^{l} m \tag{3}$$

Then we deal with the remaining summation using the rule

$$\sum_{i=1}^{p} a = a \sum_{i=1}^{p} 1 = ap \qquad (4)$$

to get

$$t_{formDictionary} = lm \qquad (5)$$

where $m \leq l$.

By characterising the above Runtime Function we find that the worst case for this algorithm is Order $nm$ or O($nm$).

---

**Algorithm 2** saveToFile(**D**)

---

**Require:** map **D**.
  1: **A** ← new collection                             ▷ *initialise collection to store strings*
  2: **for each** (k,v) pair ∈ **D do**
  3:     $a \leftarrow k$ & ”,” & $v$            ▷ *form string from k and v in requested format*
  4:     **A** ← a
  5: sort(**A**)             ▷ *ensure A is sorted alphabetically with mergesort*
  6: saveCollectionToFile(**A**,”A.csv”)       ▷ *pass collection into the provided method*

---

In analysing saveToFile, we see that the fundamental operation is either ($a \leftarrow k$ & ”,” & $v$, **A** ← a) or sort(**A**), so we will start with analysis of the former. The case we are interested in is the worst case - Order O. The Runtime Complexity function is as follows - the two fundamental operations each happen once per loop of D. Forming a string is a single operation and the complexity of adding it to a collection is dependant on the collection used, so we will assume a Linked List which inserts in constant time. Therefore if we assign the size of D to be $l$ we get

$$t_{saveToFile} = \sum_{i=1}^{l} 2 \qquad (6)$$

The constant on the right can be collapsed and using the rule

$$\sum_{i=1}^{p} 1 = p \qquad (7)$$

we get

$$t_{saveToFile} = l \qquad (8)$$

By characterising this first Runtime Function candidate, we find that the worst case for this algorithm is Order $n$ or O($n$). For the second candidate, we know that the best algorithm for sorting strings, mergesort, gives Order $n \log(n)$ worst-case performance, which is worse than Order $n$. Therefore from this analysis we can tell that the fundamental operation of the saveToFile algorithm is sort(A) and that it's characterisation is O($n \log(n)$).

3

## 2 Part 2: Implement a Trie Data Structure

Part two required the design and implementation of a Trie data structure comprised of TrieNodes to store a list of given words. Each node represents one character in a word and given words are chained together though storing lists of child nodes in each node. Methods were then written and implemented which store, split and recall words within the Trie. The TrieNode structure itself contains four variables:

1. *value* - This is a letter from a to z which this node represents.

2. *offspring* - This is a fixed list of 26 spaces, where each space can be assigned a child node representing a letter which would follow the current value in a given word.

3. *isKey* - This is a flag that is true when the current node is the end of a stored word and otherwise is false.

4. *level* - This is a number which shows how deep the current node is in the Trie structure.

The TrieNode also contains getter and setter methods, the following are the more complex examples of these:

- getOffspring(TrieNode **N**, char *c*) - We take 97 away from the ASCII value of c to get a numerical value between 0 and 25 for all letters of the alphabet, then if it exists, return the child node of **N** at that location in the offspring list.

- getAllOffspring(TrieNode **N**) - Add each node from **N**'s offspring list to a list and return it.

- setOffspring(TrieNode **N**, char *c*) - We take 97 away from the ASCII value of c and check the offspring list on **N** for an entry. If no entry exists we add a new node to the list and return true, otherwise we return false.

The Trie data structure holds only one variable, *root*, which is the address of the first TrieNode in the Trie and is assigned when the Trie is constructed. The Trie has a number of methods which control how nodes are added and and manipulated within it. I will detail these in full below.

---

**Algorithm 3** add(**T K**,*l*) return *f*

---

**Require:** Trie **T**, string **K** of length *l*.
**Ensure:** *f*, a flag indicating if **K** was added to the Trie or not.
1: TrieNode **C** ← **K**.*root*                                 ▷ *initialise node C with the value of this Trie's root*
2: *f* ← *false*
3: **for** *i* ← 1 **to** *l* **do**
4:     *k* ← **K**$_i$                                 ▷ *char k holds the value of the current character in K*
5:     TrieNode **N** ← getOffspring(**C**, *k*)
6:     **if N** := *null* **then**
7:         setOffspring(**C**, *k*)
8:         **N** ← getOffspring(**C**, *k*)
9:         setLevel(**N**, *i* + 1)                                 ▷ *call setter for level variable on node N*
10:         *f* ← *true*
11:     **C** ← **N**
12: **if C** := **not** *null* **then**
13:     setKey(**C**)        ▷ *set the isKey flag on node C to true when the last letter of K is added*
    **return** *f*

---

---

**Algorithm 4** contains(**T K**,*l*) **return** *f*

---

**Require:** Trie **T**, string **K** of length *l*.
**Ensure:** *f*, a flag indicating if **K** exists in **T** or not.
  1: TrieNode **C** ← **T**.*root*
  2: $f \leftarrow false$
  3: **for** $i \leftarrow 1$ **to** $l$ **do**
  4:      char $k \leftarrow \mathbf{K}_i$
  5:      TrieNode **N** ← getOffspring(**C**, $k$)
  6:      **if N** := *null* **then return** false
  7:      **C** ← **N**
  8:      $f \leftarrow$ getKey(**C**)            ▷ *set the value of f to the value of node C's getKey variable*
     **return** $f$

---

---

**Algorithm 5** outputBreadthFirstSearch(**T**) **return S**

---

**Require:** Trie **T**.
**Ensure: S**, a string of characters taken from reading the Trie breadth-first.
  1: **V** ← new list
  2: $l \leftarrow 1000$
  3: **S** ← new string with length of $l$
  4: $p \leftarrow 0$                                      ▷ *integer to hold position within S*
  5: **V**.add(**T**.*root*)
  6: **while V** := **not** empty **do**
  7:      TrieNode **C** ← **V**.remove()
  8:      $\mathbf{S}_p \leftarrow$ **C**.value
  9:      list **L** ← getAllOffspring(**C**)
10:      **for each** TrieNode **N** $\in$ **L do**
11:          **V**.add(**N**)

12:      $p \leftarrow p + 1$
13:      **if** $p := l$ **then**          ▷ *increase the length of S by a multiple of 10 if it's too short*
14:          String **E** ← **S**
15:          $l \leftarrow l \times 10$
16:          **S** ← new string with length of $l$
17:          **for** $i \leftarrow 1$ **to** $(l \div 10)$ **do**
18:             $\mathbf{S}_i \leftarrow \mathbf{E}_i$
     **return S**

---

---

**Algorithm 6** outputDepthFirstSearch(**T**) **return S**

---

**Require:** Trie **T**.
**Ensure:** **S**, a string of characters taken from reading the Trie breadth-first.

 1: **V** ← new stack
 2: $l ← 1000$
 3: **S** ← new string with length of $l$
 4: $p ← 0$
 5: **V**.push(**T**.*root*)
 6: **while V** := **not** empty **do**
 7:     TrieNode **C** ← **V**.pop()
 8:     list **L** ← getAllOffspring(**C**)
 9:     $m ←$ (length of **N**) $-1$
10:     **while L**:= **not** empty **do**
11:         **V**.push(**N**.get($m$))
12:         $m ← m - 1$
13:     $\mathbf{S}_p ←$ **C**.value
14:     $p ← p + 1$
15:     **if** $p := l$ **then**          ▷ *increase the length of S by a multiple of 10 if it's too short*
16:         String **E** ← **S**
17:         $l ← l × 10$
18:         **S** ← new string with length of $l$
19:         **for** $i ← 1$ **to** $(l ÷ 10)$ **do**
20:             $\mathbf{S}_i ← \mathbf{E}_i$
    **return S**

---

---

**Algorithm 7** getSubTrie(**T P**,$l$) **return U**

---

**Require:** Trie **T**, string **P** of length $l$.
**Ensure:** **U**, the Subtrie rooted at **P**.

 1: **U** ← new Trie
 2: TrieNode **S** ← **T**.*root*
 3: **if P** := empty **then return U**
 4: **for** $i ← 1$ **to** $l$ **do**
 5:     $c ← \mathbf{P}_i$
 6:     TrieNode **N** ← getOffspring(**S**, $c$)
 7:     **if N** := *null* **then return S**
 8:     **S** ← **N**
 9: **U**.*root* ← **S return U**

---

**Algorithm 8** getAllWords(**T**) **return L**

**Require:** Trie **T**
**Ensure:** **L**, a list of strings formed from **T**.
 1: **L** ← new list
 2: **V** ← new stack
 3: $l$ ← 35                                              ▷ *maximum length of an expected word*
 4: **S** ← new string of length $l$
 5: list **O** ← getAllOffspring(**T**.*root*)
 6: $m$ ← (length of **O**) −1
 7: **while O** := **not** empty **do**
 8:     **V**.push(**O**.get($m$)
 9:     $m$ ← $m − 1$
10: **while V** := **not** empty **do**
11:     TrieNode **C** ← **V**.pop()
12:     $p$ ← **C**.*level* + 1                          ▷ *position of the next character*
13:     **S**$_p$ ← **C**.*value*
14:     **if C**.*isKey* := true **then**
15:         **F** ← new string
16:         **for** $i$ ← 1 **to** $p$ **do**
17:             **F** ← **F** + **S**$_i$
18:         **L**.add(**F**)
19:     **O** ← getAllOffspring(**C**)
20:     $m$ ← (length of **O**) −1
21:     **while O** := **not** empty **do**
22:         **V**.push(**O**.get($m$))
23:         $m$ ← $m − 1$
     **return L**

# 3  Part 3: Word Auto Completion Application

The following main() algorithm relies on the data structures and algorithms explained in parts 1 and 2.

---

**Algorithm 9** main()

---

1:        ▷ *1. form a dictionary file of words and counts from the file lotr.csv.*
2: **D** ← new DictionaryFinder
3: list **L** ← readWordsFromCSV(**D**, "lotr.csv")
4: **D**.formDictionary(**L**)
5:        ▷ *2. construct a trie from the dictionary using your solution from part 2.*
6: **T** ← new Trie
7: **for each** ($k$,$v$) pair ∈ **D do**
8:      **T**.add($k$)
9:        ▷ *3. load the prefixs from lotrQueries.csv*
10: list **P** ← readWordsFromCSV("lotrQueries.csv")
11:        ▷ *4. for each prefix query:*
12: **F** ← new list       ▷ *list to store strings to be saved to file*
13: **for each** $p$ ∈ **P do**
14:        ▷ *4.1. Recover all words matching the prefix from the trie.*
15:      Trie **S** ← getSubTrie(**T**, $p$)
16:      list **M** ← getAllWords(**S**)
17:        ▷ *4.2. Choose the three most frequent words and display to standard output.*
18:      **C** ← new list       ▷ *list of lists to store frequency and position of words found*
19:      $c$ ← 0
20:      **if D**.contains($p$) **then**
21:        **M**.add("")
22:      $l$ ← size of **M**
23:      **for** $i$ ← 1 **to** $l$ **do**
24:        string $a$ ← $p$+ **M**.get($i$)
25:        $b$ ← **D**.get($a$)   ▷ *retrieve from dictionary the value for key matching the above string*
26:        $c$ ← $c + b$
27:        **E** ← new list
28:        **E**.add($b$)
29:        **E**.add($i$)
30:        **C**.add(**E**)
31:      **while C** := unsorted **do**
32:        using mergesort sort **C where C**.get($i$).get(1) > **C**.get($i + 1$).get(1)
33:      string **H** ← $p$
34:      **for** $i$ ← 1 **to** $l$ **and** $i < 4$ **do**
35:        $z$ ← **C**.get($i$).get(1) ÷ $c$
36:        **Y** ← $p$+ **M**.get(**C**.get($i$).get(2))
37:        **print Y** + " (probability " +$z$+ ")"       ▷ *see below for console output*
38:        **H** ← **H** + "," + **Y** + "," +$z$
39:      **F**.add(**H**)
40:      **print** ""
41:        ▷ *4.3. Write the results to lotrMatches.csv.*
42: saveCollectionToFile(**F**, "lotrMatches.csv")

---

## 3.1 Console Output

about (probability 0.56666666)
above (probability 0.3)
able (probability 0.1)

going (probability 0.2777778)
go (probability 0.24074075)
good (probability 0.16666667)

the (probability 0.626703)
they (probability 0.15395096)
them (probability 0.06811989)

merry (probability 0.94736844)
merely (probability 0.02631579)
merrily (probability 0.02631579)

frodo (probability 0.4909091)
from (probability 0.43636364)
front (probability 0.07272727)

great (probability 0.1969697)
ground (probability 0.18181819)
grass (probability 0.15151516)

goldberry (probability 0.6)
golden (probability 0.4)

sam (probability 1.0)


Dictionary saved to file lotrMatches.csv

# 4 Code Listing

Listing 1: DictionaryFinder.java

```
1  /***********************************************************************
2
3   Project      : CMP-5014Y - Word Auto Completion with Tries:
4                  AutoCompletion.
5
6   File         : DictionaryFinder.java
7
8   Date         : Friday 06 March 2020
9
10  Author       : Martin Siddons
11
12  Description : This class fulfills the requirements of part 1. It
       ↪ supplies
13  methods that read in a text document into a list of Strings, form a
       ↪ set of
14  words that exist in the document and count how many times each word
       ↪ occurs,
15  sort the list alphabetically and write those words and frequencies
       ↪ to file.
16
17  History      :
18  06/03/2020 - v1.0 - Initial setup, copied over readWordsFromCSV and
19  saveCollectionToFile.
20  07/03/2020 - v1.1 - Eventually settled on a implementation that
       ↪ works well in
21  Java. Editing pseudocode to fit.
22  ***********************************************************************/
23
24  import java.io.*;
25  import java.util.*;
26
27  public class DictionaryFinder {
28      HashMap<String, Integer> dict;
29
30      // Constructor
31      public DictionaryFinder() {
32          this.dict = null;
33      }
34
35      /**
36       * Static method to read in files from a CSV file and output an
           ↪ ArrayList of
37       * Strings, where each string is a line from the given file.
           ↪ Adapted from
38       * code given by AJB.
39       *
40       * @param file : File to read in.
41       * @return ArrayList of Strings read in from file.
42       */
```

```java
43    public static ArrayList<String> readWordsFromCSV(String file) {
44        ArrayList<String> words = new ArrayList<>();
45        try {
46            Scanner sc = new Scanner(new File(file));
47            sc.useDelimiter("[ ,\r]");
48            String str;
49            while (sc.hasNext()) {
50                str = sc.next();
51                str = str.trim();
52                str = str.toLowerCase();
53                words.add(str);
54            }
55        }
56        catch (FileNotFoundException e) {
57            System.out.println("Error: File not found. Aborting");
58        }
59        return words;
60    }
61
62    /**
63     * Static method to write a given collection to file.
64     * Adapted from code given by AJB.
65     *
66     * @param c    : Collection to save to file.
67     * @param file : Name of file to be created.
68     */
69    public static void saveCollectionToFile(Collection<?> c, String
        ↪ file) {
70      try {
71            FileWriter fileWriter = new FileWriter(file);
72            PrintWriter printWriter = new PrintWriter(fileWriter);
73            for (Object w : c) {
74                printWriter.println(w.toString());
75            }
76            printWriter.close();
77        }
78        catch (IOException e) {
79            System.out.println("Error: Unable to write file.
                ↪ Aborting.");
80        }
81        System.out.println("\nDictionary saved to file " + file);
82    }
83
84    /**
85     * Take an ArrayList of Strings and process each string into a
        ↪ dictionary.
86     *
87     * @param in : Arraylist to be formed into a dictionary.
88     */
89    public void formDictionary(List<String> in){
90        this.dict = new HashMap<>(in.size());
91
```

```java
 92         // For each string in the ArrayList, check if it exists in
             ↪ the
 93         // dictionary. If it doesn't, add it with a count of 1. If
             ↪ it does
 94         // exist, add 1 to the value on that entry.
 95         for (String s : in) {
 96             this.dict.put(s, this.dict.getOrDefault(s, 0) + 1);
 97         }
 98     }
 99
100     /**
101      * Call the Dictionary previously created with formDictionary
           ↪ and produce a
102      * formatted list. Ensure the list is sorted and send to
103      * saveCollectionToFile.
104      */
105     public void saveToFile() {
106         List<String> entries = new ArrayList<>();
107
108         // Take each entry from the map, format it and add it to the
             ↪ list.
109         for (Map.Entry<String, Integer> e : this.dict.entrySet()){
110             entries.add(e.getKey() + "," + e.getValue());
111         }
112         Collections.sort(entries); // Sort the list alphabetically.
113
114         // Call saveCollectionToFile to save the sorted, formatted
             ↪ list to file.
115         String filename = "output.csv";
116         saveCollectionToFile(entries, filename);
117     }
118
119     // Test Harness.
120     public static void main(String[] args) {
121         DictionaryFinder df = new DictionaryFinder();
122
123         // 1. read text document into a list of strings
124         ArrayList<String> in = readWordsFromCSV("testDocument.csv");
125
126         // 2. form a set of words that exist in the document and
             ↪ count the
127         // number of times each word occurs in a method called
             ↪ FormDictionary
128         df.formDictionary(in);
129         Collection<Integer> e = df.dict.values();
130         System.out.print("counts for words in dictionary (before
             ↪ sort): ");
131         System.out.println(e);
132
133         // 3. sort the words alphabetically; and
134         // 4. write the words and associated frequency to file.
135         df.saveToFile();
```

```
136            }
137    }
```

```
1  /************************************************************************
2
3   Project      : CMP-5014Y - Word Auto Completion with Tries :
4                     AutoCompletion.
5
6   File         : Trie.java
7
8   Date         : Thursday 12 March 2020
9
10  Author       : Martin Siddons
11
12  Description : Part 2 of the assignment is to make a Trie Data
         ↪ Structure to
13  hold strings and write methods to manipulate the trie.
14
15  History      : 12/03/2020 - Initial setup
16  14/03/2020 - Completed Q1.
17  13/05/2020 - Lost track on filling this in, just finished Q2.
18  14/06/2020 - Fixed implementation of getAllWords for part 3.
19  15/06/2020 - Finalised implementation by clearing up some methods.
20   ************************************************************************/
21
22  import java.util.ArrayList;
23  import java.util.LinkedList;
24  import java.util.List;
25  import java.util.Stack;
26
27  class TrieNode {
28      private        char       value;
29      private final TrieNode[] offspring;
30      private        boolean    isKey;
31      private        int        level; // track how deep this node is
          ↪ in the trie
32
33      // Constructors
34      public TrieNode() {
35          this.offspring = new TrieNode[26];
36      }
37
38      public TrieNode(char c) {
39          this.value     = c;
40          this.offspring = new TrieNode[26];
41      }
42
43      public char getValue()      {return value;}
44
45      public boolean getKey()     {return this.isKey;}
46
47      public int getLevel()       {return this.level;}
48
49      public void setKey()        {this.isKey = true;}
```

```java
50
51      public void setLevel(int l) {this.level = l;}
52
53
54      /**
55       * Check if a given letter exists as a child on this node.
56       *
57       * @param c : Character to find linked from this trie.
58       * @return : The TrieNode representing the given letter, or null
              ↪ if that
59       * letter is not an offspring.
60       */
61      public TrieNode getOffspring(char c) {
62          int pos = c - 97; // turn ASCII character value to numeric
                  ↪ (so a = 0)
63          if (pos < 0) {
64              System.out.println("Invalid char num: " + c);
65          }
66          return this.offspring[pos];
67      }
68
69      /**
70       * Retrieve a list of nodes of which the called-upon node is a
              ↪ parent.
71       *
72       * @return : List of TrieNodes that are offspring of the given
              ↪ node.
73       */
74      public LinkedList<TrieNode> getAllOffspring() {
75          LinkedList<TrieNode> offspring = new LinkedList<>();
76
77          for (TrieNode node : this.offspring) {
78              if (node != null) {
79                  offspring.add(node);
80              }
81          }
82          return offspring;
83      }
84
85      /**
86       * Set a node for the given character, if it doesn't exist.
87       *
88       * @param c : Character to add to the trie.
89       * @return : True if character has been added, false if not.
90       */
91      public boolean setOffspring(char c) {
92          int pos = c - 97; // turn ASCII char to numeric value (where
                  ↪ a = 0)
93          // Create the node if it doesn't exist.
94          if (this.offspring[pos] == null) {
95              this.offspring[pos] = new TrieNode(c);
96              return true;
```

```java
 97            }
 98            return false;
 99        }
100
101    }
102
103    public class Trie {
104        private TrieNode root;
105
106        // Constructor
107        public Trie() {this.root = new TrieNode();}
108
109        /**
110         * Add a given string to the Trie.
111         *
112         * @param key : String to be added to the Trie.
113         * @return : True if string has been added. If the string
                  ↪ already exists,
114         * return false.
115         */
116        public boolean add(String key) {
117            TrieNode curNode = this.root;
118            boolean wasAdded = false;
119
120            for (int i = 0; i < key.length(); i++) {
121                char curChar = key.charAt(i);
122
123                // Check if the current character is in the trie.
124                TrieNode next = curNode.getOffspring(curChar);
125
126                // If not, add it. Switch to the node for this letter.
127                if (next == null) {
128                    curNode.setOffspring(curChar);
129                    next = curNode.getOffspring(curChar);
130                    next.setLevel(i + 1);
131                    wasAdded = true;
132                }
133                curNode = next; // move to this node.
134            }
135
136            // Once at the end of the key String, set the final node to
                  ↪ be a key.
137            if (curNode != null) {
138                curNode.setKey();
139            }
140            return wasAdded;
141        }
142
143        /**
144         * Check if a given string exists in the trie.
145         *
146         * @param key : String to be checked.
```

```java
147          * @return : True if the whole string exists (and is not a
                ↪ prefix or only
148          * part of an existing word), false if not.
149          */
150         public boolean contains(String key) {
151             TrieNode curNode = this.root;
152             boolean isKey = false;
153
154             for (int i = 0; i < key.length(); i++) {
155                 char c = key.charAt(i);
156                 TrieNode next = curNode.getOffspring(c);
157                 if (next == null) {
158                     return false;
159                 }
160                 curNode = next;
161                 isKey   = curNode.getKey();
162             }
163             return isKey;
164         }
165
166         /**
167          * Traverse the trie in breadth-first order and return the
                ↪ result.
168          *
169          * @return : String of characters from reading the trie
                ↪ breadth-first.
170          */
171         public String outputBreadthFirstSearch() {
172             LinkedList<TrieNode> toVisit = new LinkedList<>();
173             int stringSize = 1000; // fixed starting size of 2000 bytes
174             char[] trieString = new char[stringSize];
175             int stringPos = 0;
176             toVisit.add(this.root);
177
178             while (! toVisit.isEmpty()) { // while the visit list is not
                    ↪ empty:
179                 TrieNode curNode = toVisit.remove(); // get the next
                        ↪ node from list
180                 trieString[stringPos] = curNode.getValue(); // add node
                        ↪ to string
181                 toVisit.addAll(curNode.getAllOffspring()); // get all
                        ↪ child nodes
182                 stringPos++;
183
184                 // If the string has reached it's limit, increase that
                        ↪ by 10 times
185                 // and copy the old string over to the new one.
186                 if (stringPos == stringSize) {
187                     char[] temp = trieString;
188                     stringSize = stringSize * 10;
189                     trieString = new char[stringSize];
190                     for (int i = 0; i < stringSize / 10; i++) {
```

```java
191                    trieString[i] = temp[i];
192                }
193            }
194        }
195        return new String(trieString);
196    }

197
198    /**
199     * Traverse the trie in depth-first order and return the result.
200     *
201     * @return String of characters from reading the trie depth-first
202     */
203    public String outputDepthFirstSearch() {
204        Stack<TrieNode> toVisit = new Stack<>();
205        int stringSize = 1000; // starting size for output string
206        char[] trieString = new char[stringSize];
207        int stringPos = 0;
208        toVisit.push(this.root);

209
210        while (! toVisit.isEmpty()) { // while the visit stack is
                ↪ not empty:
211            TrieNode curNode = toVisit.pop(); // take the node off
                    ↪ the stack

212
213            // Retrieve all the offspring from the current node and
                    ↪ add them to
214            // the toVisit stack in reverse order (due to stack
                    ↪ being FILO)
215            LinkedList<TrieNode> node = curNode.getAllOffspring();
216            while (! node.isEmpty()) {
217                toVisit.push(node.removeLast());
218            }
219            trieString[stringPos] = curNode.getValue(); // add node
                    ↪ to String
220            stringPos++;

221
222            // If the string has reached it's limit, increase that
                    ↪ by 10x and
223            // copy the old string over to the new one.
224            if (stringPos == stringSize) {
225                char[] temp = trieString;
226                stringSize = stringSize * 10;
227                trieString = new char[stringSize];
228                for (int i = 0; i < temp.length; i++) {
229                    trieString[i] = temp[i];
230                }
231            }
232        }
233        return new String(trieString);
234    }

235
236    /**
```

```
237         * returns a new trie rooted at the prefix, or null if the
              ↪ prefix is not
238         * present in this trie.
239         *
240         * @param prefix : String which denotes the root of the returned
              ↪ trie.
241         * @return : Trie structure of all branches of the given trie
              ↪ which branch
242         * from the given prefix or empty trie if prefix not found.
243         */
244        public Trie getSubTrie(String prefix) {
245            Trie subTrie = new Trie();
246            TrieNode curNode = this.root;
247
248            if (prefix.equals("")) {
249                return subTrie;
250            }
251
252            // Find the prefix within the main trie while building the
                  ↪ subtrie:
253            for (int i = 0; i < prefix.length(); i++) {
254                char c = prefix.charAt(i);
255
256                TrieNode next = curNode.getOffspring(c);
257                if (next == null) {
258                    return subTrie;
259                }
260                curNode = next;
261            }
262            subTrie.root = curNode;
263
264            return subTrie;
265        }
266
267        /**
268         * Returns all words held in the trie it is called on. This is
              ↪ written to be
269         * as time efficient as possible without regard to memory usage.
270         *
271         * @return List of strings corresponding to each word in the
              ↪ trie. Returns
272         * an empty list if there are no words in the trie it's called
              ↪ on.
273         */
274        List<String> getAllWords() {
275            List<String> words = new ArrayList<>(); // words to be
                      ↪ returned
276            Stack<TrieNode> toVisit = new Stack<>(); // stack of Nodes
                  ↪ to visit next
277            int maxWordLength = 35; // there's no words over 35 letters
                  ↪ in English
278            char[] curWord = new char[maxWordLength]; // hold the word
```

```
                          ↪ being formed
279
280        // Add the first node (root) to the toVisit stack.
281        LinkedList<TrieNode> offspring = this.root.getAllOffspring();
282        while (! offspring.isEmpty()) {
283            toVisit.push(offspring.removeLast());
284        }
285
286        while (! toVisit.isEmpty()) {
287            // Retrieve the next node and find it's position in the
                          ↪ trie
288            TrieNode curNode = toVisit.pop();
289            int charPos = curNode.getLevel() + 1; // position of the
                          ↪ next char
290            curWord[charPos] = curNode.getValue(); // add current
                          ↪ letter
291
292            // Check if the current node is the end of a word and if
                          ↪ so, add
293            // the word to our list of words. We need to ensure we
                          ↪ only copy the
294            // letters that represent this word due to a persistent
                          ↪ version of
295            // curWord being used, however this is still quicker
                          ↪ than making and
296            // storing either a deep or shallow copy of curWord in a
                          ↪ stack.
297            if (curNode.getKey()) {
298                StringBuilder foundWord = new StringBuilder();
299                for (int i = 0; i <= charPos; i++) {
300                    char curChar = curWord[i];
301                    if (curChar != '\u0000') {
302                        foundWord.append(curChar);
303                    }
304                }
305                words.add(foundWord.toString());
306            }
307
308            // Retrieve all the offspring from the current node and
                          ↪ add them to
309            // the toVisit stack in reverse order (due to stack
                          ↪ being FILO).
310            offspring = curNode.getAllOffspring();
311            while (! offspring.isEmpty()) {
312                toVisit.push(offspring.removeLast());
313            }
314        }
315        return words;
316    }
317
318    // Test Harness
319    public static void main(String[] args) {
```

```java
320         // Testing add(), should return "true" for all.
321         Trie t = new Trie();
322         System.out.println(t.add("cheers"));
323         System.out.println(t.add("cheese"));
324         System.out.println(t.add("chat"));
325         System.out.println(t.add("cat"));
326         System.out.println(t.add("can")); // REMOVE
327         System.out.println(t.add("bat") + "\n");
328
329         // Testing contains().
330         System.out.println(t.contains("cheese")); // should return
                ↪ "true"
331         System.out.println(t.contains("chose")); // should return
                ↪ "false"
332         System.out.println(t.contains("ch") + "\n"); // should
                ↪ return "false"
333
334         // Testing outputBreadthFirstSearch(), should return
                ↪ "bcaahttaetersse".
335         System.out.println(t.outputBreadthFirstSearch() + "\n");
336
337         // Testing outputDepthFirstSearch(), should return
                ↪ "batcathateersse".
338         System.out.println(t.outputDepthFirstSearch() + "\n");
339
340         // Testing getSubTrie, should return "aetersse".
341         System.out.println(t.getSubTrie("ch").outputBreadthFirstSearch()
                ↪ +
342              "\n");
343
344         // Testing getAllWords, should return the list of words from
                ↪ add() test.
345         System.out.println(t.getAllWords());
346     }
347 }
```

Listing 3: AutoCompletion.java

```java
/*****************************************************************************

 Project      : CMP-5014Y - Word Auto Completion with Tries :
                AutoCompletion.

 File         : AutoCompletion.java

 Date         : Thursday 14 May 2020

 Author       : Martin Siddons

 Description : Part 3 of the assignment is use Trie.java and
 DictionaryFinder.java to read in a large dictionary and a list of
     ↪ prefixes then
 output the three most common words on each prefix and save that to
     ↪ a file.

 History      : 14/05/2020 - Initial setup

 *****************************************************************************/

import java.util.ArrayList;
import java.util.Arrays;
import java.util.List;

public class AutoCompletion {

    public static void main(String[] args) {
        // 1. form a dictionary file of words and counts from the
            ↪ file lotr.csv.
        DictionaryFinder lotr = new DictionaryFinder();
        lotr.formDictionary(DictionaryFinder.readWordsFromCSV("lotr.csv"));

        // 2. construct a trie from the dictionary using your
            ↪ solution from
        // part 2.
        Trie lotrTrie = new Trie();
        for (String k : lotr.dict.keySet()) {
            lotrTrie.add(k);
        }

        // 3. load the prefixs from lotrQueries.csv
        ArrayList<String> prefixs =
                DictionaryFinder.readWordsFromCSV("lotrQueries.csv");

        // 4. for each prefix query:
        List<String> toFile = new ArrayList<>();
        for (String s : prefixs) {
            // 4.1. Recover all words matching the prefix from the
                ↪ trie.
            Trie subTrie = lotrTrie.getSubTrie(s);
```

```
47              List<String> matches = new
                    ↪ ArrayList<>(subTrie.getAllWords());

48
49              // 4.2. Choose the three most frequent words and display
                    ↪ to
50              // standard output.
51              // 2d array to store word frequency and position in
                    ↪ matches list.
52              // +1 ensures if we add the prefix, it won't run out of
                    ↪ bounds.
53              int[][] count = new int[matches.size() + 1][2];
54              float totalCount = 0;

55
56              if (lotr.dict.containsKey(s)) { // ensure prefix is
                    ↪ checked
57                  matches.add("");
58              }

59
60              for (int i = 0; i < matches.size(); i++) {
61                  String str = s + matches.get(i); // complete the word
62                  int wordCount = lotr.dict.get(str); // get this
                        ↪ word's count
63                  totalCount += wordCount; // inc the word counter
64                  count[i][0] = wordCount; // store the count
65                  count[i][1] = i; // store the word's position
66              }
67              Arrays.sort(count, (a, b) -> Integer.compare(b[0],
                    ↪ a[0]));
68              StringBuilder topThree = new StringBuilder();
69              topThree.append(s);

70
71              // output and save top three results
72              for (int i = 0; i < matches.size() && i < 3; i++) {
73                  float prob = count[i][0] / totalCount;
74                  String word = s + matches.get(count[i][1]);
75                  System.out.println(word + " (probability " + prob +
                        ↪ ")");
76                  topThree.append(",").append(word).append(",").append(prob);
77              }
78              toFile.add(topThree.toString());
79              System.out.println();
80          }
81      // 4.3. Write the results to lotrMatches.csv.
82      DictionaryFinder.saveCollectionToFile(toFile,
            ↪ "lotrMatches.csv");
83      }
84  }
```