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

Student number: 100246776. Blackboard ID: 100246776

Wednesday 20<sup>th</sup> May, 2020 09:55

## Contents

# 1 How to read this report?

1. Short informal descriptions accompany each algorithm. They act as reading guidelines

2. Naming conventions:

   (a) Variable names with single letter,in bold and capitalized indicate abstract data structures. For example, **M** is a map

   (b) Variables names with single letter, in bold and lowercase indicate instances(objects). For example, **n** is a node because it is derived from a user-defined class

   (c) keywords are in bold

3. Comments are either inline. If too long, they are found above a line

4. When calling a function on an object,it is passed explicitly into the function. For example, **M**.add(word,frequency) is add(**M**,word,frequency). Read as, "To **M**, add word and frequency". This separates pseudocode from most programming languages.
   However, to match function signature in specs, a trie implictly passed to every function it is called on. For e.g, outputBreadthFirstSearch() is passed the trie **t** implicitly. The rigorous way is outputBreadthFirstSearch(**t**).

5. Data types used:

   (a) unsigned numeric, boolean, string

   (b) Trie, TrieNode,AutoCompleteTrie, AutoCompleteTrieNode

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

### 2.1  Pseudocode for formDictionary

1. readFile should be a non-empty comma separated values(CSV) file. It contains words separated by commas on a single line. For example, rock,paper,scissors,rock,paper.

2. writeFile is a non-empty CSV file containing word-frequency tuples, separated by line.

3. readWordsFromCSV(readFile) returns a list of words read from readFile.

4. For each word in list, add to map. If encountering word for first time, add to map with frequency=1. If not a new word, retrieve from map and increment frequency.

5. saveToFile(writeFile) saves map to writeFile.

6. Design decisions:

    (a) Map data structures store key-value pairs.

    (b) Implementation of Map(tree map) ensure $\mathcal{O}(log(n))$ runtime complexity for lookup and insertion.

---

**Algorithm 1** formDictionary($readFile$,$writeFile$) **return**

---

**Require:** A CSV file $readFile$, an empty CSV file $writeFile$

**Ensure:** $writeFile$, contains key-frequency tuples where key is a unique word of type string in $readFile$ and frequency of type unsigned numeric, the word's frequency of occurrence in $readFile$

1: $\mathbf{D} \leftarrow Map()$           ▷ *Initialise map $\boldsymbol{D}$ to hold key-frequency tuples*

2: $\boldsymbol{L} \leftarrow readWordsFromCSV(readFile)$        ▷ $\mathbf{L}$ *is a list of words from readFile*

3: **for all** s **in L do**

4:      **if** $containsKey(\mathbf{D}, s)$ **is false then**

5:          $put(\mathbf{D}, s, 1)$          ▷ *if key not in $\boldsymbol{D}$,add it and set frequency to 1*

6:      **else**          ▷ if key in $\mathbf{D}$,retrieve and increment frequency

7:          $frequency \leftarrow get(\mathbf{D}, s) + 1$

8:          $put(\mathbf{D}, s, frequency)$

9: $saveToFile(\mathbf{D}, writeFile)$        ▷ *saveToFile writes $\boldsymbol{D}$ to a writeFile*

---

### 2.2  Algorithm Analysis for formDictionary

In a discussion board thread, Dr. Tony Bagnall stated we need to consider runtime complexity of treemap functions. Below is the runtime complexity analysis for formDictionary with such a consideration.

1. Fundamental operation:
    **if** $containsKey(\mathbf{D}, s)$ **is false then**

2. Treemap maintains $\mathcal{O}(log(n))$ lookup. Worst, average and best cases are the same. Let size of $\mathbf{L}$ to be n

3.
$$f(n) = \sum_{i=1}^{n} log(n) \tag{1}$$

$$t(n) = nlog(n) \tag{2}$$

4. The order of the runtime complexity is log linear i.e. nlog(n). The runtime complexity function is $\mathcal{O}(nlog(n))$

**Or,consider number of times fundamental operation occurs.**

1. Fundamental operation:
    **if** $containsKey(\mathbf{D}, s)$ **is false then**

2. Let size of **L** to be n

3.

$$f(n) = \sum_{i=1}^{n} 1 \tag{3}$$

$$t(n) = n \tag{4}$$

4. The order of the runtime complexity is linear i.e. n. The runtime complexity function is $\mathcal{O}(n)$

## 2.3  Pseudocode for saveToFile

1. writeFile is a non-empty CSV file containing keys-frequency tuples, separated by line.

2. D is a non-empty map from saveDictionary. It stores word-frequency pairs.

3. For all entries in D, make a keyValuePair string. For example, "brother,4".

4. On a new line in writeFile, write keyValuePair

---

**Algorithm 2** saveToFile(**D**,*writeFile*) **return**

---

**Require:** A map **D** containing word-frequency tuple where word is of type string and frequency of type unsigned numeric, an empty CSV file *writeFile*

**Ensure:** *writeFile* contains all word-frequency tuples where word is a key and frequency, the word's frequency of occurrence in *readFile*

1: **for all e in D do** ▷ *for each entry*

   ▷ for an entry, retrieve key and value i.e frequency. Format into a word-frequency tuple: "key,value"

2:     $keyValuePair \leftarrow getKey(\mathbf{e}), getValue(\mathbf{e})$

3:     $writeline(writeFile, keyValuePair)$ ▷ write keyValuePair on a newline

---

## 2.4  Algorithm Analysis for saveToFile

In a discussion board thread, Dr. Tony Bagnall stated we need to consider runtime complexity of treemap functions. Below is the runtime complexity analysis for saveToFile with such a consideration. It is also a more realistic depiction of the runtime because the algorithm's efficiency(or lack of it) depends on the data structure used.

1. Fundamental operation:
   $keyValuePair \leftarrow getKey(\mathbf{e}), getValue(\mathbf{e})$

2. Treemap maintains $\mathcal{O}(log(n))$ lookup. Worst, average and best cases are the same. Allow loop to execute n time, ie, number of entries in **D**

3.

$$f(n) = \sum_{i=1}^{n} log(n) + \sum_{i=1}^{n} log(n) \tag{5}$$

$$t(n) = 2nlog(n) \tag{6}$$

4. Ignoring constants, the order of the runtime complexity is log linear i.e. nlog(n). The runtime complexity function is $\mathcal{O}(nlog(n))$.

**Or,consider number of times fundamental operation occurs.**

1. Fundamental operation:
   $keyValuePair \leftarrow getKey(\mathbf{e}), getValue(\mathbf{e})$

2. Allow loop to execute n time, ie, number of entries in **D**

3.

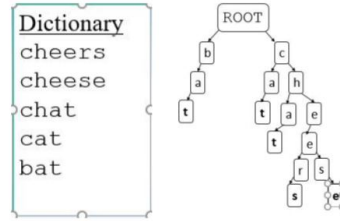$$f(n) = \sum_{i=1}^{n} 1 \tag{7}$$

$$t(n) = n \tag{8}$$

4. The order of the runtime complexity is linear i.e. n. The runtime complexity function is $\mathcal{O}(n)$

# 3   Part 2: Implement a Trie Data Structure

**TrieNode**: A TrieNode has an array of size 26. The array stores references to offspring of that node. The indexes of the array each maps to a letter. For example, 'a' maps to 0(1 in pseudocode), 'b' to 2 and so forth.A TrieNode can be a key: it represents a complete word. As such, a TrieNode has getIsKey() and setIsKey(boolean) functionalities. getIsKey() returns true if TrieNode is a key. setIsKey() sets key attribute of TrieNode to true.

**Trie**: A Trie can be viewed as containing TrieNodes with a root TrieNode at the top. The access point to the rest of the TrieNodes is through the root node.

Figure 1: Trie example



Using bat as an example, the root stores a reference to a node in it's array. The index at which node is stored 2 mapping to 'b'. A key is a node whose isKey attribute is set to true. In the diagram, a key is denoted in bold.

## 3.1   Algorithm for boolean add(String key)

A key is a unique, non-empty, not null word in the trie. If a valid word is being added and it is already a key in the trie, the algorithm returns false for not added. Otherwise, return true for added.

---

**Algorithm 3** add($word$) **return** ($added$)

---

**Require:** $word$ is of type string

**Ensure:** $added$ is a boolean, either **true** or **false**. Indicates if a word is added as key

1:   $added \leftarrow$ **false**

2:   **if r is null or** $word$ **is null then**          ▷ *r is the the root node of the trie*

3:      **return** ($added$)

4:   $asciiOffset \leftarrow 97$          ▷ *97 is 'a' in ascii. Start of small letters*

5:   $\mathbf{p} \leftarrow \mathbf{r}$          ▷ *assign r,root of trie ,to p*

6:   **for** $i \leftarrow 1$ **to** $size(word)$ **do**          ▷ for each character in word

7:      $character \leftarrow convertToChar(word_i)$    ▷ *converts word$_i$ into ascii equivalent*    ▷ *see appendix for ascii table*
     ▷ converts index to a value between 1-26, mapping it between 'a' to 'z'

8:      $index \leftarrow charToInt(character) - asciiOffset$
     ▷ *returns an array containing pointers to children of* $\mathbf{p}$

9:      $\mathbf{C} \leftarrow getOffspring(\mathbf{p})$

10:      $\mathbf{n} \leftarrow C_{index}$          ▷ *get offspring at mapped index*

11:      **if n is null then**          ▷ *If node is null, word is not in trie*

12:         $added \leftarrow$ **true**

13:         $\mathbf{t} \leftarrow TrieNode()$          ▷ *create a new trieNode*
     ▷ assign new node to this position. Index of offspring represents a character

14:         $C_{index} \leftarrow \mathbf{t}$

15:         $\mathbf{p} \leftarrow \mathbf{t}$          ▷ set parent to the new node created

16:      **else**

17:         $\mathbf{p} \leftarrow \mathbf{n}$          ▷ If node exists at index, go to it

---

18: **if** $getIsKey(\mathbf{p})$ **and** $added$ **is false then**         ▷ *If node is key and not already added*

19:      **return** ($added$)         ▷ if node already in trie, return false for not added

20: $setIsKey(\mathbf{p}, \mathbf{true})$

21: $added \leftarrow \mathbf{true}$

22: **return** ($added$)

## 3.2 Algorithm for boolean contains(String word)

A key is a unique, non-empty, not null word in the trie. The algorithm checks if the word passed as a parameter is a key in the trie. If yes, return true for contains. Return false otherwise.

---

**Algorithm 4** contains($word$) **return** ($contains$)

---

**Require:** $word$ is of type string

**Ensure:** $contains$ is a boolean, either true or false. Indicate if a word is a key

1: $contains \leftarrow \mathbf{false}$

2: **if r is null or** $word$ **is null then**         ▷ *$\mathbf{r}$ is the the root node of the trie*

3:      **return** (contains)

4: $asciiOffset \leftarrow 97$         ▷ *97 is 'a' in ascii. Start of small letters*

5: $\mathbf{p} \leftarrow \mathbf{r}$         ▷ *assign $\mathbf{r}$,root of trie ,to $\mathbf{p}$*

6: **for** $i \leftarrow 1$ **to** $size(word)$ **do**         ▷ looping character by character

7:      $character \leftarrow convertToChar(key_i)$     ▷ *converts $key_i$ into ascii equivalent*     ▷ *see appendix for ascii table*
    ▷ converts index to a value between 1-26, mapping it between 'a' to 'z'

8:      $index \leftarrow charToInt(character) - asciiOffset$
    ▷ *returns an array containing pointers to children of parent node*

9:      $\mathbf{C} \leftarrow getOffspring(\mathbf{p})$

10:      $\mathbf{n} \leftarrow C_{index}$         ▷ get offspring at mapped index

11:      **if n is null then**         ▷ *If node is null, key is not in trie*

12:         **return** (contains)

13:      $p \leftarrow n$         ▷ if a reference corresponding to the index is found, go to it.
    ▷ *If node is key,trie contains it*

14: **if** $getIsKey(\mathbf{p})$ **is true then**         ▷ getIsKey(node) returns true if node is a key

15:      $contains \leftarrow true$

16:      **return** ($contains$)
    ▷ A matching word can be found in the trie but it does not mean it is a key

17: **return** ($contains$)         ▷ trie does not contain word

---

## 3.3 Algorithm for boolean String outputBreadthFirstSearch()

This algorithm goes through a trie and outputs characters in breadth first search order.
Design decision:

1. Use a Queue because it respects first-in first-out rule. This is useful for breadth first search.

2. Underlying implementation of queue is a linkedlist. This ensures ensure O(1) for enqueue and dequeue.

---

**Algorithm 5** outputBreadthFirstSearch() **return** ($result$)

**Require:**

**Ensure:** $result$ is of type string and contains all words in breadth first search order

1: **if r is null then** ▷ $r$ is the the root node of the trie data structure
2:      **return null**
3: $asciiOffset \leftarrow 97$ ▷ 97 is 'a' in ascii. Start of small letters
4: $\mathbf{Q} \leftarrow Queue()$ ▷ holds a collection of nodes type TrieNode
5: $enqueue(\mathbf{Q}, \mathbf{r})$ ▷ Add root,$\mathbf{r}$, to front of queue,$\mathbf{Q}$
6: **while Q is not empty do** ▷ dequeue($\mathbf{Q}$) removes top most trie node and assign to temp,$\mathbf{t}$
7:      $\mathbf{t} \leftarrow dequeue(\mathbf{Q})$
     ▷ getOffspring(node) returns an array containing references to children of that node
8:      $\mathbf{C} \leftarrow getOffspring(\mathbf{t})$
9:      **for** $i \leftarrow 1$ **to** $size(\mathbf{C})$ **do**
10:         **if** $C_i$ **is not null then**
11:            $enqueue(\mathbf{q}, C_i)$ ▷ add offspring of node to the queue so it goes through them in breadth first order
12:            $result \leftarrow result + intToChar(i + asciiOffset)$ ▷ concatenate character represented by i to result
     **return** ($result$)

---

## 3.4 Algorithm for String outputDepthFirstSearch()

This algorithm is an interface method to outputDepthFirstSearch(builder,$\mathbf{r}$). It provides an empty string, builder, and the starting node of the depth first search.

---

**Algorithm 6** outputDepthSearchSearch() **return** ($result$)

**Require:**

**Ensure:** $result$, a string in depth first search order.
     ▷ $builder$ is an empty string, and $\mathbf{r}$ is the root of the Trie.
1: $builder \leftarrow ""$
2: $result \leftarrow depthFirstSearch(builder, \mathbf{r})$
3: **return** ($result$)

---

### 3.4.1 Recursive implementation of String outputDepthFirstSearch()

This algorithm goes through trie and ouputs characters in depth first search(DFS) order. DFS goes down the left-most leaf and backs up.
Design decision:

1. Recursion as opposed to a iterative version provides a more elegant and shorter solution.

2. No need to handle stack data structure.

---
**Algorithm 7** outputDepthFirstSearch(builder,$\mathbf{r}$) **return** ($result$)
---
**Require:** builder, an empty string and $\mathbf{r}$ of type TrieNode

**Ensure:** $result$ is of type string and contains all words in depth first search order

  1: **if r is null then**             $\triangleright$ *r is the the root node of the trie data structure*

  2:     **return null**

    $\triangleright$ *getOffspring(node) returns an array containing references to children of that node*

  3: $\mathbf{C} \leftarrow getOffspring(\mathbf{r})$

  4: $asciiOffset \leftarrow 97$             $\triangleright$ *97 is 'a' in ascii. Start of small letters*

  5: **for** $i \leftarrow 1$ **to** $size(\mathbf{C})$ **do**

  6:     **if** $C_i$ **is not null then**

  7:         $result \leftarrow result + intToChar(i + asciiOffset)$           $\triangleright$ *convert index to char*

  8:         $outputDepthFirstSearch(result, C_i)$        $\triangleright$ *recursive call on node's children- depth first*
    **return** ($result$)

---

## 3.5   Algorithm for getSubTrie(String prefix)

The algorithm returns a new trie rooted at the prefix, or null if the prefix is not present in this trie.

---
**Algorithm 8** getSubTrie($prefix$) **return** ($\mathbf{t}$)
---
**Require:** $\mathbf{r}$,root node of trie and $prefix$, a string representing incomplete word

**Ensure:** $\mathbf{t}$ is a trie data structure rooted at $prefix$

  1: **if r is null or** $prefix$ **is null then**

  2:     $\mathbf{r} \leftarrow \mathbf{null}$          $\triangleright$ *root is the the root node of the trie data structure*

  3: $asciiOffset \leftarrow 97$           $\triangleright$ *97 is 'a' in ascii. Start of small letters*

  4: $\mathbf{p} \leftarrow \mathbf{r}$

  5: **for** $i \leftarrow 1$ **to** $size(prefix)$ **do**

    $\triangleright$ *converts $prefix_i$ into ascii equivalent*

  6:     $character \leftarrow convertToChar(prefix_i)$

    $\triangleright$ converts index to a value between 1-26, mapping it between 'a' to 'z'

  7:     $index \leftarrow charToInt(character) - asciiOffset$        $\triangleright$ *see appendix for ascii table*

    $\triangleright$ *getOffspring(node) returns an array containing references to children of that node*

  8:     $\mathbf{C} \leftarrow getOffspring(\mathbf{p})$

  9:     **if** $C_{index}$ **is null then**           $\triangleright$ *prefix does not exist*

 10:         **return** (**null**)

 11:     $\mathbf{p} \leftarrow C_{index}$          $\triangleright$ *if node is not null, go to it*

 12: **return** $Trie(\mathbf{p})$        $\triangleright$ *return a trie rooted at the last node assigned to* $\mathbf{p}$

---

## 3.6 Algorithm for getAllWords()

This algorithm is an interface method to getAllWords(**L**,textBuilder,**r**). It provides an empty string, builder, the starting node and an empty list.

---

**Algorithm 9** getAllWords() **return** (**L**)

---

**Require:**

**Ensure: L**, a list of keys from trie.

$\quad getAllWords(\mathbf{L}, textBuilder, \mathbf{r})$             ▷ *textBuilder* is an empty string, and **r** is the root of the Trie.

$\quad$ **if** $size(\mathbf{L})$ **is null then**                           ▷ prefix does not exist

$\quad\quad$ **return** (**null**)

$\quad$ **return** (**L**)

---

This algorithm returns a list containing all words in the trie. It is called by getAllWorld().
Design decision:

1. Underlying implementation used for list, **L** is a linked list. Linked list provides O(1) runtime complexity for insertion. Linked List is always added to in the algorithm.

---

**Algorithm 10** getAllWords(**L**,**textBuilder**,**r**) **return**

---

**Require: L** is a list of words that are also keys of type string, **textBuilder**,string that can be appended to, and **r**, the root node of the trie

**Ensure:**

1: $asciiOffset \leftarrow 97$                              ▷ *97 is 'a' in ascii. Start of small letters*

2: **if** $getIsKey(\mathbf{r})$ **then**

$\quad$ ▷ *If node is a key, textBuilder contains a complete word that is also a key*

3: $\quad\quad add(\mathbf{L}, textBuilder)$

$\quad$ ▷ *getOffspring(node) returns an array of references to children of that node*

4: $\mathbf{C} \leftarrow getOffspring(\mathbf{r})$

5: **for** $i \leftarrow 1$ **to** $size(\mathbf{C})$ **do**

6: $\quad\quad \mathbf{o} \leftarrow C_i$

7: $\quad\quad$ **if o is not null then**

8: $\quad\quad\quad textBuilder \leftarrow textBuilder + intToChar(i + asciiOffset)$    ▷ *intToChar converts the int to it's char equivalent*

9: $\quad\quad\quad getAllWords(\mathbf{L}, textBuilder, \mathbf{o})$                  ▷ recursive call with updated parameters

$\quad$ ▷ *Backtracks one level up the trie. Say trie contains keys ab and ac. Once ab is appended to **L**, remove last character 'b' from 'ab' in order to append 'c'*

10: $\quad\quad\quad removeLastCharacter(textBuilder)$

# 4 Part 3: Word Auto Completion Application

**AutoCompletionTrieNode**: An AutoCompletionTrieNode is conceptually similar to a TrieNode. It is made of an array of size 26 to hold children nodes and contains a wordFrequency attribute. wordFrequency is updated whenever the node is set to a key. As such, it also contains a setter and getter for wordFrequency.

**AutoCompletionTrie**: This is same as a regular Trie.

**WordAndFrequency**: WordAndFrequency is a blueprint for an entity containing a key and frequency.

## 4.1 Algorithm for add(word,wordFrequency)

This algorithm adds a word, frequency pair to the AutoCompletionTrie. If word is not already a key in the AutoCompletionTrie, mark the node that holds last character to word as a key. Set its frequency. If a node is already a key, update its frequency.

---

**Algorithm 11** add($word$,$wordFrequency$) **return** ($added$)

---

**Require:** $word$ is a unique word and of type string,$wordFrequency$ is the frequency of occurrence of a word

**Ensure:** $added$ is a boolean, either **true** or **false**. Indicate if a word is added as key

1: $added \leftarrow$ **false**
2: **if r is null or** $word$ **is null then**                     ▷ *r is the the root node of the trie*
3:     **return** ($added$)
4: $asciiOffset \leftarrow 97$                              ▷ *97 is 'a' in ascii. Start of small letters*
5: $\mathbf{p} \leftarrow \mathbf{r}$                                  ▷ *assign $\mathbf{r}$,root of AutoCompletionTrie ,to $\mathbf{p}$*
6: **for** $i \leftarrow 1$ **to** $size(word)$ **do**                           ▷ for each character in key
7:     $character \leftarrow convertToChar(word_i)$      ▷ *converts $word_i$ into ascii equivalent*      ▷ *see appendix for ascii table*
   ▷ converts index to a value between 1-26, mapping it between 'a' to 'z'
8:     $index \leftarrow charToInt(character) - asciiOffset$
   ▷ *returns an array containing pointers to children of $\mathbf{p}$*
9:     $\mathbf{C} \leftarrow getOffspring(\mathbf{p})$
10:    $\mathbf{n} \leftarrow C_{index}$                                      ▷ get offspring at mapped index
11:    **if n is null then**                              ▷ *If node is null, word is definitely not in trie*
12:        $added \leftarrow$ **true**
13:        $t \leftarrow AutoCompletetionTrieNode()$                     ▷ create a new AutoCompletionTrieNode
   ▷ assign new node to this position. Index of offspring represents a character
14:        $C_{index} \leftarrow t$
15:        $p \leftarrow t$                                    ▷ set parent to the new node created
16:    **else**
17:        $p \leftarrow n$                                   ▷ If node exists at index, go to it
   ▷ If node is key and not already added
18: **if** $getIsKey(\mathbf{p})$ **and** $added$ **is false then**
   ▷ if word is a key, update its frequency of occurrence
19:    $newFrequency \leftarrow wordFrequency + getWordFrequency(\mathbf{p})$   ▷ if word is a key, update its frequency of occurrence
20:    $setWordFrequency(p, newFrequency)$
21:    **return** ($added$)                                 ▷ if node already in trie, return false for not added
22: $setIsKey(\mathbf{p}, \mathbf{true})$
23: $setWordFrequency(\mathbf{p}, wordFrequency)$;                      ▷ *set frequency of word when it's a key*
24: $added \leftarrow$ **true**
25: **return** ($added$)

---

## 4.2 Algorithm for addDictionaryToTrie(df)

This algorithm takes a valid CSV file as parameter. Each line in a valid file has the following format: word,frequency. The algorithm extracts word-frequency pairs from each line and calls add(word,frequency) on AutoCompletionTrie **t** to add them.

---

**Algorithm 12** addDictionaryToTrie(*df*) **return**

---

**Require:** A dictionary CSV file *df*,where each line represents a unique word and frequency of occurrence

**Ensure:**

1: **while** *df* **not EOF do**                                                        ▷ *while not end of file*
2:     *line* ← *readLine*(*df*)                                                        ▷ read line by line
3:     *word* ← *getWord*(*line*)                                    ▷ Implementation is different. Get word from line
4:     *wordFrequency* ← *getWordFrequency*(*line*)        ▷ Implementation is different. Get wordFrequency from line
        ▷ **t** is the AutoCompletionTrie being operated on. Refer to line 4 in section 1: How to read this report?
5:     *add*(**t**, *word*, *wordFrequency*)

---

**Note: Pseudocode and implementation for contains through to getSubTrie are the same as in Trie**

## 4.3 Algorithm for boolean contains(String word)

---

**Algorithm 13** contains(*key*) **return** (*contains*)

---

**Require:** *key* is a word and of type string

**Ensure:** *contains* is a boolean, either true or false. Indicate if a word is a key

1: *contains* ← **false**
2: **if r is null or** *key* **is null then**                              ▷ **r** *is the the root node of the AutoCompletionTrie*
3:     **return** (contains)
4: *asciiOffset* ← 97                                                ▷ *97 is 'a' in ascii. Start of small letters*
5: **p** ← **r**                                                              ▷ *assign* **r**,*root of trie ,to* **p**
6: **for** *i* ← 1 **to** *size*(*key*) **do**                                ▷ looping character by character
7:     *character* ← *convertToChar*(*key_i*)       ▷ *converts key_i into ascii equivalent*       ▷ *see appendix for ascii table*
        ▷ converts index to a value between 1-26, mapping it between 'a' to 'z'
8:     *index* ← *charToInt*(*character*) − *asciiOffset*
        ▷ *returns an array containing pointers to children of parent node*
9:     **C** ← *getOffspring*(**p**)
10:    **n** ← **C**_*index*                                                   ▷ get offspring at mapped index
11:    **if n is null then**                                                ▷ *If node is null, key is not in trie*
12:        *contains* ← ***false***
13:        **return** (*contains*)
14:    ***p*** ← *n*                                                ▷ if a reference corresponding to the index is found, go to it.
        ▷ *If node is key,trie contains it*
15: **if** *getIsKey*(**p**) **is true then**
16:     *contains* ← *true*
17:     **return** (*contains*)
18: **return** (*contains*)                                                ▷ trie does not contain word

---

## 4.4 Algorithm for boolean String outputBreadthFirstSearch()

---

**Algorithm 14** outputBreadthFirstSearch() **return** ($result$)

---

**Require:**

**Ensure:** $result$ is of type string and contains all words in breadth first search order

1: **if r is null then**        ▷ *r is the the root node of the trie data structure*

2:      **return null**

3: $asciiOffset \leftarrow 97$        ▷ *97 is 'a' in ascii. Start of small letters*

4: $\mathbf{Q} \leftarrow Queue()$        ▷ holds a collection of nodes type TrieNode

5: $enqueue(\mathbf{Q}, \mathbf{r})$        ▷ Add root,$\mathbf{r}$, to front of queue,$\mathbf{Q}$

6: **while Q is not empty do**

    ▷ *dequeue(**Q**) removes top most trie node and assign to temp,**t***

7:      $\mathbf{t} \leftarrow dequeue(\mathbf{Q})$

    ▷ *getOffspring(node) returns an array containing references to children of that node*

8:      $\mathbf{C} \leftarrow getOffspring(\mathbf{t})$

9:      **for** $i \leftarrow 1$ **to** $size(\mathbf{C})$ **do**

10:        **if** $C_i$ **is not null then**

11:          $add(\mathbf{q}, C_i)$      ▷ add offspring of node to the queue so it goes thorough them in breadth first order

    ▷ *concatenate character represented by i to result*

12:          $result \leftarrow result + intToChar(i + asciiOffset)$

    **return** ($result$)

---

## 4.5 Algorithm for String outputDepthFirstSearch()

---

**Algorithm 15** outputDepthSearchSearch() **return** ($result$)

---

**Require:**

**Ensure:** $result$, a string in depth first search order.      ▷ *builder* is an empty string, and $\mathbf{r}$ is the root of the Trie.

1: $result \leftarrow depthFirstSearch(builder, \mathbf{r})$

2: **return** ($result$)

---

### 4.5.1 Recursive implementation of String outputDepthFirstSearch()

---

**Algorithm 16** outputDepthFirstSearch(builder,$\mathbf{r}$) **return** ($result$)

---

**Require:** builder, an empty string and $\mathbf{r}$ of type TrieNode

**Ensure:** $result$ is of type string and contains all words in depth first search order

1: **if r is null then**        ▷ *r is the the root node of the trie data structure*

2:      **return null**

    ▷ *getOffspring(node) returns an array containing references to children of that node*

3: $\mathbf{C} \leftarrow getOffspring(\mathbf{r})$

4: $asciiOffset \leftarrow 97$        ▷ *97 is 'a' in ascii. Start of small letters*

5: **for** $i \leftarrow 1$ **to** $size(\mathbf{C})$ **do**

6:      **if** $C_i$ **is not null then**

7:        $result \leftarrow result + intToChar(i + asciiOffset)$        ▷ *convert index to char(0 maps to 'a')*

8:        $outputDepthFirstSearch(result, C_i)$        ▷ *recursive call on node's children- depth first*

    **return** ($result$)

---

## 4.6 Algorithm for getSubTrie(String prefix)

---

**Algorithm 17** getSubTrie($prefix$) **return** (**t**)

---

**Require:** **r**, root node of AutoCompletionTrie and $prefix$, a string representing part of a word

**Ensure:** **t** is a AutoCompletionTrie data structure rooted at $prefix$

1: **if r is null or** $prefix$ **is null then**

2:      **r** ← **null**                 ▷ *root is the the root node of the trie data structure*

3: $asciiOffset \leftarrow 97$               ▷ *97 is 'a' in ascii. Start of small letters*

4: **p** ← **r**

5: **for** $i \leftarrow 1$ **to** $size(prefix)$ **do**

     ▷ *converts $prefix_i$ into ascii equivalent*

6:      $character \leftarrow convertToChar(prefix_i)$

     ▷ converts index to a value between 1-26, mapping it between 'a' to 'z'

7:      $index \leftarrow charToInt(character) - asciiOffset$        ▷ see appendix for ascii table

     ▷ *getOffspring(node) returns an array containing references to children of that node*

8:      **C** ← $getOffspring($**p**$)$

9:      **if** $C_i$ **is null then**                    ▷ prefix does not exist

10:        **return** (**null**)

11:      **p** ← $C_{index}$

12: **return** ($AutoCompletionTrie($**p**$)$)      ▷ return an AutoCompletionTrie rooted at the last node assigned to **p**

---

## 4.7 Algorithm for wordFrequencyInfo()

wordFrequencyInfo() is an interface function. It passes an empty queue, empty string and starting node(usually root) to wordFrequencyInfo(**Q**, $textBuilder$, **r**). The queue returned should contain WordAndFrequency entities(For definition,see part 3 introduction).

---

**Algorithm 18** wordFrequencyInfo() **return** (**Q**)

---

**Require:**

**Ensure: Q**, a queue that guarantees head will contain word with highest frequency or if all elements have same frequency, shortest key length first.

     ▷ **Q** *guarantees sorted order for head of queue*

1: **Q** ← $Queue()$              ▷ *textBuilder is an empty string, and **r** is the root of the AutoCompletetionTrie*

2: $wordFrequencyInfo($**Q**$, textBuilder, $**r**$)$

3: **return** (**Q**)

---

## 4.8 Algorithm for wordFrequencyInfo(Q, $textBuilder$, r)

wordFrequencyInfo(**Q**, $textBuilder$, **r**) is called by wordFrequencyInfo(). The algorithm goes through the AutoCompletion-Trie recursively. When it encounters a node that is a key, it means that the textBuilder contains a complete word. A WordAndFrequency entity is then made and added to the queue.

     The implementation of this pseudocode uses a priority queue.

1. Ability to specify priority of an element. As such, no need to sort

2. Adding to priority queue is log(n) which is fairly fast

3. Polling/removing from priority queue is O(1)

---

**Algorithm 19** wordFrequencyInfo($\mathbf{Q}$, *textBuilder*, $\mathbf{r}$) **return**

---

**Require:** $\mathbf{Q}$, an empty queue,*textBuilder*, a string to append to, and $\mathbf{r}$, the root of the AutoCompletionTrie
**Ensure:** $\mathbf{Q}$ contains all WordFrequency entities from AutoCompletionTrie. Head guaranteed to be in order
  1: **if** $getIsKey(\mathbf{r})$ **then**                           ▷ Make a WordAndFrequency entity with word and frequency of occurrence
  2:    $\mathbf{t} \leftarrow WordAndFrequency(textBuilder, getFrequency(\mathbf{r}))$
  3:    $add(\mathbf{Q}, t)$
                                  ▷ *getOffspring(node) returns an array containing references to children of that node*
  4: $\mathbf{C} \leftarrow getOffspring(\mathbf{r})$
  5: **for** $i \leftarrow 1$ to $size(C)$ **do**
  6:    **if** $C_i$ *is* **not null then**
  7:        $textBuilder \leftarrow textBuilder + intToChar(97 + i)$        ▷ *intToChar converts an int to it's char equivalent*
  8:        $wordFrequencyInfo(\mathbf{Q}, textBuilder, C_i)$
  9:        $removeLastCharacter(textBuilder)$               ▷ *backtracks one level up the AutoCompletionTrie*

---

## 4.9 Algorithm for prefixMatches(P)

This function finds the most probable matches for each prefix in the list passed to it and writes them to a CSV file.

---

**Algorithm 20** prefixMatches($\mathbf{P}$) **return**

---

**Require:** $\mathbf{P}$, a list of prefixes
**Ensure:** Adding top three most probable words for each prefix in $\mathbf{P}$. Prefix, words and probability should be written to a CSV file.
  $\mathbf{T} \leftarrow List()$
  **for all** prefix in $\mathbf{P}$ **do**
    $add(\mathbf{T}, prefix)$                     ▷ $\mathbf{T}$ is will be saved to a CSV lotrMatches. To respect format of file, add prefix
    $\mathbf{s} \leftarrow getSubTrie(prefix)$
    $\mathbf{Q} \leftarrow wordFrequencyInfo()$                 ▷ $\mathbf{Q}$ contains WordAndFrequency entities for that prefix
    $totalFreq \leftarrow 0$
    **for all** $w$  in  $\mathbf{Q}$ **do**                  ▷ For each WordAndFrequency entity in $\mathbf{Q}$
        $totalFreq = totalFreq + getFreq(\mathbf{w})$        ▷ Add frequency of WordAndFrequency entity to totalFreq
    **if** $size(\mathbf{Q}) < 3$ **then**                 ▷ if queue has less than 3 elements, set limit to size of $\mathbf{Q}$
        $limit \leftarrow size(\mathbf{Q})$
    **else**
        $limit \leftarrow 3$                        ▷ else we want top 3 most probable words
    **for** $i \leftarrow 1$ to $limit$ **do**
        $\mathbf{w} \leftarrow dequeue(\mathbf{Q})$                  ▷ remove top element
        **if** $\mathbf{w}$ *is not null* **then**
            $prob \leftarrow getFreq(\mathbf{w})/totalFreq$          ▷ find probability of given word
            $completeKey \leftarrow prefix + \mathbf{w}$    ▷ Need to add prefix to make a complete word. This is due to getSubTrie()
            $\mathbf{T} \leftarrow add(\mathbf{T}, completeKey)$
            $\mathbf{T} \leftarrow add(\mathbf{T}, prob)$
  $saveToFile("lotrMatches.csv", \mathbf{T})$
  clear($\mathbf{T}$)                            ▷ clear contents of $\mathbf{T}$ to add information for next prefix

---

Figure 2: lotrMatches format

```
ab,about,0.566667,above,0.3,able,0.1
go,going,0.277778,go,0.240741,good,0.166667
the,the,0.626703,they,0.153951,them,0.06812
mer,merry,0.947368,merely,0.026316,merrily,0.026316
fro,frodo,0.490909,from,0.436364,front,0.072727
gr,great,0.19697,ground,0.181818,grass,0.151515
gol,goldberry,0.6,golden,0.4
sam,sam,1
```

**Note: pseudocode for saveToFile("lotrMatches.csv",T) is not very relevant**

## 5 Code Listing

Listing 1: DictionaryMaker.java

```
1  /*****************************************************************************
2
3   File        : DictionaryMaker.java
4
5   Date        : Wednesday 18th March 2020
6
7   Author      : 100246776
8
9   Description : Provide implementation for formDictionary and saveToFile in part 1
10
11  Last update : 18th May 2020
12
13  *****************************************************************************/
14
15  package dsacoursework2;
16  import java.io.*;
17  import java.util.*;
18  public class DictionaryMaker {
19      /**
20       * Reads all the words in a comma separated text document into an Array
21       *
22       * @param
23       */
24      public static ArrayList<String> readWordsFromCSV(String file, String delimiter)
              ↪ throws FileNotFoundException {
25          Scanner sc = new Scanner(new File(file));
26          sc.useDelimiter(delimiter);
27          ArrayList<String> words = new ArrayList<>();
28          String str;
29          while (sc.hasNext()) {
30              str = sc.next();
31              str = str.trim();
32              str = str.toLowerCase();
33              words.add(str);
34          }
35          return words;
36      }
37
38      /**
39       * Forms a dictionary of word,frequency tuples. A word should be unique
40       * @param readFile, the file to be read
41       * @param writeFile, the file to be written to
42       * @throws IOException
43       */
44      public void formDictionary(String readFile, String writeFile) throws IOException {
45
46          ArrayList<String> listOfWords = readWordsFromCSV(readFile, ",");
47          TreeMap<String, Integer> frequencyDictionary = new TreeMap<>();//treeMap
                  ↪ maintains sorted order
48          for (String s : listOfWords) {
49              if (!frequencyDictionary.containsKey(s)) {//if treemap does not contain
                      ↪ the key
50                  frequencyDictionary.put(s, 1);//if key not already in map, set
                          ↪ frequency to 1
51              } else {
52                  int frequency=frequencyDictionary.get(s) + 1;
```

16

```
53              frequencyDictionary.put(s,frequency);//if key already in map,
                    ↪ increase its frequency by 1
54            }
55        }
56        System.out.println(frequencyDictionary);//testing purposes
57        saveToFile(frequencyDictionary, writeFile);
58        System.out.println("file saved");
59    }
60
61    /**
62     * Writes a map to a CSV file
63     * @param frequencyDictionary, a map of word,frequency pair from formDictionary
64     * @param writeFile, an empty CSV file to write to
65     * @throws IOException
66     */
67    public static void saveToFile(TreeMap<String, Integer> frequencyDictionary,
          ↪ String writeFile) throws IOException {
68        FileWriter fileWriter = new FileWriter(writeFile);
69        PrintWriter printWriter = new PrintWriter(fileWriter);
70        for (Map.Entry<String, Integer> entry : frequencyDictionary.entrySet())//loop
              ↪ through TreeMap
71        {
72            printWriter.println(entry.getKey() + "," + entry.getValue());
73        }
74        printWriter.close();
75    }
76
77    public static void main(String[] args) throws Exception {
78        DictionaryMaker df = new DictionaryMaker();
79        /*
80        Testing done by comparing test.csv to testDictionary, provided in zip file
81         */
82        df.formDictionary("testDocument.csv", "test.csv");
83    }
84 }
```

Listing 2: TrieNode.java

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

  File        : TrieNode.java

  Date        : Wednesday 18th March 2020

  Author      : 100246776

  Description : Implementation of TrieNode

  Last update : 18th May 2020

  ******************************************************************************/
package dsacoursework2;

public class TrieNode {
    private boolean isKey;//isKey=true if node completes a word
    private final TrieNode[] offspring;

    public TrieNode() {
        this.offspring = new TrieNode[26];
    }//An array to hold the references to offspring

    /**
     *
     * @return an array of reference to children nodes
     */
    public TrieNode[] getOffspring() {
        return offspring;
    }

    /**
     *
     * @return either true of false if node is a key
     */
    public boolean getIsKey() {
        return this.isKey;
    }

    /**
     * Setter method to set isKey to true or false
     * @param bool
     */
    public void setIsKey(boolean bool) {
        this.isKey = bool;
    }//set isKey attribute of the trie node
}
```

Listing 3: Trie.java

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

 File        : Trie.java

 Date        : Wednesday 25th March 2020

 Author      : 100246776

 Description : Implementation of Trie, Solutions to part 1, Testing

 Last update : 18th May 2020

 **************************************************************************/
package dsacoursework2;
import java.util.LinkedList;
import java.util.Queue;

public class Trie {
    private static final int ASCII_OFFSET = 97;
    private TrieNode root;

    public Trie() {
        root = new TrieNode();
    }//A trie is constructed with a root node

    public Trie(TrieNode trieNode) {
        root = trieNode;
    }

    /**
     * Adds a word to a trie if word not already a key in the trie.
     * @param word should be valie i.e not null or empty
     * @return True for added. False for not added
     */
    boolean add(String word) {
        boolean added = false;
        if (root == null || word==null || word.equals(""))
        {
            System.out.println("Some null reference");
            return added;
        }
        TrieNode parent = root;
        for (int i = 0; i < word.length(); i++) //for each character in the key
        {
            char character = word.charAt(i);//converts i to ascii equivalent
            int index = character - ASCII_OFFSET;//b being 98 and a being 97, then
                ↪ the index of b maps to 1(98-97)
            TrieNode nodeReference=parent.getOffspring()[index];
            if (nodeReference == null)//if no reference to any node
            {
                added = true;
                TrieNode temp = new TrieNode();//create a new node
                parent.getOffspring()[index] = temp; //add temp to this index position
                parent = temp;//set parent to the new node created
            } else {
                parent = nodeReference;//if node exists at particular index, go to it
            }
        }
        if (parent.getIsKey() && !added) {//if node is a key, return false for not
```

```java
57 │                ↪ added
58 │
59 │                System.out.println(word+" already in trie");
60 │                return added;//if node already in trie, return false for not added
61 │            }
62 │            parent.setIsKey(true);
63 │            added=true;
64 │            System.out.println(word+" added");
65 │            return added;
66 │        }
67 │
68 │        /**
69 │         * If word is a key in trie, return true. Otherwise return false.
70 │         * @param word should be valid i.e not null and not empty
71 │         * @return true for word in trie. False for word not in trie
72 │         */
73 │        public boolean contains(String word) {
74 │            boolean contains=false;
75 │            if (root == null  || word ==null || word.equals("")) {
76 │                System.out.println("Some null reference");
77 │                return contains;
78 │            }
79 │            TrieNode p = root;
80 │            for (int i = 0; i < word.length(); i++) {
81 │                char c = word.charAt(i);
82 │                int index = c - ASCII_OFFSET;
83 │                TrieNode[] children=p.getOffspring();
84 │                TrieNode n=children[index];
85 │                if (n== null) {
86 │                    /*
87 │                    Say, we look for word cat in trie. In the offspring array of root,
88 │                        ↪ there is no reference at index 2
88 │                    representing character 'c', this means word 'cat' is not present in
89 │                        ↪ trie. Return false.
89 │                    Same applies for trie nodes further down the tree
90 │                     */
91 │                    System.out.println(word+ " is not in trie");
92 │                    return contains;
93 │                }
94 │                p = n;//if a reference corresponding to the index is found, go to it.
95 │            }
96 │            /*
97 │            Return true only when word is a key
98 │             */
99 │            if (p.getIsKey()) {
100│                System.out.println(word+ " is in trie");
101│                contains= true;
102│                return contains;
103│            }
104│            /*
105│            Say,we look for word 'cat'.The characters 'c','a','t' are in the trie.
105│                ↪ However, node that holds 't' is not a
106│            key. This means 'cat' in the trie is part of another word that is a key. For
106│                ↪ e.g catalogue. In this case, return
107│            false
108│             */
109│            System.out.println(word+ " is not in trie");
110│            return contains;//trie does not contain the key
111│        }
112│
113│        /**
```

```java
114        * Goes through trie and outputs each character in breadth first search order
115        * @return A string containing characters in that order
116        * Design decision:
117        * Use a Queue because it respects first-in first-out rule. This is useful for
              ↪ breadth first search.
118        * Underlying implementation of queue is a linkedlist. This ensures ensure O(1)
              ↪ for enqueue and dequeue.
119        */
120       public String outputBreadthFirstSearch() {
121           if (root == null)
122               return null;
123           StringBuilder sb = new StringBuilder();
124           Queue<TrieNode> queue = new LinkedList<>();
125           queue.add(root);
126           while (!queue.isEmpty()) {
127               TrieNode temp = queue.remove();//remove oldest node added
128               TrieNode[] c=temp.getOffspring();
129               for (int i = 0; i < c.length; i++) {
130                   if (c[i] != null) {
131                       queue.add(c[i]);//add offspring of node to the queue so it goes
                              ↪ through them in breadth first order
132                       sb.append((char) (i + ASCII_OFFSET));//index of offspring
                              ↪ reference corresponds to character
133                   }
134               }
135           }
136           System.out.println(sb.toString());//testing purposes
137           return sb.toString();//toString method returns a string object
138       }
139
140       /**
141        * Interface method to outputDepthFirstSearch(StringBuilder sb, TrieNode trieNode)
142        * @return A string of characters in depth first search order
143        */
144       public String outputDepthFirstSearch()
145       {
146           StringBuilder sb=new StringBuilder();
147           String result=outputDepthFirstSearch(sb,root);
148           return result;
149       }
150
151       /**
152        * Goes through trie and outputs character in depth first search order
153        * @param sb An empty stringbuilder to store characters
154        * @param trieNode Start point of depth first Search. Can specify any node as
              ↪ returned by getSubTrie(String prefix)
155        * @return A string of characters in depth first search order
156        * Design decision:
157        * Recursion as opposed to a iterative version provides a more elegant and
              ↪ shorter solution
158        * No need to handle stack data structure
159        */
160       private String outputDepthFirstSearch(StringBuilder sb, TrieNode trieNode) {
161           if (trieNode == null) {//base case for non=existing trie
162               return null;
163           }
164           TrieNode[] children = trieNode.getOffspring();
165           for (int i = 0; i < children.length; i++) {
166               if (children[i] != null)//ignore nulls
167               {
```

```
168             sb.append((char) (i + ASCII_OFFSET));//convert index to char
169             outputDepthFirstSearch(sb, children[i]);//recursive call on node's
                    ↪ children- depth first
170         }
171     }
172     return sb.toString();
173 }
174
175 /**
176  * @param prefix, an incomplete word
177  * @return A Trie rooted at prefix
178  */
179 public Trie getSubTrie(String prefix) {
180     int index;
181     if (root == null || prefix==null || prefix.equals("")) {
182         return null;
183     }
184     TrieNode parent = root;
185     for (int i = 0; i < prefix.length(); i++) {
186         index = prefix.charAt(i) - ASCII_OFFSET;
187         if (parent.getOffspring()[index] == null) {//prefix does not exist
188             return null;
189         }
190         parent = parent.getOffspring()[index];//if node is not null, go to it
191     }
192     return new Trie(parent);//return a trie rooted at the last node assigned to
            ↪ parent
193 }
194
195 /**
196  * Interface method to getAllWords(LinkedList listOfWords, StringBuilder sb,
        ↪ TrieNode root)
197  * @return A LinkedList containing all keys in the trie
198  */
199 public LinkedList<String> getAllWords()//interface method for private getAllWords
        ↪ method
200 {
201     if(root==null){
202         return null;
203     }
204     LinkedList<String> listOfWords = new LinkedList<>();
205     getAllWords(listOfWords, new StringBuilder(), root);
206     if(listOfWords.size()==0){
207         return null;
208     }
209     return listOfWords;
210 }
211
212 /**
213  * Populates a LinkedList with keys from the trie.
214  * @param listOfWords A linkedlist
215  * @param sb An empty StringBuilder
216  * @param root Root node of Trie
217  */
218 private void getAllWords(LinkedList listOfWords, StringBuilder sb, TrieNode root)
        ↪ {
219     if (root.getIsKey()) {
220         System.out.println(sb.toString());//for test purposes
221         listOfWords.add(sb.toString());//if key is found, add to listOfWords
222     }
```

```java
223            TrieNode[] offspring = root.getOffspring();
224            for (int i = 0; i < offspring.length; i++) {
225                TrieNode o=offspring[i];
226                if (o != null) {
227                    getAllWords(listOfWords, sb.append((char) (ASCII_OFFSET + i)), o);
228                    /*
229                    Say, the key 'bat'is added to stringBuilder, there might still be
                         ↪ other keys to find such as bass
230                    Once function returns, we continue adding to the string "ba" if there
                         ↪ any other keys.
231                     */
232                    sb.setLength(sb.length() - 1);//backtracks one level
233                }
234            }
235        }
236
237
238    public static void main(String[] args) {
239        Trie trie = new Trie();
240
241        System.out.println("Testing add(key)...");
242        if (!trie.add("bat")) throw new AssertionError();
243        if (!trie.add("cat")) throw new AssertionError();
244        if (!trie.add("chat")) throw new AssertionError();
245        if (!trie.add("cheese")) throw new AssertionError();
246        if (!trie.add("cheers")) throw new AssertionError();
247        if (!trie.add("cheer")) throw new AssertionError();
248        if (trie.add("bat")) throw new AssertionError();
249        if (trie.add("cat")) throw new AssertionError();
250        if (trie.add("chat")) throw new AssertionError();
251        if (trie.add("cheese")) throw new AssertionError();
252        if (trie.add("cheers")) throw new AssertionError();
253        if (trie.add(null)) throw new AssertionError();
254        if (trie.add("")) throw new AssertionError();
255
256        System.out.println("\nTesting contains(key)...");
257        if (!trie.contains("bat")) throw new AssertionError();
258        if (!trie.contains("cat")) throw new AssertionError();
259        if (!trie.contains("cat")) throw new AssertionError();
260        if (!trie.contains("cheer")) throw new AssertionError();
261        if (!trie.contains("cheers")) throw new AssertionError();
262        if (trie.contains("coronavirus")) throw new AssertionError();
263        if (trie.contains("china")) throw new AssertionError();
264        if (trie.contains("ch")) throw new AssertionError();
265        if (trie.contains("")) throw new AssertionError();
266        if (trie.contains(null)) throw new AssertionError();
267
268        System.out.println("\nTesting outputBreadthFirstSearch()");
269        if(!trie.outputBreadthFirstSearch().equals("bcaahttaetersse")) throw new
               ↪ AssertionError();
270
271        System.out.println("\nTesting outputDepthFirstSearch()");
272        System.out.println(trie.outputDepthFirstSearch());
273        if(!trie.outputDepthFirstSearch().equals("batcathateersse")) throw new
               ↪ AssertionError();
274
275        System.out.println("\nTesting getAllWords()");
276        trie.getAllWords();
277        System.out.println("Test passed. If adding more keys, check if getAllWords()
               ↪ finds all keys");
```

```java
278
279          System.out.println("\nTesting getSubTrie(String prefix)");
280          Trie subtrie=trie.getSubTrie("ch");
281          subtrie.getAllWords();
282          System.out.println("Test passed. If adding more keys, check if getAllWords()
                  ↪ finds all keys");
283
284          System.out.println("\nTesting getSubTrie(String prefix)");
285          Trie subtrie1=trie.getSubTrie("coron");
286          if(subtrie1==null){
287              System.out.println("Invalid prefix");
288          }
289
290      }
291  }
```

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

 File        : AutoCompletionTrieNode.java

 Date        : Monday 30th March 2020

 Author      : 100246776

 Description : Implementation of AutoCompletionTrieNode

 Last update : 18th May 2020

 ***************************************************************************/
package dsacoursework2;

public class AutoCompletionTrieNode {
    private boolean isKey;
    private final AutoCompletionTrieNode[] offspring;
    private int wordFrequency;

    /**
     * wordFrequency is only set when the node is a key
     */
    public AutoCompletionTrieNode()
    {
        this.offspring = new AutoCompletionTrieNode[26];
        this.wordFrequency = 0;
    }

    /**
     *
     * @return an array of offspring
     */
    public AutoCompletionTrieNode[] getOffSpring() {
        return offspring;
    }

    /**
     *
     * @return true if isKey is set to true. False otherwise
     */
    public boolean getIsKey() {
        return this.isKey;
    }

    /**
     *
     * @param bool, true when node needs to be a key. False otherwise.
     */
    public void setIsKey(boolean bool) {
        this.isKey = bool;
    }

    /**
     *
     * @param wordFrequency an integer that specifies the frequency of occurrence of
     *         ↪ word
     */
    public void setWordFrequency(int wordFrequency) {
```

```java
59            this.wordFrequency = wordFrequency;
60        }
61
62
63        public int getWordFrequency()
64        {
65            return wordFrequency;
66        }
67  }
```

Listing 5: AutoCompletionTrie.java

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

   File        : AutoCompletionTrie.java

   Date        : Wednesday 18th March 2020

   Author      : 100246776

   Description : Solutions for part 3

   Last update : 18th May 2020

 ******************************************************************************/
package dsacoursework2;

import java.io.*;
import java.text.DecimalFormat;
import java.util.*;

public class AutoCompletionTrie {
    private static final int ASCII_OFFSET = 97;
    private AutoCompletionTrieNode root;

    public AutoCompletionTrie() {
        root = new AutoCompletionTrieNode();
    }

    public AutoCompletionTrie(AutoCompletionTrieNode trieNode) {
        root = trieNode;
    }

    /**
     * Extracts word-frequency pairs for each line of a CSV file and call add method
     *      ↪ to add to AutoCompletionTrie
     * @param dictionary, a CSV file with each line of the format word,frequency
     * @throws IOException
     */
    public void addDictionaryToTrie(String dictionary) throws IOException {
        BufferedReader csvReader = new BufferedReader(new FileReader(dictionary));
        String line;
        String[] dictionaryEntry;
        while ((line = csvReader.readLine()) != null) {//reads line by line until end
             ↪ of document
            try {
                /*
                E.g. Alive,2 becomes dictionaryEntry[0]="Alive", dictionaryEntry[1]=2
                 */
                dictionaryEntry = line.split(",");
                String word = dictionaryEntry[0];
                int wordFrequency = Integer.parseInt(dictionaryEntry[1]);
                boolean isAdded = this.add(word, wordFrequency);//add word and
                     ↪ frequency to trie
                if (isAdded) {
                    System.out.printf(" %s added to trie\n",
                         ↪ dictionaryEntry[0]);//test purposes
                }
            } catch (Exception e) {
                e.printStackTrace();
            }
```

```java
56              }
57              csvReader.close();
58          }
59
60
61          /**
62           * If word is a key , update its wordFrequency. Otherwise , add the word and
                ↪ wordFrequency to the trie
63           * @param word
64           * @param wordFrequency
65           * @return True for added and false for not added
66           */
67          public boolean add(String word, int wordFrequency) {
68              boolean added = false;
69              if (root == null || word == null) {
70                  return added;
71              }
72              AutoCompletionTrieNode p = root; //the root
73              for (int i = 0; i < word.length(); i++) //for each character in the key
74              {
75                  char c = word.charAt(i);//converts i to char
76                  int index = c - ASCII_OFFSET;//b being 98 and a being 97, then the index
                    ↪ of b maps to 1(98-97)
77                  AutoCompletionTrieNode[] offSpring = p.getOffSpring();
78                  AutoCompletionTrieNode n = offSpring[index];
79                  if (n == null)//if no reference to any node
80                  {
81                      added = true;
82                      AutoCompletionTrieNode temp = new AutoCompletionTrieNode();//create a
                        ↪ new node
83                      p.getOffSpring()[index] = temp;
84                      p = temp;//set p to the new node created
85                  } else {
86                      p = n;//if node exists at particular index , go to it
87                  }
88              }
89              if (p.getIsKey() && !added) {
90                  //if word is a key, update its frequency of occurrence
91                  int newFrequency = wordFrequency + p.getWordFrequency();
92                  p.setWordFrequency(newFrequency);
93                  System.out.printf("%s is a key in trie. Frequency is now %d\n", word,
                    ↪ newFrequency);
94                  return added;//added is still false. No new key added
95              }
96              p.setIsKey(true);
97              added = true;
98              p.setWordFrequency(wordFrequency);//set frequency of word when it is a key
99              return added;
100         }
101
102
103         public boolean contains(String key) {
104             boolean contains = false;
105             if (root == null || key == null || key.equals("")) {
106                 System.out.println("Some null reference");
107                 return contains;
108             }
109             AutoCompletionTrieNode p = root;//starts at the root
110             for (int i = 0; i < key.length(); i++) {
111                 char c = key.charAt(i);//converts i to char
```

```java
112             int index = c - ASCII_OFFSET;//b being 98 and a being 97, then the index
                    ↪ of b maps to 1(98-97)
113             AutoCompletionTrieNode[] children = p.getOffSpring();
114             AutoCompletionTrieNode n = children[index];
115             if (n == null) {//if a index that maps to a character is not found, key
                    ↪ is not in trie
116                 /*
117                 Say, we look for word cat in trie. In the offspring array of root,
                        ↪ there is no reference at index 2
118                 representing character 'c', this means word 'cat' is not present in
                        ↪ trie. Return false.
119                 Same applies for trie nodes further down the tree
120                  */
121                 System.out.println(key + " is not in trie");
122                 return contains;
123             }
124             p = n;//if index exists, go to the node
125         }
126         /*
127         Return true only when word is a key
128          */
129         if (p.getIsKey()) {
130             System.out.println(key + " is in trie");
131             contains = true;
132             return contains;
133         }
134         /*
135         Say,we look for word 'cat'.The characters 'c','a','t' are in the trie.
                ↪ However, node that holds 't' is not a
136         key. This means 'cat' in the trie is part of another word that is a key. For
                ↪ e.g catalogue. In this case, return
137         false
138          */
139         System.out.println(key + " is not in trie");
140         return contains;//trie does not contain the key
141     }
142
143
144     public String outputBreadthFirstSearch() {
145         if (root == null)
146             return null;
147         StringBuilder sb = new StringBuilder();
148         Queue<AutoCompletionTrieNode> queue = new LinkedList<>();
149         queue.add(root);
150         while (!queue.isEmpty()) {
151             AutoCompletionTrieNode temp = queue.remove();//remove oldest node added
152             AutoCompletionTrieNode[] c = temp.getOffSpring();
153             for (int i = 0; i < c.length; i++) {
154                 if (c[i] != null) {
155                     queue.add(c[i]);//add offspring of node to the queue so it goes
                            ↪ through them in breadth first order
156                     sb.append((char) (i + ASCII_OFFSET));//index of offspring
                            ↪ corresponds to character
157                 }
158             }
159         }
160         System.out.println(sb.toString());//for testing purposes
161         return sb.toString();
162     }
163
```

```java
164      public String outputDepthFirstSearch()//interface function
165      {
166          StringBuilder sb = new StringBuilder();
167          String result = outputDepthFirstSearch(sb, root);
168          return result;
169      }
170
171
172      private String outputDepthFirstSearch(StringBuilder sb, AutoCompletionTrieNode
            ↪ trieNode) {
173          if (trieNode == null) {//base case for non=existing trie
174              return null;
175          }
176          AutoCompletionTrieNode[] children = trieNode.getOffSpring();
177          for (int i = 0; i < children.length; i++) {
178              if (children[i] != null)//ignore nulls
179              {
180                  sb.append((char) (i + ASCII_OFFSET));//convert index to char(0 maps
                        ↪ to 'a')
181                  outputDepthFirstSearch(sb, children[i]);//recursive call on node's
                        ↪ children- depth first
182              }
183          }
184          return sb.toString();
185      }
186
187
188      public AutoCompletionTrie getSubTrie(String prefix) {
189          int index;
190          if (root == null || prefix == null || prefix.equals("")) {
191              return null;
192          }
193          AutoCompletionTrieNode p = root;
194          for (int i = 0; i < prefix.length(); i++) {
195              index = prefix.charAt(i) - ASCII_OFFSET;
196              if (p.getOffSpring()[index] == null) {//prefix does not exist
197                  return null;
198              }
199              p = p.getOffSpring()[index];
200          }
201          return new AutoCompletionTrie(p);
202      }
203
204
205      private void getAllWords(ArrayList listOfWords, StringBuilder sb,
            ↪ AutoCompletionTrieNode root) {
206          if (root.getIsKey()) {
207              System.out.println(sb.toString());
208              listOfWords.add(sb.toString());
209          }
210          AutoCompletionTrieNode[] children = root.getOffSpring();
211          for (int i = 0; i < children.length; i++) {
212              if (children[i] != null) {
213                  getAllWords(listOfWords, sb.append((char) (ASCII_OFFSET + i)),
                        ↪ children[i]);
214                  /*
215                  Say, the key 'bat'is added to stringBuilder, there might still be
                        ↪ other keys to find such as bass
216                  Once function returns, we continue adding to the string "ba" if there
                        ↪ any other keys.
```

```
217                    */
218                    sb.setLength(sb.length() - 1);//backtracks one level
219                }
220            }
221        }
222
223        public ArrayList<String> getAllWords()
224        {
225            ArrayList<String> listOfWords = new ArrayList<>();
226            getAllWords(listOfWords, new StringBuilder(), root);
227            return listOfWords;
228        }
229
230
231        public static class WordAndFrequency {
232            /*
233            Objects of WordAndFrequency hold a word string and frequency of that key
234             */
235            private String key;
236            private int freq;
237
238            WordAndFrequency(String key, int freq)
239            {
240                this.key = key;
241                this.freq = freq;
242            }
243            WordAndFrequency(String key, String freq) {
244                this.key = key;
245                this.freq = Integer.parseInt(freq);
246            }
247
248            int getFreq() {
249                return freq;
250            }
251
252            String getKey() {
253                return key;
254            }
255        }
256
257        /*
258        Design decision for priority queue:
259        1. Ability to specify priority of an element. As such, no need to sort
260        2. Adding to priority queue is log(n) which is fairly fast
261        3. Poll() is O(1) time complexity
262        4. Better space complexity than using an arraylist that doubles in size to
263            ↪ accommodate more elements
264        5. Priority queue guarantees maintaining order for highest priority value. It
265            ↪ does not necessarily maintain sorted order.
266            As we only ever use it to find highest priorities, it is better to use it
267                ↪ instead of a treeset,
268             which maintains sorted order throughout
269            It is also more computationally expensive to maintain trees
270         */
271        /**
272         * An interface function
273         * @return A queue, q, containing WordAndFrequency objects
274         */
275        private PriorityQueue<WordAndFrequency> wordFrequencyInfo() {
276            //Each WordAndFrequency object has a key and its frequency
```

```
217                    */
218                    sb.setLength(sb.length() - 1);//backtracks one level
219                }
220            }
221        }
222
223        public ArrayList<String> getAllWords()
224        {
225            ArrayList<String> listOfWords = new ArrayList<>();
226            getAllWords(listOfWords, new StringBuilder(), root);
227            return listOfWords;
228        }
229
230
231        public static class WordAndFrequency {
232            /*
233            Objects of WordAndFrequency hold a word string and frequency of that key
234             */
235            private String key;
236            private int freq;
237
238            WordAndFrequency(String key, int freq)
239            {
240                this.key = key;
241                this.freq = freq;
242            }
243            WordAndFrequency(String key, String freq) {
244                this.key = key;
245                this.freq = Integer.parseInt(freq);
246            }
247
248            int getFreq() {
249                return freq;
250            }
251
252            String getKey() {
253                return key;
254            }
255        }
256
257        /*
258        Design decision for priority queue:
259        1. Ability to specify priority of an element. As such, no need to sort
260        2. Adding to priority queue is log(n) which is fairly fast
261        3. Poll() is O(1) time complexity
262        4. Better space complexity than using an arraylist that doubles in size to
            ↪ accommodate more elements
263        5. Priority queue guarantees maintaining order for highest priority value. It
            ↪ does not necessarily maintain sorted order.
264            As we only ever use it to find highest priorities, it is better to use it
                ↪ instead of a treeset,
265             which maintains sorted order throughout
266            It is also more computationally expensive to maintain trees
267         */
268        /**
269         * An interface function
270         * @return A queue, q, containing WordAndFrequency objects
271         */
272        private PriorityQueue<WordAndFrequency> wordFrequencyInfo() {
273            //Each WordAndFrequency object has a key and its frequency
```

```java
274            PriorityQueue<WordAndFrequency> q = new PriorityQueue<>(new
                  ↪ CompareByFreqAndAlphabet());
275            wordFrequencyInfo(q, new StringBuilder(), root);//get all words in the trie
                  ↪ and their associated frequency
276            return q;
277        }
278
279        /**
280         * Comparator logic for priority queue
281         */
282        static class CompareByFreqAndAlphabet implements Comparator<WordAndFrequency> {
283            public int compare(WordAndFrequency m1, WordAndFrequency m2) {
284                if (m1.getFreq() < m2.getFreq()) {
285                    return 1;
286                } else if (m1.getFreq() > m2.getFreq()) {
287                    return -1;
288                } else {
289                    return m1.getKey().compareTo(m2.getKey());//if same frequency, sort
                          ↪ by key length i.e choose shorter word
290                }
291            }
292        }
293
294
295        /**
296         * The algorithm goes through the AutoCompletionTrie recursively.
297         * When it encounters a node that is a key, it means that the stringbuilder
                  ↪ contains a complete word.
298         * A WordAndFrequency entity is then made and added to the priority queue.
299         * @param queue, an empty priority queue
300         * @param sb, an empty stringbuilder
301         * @param root, starting node in an AutoCompletionTrie
302         */
303        private void wordFrequencyInfo(PriorityQueue queue, StringBuilder sb,
              ↪ AutoCompletionTrieNode root) {
304            if (root.getIsKey()) {
305                queue.add(new WordAndFrequency(sb.toString(), root.getWordFrequency()));
306            }
307            AutoCompletionTrieNode[] offspring = root.getOffSpring();
308            for (int i = 0; i < offspring.length; i++) {
309                if (offspring[i] != null) {
310                    wordFrequencyInfo(queue, sb.append((char) (97 + i)), offspring[i]);
311                    sb.setLength(sb.length() - 1);//backtracks one level
312                }
313            }
314        }
315
316        /**
317         * This function finds the most probable matches for each prefix in the list
318         * Writes matches to a CSV file.
319         * @param prefixes, an ArrayList of prefixes
320         * @throws IOException
321         */
322        public void prefixMatches(ArrayList<String> prefixes) throws IOException {
323            final DecimalFormat df = new DecimalFormat("#.######");
324            ArrayList<String> toWrite=new ArrayList<>();
325            for (String prefix : prefixes) {
326                toWrite.add(prefix);
327                AutoCompletionTrie currentSubTrie = this.getSubTrie(prefix);//get subtrie
                      ↪ rooted at prefix
```

```java
328            /*
329            temp contains a WordAndFrequency objects. The object contains a unique
                   ↪ word and frequency of occurrence
330             */
331            PriorityQueue<WordAndFrequency> temp = currentSubTrie.wordFrequencyInfo();
332            int totalFreq = 0;
333            for (WordAndFrequency wordAndFrequency : temp) {
334                /*
335                adding frequency attribute of each object to find totalFrequency
336                 */
337                totalFreq += wordAndFrequency.getFreq();
338            }
339
340            int limit = Math.min(temp.size(), 3);
341
342            for (int i = 0; i < limit; i++)
343            {
344                WordAndFrequency wf=temp.poll();//removes and return top most object
                       ↪ in priority queue
345                if(wf !=null){
346                    toWrite.add(prefix+wf.getKey());
347                    toWrite.add(df.format((double)wf.getFreq()/totalFreq));
348                }
349            }
350            saveToFile("lotrMatches.csv", toWrite);
351            toWrite.clear();
352        }
353    }
354
355    /**
356     * Writes ArrayList to a file in a specific format
357     * @param writeFile, name of file to write to
358     * @param toWrite, an Arraylist to write to file
359     * @throws IOException
360     */
361    private static void saveToFile(String writeFile, ArrayList<String> toWrite)
          ↪ throws IOException
362    {
363        BufferedWriter writer = new BufferedWriter(
364                new FileWriter(writeFile, true));  //with append set to true, file
                       ↪ data is not overwritten
365        for(int i=0;i<toWrite.size();i++)
366        {
367            if(i != toWrite.size()-1)
368            {
369                writer.write(toWrite.get(i)+",");
370            }
371            else {
372                writer.write(toWrite.get(i)+System.lineSeparator());
373            }
374        }
375        writer.close();
376    }
377
378
379    public static void main(String[] args) throws IOException {
380
381        AutoCompletionTrie trie = new AutoCompletionTrie();
382        DictionaryMaker df = new DictionaryMaker();
383        /*
```

```
384         Writes lotr as a dictionary in csv format
385          */
386         df.formDictionary("lotr.csv", "dictionaryOfWords.csv");
387
388         trie.addDictionaryToTrie("dictionaryOfWords.csv");
389         /*
390         Load queries into an ArrayList
391          */
392         ArrayList<String> lotrQueries =
              ↪ DictionaryMaker.readWordsFromCSV("lotrQueries.csv", "\n");
393
394         trie.prefixMatches(lotrQueries);
395
396     }
397 }
```

# 6 Appendix

[b]

Figure 3: Ascii Table

## ASCII Table

| Dec | Hex | Oct | Char | Dec | Hex | Oct | Char | Dec | Hex | Oct | Char | Dec | Hex | Oct | Char |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | | 32 | 20 | 40 | [space] | 64 | 40 | 100 | @ | 96 | 60 | 140 | ` |
| 1 | 1 | 1 | | 33 | 21 | 41 | ! | 65 | 41 | 101 | A | 97 | 61 | 141 | a |
| 2 | 2 | 2 | | 34 | 22 | 42 | " | 66 | 42 | 102 | B | 98 | 62 | 142 | b |
| 3 | 3 | 3 | | 35 | 23 | 43 | # | 67 | 43 | 103 | C | 99 | 63 | 143 | c |
| 4 | 4 | 4 | | 36 | 24 | 44 | $ | 68 | 44 | 104 | D | 100 | 64 | 144 | d |
| 5 | 5 | 5 | | 37 | 25 | 45 | % | 69 | 45 | 105 | E | 101 | 65 | 145 | e |
| 6 | 6 | 6 | | 38 | 26 | 46 | & | 70 | 46 | 106 | F | 102 | 66 | 146 | f |
| 7 | 7 | 7 | | 39 | 27 | 47 | ' | 71 | 47 | 107 | G | 103 | 67 | 147 | g |
| 8 | 8 | 10 | | 40 | 28 | 50 | ( | 72 | 48 | 110 | H | 104 | 68 | 150 | h |
| 9 | 9 | 11 | | 41 | 29 | 51 | ) | 73 | 49 | 111 | I | 105 | 69 | 151 | i |
| 10 | A | 12 | | 42 | 2A | 52 | * | 74 | 4A | 112 | J | 106 | 6A | 152 | j |
| 11 | B | 13 | | 43 | 2B | 53 | + | 75 | 4B | 113 | K | 107 | 6B | 153 | k |
| 12 | C | 14 | | 44 | 2C | 54 | , | 76 | 4C | 114 | L | 108 | 6C | 154 | l |
| 13 | D | 15 | | 45 | 2D | 55 | - | 77 | 4D | 115 | M | 109 | 6D | 155 | m |
| 14 | E | 16 | | 46 | 2E | 56 | . | 78 | 4E | 116 | N | 110 | 6E | 156 | n |
| 15 | F | 17 | | 47 | 2F | 57 | / | 79 | 4F | 117 | O | 111 | 6F | 157 | o |
| 16 | 10 | 20 | | 48 | 30 | 60 | 0 | 80 | 50 | 120 | P | 112 | 70 | 160 | p |
| 17 | 11 | 21 | | 49 | 31 | 61 | 1 | 81 | 51 | 121 | Q | 113 | 71 | 161 | q |
| 18 | 12 | 22 | | 50 | 32 | 62 | 2 | 82 | 52 | 122 | R | 114 | 72 | 162 | r |
| 19 | 13 | 23 | | 51 | 33 | 63 | 3 | 83 | 53 | 123 | S | 115 | 73 | 163 | s |
| 20 | 14 | 24 | | 52 | 34 | 64 | 4 | 84 | 54 | 124 | T | 116 | 74 | 164 | t |
| 21 | 15 | 25 | | 53 | 35 | 65 | 5 | 85 | 55 | 125 | U | 117 | 75 | 165 | u |
| 22 | 16 | 26 | | 54 | 36 | 66 | 6 | 86 | 56 | 126 | V | 118 | 76 | 166 | v |
| 23 | 17 | 27 | | 55 | 37 | 67 | 7 | 87 | 57 | 127 | W | 119 | 77 | 167 | w |
| 24 | 18 | 30 | | 56 | 38 | 70 | 8 | 88 | 58 | 130 | X | 120 | 78 | 170 | x |
| 25 | 19 | 31 | | 57 | 39 | 71 | 9 | 89 | 59 | 131 | Y | 121 | 79 | 171 | y |
| 26 | 1A | 32 | | 58 | 3A | 72 | : | 90 | 5A | 132 | Z | 122 | 7A | 172 | z |
| 27 | 1B | 33 | | 59 | 3B | 73 | ; | 91 | 5B | 133 | [ | 123 | 7B | 173 | { |
| 28 | 1C | 34 | | 60 | 3C | 74 | < | 92 | 5C | 134 | \ | 124 | 7C | 174 | | |
| 29 | 1D | 35 | | 61 | 3D | 75 | = | 93 | 5D | 135 | ] | 125 | 7D | 175 | } |
| 30 | 1E | 36 | | 62 | 3E | 76 | > | 94 | 5E | 136 | ^ | 126 | 7E | 176 | ~ |
| 31 | 1F | 37 | | 63 | 3F | 77 | ? | 95 | 5F | 137 | _ | 127 | 7F | 177 | |

# References

[1]  Dr. Anthony Bagnall. *Tries*. Lecture notes. University of East Anglia,Norwich

[2]  Dr. Anthony Bagnall. *Trees*. Lecture notes. University of East Anglia,Norwich

[3]  Dr. Anthony Bagnall,Jason Lines. *Algorithm design and recursion*. Labsheet. University of East Anglia,Norwich

[4]  Dr. Anthony Bagnall. *Lab 9: Trees* . Labsheet. University of East Anglia,Norwich

[5]  Dr. Anthony Bagnall. *Lab 10: Trie* . Labsheet. University of East Anglia,Norwich