COMP 3270 FALL 2018
**Programming Project: Autocomplete**

Name: Zachary Parker                    Date Submitted: 11-26-18

1. **Pseudocode**: Understand the strategy provided for *TrieAutoComplete*. State the algorithm for the functions precisely using numbered steps that follow the pseudocode conventions that we use. Provide an approximate efficiency analysis by filling the table given below, for your algorithm.

*Add*
- Pseudocode:

```
Add(word: string, weight: double)
1       if word is null
2               throw an exception
3       if weight < 0
4               throw an exception
5       Node curr = root of tree
6       Node nextNode = null
7       char index = '-'
8       for i to word.length
9               index = the character in word at index i
10              nextNode = the child of curr with character index
11              if nextNode is null
12                      instantiate nextNode with character index,
                        parent curr, and given weight
13                      add new node to children map
14              if curr's subtree max weight < weight
15                      curr's subtree max weight = weight
16              point curr to nextNode
17      set curr's word value
19      set curr's weight
20      set that curr is a word
```

- Complexity analysis:

| Step # | Complexity stated as O(_) |
|--------|---------------------------|
| 1 | O(1) |
| 2 | O(1) |
| 3 | O(1) |
| 4 | O(1) |
| 5 | O(1) |
| 6 | O(1) |
| 7 | O(1) |
| 8 | O(n) |
| 9 | O(n) |
| 10 | O(n) |
| 11 | O(n) |
| 12 | O(n) |
| 13 | O(n) |
| 14 | O(n) |
| 15 | O(n) |
| 16 | O(n) |
| 17 | O(1) |
| 18 | O(1) |
| 19 | O(1) |

Complexity of the algorithm = O(n)

*topMatch*

- Pseudocode:

```
topMatch(prefix: string)
1      if prefix is null
2            throw an exception
3      nodeQueue = new priority queue of Nodes using a
       reverse weight comparator
4      Node curr = root of tree
5      char index = '-'
6      boolean wordExists = true
7      string empty = ""
8      String topMatch = ""
9      for i to prefix.length
10           index = the character in the word at index i
11           if a child of curr contains the value of index
12                 point curr to that child
13           else
14                 wordExists = false
15                 break
16     if (wordExists is false)
17           return empty
18     while (curr is a word and curr's weight < curr's
       subtree max weight)
19           for every node in curr's children
20                 add that node to nodeQueue
21           point curr to the head of queue and remove
             element from queue
22     topMatch = curr's word value
23     return topMatch
```

- Complexity analysis:

| Step # | Complexity stated as O(_) |
|--------|---------------------------|
| 1 | O(1) |
| 2 | O(1) |
| 3 | O(1) |
| 4 | O(1) |
| 5 | O(1) |
| 6 | O(1) |
| 7 | O(1) |
| 8 | O(1) |
| 9 | O(n) |
| 10 | O(n) |
| 11 | O(n) |
| 12 | O(n) |
| 13 | O(n) |
| 14 | O(n) |
| 15 | O(n) |
| 16 | O(1) |
| 17 | O(1) |
| 18 | O(n) |
| 19 | $O(n^2)$ |
| 20 | $O(n^2)$ |
| 21 | O(n) |
| 22 | O(1) |
| 23 | O(1) |

Complexity of the algorithm = $O(n^2)$

*topMatches*

Pseudocode:

**topMatches**(prefix: string, k: integer)
```
1       if prefix is null
2               throw an exception
3       Node curr = root of tree
4       char index = '-'
4       boolean matchExists = true
5       nodeQueue = new priority queue of Nodes using a
        reverse weight comparator
6       emptyList = new ArrayList of strings
7       topMatches = new ArrayList of Nodes
8       topMatchesString = new ArrayList of strings
9       if k = 0
10              return emptyList
11      for i to prefix.length
12              index = the character in the word at index i
13              point curr to its child with the value of index if
                possible
14              if curr is null
15                      return emptyList
16      if curr is not null
17              add curr to nodeQueue
18      while nodeQueue is not empty
19              point curr to head of nodeQueue and remove that
                node from queue
20              for every node in curr's children
21                      add that node to nodeQueue
22              if the size of topMatches >= k
23                      sort topMatches in descending order of
                        value
24                      break
25              if curr points to a word
26                      add curr to topMatches
27      for every node in topMatches
28              add that node's word to topMatchesString
29      return topMatchesString
```

- Complexity analysis:

| Step # | Complexity stated as O(_) |
|--------|---------------------------|
| 1 | O(1) |
| 2 | O(1) |
| 3 | O(1) |
| 4 | O(1) |
| 5 | O(1) |
| 6 | O(1) |
| 7 | O(1) |
| 8 | O(1) |
| 9 | O(1) |
| 10 | O(1) |
| 11 | O(n) |
| 12 | O(n) |
| 13 | O(n) |
| 14 | O(n) |
| 15 | O(n) |
| 16 | O(1) |
| 17 | O(1) |
| 18 | O(n) |
| 19 | O(n) |
| 20 | $O(n^2)$ |
| 21 | $O(n^2)$ |
| 22 | O(n) |
| 23 | O(n log(n)) |
| 24 | O(1) |
| 25 | O(n) |
| 26 | O(n) |
| 27 | O(n) |
| 28 | O(n) |
| 29 | O(1) |

Complexity of the algorithm = $O(n^2)$

2.**Testing**: Complete your test cases to test the *TrieAutoComplete* functions based upon the criteria mentioned below.

### Test of correctness:

Assuming the trie already contains the terms {"ape, 6", "app, 4", "ban, 2", "bat, 3", "bee, 5", "car, 7", "cat, 1"}, you would expect results based on the following table:

| Query | k | Result |
|-------|-----|--------|
| "" | - | Car |
| "a" | - | Ape |
| "ap" | - | Ape |
| "b" | - | Bee |
| "ba" | - | Bat |
| "c" | - | Car |
| "ca" | - | Car |
| "cat" | - | Cat |
| "d" | - | "" |
| " " | - | "" |
| "" | 8 | {"car", "ape", "bee", "app", "bat", "ban", "cat"} |
| "" | 1 | {"car"} |
| "" | 2 | {"car", "ape"} |
| "" | 3 | {"car", "ape", "bee"} |
| "a" | 1 | {"ape"} |
| "ap" | 1 | {"ape"} |
| "b" | 2 | {"bee", "bat"} |
| "ba" | 2 | {"bee", "bat"} |
| "d" | 100 | {} |

3.**Analysis**: Answer the following questions. Use data wherever possible to justify your answers, and keep explanations brief but accurate:

i. What is the order of growth (big-Oh) of the number of compares (in the worst case) that each of the operations in the *Autocompletor* data type make?
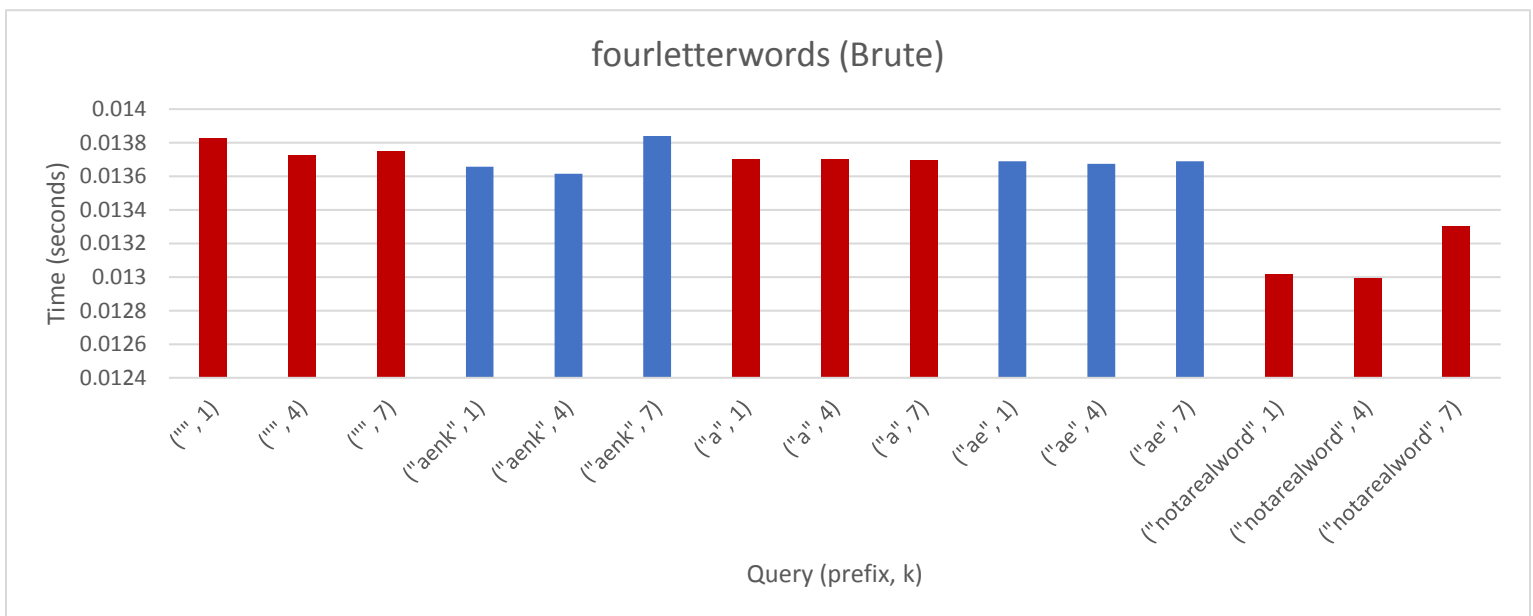
add: O(n)

topMatch: $O(n^2)$

topMatches: $O(n^2)$

This information is presented above in Problem 1.
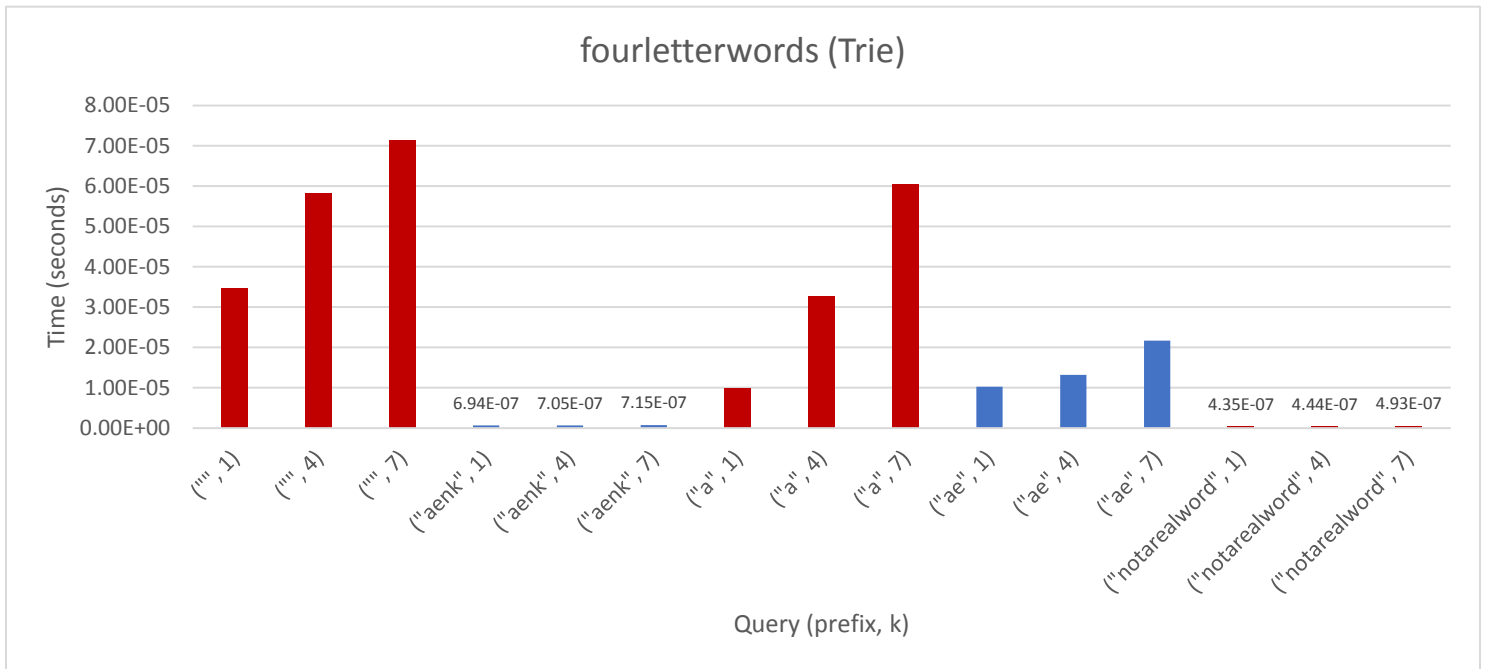
ii. How does the runtime of *topMatches()* vary with k, assuming a fixed prefix and set of terms? Provide answers for *BruteAutocomplete* and *TrieAutocomplete*. Justify your answer, with both data and algorithmic analysis.

The following observations were made using the given AutocompletorBenchmark class and fourletterwords text file.

BruteAutocomplete: The runtime of topMatches() is nearly the exact same for a fixed prefix and set of terms. This suggests that the runtime of topMatches() in this given implementation is not actually affected the value of k.
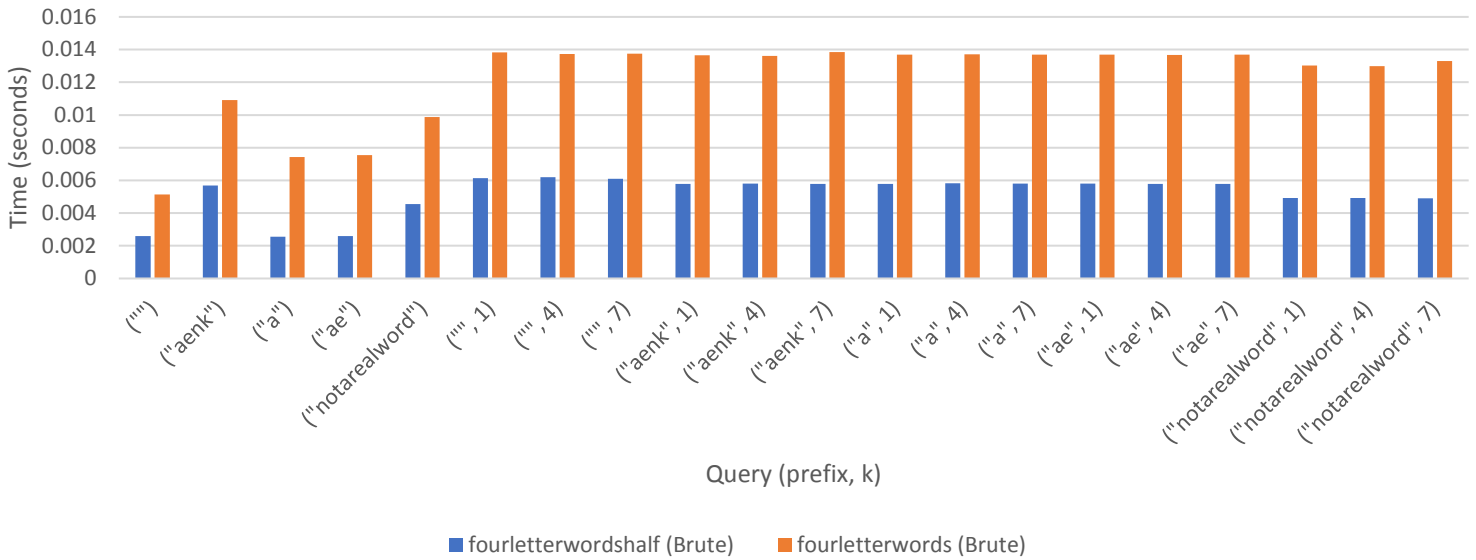


fourletterwords (Brute)

**fourletterwords (Trie)**



iii. How does increasing the size of the source and increasing the size of the prefix argument affect the runtime of *topMatch* and *topMatches*? (Tip: Benchmark each implementation using fourletterwords.txt, which has all four-letter combinations from aaaa to zzzz, and fourletterwordshalf.txt, which has all four-letter word combinations from aaaa to mzzz. These datasets provide a very clean distribution of words and an exact 1-to-2 ratio of words in source files.)
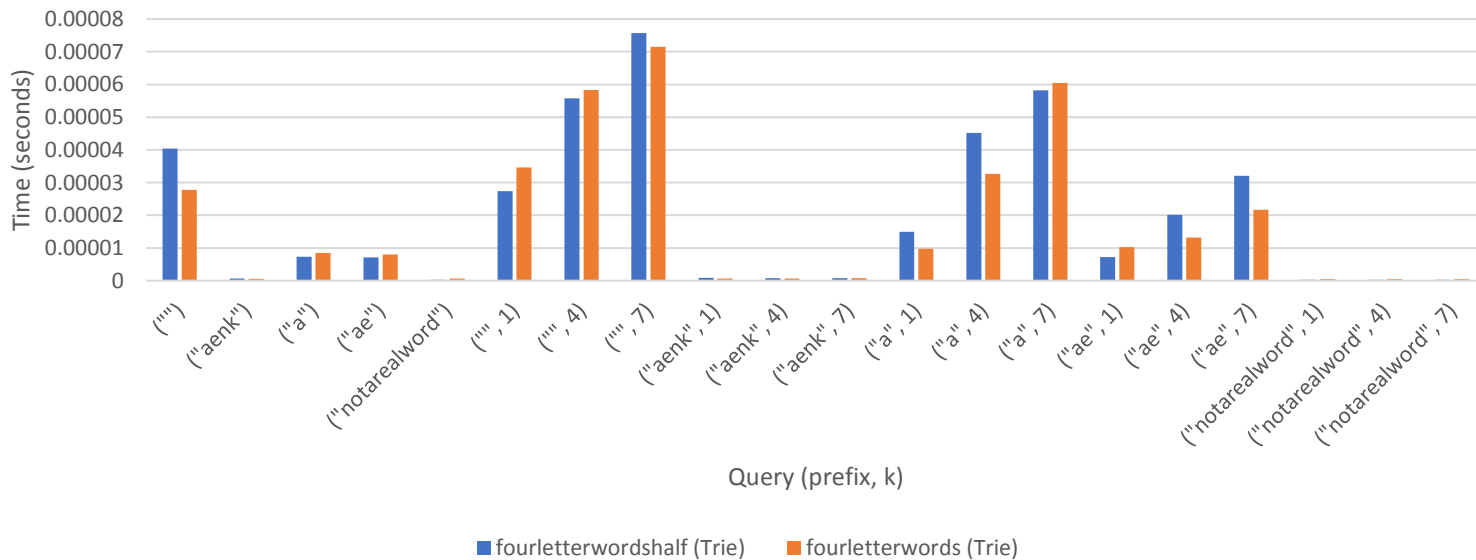
**9**

**topMatch and topMatches: Source Size Comparison (Brute)**

For TrieAutocomplete, increasing the size of the source does not discernibly alter the runtime of topMatch and topMatches. This means that this implementation is far more scalable in comparison to BruteAutocomplete.



**topMatch and topMatches: Source Size Comparison (Trie)**

4. Graphical Analysis: Provide a graphical analysis by comparing the following:

    i.    The big-Oh for *TrieAutoComplete* after analyzing the pseudocode and big-Oh for *TrieAutoComplete* after the implementation.

Pseudocode:

add: O(n)
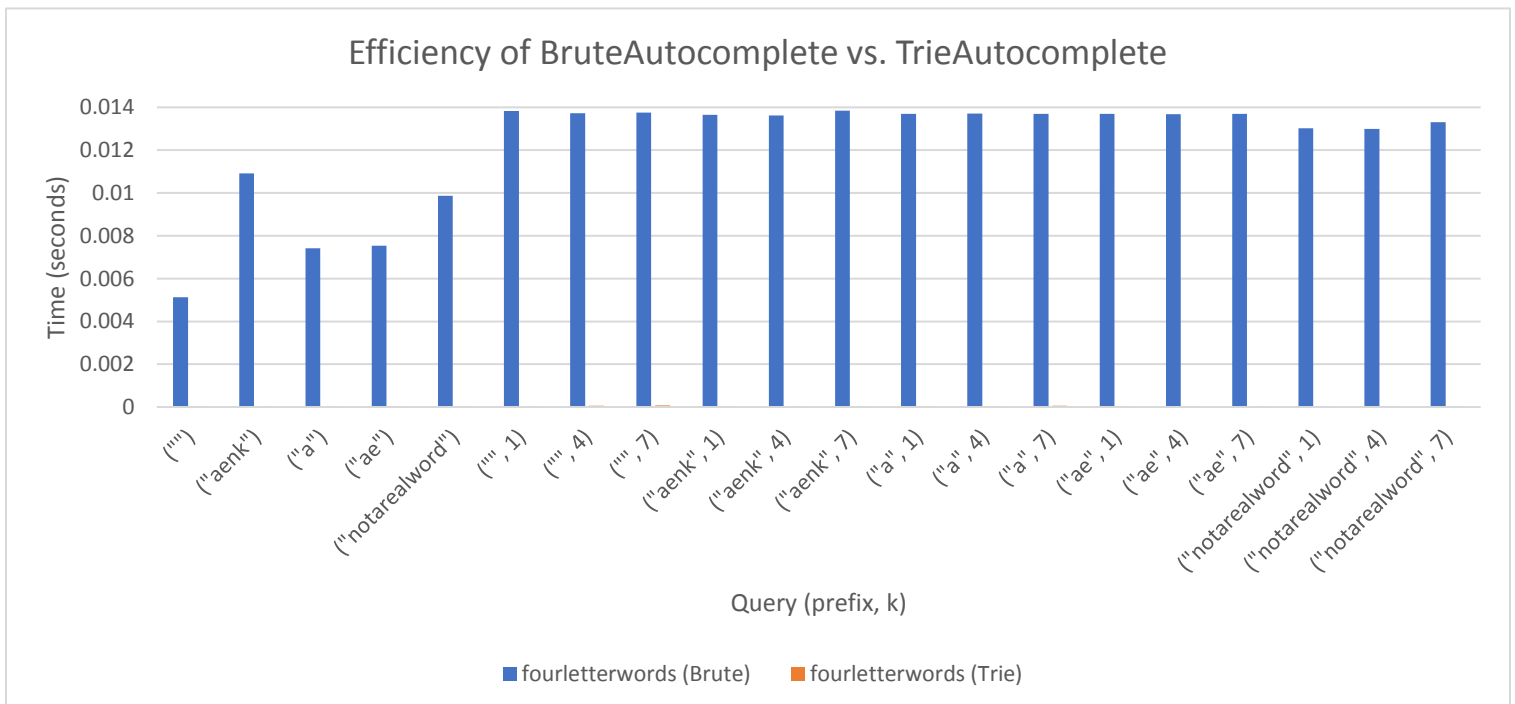
topMatch: $O(n^2)$

topMatches: $O(n^2)$


Implementation:

add: O(n)

topMatch: $O(n^2)$

topMatches: $O(n^2)$


    ii.    Compare the *TrieAutoComplete* with *BruteAutoComplete*.

In all given benchmark cases, TrieAutocomplete is orders of magnitude more efficient than BruteAutocomplete. Note that, when compared to the brute force implementation, the allotted time of the trie strategy is hardly visible in the graph.



Efficiency of BruteAutocomplete vs. TrieAutocomplete

Efficiency of BruteAutocomplete vs. TrieAutocomplete