

ASSIGNMENT 4: HASHING for STRINGS

Due Date: Friday April 15th, 2022. 11:59 pm.

Note: A maximum of one late day will be allowed for this assignment (with 10% late penalty)

Goal: The goal of this assignment is to get some practice with collision resolution and hash functions. On the side you will also learn basic string manipulations. It's a fun assignment where the task is to find valid anagrams of a given input.

Problem Statement: You are given a vocabulary *V* of (lowercase) English words. It uses letters of the English alphabet [a-z], digits [0-9], and the apostrophe symbol [']. No other characters are used in the vocabulary *V*. Your goal is to print out all valid anagrams of an input string present in the vocabulary. The input string will be a sequence of at most 12 characters.

Anagram: Two strings are anagrams of each other if by rearranging letters of one string you can obtain the other. For example, "a bit" is an anagram of "bait", and "super" is an anagram of "purse". Note that we can add spaces at will, i.e., we won't count spaces when matching characters for checking anagrams. In this assignment, you will load *V* (Vocabulary) from the text file and then be ready to compute anagrams. You will be provided an input file also in the text format. In both vocabulary and input files there will be one string written per line. Your goal will be to compute all valid anagrams (i.e., each word within your anagram must be present in *V*) of all input strings. After computing all valid anagrams of one string you must output a '-1' to indicate that you are done computing anagrams of this string. For the purpose of this assignment, you only have to compute anagrams with a maximum of 2 spaces in them (i.e., three words at most). However, each permutation of these words will also be a valid anagram. This is the first assignment in the course where you will be evaluated not only on the correctness and complexity of your code, but also on the runtime efficiency of the code. You can compute the time taken for your code to run using the built-in `getTimeMillis()` command.

Vocabulary File: The vocabulary file (`vocabulary.txt`) will be provided in the resources of the assignment. The first line of the Vocabulary will indicate the number of words in the Vocabulary (*V*), followed by one word per line (all lowercase and no spaces). A sample `vocabulary.txt` is given below:

```
6
a
it
bit
bat
tab
```

i

Input File: The input file (input.txt) will be an input to the code at runtime. The first line will have the number (K) of input strings. This will be followed by K lines, with one string per line. It will have only lowercase letters, digits, and an apostrophe. It will not have a space. A sample input.txt is given as under:

2

bait

bb

Output File: You will produce all valid anagrams of each input string that can be constructed using the vocabulary and output -1 after finishing with one input and moving onto the next. The output for a particular string should be in lexicographic order. Lexicographic ordering is done based on ASCII codes: i.e., lowercase>digits>apostrophe>space. For example, for the input file above you will output:

a bit

bat i

bit a

i bat

i tab

tab i

-1

-1

Note that for the second input word there were no valid anagrams found. Also note that the number of '-1's in the output should be exactly the same as the number of input words in input.txt. Your output must be produced on stdout (without any other extra information). Further note that output anagrams should not have contiguous spaces. They should not start with a space, or end with a space. These will be required for correctly autograding your assignment.

Hashing: The main purpose of the assignment is to have you store vocabulary appropriately and have you check for anagrams efficiently. There may be many ways to store the vocabulary, but in this assignment you must hash each valid word. You will have to implement your own hash function and your own collision resolution. You should implement quadratic probing as discussed in the class. The goal is that your anagram computation should be as efficient as possible (given quadratic probing). You may use any function within Java built-in String class, **except hashCode() or any other inbuilt hash functions.**

Tip: To compute better time efficiency not only will you have to implement a good hashing mechanism, you will also have to create an optimized approach to search through the space of anagrams. This may take some trial and error, so start early!

For Fun: You must find some friend of yours (either in the class or otherwise) and output a funny anagram of their name. Share this anagram with TA at the time of demo. There are no points for a more humorous anagram, although there are points for completing this part of the task.

Code: Your code will be run using the following command:

```
javac Anagram.java
```

```
java Anagram vocabulary.txt input.txt
```

This implies that we can change the vocabulary.txt at the time of final evaluation. However, its size will be in the range of the size of the vocabulary.txt we are providing with the assignment. Also, there will be no words in the vocabulary that have sizes 1 or 2. That is, all valid words will be at least three characters long.

What to submit?

1. Submit your code in a .zip file named in the format **<EntryNo>.zip** e.g., **2018ME10000.zip**. Make sure that when we run “unzip yourfile.zip”, there should be a directory <EntryNo> created which should contain all your code files including Anagram.java. (do not change the name of this file since this is required for autograding). In addition to your code, “writeup.txt” should also be present in the directory (details below). Thus the directory structure should be as follows after unzipping:

<EntryNo>

- Anagram.java
- Any other files needed for implementation
- writeup.txt

2. You will be penalized for any submissions that do not conform to this requirement.
3. The writeup.txt should have a line that lists names of all students you discussed/collaborated with (see guidelines on collaboration vs. cheating on the course home page). If you never discussed the assignment with anyone say None. After this line, you are welcome to write something about your code, though this is not necessary.

What is allowed? What is not?

1. This is an individual assignment.
2. Your code must be your own. You can browse online resources for any general ideas/concepts, but you are supposed to search for/look at specific code meant to solve these or related problems.

3. You should develop your algorithm using your own efforts. You should not Google search for direct solutions to this assignment. However, you are welcome to Google search for generic Java-related syntax.
4. You must not discuss this assignment with anyone outside the class. **Make sure you mention the names in your write-up in case you discuss with anyone from within the class.** Please refer to the plagiarism related guidelines covered in the first lecture and follow them carefully. In case of any doubts, you are free to contact the TAs or the instructors.
5. You are not allowed to use built-in (or anyone else's) implementations of stacks, queues, vectors, growable arrays and/or other similar data structures - it is ok to use fixed size arrays as covered in the class. As stated above, you can use any functions from the java in-built String class, **except for hashCode(), or any other in-built hash functions.** A key aspect of the course is to have you learn how to implement these data structures. You are free to use your own implementation of any data structures from one of the earlier assignments.
6. Your submitted code will be automatically evaluated against another set of benchmark problems. You get a significant penalty if your output is not automatically parsable and does not follow input-guidelines.
7. We will run plagiarism detection software. Anyone found guilty will be awarded a suitable penalty as per IIT rules.

Evaluation Criteria

The assignment is worth 12 points. Your code will be autograded at the demo time against a series of tests. A separate demo will be taken and points will be reserved for correctness of the code (in terms of whether any built-in functions are being used or not), efficiency as well as your ability to answer questions related to your own code. Your code should be as efficient as possible (primarily think about efficiency in terms of the time and memory complexity, and removing any obvious inefficiencies/redundant operations resulting in very slow implementations). Marks will be deducted for inefficient code/implementations.