

This assignment has several parts -- a comparison of dictionary ADTs (section 19.2.2 of the text) and a concordance-production application using the dictionary ADTs. A Webster's dictionary definition of concordance is: "an alphabetical list of the main words in a work." In addition to the main words, I want you to keep track of all the line numbers where these words occur.

WORD & LINE CONCORDANCE

The goal of this assignment is to process a textual, data file to generate a word concordance with line numbers for each word. A dictionary ADT is perfect to store the word concordance with the word being the dictionary *key* and a list of its line numbers being the associated *value* with the key. Since the concordance should only keep track of the "main" words, there will actually be two files of textual information: the data file to be processed and a stop-word file. The stop-word file will contain a list of stop words (e.g., "a", "the", etc.) -- these words will **not** be included in the concordance even if they do appear in the data file. The stop words are to be stored in a second dictionary which will be accessed to ensure its contents do not appear in the concordance. Words in the data file are to be extracted, compared with the stop words, and, when appropriate, added to the concordance list along with the line number of the current occurrence of the word. Often, a word will be encountered several times--each line of encounter is to be recorded in the concordance list, but each word is to appear only once. The output of the program should be a text file containing the concordance words printed out in alphabetical order along with their corresponding line numbers.

NOTES:

- a) Words are defined to be sequences of letters that are delimited by any non-letter (e.g., white space, punctuation, brackets, parentheses, dashes, double quotes, etc.)
- b) There is to be no distinction made between upper and lower case characters, i.e., "CAT" is the same word as "cat".
- c) Blank lines are to be counted in the line numbering.
- d) The word-concordance application should not reference the dictionary ADT implementation, except through the dictionary operations.
- e) It is strongly suggested that the logic for reading words and assigning line numbers to them be developed and tested separately from other aspects of the program. This could be accomplished by reading a sample file and printing out the words recognized and the lines they appeared on with no effort to avoid duplicates or associate words with more than one line.

DICTIONARY ADT

I want you to implement the dictionary ADT three ways:

- 0) use the HashDict implementation from the textbook and lab10. (should not be anything to do here)
- 1) use a binary-search tree (BST) implementation as the underlying data structure, and
- 2) use an AVL tree implementation as the underlying data structure. (you do not need to implement the dictionary `pop` operation since that would require an AVL tree delete operation. That means you need to implement the ADT `add` operation so that it replaces if the item already exists in the tree.)

Time the word-concordance application with all three dictionary ADT implementations, and complete the following table:

Word-concordance Program	Execution Time (seconds)
Dictionary ADT implemented using HashDict	
Dictionary ADT implemented using a single BST	
Dictionary ADT implemented using a single AVL tree	

DATA FILES

The stop words are in the file

http://www.cs.uni.edu/~fienup/cs052s11/homework/stop_words.txt

The textual information to be examined is in the file

<http://www.cs.uni.edu/~fienup/cs052s11/homework/hw6data.txt>

HINTS:

- 0) The word-concordance application code and each dictionary ADT should be developed and tested separately.
- 1) A built-in Python list or a linked implementation of the queue ADT may be useful to store the lines-of-occurrence for each word.

EXTRA CREDIT POSSIBILITIES:

- 0) Use a better definition of a word that allows words to contain an apostrophe or single hyphens. For example, "it's" and "end-of-line-characters" should be considered words.
- 1) implement the dictionary ADT using the a closed-address hash table (like HashDict) with each slot contains a BST instead of a linked-list
- 2) implement the dictionary ADT using the an open-address hash table (like HashTable from lab10) with linear probing.
- 3) implement the dictionary ADT using the an open-address hash table (like HashTable from lab10) with quadratic probing.

SUBMISSION

Submit ALL necessary files to run your concordance-production application using the dictionary ADTs as a single zipped file (called hw6.zip) electronically at

https://www.cs.uni.edu/~schafer/submit/which_course.cgi

Include in your zip file a "results" file containing your timing results for your all three dictionary ADTs.