

# Cosc 242 Assignment

**Due: 4pm Friday September 16<sup>th</sup> 2016**

## Overview

In this assignment you will expand and modify code written during the labs to produce two separate programs, one based around a hash table and one based around a tree data structure. These programs can be used to process two groups of words. The first group of words will be read from *stdin* and will be inserted into the data structure. The second group of words will be read from a file specified on the command line. If any word read from the file is not contained in the data structure then it should get printed to *stdout*. A moment of reflection will confirm that these programs can be used to perform a basic spell check. You can begin working on this assignment while you are completing the labs it is based on.

## Assignment groups

Everyone in the class should form groups of three to work on the this assignment. You have until the end of Wednesday August 24<sup>th</sup> to select your groups. Please send an email to [ihewson@cs.otago.ac.nz](mailto:ihewson@cs.otago.ac.nz) with the names and University user codes of your group members (e.g. stuad123, brofr456, jonsu789). If you don't select your own group then one will be chosen for you containing students who have completed a similar amount of internal assessment. You will be informed of your group members via email, so please check your University email regularly and let Iain know if there are any problems. The related labs, and all of the programming for this assignment should be done working in your groups.

You must not collaborate with, or discuss issues related to the assignment with anyone who is not a member of your group.

## Provided files

We have provided some files for you to use when completing this assignment. Apart from the two executables, they all contain comments which explain what they do and how to use them. The files can be found in the directory `/home/cshome/coursework/242/asgn-files/` and are as follows:

- `print-stats.txt` — contains extra functions which you should add to your `htable` files when completing the hash table part.

- `output-dot.txt` — contains extra functions which you should add to your `tree` files when completing the `bst/rbt` part.
- `sample-htable` and `sample-tree` — are two executables (compiled on Linux) which you can use to see if your programs are working correctly. If your programs are correct then they should behave exactly the same as the corresponding sample programs.

## Part I — Comparing hashing strategies

For this part of the assignment you will modify the hash table implementation that you developed during labs and add some new features.

Your program needs to meet the following requirements:

- The default size of the hash table should be 113.
- Add a `htable_print_entire_table()` function to print the entire contents of the hash table like `sample-htable` using the format string `"\%5d_\%5d_\%5d_\%s\n"` to print each line (spaces have been made visible so you can count them).
- You should be able to use either *linear probing* or *double hashing* as a collision resolution strategy, based on an enumerated type which gets passed to `htable_new()`. You can find more information about enumerated types in the `rbt` labs. The definition of the enumerated type in `htable.h` should look like this:

```
typedef enum hashing_e { LINEAR_P, DOUBLE_H } hashing_t;
```

- Information should be collected about the state of a hash table as it is being built. You can use this information to see the difference between the two collision resolution strategies. There are a number of different attributes of a hash table which could give us an idea of how effective our hashing strategy is. We are going to look at three.
  1. The proportion of keys which were placed in their home cells. i.e. the percentage of keys that were placed into the hash table without colliding with another key.
  2. The maximum number of collisions that occurred before placing a key into the hash table.
  3. The average number of collisions which occurred before placing a key into the hash table.

There is an easy, but ingenious, way that we can record all of these statistics as the table builds. All we need to do is add a `stats` array to our hash-table struct. Whenever we add a new key to our hash-table we just record how many collisions occurred before the key found an empty cell. (Note that unlike our `keys` and `freqs` arrays this one gets filled in sequential order, indexed by `num_keys`).

Your hash-table struct in `htable.c` should now look like this:

```
struct htablerec {
    int capacity;
    int num_keys;
```

```

    char **keys;
    int *freqs;
    int *stats;
    hashing_t method;
};

```

We have provided a couple of functions in the file `print-stats.txt` for you to use when displaying the information collected while building your hash table. Please don't alter the output of these two functions, so that your program's output is consistent with `sample-htable`.

## Part II — Comparing BSTs and RBTs

For this part of the assignment you will combine the code you have written for binary search trees and red-black trees in the labs. Once you have done this, the only difference between a bst and an rbt is that the rbt calls a `tree_fix()` function after every insertion to make sure that the tree is balanced.

This program needs to meet the following requirements:

- Create a combination `tree` ADT which can be either an ordinary bst or a balanced rbt depending on an enumerated type which gets passed to `tree_new()`.
- Make a static `type_t` variable called `tree_type` which can hold the value passed to the tree 'constructor'. The definition of the enumerated type in `tree.h` should look like this:

```
typedef enum tree_e { BST, RBT } tree_t;
```

- You don't need to implement a `tree_remove()` function since removing nodes from an rbt can be a bit tricky. In fact you might as well remove your old `bst_remove()` code since using it would break an rbt.
- Add a `tree_depth()` function which should return the length of the longest path between the root node and the furthest leaf node.
- Add a `frequency` field to your `struct tree_node` and update the frequency whenever a duplicate item is added to the tree.
- Add the two `output_dot` graph printing functions which we have provided to your tree ADT. These will enable you to visualise what your tree looks like – including the colours and frequencies. You might find it useful to use these functions when completing your tree labs. When you run your tree program with the `-o` option it should produce a file (`tree-view.dot` by default) containing a representation of your tree using the "*dot*" language<sup>1</sup>. You can produce a nice image of your tree by running the command

```
dot -Tpdf < tree-view.dot >| tree-view.pdf
```

in a terminal.

- You may have noticed that the rbt you implemented in labs doesn't ensure that the root is always black. This doesn't affect the structure of the tree at all, but you should try to fix this problem in your program.

---

<sup>1</sup>Dot is a plain text graph description language. For more information see <http://www.graphviz.org>.

## The `htable-main.c` and `tree-main.c` files

For both parts of this assignment you will need to create a main file (`htable-main.c` and `tree-main.c`) which use the corresponding data structure to perform a number of tasks. By default, words are read from `stdin` and added to your data structure before being printed out alongside their frequencies to `stdout`. This should be done by passing a `print_info` function, shown below, to either `tree_preorder` or `htable_print`.

```
static void print_info(int freq, char *word) {
    printf("%-4d %s\n", freq, word);
}
```

All memory allocated should be deallocated before your program finishes.

All words should be read using the `getword()` function from the lab book. Helper functions like `getword` and `emalloc` should be in your `mylib.c` file.

You should use the `getopt` library to help you process options given on the command line. Here is an example of how to use it:

```
const char *optstring = "ab:c";
char option;

while ((option = getopt(argc, argv, optstring)) != EOF) {
    switch (option) {
        case 'a':
            /* do something */
        case 'b':
            /* the argument after the -b is available
               in the global variable 'optarg' */
        case 'c':
            /* do something else */
        default:
            /* if an unknown option is given */
    }
}
```

You need to include `getopt.h` to use the `getopt` library. The letters listed in `optstring` are possible valid options. The colon following the letter `b` indicates that `b` takes an argument. As the options are being processed by `getopt`, they get shifted to the front of the `argv` array. After processing, the index of the first non-option argument is available in the global variable `optind`. For more information have a look at the man page for `getopt.h` (type `man 3 getopt`).

When given the command line option `-c filename`, your program will be used to process two groups of words. The first group will be read from *stdin* and put into the hash table or tree as usual. These words will function as a dictionary. The second group of words will be read from the file specified on the command line. If any word read from the file is not contained in the dictionary then it should get printed to *stdout*. Running your program with a command like this

```
./htable-asgn < dictionary.txt -c document.txt
```

or

```
./tree-asgn < dictionary.txt -c document.txt
```

should print out a list of every word from `document.txt` which is not found in `dictionary.txt`. If there is no output then `document.txt` has no misspelled words (as defined in `dictionary.txt`).

The exact behaviour of your program when given the `-c filename` option should be as follows:

1. Take `filename` to be the plain text file that we want to check the spelling in.
2. Read words from `stdin` (using the `getword()` function) and put them into our hash-table or tree. This data structure will now function as our dictionary.
3. For each word we read from `filename` (using `getword()`) check to see if it is in our dictionary. If it is then don't do anything. If it is not then print the word to `stdout`.
4. When finished checking `filename` for unknown words print timing information and unknown word count to `stderr` like this:

```
Fill time      : 0.320000
Search time    : 0.180000
Unknown words = 8690
```

When your program is given the `-h` option, or an invalid option is given, then a usage message should be printed and your program should exit.

- The hash table program that you write should respond to command line arguments as specified in this table:

| Option                    | Action performed  |
|---------------------------|---|
| <code>-c filename</code>  | Check the spelling of words in <i>filename</i> using words read from <code>stdin</code> as the dictionary. Print all unknown words to <code>stdout</code> . Print timing information and unknown word count to <code>stderr</code> . When this option is given then the <code>-p</code> option should be ignored. |
| <code>-d</code>           | Use double hashing as the collision resolution strategy (linear probing is the default).  |
| <code>-e</code>           | Display entire contents of hash table on <code>stderr</code> using the format string <code>"%5d_%5d_%5d_%s\n"</code> to display the index, frequency, stats, and the key if it exists. (Note that spaces have been made visible in the format string so you can see how many there are).                          |
| <code>-p</code>           | Print stats info using the functions provided in <code>print-stats.txt</code> instead of printing the frequencies and words   |
| <code>-s snapshots</code> | Display up to the given number of stats snapshots when given <code>-p</code> as an argument. If the table is not full then fewer snapshots will be displayed. Snapshots with 0 entries are not shown.   |
| <code>-t tablesize</code> | Use the first prime $\geq$ <i>tablesize</i> as the size of your hash table. You can assume that <i>tablesize</i> will be a number greater than 0.   |
| <code>-h</code>           | Print a help message describing how to use the program  |

- The tree version of this assignment should respond to command line arguments as specified in this table:

| Option                   | Action performed  |
|--------------------------|---|
| <code>-c filename</code> | Check the spelling of words in <i>filename</i> using words read from stdin as the dictionary. Print all unknown words to stdout. Print timing information and unknown word count to stderr. When this option is given then the <code>-d</code> and <code>-o</code> options should be ignored. |
| <code>-d</code>          | Print the depth of the tree to stdout and don't do anything else  |
| <code>-p</code>          | Print stats info using the functions provided in <code>print-stats.txt</code> instead of printing the frequencies and words   |
| <code>-f filename</code> | Write the "dot" output to <i>filename</i> instead of the default file name if <code>-o</code> is also given.  |
| <code>-o</code>          | Output a representation of the tree in "dot" form to the file 'tree-view.dot' using the functions given in <code>output-dot.txt</code> .  |
| <code>-r</code>          | Make the tree an rbt instead of the default bst.  |
| <code>-h</code>          | Print a help message describing how to use the program  |

## Submission

In order to submit your assignment files open a terminal and change into a directory which contains all of your files including the report. Type the command `asgn-submit` and press return. You should see a list of all your files printed like this:

```
htable-main.c      mylib.c           tree.c
htable.c           mylib.h           tree.h
htable.h           tree-main.c
```

If the file list looks correct then just press return, and you should see the message:

```
Submission complete.
```

- Your assignment should be submitted before 4pm on the due date.
- All group members must submit a copy of the assignment.

## Marking

This assignment is worth 15% of your final mark for Cosc 242. It is possible to get full marks. In order to do this you must write correct, well-commented code which meets the specifications.

Program marks are awarded for both implementation and style (although it should be noted that it is very bad style to have an implementation that doesn't work).

| <b>Allocation of marks</b> |    |
|----------------------------|----|
| Implementation             | 10 |
| Style/Readability          | 5  |
| Total                      | 15 |

In order to maximise your marks please take note of the following points:

- Your code should compile without warnings on the Linux lab machines using these commands:

```
gcc -W -Wall -ansi -pedantic -lm htable*.c mylib.c -o htable-asgn
gcc -W -Wall -ansi -pedantic -lm tree*.c mylib.c -o tree-asgn
```

If your code does not compile, it is considered to be a very, very bad thing!

- Your program should use good C layout as demonstrated in the lab book.
- No line should be more than 80 characters long.
- Most of your comments should be in your function headers. A function header should include:
  - A description of what the function does.
  - A description of all the parameters passed to it.
  - A description of the return value if there is one.
  - Any special notes.

Any assignments submitted after the due date and time will lose marks at a rate of 10% per day late.

You should not discuss issues pertaining to the assignment with anyone not in your group. All programs will be checked for similarity.

Part of this assignment involves you clarifying exactly what your programs are required to do. Don't make assumptions, only to find out that they were incorrect when your assignment gets marked.

You should check your University email regularly, since email clarifications may be sent to the class email list.

If you have any questions about this assignment, or the way it will be assessed, please see Iain or send an email to [ihewson@cs.otago.ac.nz](mailto:ihewson@cs.otago.ac.nz).