

COMP20003 Algorithms and Data Structures Second (Spring) Semester 2020

[Assignment 1]

Melbourne Census Dataset Information Retrieval using a Linked List

Handed out: Monday, 17 of August

Due: 11:59 PM, Sunday, 30 of August

Marks: 10 (10% of total mark)

Purpose

The purpose of this assignment is for you to:

- Improve your proficiency in C programming and your dexterity with dynamic memory allocation.
- Demonstrate understanding of a concrete data structure (linked list).
- Practice multi-file programming and improve your proficiency in using UNIX utilities.

Background

A dictionary is an abstract data type that stores and supports lookup of key, value pairs. For example, in a telephone directory, the (string) key is a person or company name, and the value is the phone number. In a student record lookup, the key would be a student ID number and the value would be a complex structure containing all the other information about the student.

Your task

In this assignment, you will create a simple dictionary based on a linked list to store information from the City of Melbourne Census of Land Use and Employment (CLUE). A user will be able to search this dictionary to retrieve information about businesses in Melbourne using the business name (key).

Your implementation will build the dictionary by reading census data from a file and inserting each property record as a node in a linked list. You will also implement a method to search for a key in the list, outputting any records that match the key. Note that keys are not guaranteed to be unique!

Dataset

The dataset comes from the City of Melbourne Open Data website, which provides a variety of data about Melbourne that you can explore and visualize online:

<https://data.melbourne.vic.gov.au/>

The dataset used in this project is a subset of the Business establishment trading name and industry classification 2018 dataset, accessed from:

<https://data.melbourne.vic.gov.au/Business/Business-establishment-trading-name-and-industry-c/vesm-c7r2>

Census year - the year in which surveying was completed (2018)

Block ID - an ID number to identify city blocks (about 606 in total)

Property ID - an ID number to identify an individual property

Base property ID - an ID number to identify a parcel of land (which may contain

multiple properties)

CLUE small area - city area name (e.g., Melbourne CBD)

Trading name - name of the business located at this property

Industry (ANZSIC4) code - numeric code to describe the industry in which the business operates

Industry (ANZSIC4) description - name of the industry corresponding to the code

x coordinate - longitude of the establishment

y coordinate - latitude of the establishment

Location - location as a (lat,long) pair (used for visualization)

The fields <CLUE small area>, <Trading name>, <Industry (ANZSIC4) description>, and <Location> are alphabetic strings of varying length. You may assume that none of these fields are more than 128 characters. The dataset is in csv format, with each field separated by a comma. Note that string fields may contain commas (the <Location> field always contains a comma); in these cases the string is enclosed in quotation marks.

For the purposes of this assignment, you may assume that the input data is well-formatted, that the input file is not empty, and that the maximum length of an input record (a single full line of the csv file) is 512 characters. This number could help you choose a reading buffer size.

The <Trading name> should be used as the key in your dictionary implementation. The other columns can be treated as the associated <data>.

Implementation Details

Your `Makefile` should produce an executable program called `dict`. This program should take two command line arguments: (1) the name of the data file used to build the dictionary, and (2) the name of an output file.

Your `dict` program should:

- Construct a linked list to store the information contained in the data file specified in the command line argument. Each record (row) should be stored in a separate Node.
- Search the linked list for records, based on keys. The keys will be read in from `stdin`, i.e. from the screen. Remember that the entries in the file do not necessarily have unique keys, so your search must locate *all* keys matching the search key, and output all the data found.
- Your program will look up each key and output the information (the data found) to the output file specified by the second command line parameter. If the key is not found in the tree, you must output the word `NOTFOUND`.

For testing, it may be convenient to create a file of keys to be searched, one per line, and redirect the input from this file. Use the UNIX operator `<` to redirect input from a file.

Examples of use:

- `dict datafile outputfile` then type in keys; or
- `dict datafile outputfile < keyfile`

Example output

This is an example of what might be output to the file after searching for two keys:

```
In a Rush Espresso -- > Census year: 2018 || Block ID: 44 || Property ID: 105956 || Base property
ID: 105956 || CLUE small area: Melbourne (CBD) || Industry (ANZSIC4) code: 4511 || Industry (ANZSIC4)
description: Cafes and Restaurants || x coordinate: 144.96174 || y coordinate: -37.81561 || Location:
(-37.81560561, 144.9617411) ||
```

```
In a Rush Espresso -- > Census year: 2018 || Block ID: 1101 || Property ID: 108973 || Base
property ID: 108973 || CLUE small area: Docklands || Industry (ANZSIC4) code: 4511 || Industry (ANZSIC4)
description: Cafes and Restaurants || x coordinate: 144.95223 || y coordinate: -37.81761 || Location:
(-37.81761044, 144.9522269) ||
```

```
Tim Hortons -- > NOTFOUND
```

The format need not be exactly as above. Only variations in whitespace/tabs are permitted.

Requirements

The following implementation requirements must be adhered to:

- You *must* write your implementation in the C programming language.
- You *must* write your code in a modular way, so that your implementation could be used in another program without extensive rewriting or copying. This means that the linked list operations are kept together in a separate .c file, with its own header (.h) file, separate from the main program.
- Your code should be easily extensible to different dictionaries. This means that the functions for insertion, search, and deletion take as arguments not only the item being inserted or a key for searching and deleting, *but also a pointer to a particular dictionary*, e.g. `insert(dict, item)`.
- Your implementation must read the input file *once only*.
- Your program should store strings in a space-efficient manner. If you are using `malloc()` to create the space for a string, remember to allow space for the final end of string `'\0'` (NULL).
- A Makefile is *not* provided for you. The Makefile should direct the compilation of your program. To use the Makefile, make sure it is in the same directory of your code, and type `make dict` to make the dictionary. You must submit your makefile with your assignment.

Hint: If you haven't used make before, try it on simple programs first. If it doesn't work, read the error messages carefully. A common problem in compiling multifile executables is in the included header files. Note also that the whitespace before the command is a tab, and not multiple spaces. It is *not* a good idea to code your program as a single file and then try to break it down into multiple files. Start by using multiple files, with minimal content, and make sure they are communicating with each other before starting more serious coding.

Resources: Programming Style (2 Marks)

Two locally-written papers containing useful guidelines on coding style and structure can be found on the *LMS Resources* → *Project Coding Guidelines*, by Peter Schachte, and below and adapted version of the *LMS Resources* → *C Programming Style*, written for Engineering Computation COMP20005 by Aidan Nagorcka-Smith. *Be aware that your programming style will be judged with 2 marks.*

```

1  /** *****
2  * C Programming Style for Engineering Computation
3  * Created by Aidan Nagorcka-Smith (aidann@student.unimelb.edu.au) 13/03/2011
4  * Definitions and includes
5  * Definitions are in UPPER_CASE
6  * Includes go before definitions
7  * Space between includes, definitions and the main function.
8  * Use definitions for any constants in your program, do not just write them
9  * in.
10 *
11 * Tabs may be set to 4-spaces or 8-spaces, depending on your editor. The code
12 * Below is ``gnu'' style. If your editor has ``bsd'' it will follow the 8-space
13 * style. Both are very standard.
14 */
15
16 /**
17 * GOOD:
18 */
19
20 #include <stdio.h>
21 #include <stdlib.h>
22 #define MAX_STRING_SIZE 1000
23 #define DEBUG 0
24 int main(int argc, char **argv) {
25     ...
26
27 /**
28 * BAD:
29 */
30
31 /* Definitions and includes are mixed up */
32 #include <stdlib.h>
33 #define MAX_STRING_SIZE 1000
34 /* Definitions are given names like variables */
35 #define debug 0
36 #include <stdio.h>
37 /* No spacing between includes, definitions and main function*/
38 int main(int argc, char **argv) {
39     ...
40
41 /** *****
42 * Variables
43 * Give them useful lower_case names or camelCase. Either is fine,
44 * as long as you are consistent and apply always the same style.
45 * Initialise them to something that makes sense.
46 */
47
48 /**
49 * GOOD: lower_case
50 */
51
52 int main(int argc, char **argv) {
53
54     int i = 0;
55     int num_fifties = 0;
56     int num_twenties = 0;
57     int num_tens = 0;
58
59     ...
60 /**
61 * GOOD: camelCase
62 */
63
64 int main(int argc, char **argv) {
65

```

```

66  int i = 0;
67  int numFifties = 0;
68  int numTwenties = 0;
69  int numTens = 0;
70
71  ...
72  /**
73  * BAD:
74  */
75
76  int main(int argc, char **argv) {
77
78      /* Variable not initialised – causes a bug because we didn't remember to
79       * set it before the loop */
80      int i;
81      /* Variable in all caps – we'll get confused between this and constants
82       */
83      int NUM_FIFTIES = 0;
84      /* Overly abbreviated variable names make things hard. */
85      int nt = 0
86
87      while (i < 10) {
88          ...
89          i++;
90      }
91
92      ...
93
94      /** *****
95       * Spacing:
96       * Space intelligently, vertically to group blocks of code that are doing a
97       * specific operation, or to separate variable declarations from other code.
98       * One tab of indentation within either a function or a loop.
99       * Spaces after commas.
100      * Space between ) and {.
101      * No space between the ** and the argv in the definition of the main
102      * function.
103      * When declaring a pointer variable or argument, you may place the asterisk
104      * adjacent to either the type or to the variable name.
105      * Lines at most 80 characters long.
106      * Closing brace goes on its own line
107      */
108
109      /**
110      * GOOD:
111      */
112
113      int main(int argc, char **argv) {
114
115          int i = 0;
116
117          for(i = 100; i >= 0; i--) {
118              if (i > 0) {
119                  printf("%d bottles of beer, take one down and pass it around,"
120                      " %d bottles of beer.\n", i, i - 1);
121              } else {
122                  printf("%d bottles of beer, take one down and pass it around."
123                      " We're empty.\n", i);
124              }
125          }
126
127          return 0;
128      }
129
130      /**
131      * BAD:

```

```

132 */
133
134 /* No space after commas
135 * Space between the ** and argv in the main function definition
136 * No space between the ) and { at the start of a function */
137 int main(int argc, char ** argv){
138     int i = 0;
139     /* No space between variable declarations and the rest of the function.
140     * No spaces around the boolean operators */
141     for(i=100;i>=0;i--) {
142         /* No indentation */
143         if (i > 0) {
144             /* Line too long */
145             printf("%d bottles of beer, take one down and pass it around, %d
146 bottles of beer.\n", i, i - 1);
147         } else {
148             /* Spacing for no good reason. */
149
150             printf("%d bottles of beer, take one down and pass it around."
151 " We're empty.\n", i);
152
153         }
154     }
155     /* Closing brace not on its own line */
156     return 0;}
157
158 /** *****
159 * Braces:
160 * Opening braces go on the same line as the loop or function name
161 * Closing braces go on their own line
162 * Closing braces go at the same indentation level as the thing they are
163 * closing
164 */
165
166 /**
167 * GOOD:
168 */
169
170 int main(int argc, char **argv) {
171
172     ...
173
174     for(...) {
175         ...
176     }
177
178     return 0;
179 }
180
181 /**
182 * BAD:
183 */
184
185 int main(int argc, char **argv) {
186
187     ...
188
189     /* Opening brace on a different line to the for loop open */
190     for(...)
191     {
192         ...
193         /* Closing brace at a different indentation to the thing it's
194         closing
195         */
196     }
197

```

```

198 /* Closing brace not on its own line. */
199     return 0;}
200
201 /** *****
202 * Commenting:
203 * Each program should have a comment explaining what it does and who created
204 * it.
205 * Also comment how to run the program, including optional command line
206 * parameters.
207 * Any interesting code should have a comment to explain itself.
208 * We should not comment obvious things – write code that documents itself
209 */
210
211 /**
212 * GOOD:
213 */
214
215 /* change.c
216 *
217 * Created by Aidan Nagorcka-Smith (aidann@student.unimelb.edu.au)
218 13/03/2011
219 *
220 * Print the number of each coin that would be needed to make up some
221 change
222 * that is input by the user
223 *
224 * To run the program type:
225 * ./coins --num_coins 5 --shape_coins trapezoid --output blabla.txt
226 *
227 * To see all the input parameters, type:
228 * ./coins --help
229 * Options::
230 * --help                Show help message
231 * --num_coins arg       Input number of coins
232 * --shape_coins arg     Input coins shape
233 * --bound arg (=1)      Max bound on xxx, default value 1
234 * --output arg          Output solution file
235 *
236 */
237
238 int main(int argc, char **argv) {
239
240     int input_change = 0;
241
242     printf("Please input the value of the change (0–99 cents
243 inclusive):\n");
244     scanf("%d", &input_change);
245     printf("\n");
246
247     // Valid change values are 0–99 inclusive.
248     if(input_change < 0 || input_change > 99) {
249         printf("Input not in the range 0–99.\n")
250     }
251
252     ...
253
254 /**
255 * BAD:
256 */
257
258 /* No explanation of what the program is doing */
259 int main(int argc, char **argv) {
260
261     /* Commenting obvious things */
262     /* Create a int variable called input_change to store the input from
263 the

```

```

264  * user. */
265  int input_change;
266
267  ...
268
269  /** *****
270  * Code structure:
271  * Fail fast – input checks should happen first, then do the computation.
272  * Structure the code so that all error handling happens in an easy to read
273  * location
274  */
275
276  /**
277  * GOOD:
278  */
279  if (input_is_bad) {
280      printf("Error: Input was not valid. Exiting.\n");
281      exit(EXIT_FAILURE);
282  }
283
284  /* Do computations here */
285  ...
286
287  /**
288  * BAD:
289  */
290
291  if (input_is_good) {
292      /* lots of computation here, pushing the else part off the screen.
293      */
294      ...
295  } else {
296      fprintf(stderr, "Error: Input was not valid. Exiting.\n");
297      exit(EXIT_FAILURE);
298  }

```

Additional Support

Your tutors will be available to help with your assignment during the scheduled workshop times. Questions related to the assignment may be posted on the Piazza forum, using the folder tag *assignment1* for new posts. You should feel free to answer other students' questions if you are confident of your skills.

A tutor will check the Discussion Forum regularly, and answer some questions, but be aware that for some questions you will just need to use your judgment and document your thinking. For example, a question like, "How much data should I use for the experiments?", will not be answered; you must try out different data and see what makes sense.

Submission

Your C code files (including your `Makefile` and any other files needed to run your code) should be submitted through the LMS under **Assignment 1: Code** in the **Assignments** tab.

Your programs *must* compile and run correctly on JupyterHub. You may have developed your program in another environment, but it still *must* run on JupyterHub at submission time. For this reason, and because there are often small, but significant, differences between compilers, it is suggested that if you are working in a different environment, you upload and test your code on JupyterHub at reasonably frequent intervals.

A common reason for programs not to compile is that a file has been inadvertently omitted from the submission. Please check your submission, and resubmit all files if necessary.

Assessment

There are a total of 10 marks given for this assignment.

Your C program will be marked on the basis of accuracy, readability, and good C programming structure, safety and style, including documentation (2 marks). Safety refers to checking whether opening a file returns something, whether mallocs do their job, etc. The documentation should explain all major design decisions, and should be formatted so that it does not interfere with reading the code. As much as possible, try to make your code self-documenting, by choosing descriptive variable names. The remainder of the marks will be based on the correct functioning of your submission.

Plagiarism

This is an individual assignment. The work must be your own.

While you may discuss your program development, coding problems and experimentation with your classmates, you must not share files, as this is considered plagiarism.

If you refer to published work in the discussion of your experiments, be sure to include a citation to the publication or the web link.

“Borrowing” of someone else’s code without acknowledgment is plagiarism. Plagiarism is considered a serious offense at the University of Melbourne. You should read the University code on Academic integrity and details on plagiarism. Make sure you are not plagiarizing, intentionally or unintentionally.

You are also advised that there will be a C programming component (on paper, not on a computer) in the final examination. Students who do not program their own assignments will be at a disadvantage for this part of the examination.

Late policy

The late penalty is 10% of the available marks for that project for each day (or part thereof) overdue. Requests for extensions on medical grounds will need to be supported by a medical certificate. Any request received less than 48 hours before the assessment date (or after the date!) will generally not be accepted except in the most extreme circumstances. In general, extensions will not be granted if the interruption covers less than 10% of the project duration. Remember that departmental servers are often heavily loaded near project deadlines, and unexpected outages can occur; these will not be considered as grounds for an extension.

Students who experience difficulties due to personal circumstances are encouraged to make use of the appropriate University student support services, and to contact the lecturer, at the earliest opportunity.

Finally, **we are here to help!** There is information about getting help in this subject on the LMS. Frequently asked questions about the project will be answered in Piazza.