

Deep Shah
Shashank Seeram
CS 214
PA 4: Indexer

Program Design:

The index uses tokenizer, sorted-list and sorted-list main to index all files given as input. The tokenizer is responsible for creating tokens from a file stream, and also making sure that every token is correctly formatted (uppercase converted to lowercase) and doesn't contain separator characters. This is done by using a buffer to read in large chunks of the file at once and the parsing through the buffer to get the token when asked.

The sorted-list main has nodes with a token and the corresponding sorted list, which will store the files the token appears in and that token's frequency. The sorted-list contains nodes which hold a filename (path from input directory) and the frequency. Both of the lists are linked-lists with a front and each of the nodes are holding the next value. The front is stored in the sorted-list main or sorted-list structs, which allows the program to perform various search and reorganizing functions.

The Index, uses the other 3 programs to create the lists. It created the main sorted list by creating tokens using tokenizer and then the inserting (or searching) for these tokens in the list. When the node is created or found it searches through the normal sorted list associated with that main sorted list node for the filename, and if found it will increment the frequency count. Otherwise it will create the proper node and set 1 for the frequency.

The indexer contains the function indexDir which can recursively index directories if they happen to have subdirectories of their own. It also calls indexFile if it happens across a file. IndexFile uses the tokenizer to actually build the required parts of the main sorted list tree, which indexDir is more responsible for the file input/output and maintaining path names, all the while checking for errors. When everything is created, a function names output in indexer creates the output file, after which the structure is freed from memory, and the program closes.

Running time:

The worst case running time for this program is fairly simple if one is to ignore all the details that only lead to coefficients that are removed anyways. When you consider the main sorted list structure, there is a node for every unique token. Each token then has a sorted list which has a node

for each file that token appears in. Assume you have m unique tokens and k files. Also assume there are n tokens altogether. The program searches through the list first to find the token, then through the token's sorted list to find the file entry. With that in mind, in the worst case the program must run through the entire main list of m nodes, and then run through the sub-list of whatever token it finds for at most k nodes, since there are only k files. In between the program performs memory allocations and pointer reassignments, and these will amount to cn , where c is some constant. So in total, the program, in the worst case, performs $m \cdot k$ operations n times, plus $cn \Rightarrow nmk + cn \Rightarrow n(mk + c)$. Since this is the worst case and we are computing $O()$, the cn is irrelevant in the face of the nmk , and so $O(n)$ in the absolute worst case, since mk is also a constant. On the other hand, this worst case is extremely unlikely to happen, as the computations overestimate by a lot in order to achieve $O()$. One thing to note, is that for enormous files, n will be much larger than $m \cdot k$, so $O(n)$ is reasonable, but if somehow $m \cdot k$ is larger than n then $O()$ will be closer to $O(n^2)$ or simply $O(nmk)$.

Memory Usage:

Using the variables from the running time analysis, n tokens, m unique tokens, and k files, there will be at most m allocated nodes for the main sorted list, each of which can have a max of k allocated nodes, one for each file. Counting each node as a single unit of memory (not an actual bit or byte, just a unit for the sake of the analysis), this will take at most $m \cdot k$ nodes. Each node holds a certain amount of data, so the it would be $(m \cdot c)(k \cdot i)$ where c is how many bytes a single main sorted list node has, and i how many bytes a sorted list node has. So the final memory usage is $mkci$ bytes of memory.