

CSCI 2270 Data Structures Recitation 14

Instructors: Hoenigman/Zagrodzki/Zietz

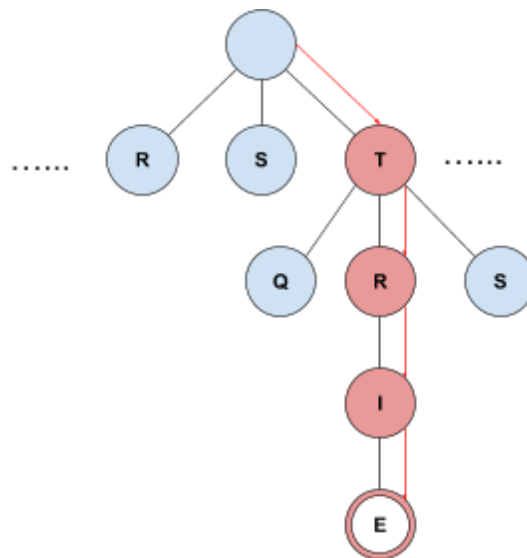
Trie

Learning Objectives:

- Understanding a new data structure Trie and its applications.

1. Trie

Recall how you look up for a word in an English dictionary (of course the paperback). Say you wanna find a word “trie”. You first go to the page where the first letter, “t”, starts. Among those words beginning with “t”, you look for those who have a letter “r” following “t”, and so forth. If we visualize this process, we can present this as a “tree” structure as follows:



We notice some properties of this tree --- the root is empty, and all the other nodes are represented by a letter. Thus, looking for a word is simply a search problem of the tree. Starting from the root of the tree, we keep going down until we reach the end of the word. Then, all the nodes along the path from root to the destination become the spelling of the word.

One thing nice about trie is that we learned about BST of general tree structures in the class, but we haven't seen how this structure can be used to useful applications. Trie, in this sense, uses a tree structure, and can complete simple information retrieval tasks.

1.1 Structure

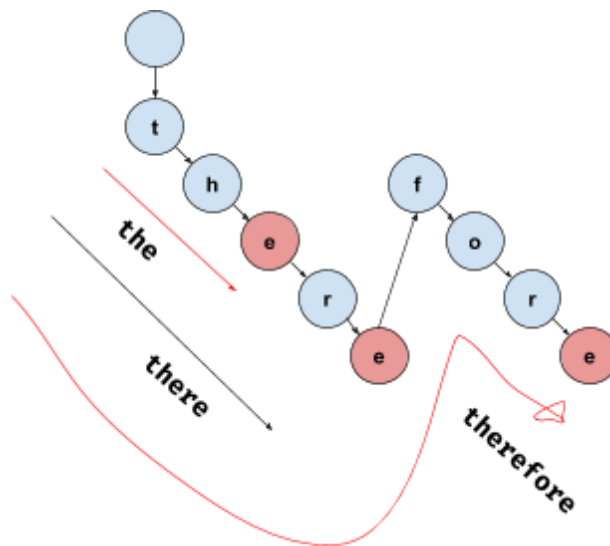
As usual, we first define the structure of each node in this tree. Recall that in BST, we have a field called value, and pointers to its left and right children. For trie, there are more than two possible children (26 at most in English). So we arrange them into an array, and define the struct as follows:

```
struct TrieNode
{
    struct TrieNode *children[ALPHABET_SIZE];
    bool isEndOfWord;
};
```

Note that in the structure, `ALPHABET_SIZE` can be different. For English, we set it to 26. Another field we need in this structure is a boolean called `isEndOfWord` (not “isEndofWorld”). Here, we present a simple example of the usage of this variable. Consider that we have three words:

the, there, therefore

All these three words are in the path from the root to the leaf which ends with “therefore”. If we want to find the word “the”, we can simply stop at the first “e” along the path. At that point, all the letters from the root to this node become a word “the”. Thus, the node “e” could mark as the end of the word, and that’s why we need a boolean variable.



Except the basic definition above, we can also associate each word with a count, indicating how many times it appears given a context. **The first task of this recitation is to modify the structure above so that the number of a word is recorded.**

1.2 Inserting a word

Before we read in any word, the trie has only one empty node, which is also the root. As we read in more words, we will probably need to add new nodes. If the word is already in the trie, we should record the time this word appears in the context.

Here's a start code, and you need to finish this code in the recitation.

```
void insert(struct TrieNode* root, string word) {  
  
    /* Start from the root */  
    struct TrieNode *pt = root;  
  
    /* Iterate over all letters in the word */  
    for (int i = 0; i < word.length(); i++) {  
        /* If the letter is in the trie as a node,  
        we can keep going down;  
        otherwise, we need to create a node for  
        this letter.  
        */  
    }  
  
    /* Once we finish iterating over the word,  
    we can mark the node of the last letter  
    as "end of the word".  
    */  
    pt->isEndOfWord = true;  
}
```

1.3 Search for a word

Searching for a word is similar to inserting a word. In fact, the common part in both searching and inserting is that at each node `i`, we need to look at the next letter `j` from node `i`'s children, to see if a node `j` exists in that pointer list to form a bi-gram "`ij`".

Following is a starter code.

```
int search(struct TrieNode* root, string word) {
    struct TrieNode* pt = root;

    /*
     * Starting from the root,
     * we iterate over all the letters of this word.
     */
    for (int i = 0; i < word.length(); i++){

    }

    /*
     * At this point, we finish iterating the word.
     * Does that mean we've found this word?
     * What additional condition should we check?
     */

}
```

2. Recitation Task

In summary, you need to finish the following tasks:

Programming Exercise:

1. Modify the structure of trie nodes so that we can record the number of times a word appear.
2. Complete inserting and search function.
3. Read in a text file provided on Moodle, and create a trie. Then when you search for a word, it should output the number of times that word appeared in the text file. Otherwise, output -1.

Analysis:

1. Implement a simple solution where we use array to store words. Then compare the difference of running time of searching a word using a trie and an array.

