

Name: \_\_\_\_\_

Name: \_\_\_\_\_

## Lab - Huffman Encoding for Text Compression

A Huffman code is a variable length code that minimizes the length of the message. Given a string  $X$  defined over some alphabet, such as the Unicode character set, the goal is to efficiently encode  $X$  into a small binary string (using only the characters 0 and 1). The basic idea is that symbols that are frequently encountered in the message will be represented by short codes. On the other hand, infrequent symbols will have longer codes.

Text compression is useful in any situation where we wish to reduce bandwidth for digital communications, to minimize the time needed to transmit the text. Likewise, text compression is useful for storing large documents more efficiently, to allow a fixed-capacity storage device to contain as many documents as possible.

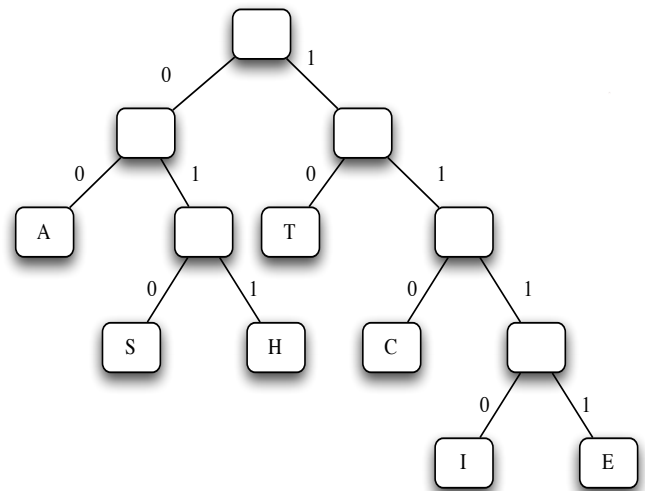
A binary tree represented the encoding of a string. Each edge of the tree represents a bit in a code-word, with an edge to a left child representing a “0” and an edge to a right child representing a ‘1’. Each leaf is associated with a specific character., and the code-word for that character is defined by the sequence of bits associated with the edges in the path from the root of the tree to the leaf of the character.

The tree shown here was constructed from the message:

THE CAT IS THAT CAT

The code for each letter is read by Write the code for each letter in the message:

SYMBOL	CODE
A	00
C	
E	
H	
I	
S	
T	



## Getting the Frequencies

The structure of the Huffman tree must depend on the *frequencies* of the letters, where we have, for each character  $c$ , a count  $f(c)$  of the number of times  $c$  appears in the string  $X$ . Therefore, the first task is to compile a count of each letter in the message. Consider the following message:

CHEESES CHEESES

T CLEESE SELLS THE CHEESES

How many times does each character appear in the message?

Complete the table.

SYMBOL	COUNT
L	3
C	
E	
H	
S	
T	

Design an **algorithm** for a method that computes the frequency count for the letters in a message. The message will be stored in a file, so the method will have as input the name of the file. It will return an integer array of size 256 containing the frequencies of each letter of the alphabet. The method will iterate to get the characters one at a time from the file. **While the code is given here for the method header, for this lab you are to only provide the algorithm, not the code.**

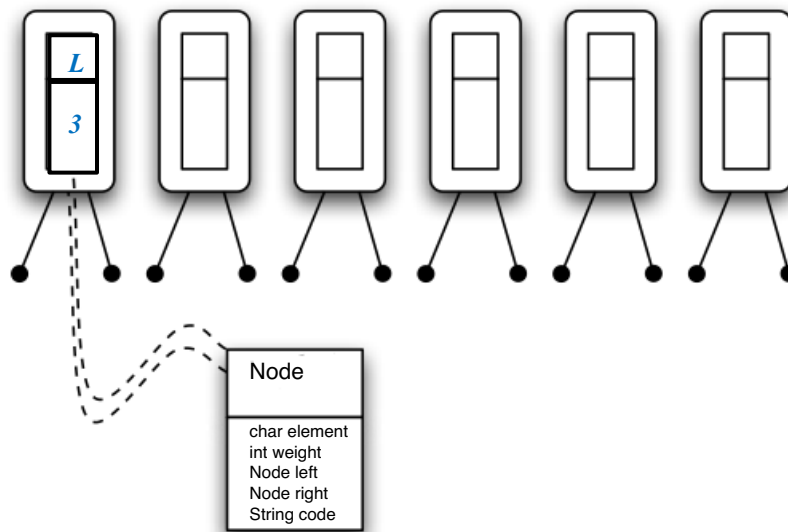
```
/** Get the frequency of the characters */  
public static int[] getCharacterFrequency(String filename) {
```

## Creating the Initial Forest of Trees

Once the counts have been done, the next step is to create a collection (forest) of Huffman trees. To start there will be one tree for each possible symbol. Each of the trees will have a single node, where the data portion contains the character and the frequency (designated as “weight”) of the character. The frequency will be the count for that character found in the character frequency array.

Fill in the initial trees for the characters and counts that were found in the previous section.

Forest of Huffman trees



```
public class Node {  
    char element;        // the character  
    int weight;          // weight of this subtree  
    Node left;           // left subtree  
    Node right;          // right subtree  
    String code = "";    // code from the root  
  
    /** Create an empty node */  
    public Node() {  
    }  
  
    /** Create a node with the weight and character */  
    public Node(int weight, char element) {  
        this.weight = weight;  
        this.element = element;  
    }  
}
```

Design the ***algorithm*** for a method that creates the forest of trees given an array of the counts. These will be the leaves in the Huffman tree. The method header is:

```
/** Get a Huffman tree from the codes */  
public static Comparable<T> getHuffmanTree(int[] counts) {
```

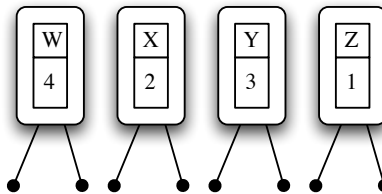
## Building the Code Tree

At this point an iterative process will be applied to construct a single Huffman tree out of the forest of trees.

Let's look at a small example. At each step the two trees with the smallest frequency will be removed from the forest and replaced by a single tree. The new tree will have the tree with the smallest frequency on the left and the second smallest tree will be on the right. The **code** instance variable in the root will be the concatenation of the **code** from the two smallest trees. The **weight** in the root will be the sum of the **weights** of the two smaller trees.

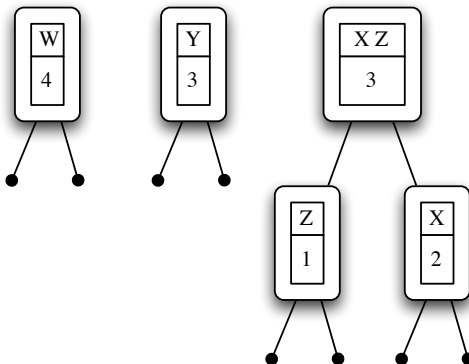
Step 1: In the following example, the two smallest trees in the forest are the Z and the X.

Initial forest of Huffman trees



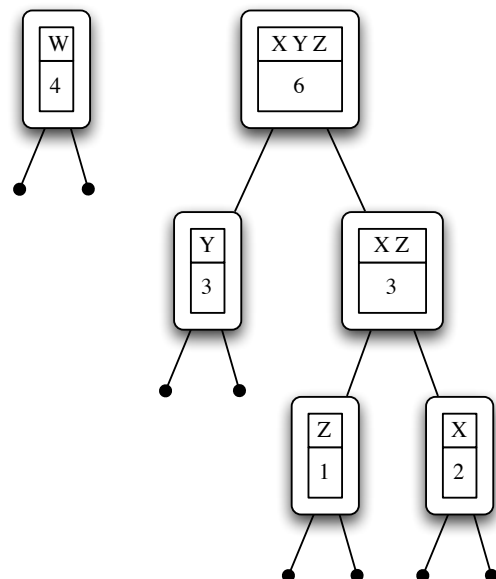
Step 2: Those two trees will be removed and then used to form the tree XZ, shown below. Since the count for Z is less than that of X, Z will go on the left. (This is an arbitrary choice that will make no difference in the effectiveness of the code, but it can make it easier to check that the construction process is working properly.)

Forest of Huffman trees after combining X and Z



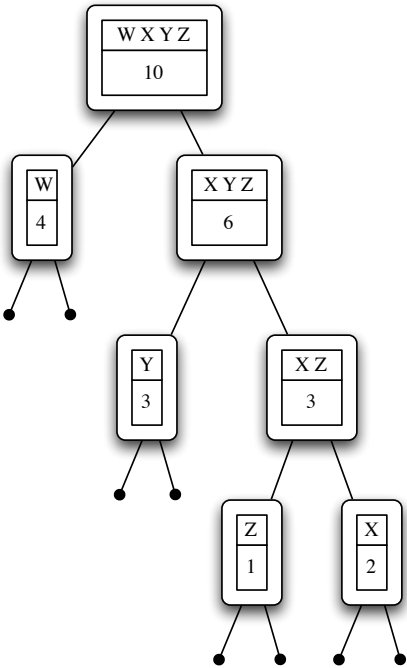
Step 3: Now the two smallest trees are the Y and XZ, both with a count of 3. Which one should go on the left? Pick either one. If you need to pick among several trees of the same frequency, any one of them can be safely chosen.

Forest of Huffman trees after combining XZ and Y



Step 4: The final step is to combine the last two trees in the forest.

Final forest of Huffman trees



**Practice**

Looking back at the example forest of trees constructed on page 3, show the trees combined at each step and the final Huffman tree resulting from the initial forest that was constructed in the previous section.

Trees combined in step 1:

Trees combined in step 2:



Trees combined in step 3:



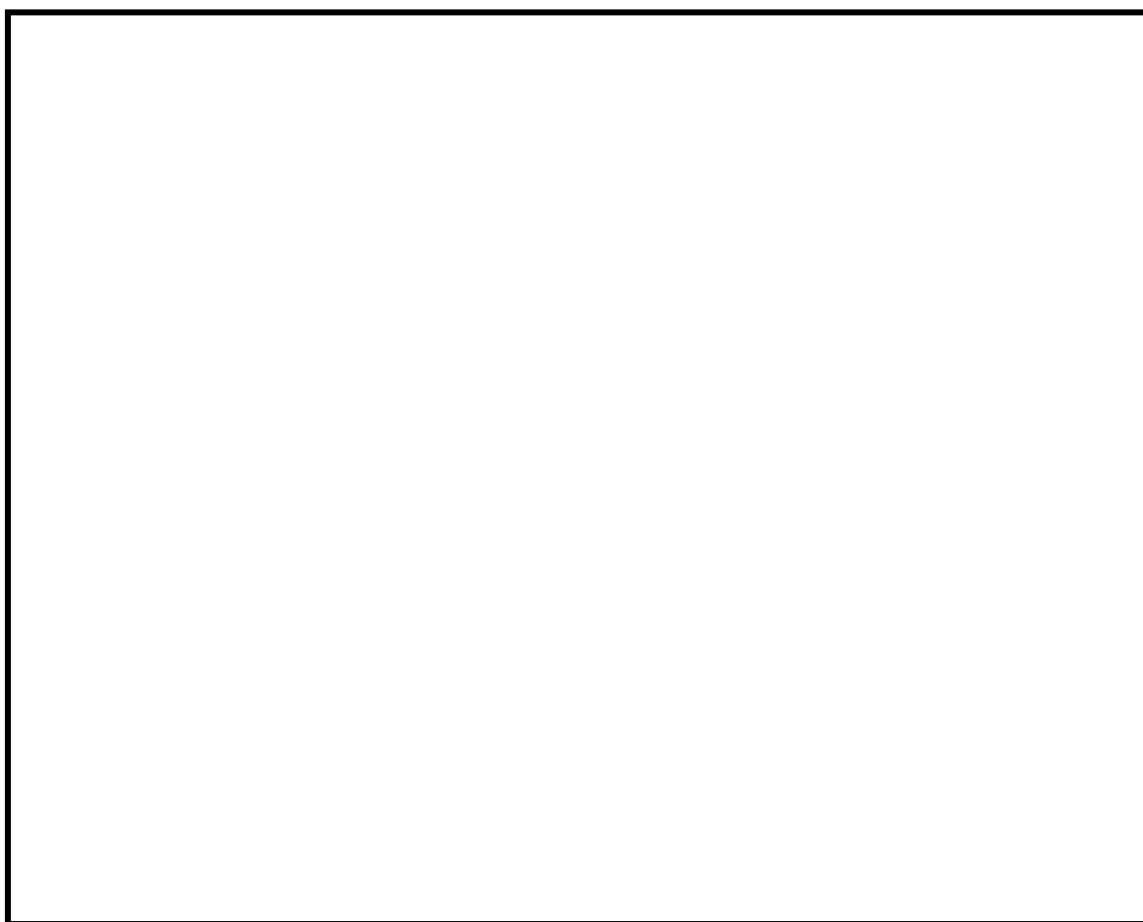
Trees combined in step 4:



Trees combined in step 5:



Final tree:



If it is correct, the message 111011110101100 will either decode to CHEST or HCEST.

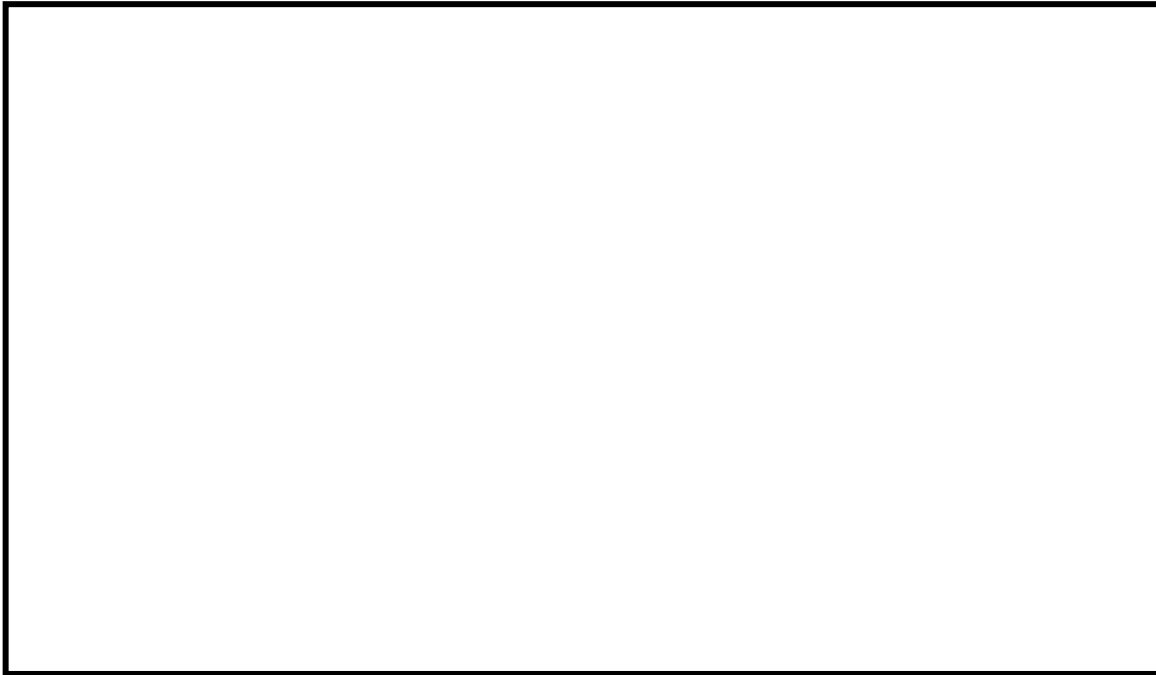


## Encoding the Message

Encode the message, CHEESE, using the final Huffman tree from the previous section.



Design an algorithm for encoding a single character.



Develop an algorithm for encoding a file and displaying the result.

