

CSE 143, Summer 2012

Programming Assignment #7: Huffman Coding (40 points)

Due Friday, August 10, 2012, 11:30 PM

No submissions for this assignment will be accepted after Monday, August 13, at 11:30pm.

This program focuses on binary trees and priority queues. Turn in files named `HuffmanTree.java`, `HuffmanNode.java`, `secretmessage.huf`, and `secretmessage.huf.counts` from the Homework section of the web site. You will need support files `HuffMain.java`, `Bit*.java`, and input files from the course web page.

Huffman Coding:

Huffman coding is an algorithm devised by David A. Huffman of MIT in 1952 for compressing text data to make a file occupy a smaller number of bytes. This relatively simple compression algorithm is powerful enough that variations of it are still used today in computer networks, fax machines, modems, HDTV, and other areas.

Normally text data is stored in a standard format of 8 bits per character, commonly using an encoding called ASCII that maps every character to a binary integer value from 0-255. The idea of Huffman coding is to abandon the rigid 8-bits-per-character requirement and use different-length binary encodings for different characters. The advantage of doing this is that if a character occurs frequently in the file, such as the letter 'e', it could be given a shorter encoding (fewer bits), making the file smaller. The tradeoff is that some characters may need to use encodings that are longer than 8 bits, but this is reserved for characters that occur infrequently, so the extra cost is worth it.

The table below compares ASCII values of various characters to possible Huffman encodings for the text of Shakespeare's *Hamlet*. Frequent characters such as space and 'e' have short encodings, while rarer ones like 'z' have longer ones.

Character	ASCII value	ASCII (binary)	Huffman (binary)
' '	32	00100000	10
'a'	97	01100001	0001
'b'	98	01100010	0111010
'c'	99	01100011	001100
'e'	101	01100101	1100
'z'	122	01111010	00100011010

The steps involved in Huffman coding a given text source file into a destination compressed file are the following:

1. Examine the source file's contents and count the number of occurrences of each character.
2. Place each character and its frequency (count of occurrences) into a sorted "priority" queue.
3. Convert the contents of this priority queue into a binary tree with a particular structure.
4. Traverse the tree to discover the binary encodings of each character.
5. Re-examine the source file's contents, and for each character, output the encoded binary version of that character to the destination file.

Encoding a File:

For example, suppose we have a file named `example.txt` with the following contents:

ab ab cab

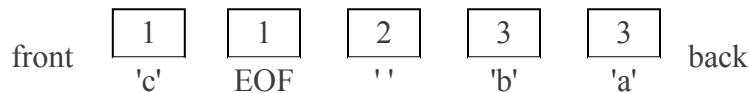
In the original file, this text occupies 10 bytes (80 bits) of data. The 10th is a special "end-of-file" (EOF) byte.

byte	1	2	3	4	5	6	7	8	9	10
char	'a'	'b'	' '	'a'	'b'	' '	'c'	'a'	'b'	EOF
ASCII	97	98	32	97	98	32	99	97	98	256
binary	0110000 1	0110001 0	0010000 0	0110000 1	0110001 0	0010000 0	0110001 1	0110000 1	0110001 0	N/A

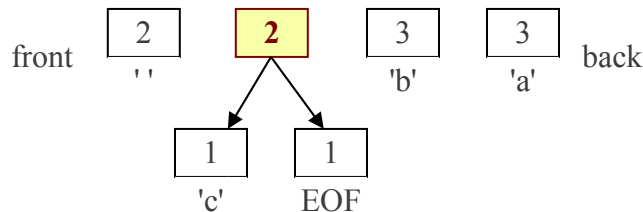
In Step 1 of Huffman's algorithm, a count of each character is computed. (In this assignment, our provided client program does this part for you, so you don't need to do it yourself.) The counts are represented as a map:

```
{ ' ' =2, 'a'=3, 'b'=3, 'c'=1, EOF=1 }
```

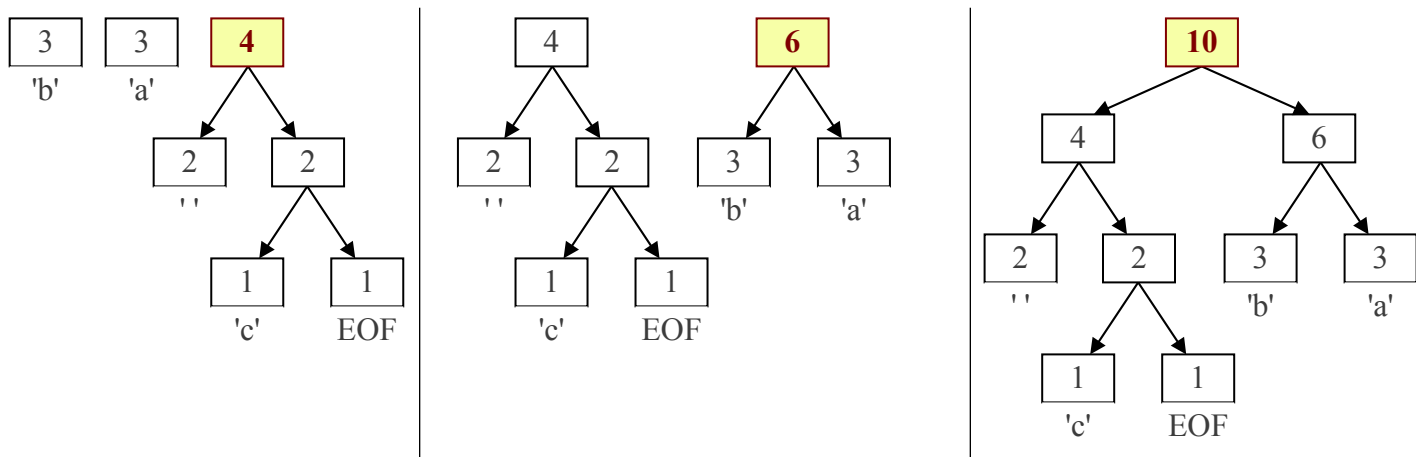
Step 2 of the algorithm places these counts into binary tree nodes, each storing a character and a count of its occurrences. The nodes are put into a priority queue, which keeps them in sorted order with smaller counts at the front of the queue. (The priority queue is somewhat arbitrary in how it breaks ties, such as 'c' being before EOF and 'b' being before 'a').



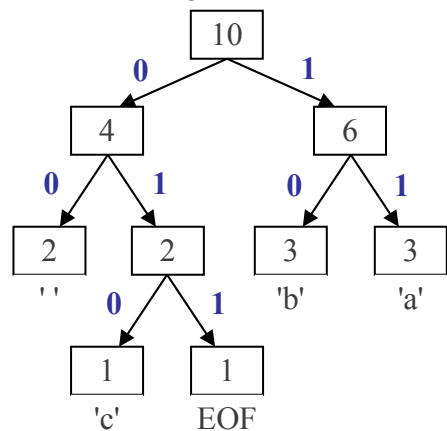
Now the algorithm repeatedly removes the two nodes from the front of the queue (the two with the smallest frequencies) and joins them into a new node whose frequency is their sum. The two nodes are placed as children of the new node; the first removed becomes the left child, and the second the right. The new node is re-inserted into the queue in sorted order:



This process is repeated until the queue contains only one binary tree node with all the others as its children. This will be the root of our finished Huffman tree. The following diagram shows this process:



Notice that the nodes with low frequencies end up far down in the tree, and nodes with high frequencies end up near the root of the tree. This structure can be used to create an efficient encoding. The Huffman code is derived from this tree by thinking of each left branch as a bit value of 0 and each right branch as a bit value of 1:



The code for each character can be determined by traversing the tree. To reach ' ' we go left twice from the root, so the code for ' ' is 00. The code for 'c' is 010, the code for EOF is 011, the code for 'b' is 10 and the code for 'a' is 11. By traversing the tree, we can produce a map from characters to their binary representations. For this tree, it would be: {' '=00, 'a'=11, 'b'=10, 'c'=010, EOF=011}

Using this map, we can encode the file into a shorter binary representation. The text ab ab cab would be encoded as:

char	'a'	'b'	' '	'a'	'b'	' '	'c'	'a'	'b'	EOF
------	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

binary	11	10	00	11	10	00	010	11	10	011
---------------	----	----	----	----	----	----	-----	----	----	-----

The overall encoded contents of the file are 1110001110000101110011, which is 22 bits, or almost 3 bytes, compared to the original file which was 10 bytes. (Many Huffman-encoded text files compress to about half their original size.)

byte	1		2		3	
char	a	b	a	b	c	a
binary	11	10	00	11	10	00

Since the character encodings have different lengths, often the length of a Huffman-encoded file does not come out to an exact multiple of 8 bits. Files are stored as sequences of whole bytes, so in cases like this the remaining digits of the last bit are filled with 0s. You do not need to worry about this in the assignment; it is part of the underlying file system.

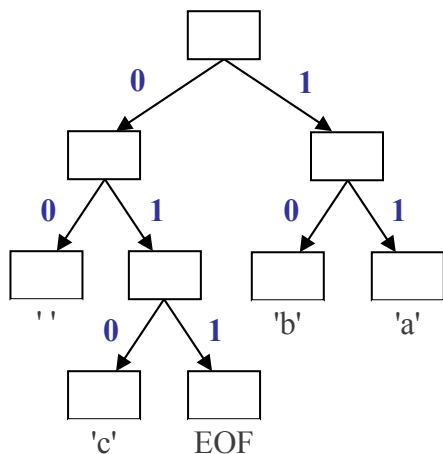
It might worry you that the characters are stored without any delimiters between them, since their encodings can be different lengths and characters can cross byte boundaries, as with 'a' at the end of the second byte above. But this will not cause problems in decoding the compressed file, because Huffman encodings have a *prefix property* where no character's encoding can ever occur as the start of another's encoding. This is important when you decode the file later.

Decoding a File:

You can use a Huffman tree to decode text that was compressed with its encodings. The decoding algorithm is to read each bit from the file, one at a time, and use this bit to traverse the Huffman tree. If the bit is a 0, you move left in the tree. If the bit is 1, you move right. You do this until you hit a leaf node. Leaf nodes represent characters, so once you reach a leaf, you output that character. For example, suppose we are asked to decode a file containing the following bits:

1011010001101011011

Using the Huffman tree, we walk from the root until we find characters, then we output them and go back to the root.



Implementation Details:

In this assignment you will create a class `HuffmanTree` to represent the overall tree of character frequencies drawn on the previous page. You will also create a class `HuffmanNode` where each node stores information about one character. The contents of the `HuffmanNode` class are up to you, but it should not perform a large share of the overall algorithm.

Your `HuffmanTree` class must have the following public constructor and methods:

`public HuffmanTree(Map<Character, Integer> counts)`

In this constructor you are passed a map from characters (`char`) to the number of occurrences of that character (`int`). You should use this map to build your Huffman tree using a priority queue (`PriorityQueue`) as previously described.

`public Map<Character, String> createEncodings()`

In this method you should traverse your Huffman tree and produce a mapping from each character in the tree to its encoded binary representation as a `String`. For the example shown on the previous pages, the map is the following:

```
{' '=00, 'a'=11, 'b'=10, 'c'=010, EOF=011}
```

The client may want to modify the map you return, so you shouldn't return any map that you expect to remain untouched.

`public void compress(List<Character> input, BitOutputStream output)`

In this method you should read the text data from the given list of characters and use your Huffman encodings to write a Huffman-compressed version of this data to the given bit output file stream. You will use a `BitOutputStream` object to help you write the binary output one bit at a time, as described below.

`public void decompress(BitInputStream input, PrintStream output)`

In this method you should read the compressed binary data from the given bit input file stream and use your Huffman tree to write a decompressed text version of this data to the given output file stream. You may assume that all characters in the input file were represented in the map of counts passed to your tree's constructor. You will use a `BitInputStream` object to help you read the binary input one bit at a time, as described below.

You may have additional methods, so long as they are `private`. Note that the methods might be called in any order.

If a parameter passed to any method above is `null`, you should throw an `IllegalArgumentException`. If your object is asked to compress/decompress an empty file, the result should also be an empty file. Methods that traverse your tree should be implemented recursively whenever practical.

Your classes will interact with a provided `HuffMain` client program that prompts the user for a file name to compress. When compressing a file, the `HuffMain` program also saves a separate file with a `.count` extension that stores the counts of every character in the original file. This is so that the `HuffmanTree` can be reconstructed later when decompressing.

BitOutputStream and BitInputStream:

To compress/decompress files, you will want to read and write binary data one bit at a time. Java's built-in input/output streams read an entire byte at a time, which makes it difficult to examine each bit. Therefore we are providing you with `BitOutputStream` and `BitInputStream` classes with `writeBit` and `readBit` methods to make it easier. The constructors and methods of these classes throw exceptions if something fails during the input/output process.


BitOutputStream Method	Description
<code>public void writeBit(int bit)</code>	writes a single 0 or 1 bit to the output
<code>public void writeBits(String bits)</code>	treats each character of the given string as a bit ('0' or '1') and writes each of those bits to the output
<code>public void close()</code>	stops writing (important to call this to ensure data is saved)
BitInputStream Method	Description
<code>public int readBit()</code>	reads a single 0 or 1 bit from input; returns -1 at end of file
<code>public boolean hasNextBit()</code>	returns <code>true</code> if more bits remain to be read, else <code>false</code>
<code>public void close()</code>	stops reading

Development Strategy and Hints:

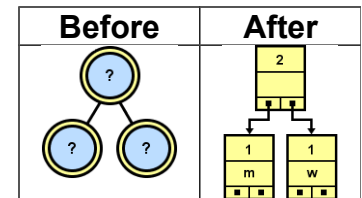
We suggest that you first focus on building your Huffman tree properly from the given map of character counts. Then work on creating the map of `char` \rightarrow `String` encodings from your tree. Then work on using your encodings to compress files, and lastly work on trying to decompress a file that you have previously compressed.

For your nodes to be able to be stored in a priority queue, the queue needs to know how to sort them. Therefore your node class must implement the `Comparable` interface as discussed in lecture and section. Nodes should be compared by character frequency, where a character that occurs fewer times is "less than" one that occurs more often. If two nodes have the same number of occurrences, they are considered "equal" in this context.

Consider writing a `toString` method in your `HuffmanNode` so you can easily print nodes or priority queues of them. Note that if you `println` a priority queue, it does not necessarily show its elements in the order they would be returned.

You can examine your binary tree in jGRASP's debugger. Set a breakpoint and drag your tree from the left to the right side of the program. The debugger initially will not know how to display your node data. To fix this, from the "Viewer" window, click the "wrench" icon . In the "Presentation View Configuration" box, type an expression into the "Value Expressions" box. To see multiple fields, use the following pattern:

```
_node_.field1#_node_.field2
```



It can be difficult to tell whether you have compressed/decompressed a file correctly. If you open a Huffman-compressed binary file in a text editor, the appearance will be gibberish (because the text editor will try to interpret the bytes as ASCII encodings, which is not the way the data is stored). While developing your program, it can be helpful to write out each 0 or 1 as an entire character (byte) rather than as a bit. This defeats the purpose of compression, because the "compressed" file is actually larger than the original, but it can help you see whether the 0s and 1s are what you expect.

To write out your 0s and 1s as entire bytes instead of as bits, you can simply set the `DEBUG` flag to true at the top of the `HuffMain` client.

The provided `HuffMain` client program can compress any text file. We suggest you start with a very small input file such as the example shown in this document, and work your way up to larger files once that works.

A common error involves not closing the `BitOutputStream` when compressing a file. If you see extra characters displayed at the end of a decompressed file, make sure that you have closed your `BitOutputStream`.

Creative Aspect (`secretmessage.huf` and `secretmessage.huf.counts`):

Along with your program you should turn in files named `secretmessage.huf` and `secretmessage.huf.counts` that represent a "secret" compressed message from you to your TA, and its counts file. The message can be anything you want, as long as it is not offensive. Your TA will decompress your message with your tree and read it while grading.

Style Guidelines and Grading:

Part of your grade will come from appropriately utilizing binary trees and recursion to implement your Huffman tree as described previously. We will also grade on the elegance of your recursive algorithms; don't create special cases in your recursive code if they are not necessary or repeat cases already handled. Redundancy is another major grading focus; some methods are similar in behavior or based off of each other's behavior. You should avoid repeated logic as much as possible. Your class may have other methods besides those specified, but any other methods you add should be `private`.

You should follow good general style guidelines such as: making fields `private` and avoiding unnecessary fields; declaring collection variables using interface types; appropriately using control structures like loops and `if/else`; properly using indentation, good variable names and types; and not having any lines of code longer than 100 characters.

Comment your code descriptively in your own words at the top of your class, each method, and on complex sections of your code. Comments should explain each method's behavior, parameters, return, pre/post-conditions, and exceptions. For reference, our `HuffManTree` class is around 125 lines long (60 "substantive") including comments and blank lines.