

Thanks to Owen Astrachan (Duke) and Julie Zelenski for creating this. Modifications by Keith Schwarz, Stuart Reges, Marty Stepp.

感谢 Owen Astrachan (杜克大学) 和 Julie Zelenski 制作此文档。修改: Keith Schwarz, Stuart Reges, Marty Stepp.

This assignment focuses on binary trees and priority queues. We provide you with several other support files, but you should not modify them. For example, we provide you with a huffmanmain.cpp that contains the program's overall text menu system; you must implement the functions it calls to perform various file compression / decompression operations.

本作业的重点是二叉树和优先队列。我们为您提供了其他几个支持文件，但您不应修改它们。例如，我们提供的 huffmanmain.cpp 包含程序的整个文本菜单系统；您必须实现它调用的函数，以执行各种文件压缩/解压缩操作。

Due: Monday, November 11th, 2019 at 11:59pm

截止时间: 11 月 11th, 2019 日星期一 11:59pm

Huffman Encoding 赫夫曼编码

Huffman encoding 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, HDTV, and other applications. Normally text data is stored in a standard format of 8 bits per character using an encoding called ASCII that maps every character to a binary integer value from 0-255.

哈夫曼编码是麻省理工学院的戴维·A·哈夫曼 (David A. Huffman) 于 1952 年设计的一种算法，用于压缩文本数据，使文件占用的字节数更少。这种相对简单的压缩算法功能强大。其变体至今仍用于计算机网络、高清电视和其他应用中。通常，文本数据以每个字符 8 位的标准格式存储，使用一种称为 ASCII 的编码。该编码将每个字符映射为 0-255 之间的二进制整数值。

The idea of Huffman encoding is to abandon the rigid 8-bits-per-character requirement and use variable-length binary encodings for different characters. The advantage of doing this is that if a character occurs frequently in the file, such as the common letter 'e', it could be given a

哈夫曼编码的理念是放弃每个字符 8 位的硬性规定，对不同字符使用长度可变的二进制编码。这样做的好处是，如果一个字符在文件中频繁出现，例如常见的字母“e”，就可以给它一个

```
Welcome to CS 106X Shrink-It!
...
1) build character frequency table
2) build encoding tree
3) build encoding map
4) encode data
5) decode data
C) compress file
D) decompress file
F) free tree memory
B) binary file viewer
T) text file viewer
S) side-by-side file comparison
Q) quit
Your choice? c
Input file name: large.txt
Output file name (Enter for large.huf):
Reading 9768 uncompressed bytes.
Compressing ...
Wrote 5921 compressed bytes.
```

example output from provided HuffmanMain client

所提供的 HuffmanMain 客户端的输出示例

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 so infrequently that the extra cost is worth it.

更短的编码 (更少的比特)，使文件更小。这样做的代价是，有些字符可能需要使用比 8 位更长的编码，但这只适用于不常出现的字符，额外的成本是值得的。

The table below compares ASCII values of various characters to possible Huffman encodings for some English text. Frequent characters such as space and 'e' have short encodings, while rare ones like 'x' and 'z' have longer ones.

下表比较了一些英文文本中各种字符的 ASCII 值和可能的哈夫曼编码。空格和“e”等常见字符的编码较短，而“x”“z”等罕见字符的编码较长。

Character 人物	ASCII value ASCII 值	ASCII (binary) ASCII (二进制)	Huffman (binary) 赫夫曼 (二进制)
' '	32	00100000	10
'a'	97	01100001	0001
'b'	98	01100010	011010
'c'	99	01100011	001100
'e'	101	01100101	1100
'z' 'Z'	122	01111010	00100011010

The steps involved in Huffman encoding a given text source file into a destination compressed file are:

将给定的文本源文件转换为目标压缩文件的哈夫曼编码步骤如下：

count frequencies: Examine a source file's contents and count the number of occurrences of each character.

计算频率检查源文件的内容并计算每个字符出现的次数。

build encoding tree: Build a binary tree with a particular structure, where each leaf node stores a character and its frequency count. A priority queue is used to help build the tree along the way.

构建编码树：构建具有特定结构的二叉树，其中每个叶节点存储一个字符及其频率计数。在此过程中，优先队列会帮助建立编码树。

build encoding map: Traverse the binary tree to discover the binary encodings of each character.

构建编码图：遍历二叉树，找出每个字符的二进制编码。

encode data: Re-examine the source file's contents, and for each character, output the encoded binary version of that character to the destination file.

编码数据：重新检查源文件的内容，并为每个字符向目标文件输出该字符的二进制编码版本。

Encoding a File, Step 1: Counting Frequencies

为文件编码。第 1 步：计算频率

For example, suppose we have a file example.txt whose contents are **ababca**b**** In the original file, this text occupies 10 bytes, or 80 bits, of data. The 10th is a special "end-of-file" (EOF) byte.

例如，假设我们有一个文件 example.txt，其内容为 **ababca**b**** 在原文件中，这段文字占 10 个字节，即 80 位数据。10th 是一个特殊的“文件结束”（EOF）字节。

byte 字节	1	2	3	4	5	6	7	8	9	10
char 字符	'a'	'b'	' '	'a'	'b'	'	'c'	'a'	'b'	EOF
ASCII	97	98	32	97	98	32	99	97	98	256
binary 二进制	01100001	01100010	00100000	01100001	01100010	00100000	01100011	01100001	01100010	N/A 不适用

In Step 1 of Huffman's algorithm, a count for each character is computed. The counts are represented as a map:

在哈夫曼算法的第 1 步，计算每个字符的计数。计数以映射的形式表示：

{ ' ':2, 'a':3, 'b':3, 'c':1, EOF:1 }

Encoding a File, Step 2: Building an Encoding Tree

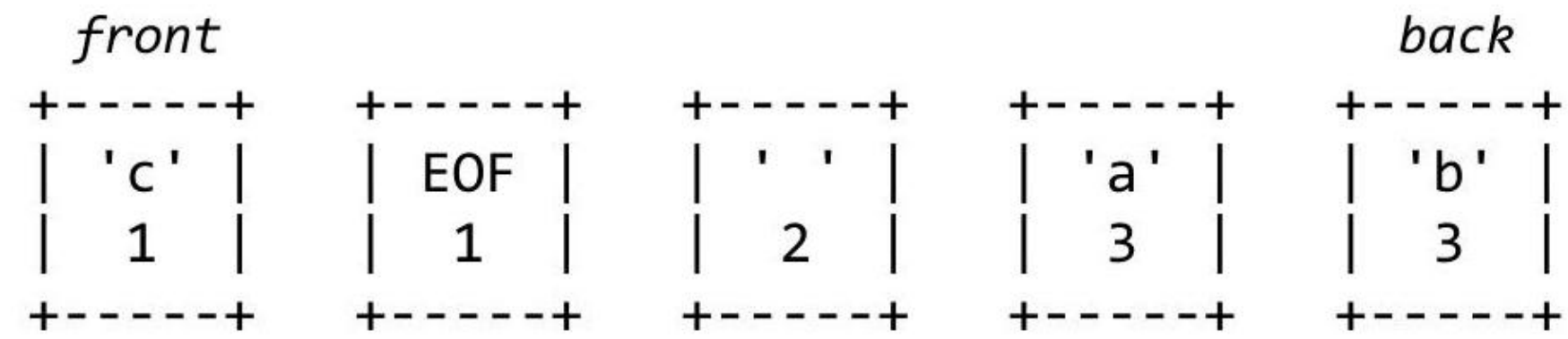
编码文件。第 2 步：构建编码树

Step 2 of Huffman's algorithm places our counts into binary tree nodes, with each node storing a character and a count of its occurrences. The nodes are then put into a priority queue, which keeps them in prioritized order with smaller counts having higher priority, so that characters with lower counts will come out of the queue sooner. The priority queue is

哈夫曼算法的第 2 步将计数放入二叉树节点。每个节点存储一个字符及其出现次数。然后将这些节点放入优先队列，按照优先级顺序排列。计数越小优先级越高，这样计数较低的字符会更快从队列中出来。优先队列为

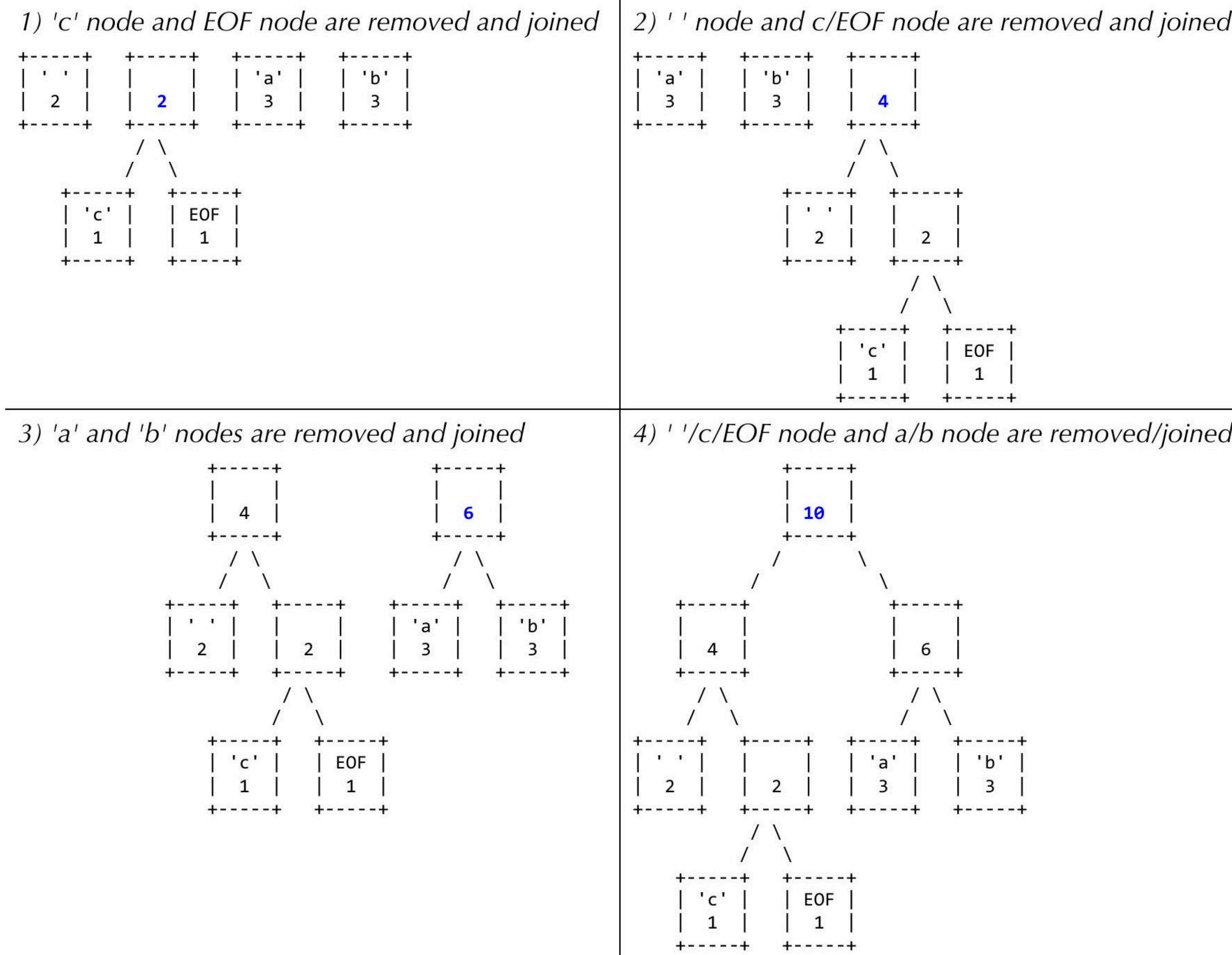
somewhat arbitrary in how it breaks ties, such as 'c' being before EOF and 'a' being before 'b'.

在如何打破并列关系方面有些武断，例如“c”位于 EOF 之前，而“a”位于“b”之前。



Now the algorithm repeatedly removes the two nodes from the front of the queue—the two with the smallest frequencies—and marries them into a new node whose frequency is their sum. The two nodes are wired in as children of the new node; the first removed becomes the left child, and the second becomes 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 final Huffman tree. The following diagram shows this process:

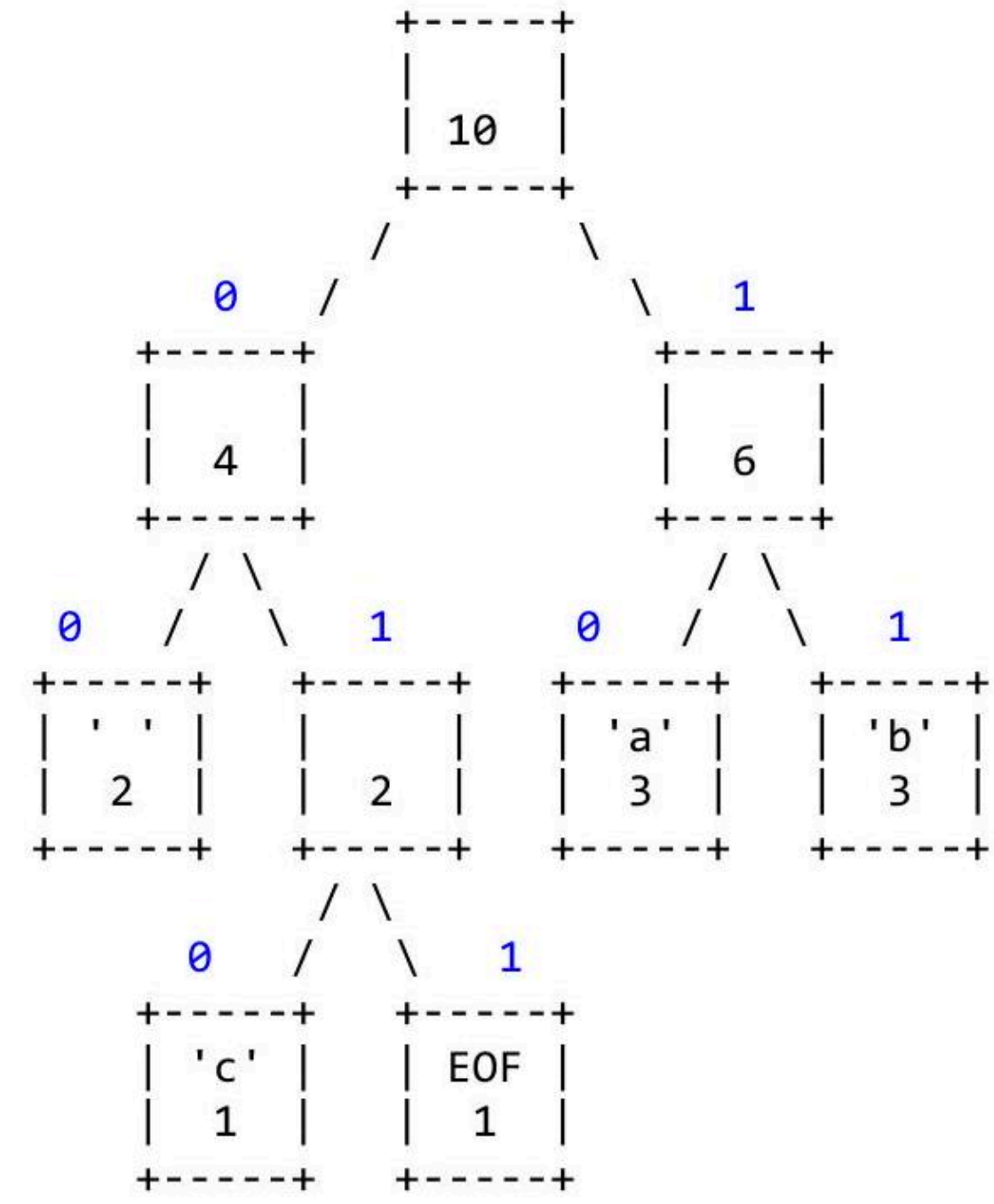
现在，算法会反复从队列前方移除两个节点—频率最小的两个节点—并将它们连接到一个新节点中，新节点的频率就是这两个节点的频率之和。这两个节点作为新节点的子节点被连接进来：第一个被移除的节点成为左节点，第二个节点成为右节点。新节点按排序顺序重新插入队列。这个过程不断重复，直到队列中只有一个二叉树节点，其他节点都是它的子节点。这就是我们最终的哈夫曼树的根节点。下图展示了这一过程：



编码文件，第 3 步：构建编码图

The Huffman code for each character is derived from your binary tree by thinking of each left branch as a bit value of 0 and each right branch as a bit value of 1, as shown in the diagram at right. 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 ' a ' is **10** and the code for ' b ' is **11**. By traversing the tree, we can produce a map from characters to their binary representations.

如右图所示，每个字符的哈夫曼代码都是从二进制树中导出的，方法是将每个左分支视为位值 0，而将每个右分支视为位值 1，每个字符的代码可以通过遍历树来确定。要到达""，我们要从根部向左走两次，因此""的代码为 **00**。c 的代码是 **010**，EOF 的代码是 **011**，a 的代码是 **10**，b 的代码是 **11**。通过遍历这棵树，我们可以生成从字符到其二进制表示的映射。



Though the binary representations are integers, since they consist of binary digits and can be arbitrary length, we will store them as strings. For this tree, it would be:

虽然二进制表示法是整数，但由于它们由二进制数字组成，长度也可以是任意的，因此我们将它们存储为字符串。对于这棵树来说，它应该是

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

Encoding a File, Step 4: Encoding the Text Data:

编码文件，第 4 步：编码文本数据：

Using the encoding map, we can encode the file's text into a shorter binary representation. Using the preceding encoding map, the text **abacab** would be encoded as:

使用编码映射，我们可以将文件文本编码为更短的二进制表示法。使用前面的编码映射，文本 **abacab** 将被编码为

101100101100010101011

The following table details the char-to-binary mapping in more detail. The overall encoded contents of the file require 22 bits, or almost 3 bytes, compared to the original file of 10 bytes.

下表详细介绍了字符到二进制的映射。与原始文件的 10 字节相比，文件的总体编码内容需要 22 位，即近 3 字节。

char 聚焦	'a'	'b'	' '	'a'	'b'	' '	'a'	'b'	EOF
binary 二进制	10	11	00	10	11	00	010	10	011

Since the character encodings have different lengths, often the length of a Huffmanencoded 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 byte are filled with zeroes. You do not need to worry about this; it is managed for you as part of the underlying file system.

由于字符编码的长度不同，哈夫曼编码文件的长度往往不是 8 位的精确倍数。文件是以整字节序列的形式存储的，因此在这种情况下，最后一个字节的剩余数字会被填为 0。你无需担心这个问题，底层文件系统会为你管理这个问题。

byte	1	2	3
char	a b a	b c a	b EOF
binary	<u>10 11 00 10</u>	<u>11 00 010 1</u>	<u>0 11 011 00</u>

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. But this will not cause problems in decoding the file, because Huffman encodings by definition have a useful prefix property where no character's encoding can be the prefix of another's.

由于字符的编码长度可能不同，而且字符可能跨越字节边界，例如第二个字节末尾的" a "，因此在存储字符时它们之间没有任何分隔符，这可能会让你担心。但这不会给文件解码带来问题，因为根据定义，哈夫曼编码有一个有用的前缀属性，即任何字符的编码都不能是另一个字符的前缀。

Decoding a File 解码文件

You can use a Huffman tree to decode text that was previously encoded with its binary patterns. 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 given the same encoding tree above, and we are asked to decode a second file containing the following bits:

您可以使用哈夫曼树来解码之前用二进制模式编码的文本。解码算法是从文件中读取每一位，每次读取一位，并使用该位遍历哈夫曼树。如果该位为 0，则在树中向左移动。如果该位为 1，则向右移动。这样做直到遇到叶节点为止。叶节点代表字符，因此一旦到达叶节点，就会输出该字符。例如，假设我们得到了上述相同的编码树，要求我们解码包含以下位的第二个文件：

11100100010101010011

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

使用哈夫曼树，我们从树根开始走，直到找到字符，然后输出字符并返回树根。

We read a 1 (right), then a 1 (right). We reach ' b ' and output b. Back to root.

我们读取一个 1（右），然后读取一个 1（右）。读到" b "，输出 b。

1110010001001010011

We read a 1 (right), then a 0 (left). We reach ' a ' and output a. (Back to root.)

我们读取一个 1（右），然后读取一个 0（左）。读到" a "，输出 a。

1110010001001010011

We read a 0 (left), then a 1 (right), then a 0 (left). We reach ' c ' and output c.

我们读取 0（左），然后是 1（右），接着是 0（左），读到" c "，输出 c。

1110010001001010011

We read a 0 (left), then a 0 (left). We reach ' ' and output a space.

我们读取一个 0（左），然后读取一个 0（左）。我们读到""，然后输出一个空格。

1110010001001010011

We read a 1 (right), then a 0 (left). We reach ' a ' and output a.

我们读取一个 1（右侧），然后读取一个 0（左侧）。读到" a "，输出 a。

1110010001001010011

We read a 0 (left), then a 1 (right), then a 0 (left). We reach ' c ' and output c.

我们读取 0（左），然后是 1（右），接着是 0（左），我们读到" c "，然后输出 c。

1110010001001010011

We read a 1 (right), then a 0 (left). We reach ' a ' and output a.

我们读取一个 1（右侧），然后读取一个 0（左侧）。读到" a "，输出 a。

1110010001001010011

We read a 0, 1, 1. This is our EOF encoding pattern, so we stop. The overall decoded text is bac ac.

我们读取一个 0, 1, 1。这是我们的 EOF 编码模式，因此我们停止。整个解码文本为 bac ac。

Provided Code 提供代码

We provide you with a file `HuffmanNode.h` that declares some useful support code, including the `HuffmanNode` structure, which represents a node in a Huffman encoding tree.

我们为您提供了一个文件 `HuffmanNode.h`，该文件声明了一些有用的支持代码，包括 `HuffmanNode` 结构，该结构表示哈夫曼编码树中的一个节点。

```
struct HuffmanNode {
    int character; // character being represented by this node
    int count; // number of occurrences of that character
    HuffmanNode* zero; // 0 (left) subtree (NULL if empty)
    HuffmanNode* one; // 1 (right) subtree (NULL if empty)
    ...
};
```

The character field is declared as type `int`, but you should think of it as a `char`. (Types `char` and `int` are largely interchangeable in C++, but using `int` here allows us to

字符字段声明为 `int` 类型，但您应该将其视为 `char`。（在 C++ 中，`char` 和 `int` 类型在很大程度上是可以互换的，但在这里使用 `int` 可以让我们

sometimes use character to store values outside the normal range of `char`, for use as special flags.) The character field can take one of three types of values:

有时使用字符来存储 `char` 正常范围之外的值，以用作特殊标记）。字符字段可存储三种值之一：

an actual char value: 一个实际字符值；

the constant `PSEUDO_EOF` (defined in `bitstream.h` in the Stanford library), which represents the pseudo-EOF value that you will need to place at the end of an encoded stream; or

常量 `PSEUDO_EOF`（定义在 Stanford 库的 `bitstream.h`），它表示需要放在编码流末尾的伪 EOF 值；或

the constant `NOT_A_CHAR` (defined in `bitstream.h` in the Stanford library), which represents something that isn't actually a character. This can be stored in interior nodes of the Huffman encoding tree, because such nodes do not represent any one individual character.

常量 `NOT_A_CHAR`（定义于斯坦福库中的 `bitstream`。这可以存储在哈夫曼编码树的内部节点中，因为这些节点并不代表任何一个单独的字符。

Bit Input/Output Streams 比特输入/输出流

In parts of this program you will need to read and write bits to files. In the past we have wanted to read input an entire line or word at a time, but for this program, it's much better to read one single character (byte) at a time. You should use the following in/output stream functions:

在本程序的部分内容中，您需要读取和写入文件的位。过去，我们希望一次读取整行或整字的输入，但在本程序中，一次读取一个字符（字节）会更好。您应该使用以下输入/输出流函数：

ostream (output stream) member 输出流	Description 说明
void put(int byte)	writes a single byte (character, 8 bits) to the output stream 将单字节（字符，8 位）写入输出流

istream (input stream) member istream（输入流）成员	Description 说明
int get()	reads a single byte (character, 8 bits) from input (or EOF) 从输入（或 EOF）读取单字节（字符，8 位）

You might also find that you want to read an input stream, then "rewind" it back to the start and read it again. To do this on an input stream variable named `input`, you can use the `rewindStream` function from `filelib.h`:

您可能还会发现，您想读取一个输入流，然后将其“倒带”回起点并再次读取。要在名为 `input` 的输入流变量上实现这一功能，可以使用 `filelib.h` 中的 `rewindStream` 函数：

```
rewindStream(input); // tells the stream to seek back to the beginning
```

To read or write a compressed file, even a whole byte is too much! You will want to read and write binary data one single bit at a time, which is not directly supported by the default in/output streams. Fortunately, the Stanford C++ library provides `obitstream` and `ibitstream` classes with `writeBit` and `readBit` members to make it easier.

要读取或写入压缩文件，即使是整个字节也太多了！您需要一次读写一个比特的二进制数据，而默认的输入/输出流并不直接支持这一点。幸运的是，斯坦福 C++ 程序库提供了带有 `writeBit` 和 `readBit` 成员的 `obitstream` 和 `ibitstream` 类，让这一切变得更容易。

