

Assignment 3

cpe 357 Winter 2018

"I'd crawl over an acre of 'Visual This++' and 'Integrated Development That' to get to gcc, Emacs, and gdb. Thank you."
(By Vance Petree, Virginia Power)

— /usr/games/fortune

Due by 11:59:59pm, Monday, February 12th.
This assignment is to be done individually.

Programs: hencode and hdecode

This assignment is to build a file compression tool, **hencode**, that will use Huffman codes to compress a given file, and a file decompression tool, **hdecode**, that will uncompress a file compressed with **hencode**.

Usage:

```
hencode infile [ outfile ]
hdecode [ ( infile | - ) [ outfile ] ]
```

hencode must:

- take a command line argument that is the name of the file to be compressed; and
- take a second, optional, command line argument that is the name of the outfile. (If this argument is missing, the output should go to stdout.)
- properly encode its input into a binary Huffman-coded file according to the algorithm below.

hdecode must:

- Take an optional command line argument that is the name of the file to be uncompressed. If this argument is missing, or if the input file name is “-”, **hdecode** will take its input from standard input.
- Take a second, optional, command line argument that is the name of the output file. (If this argument is missing, the output should go to stdout.)
- Properly decode its input and restore the original data.

Both must:

- use only UNIX unbuffered IO (**read(2)** and **write(2)**) for reading and writing the files¹; and
- share as much of the code base as is reasonable between the two programs; and
- demonstrate robust programming techniques including proper handling of errors; and
- adhere to a “reasonableness” standard with respect to performance.

See also the “Pesky Details” section below.

¹That means go ahead and use **stdio** for error messages.

Huffman Codes

In 1952, David A. Huffman proposed a method for compressing files by generating a variable length encoding of characters based on the relative frequency with which each character occurs in a file². This encoding is what is known as *prefixless code* because no letter encoding is a prefix of any other encoding. In practice, this means that the coding can be stored in a binary tree with the characters at the leaves. Each character's encoding can be determined by looking at the sequence of right or left child node transitions and representing those as either zero bits or one bits.

Generating the Tree

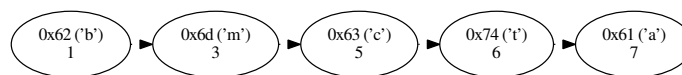
Note: Be sure to follow the tiebreaker conventions outlined in this document. If you do something different, you will still get a valid encoding of your file, but will be different from the reference one. That means that both programs will have to work perfectly all the way through in order to be tested. If you follow the tiebreakers partial credit becomes a possibility.

To build the tree, it is first necessary to generate a histogram of all the characters in the file, remembering that any value from 0 to 255 represents a valid byte. We will demonstrate the algorithm through the following encoding example.

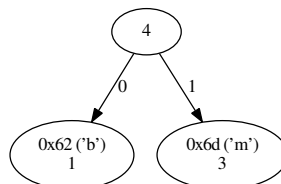
If the file you are encoding consists of the string “bmmcccccttttttaaaaaa” the histogram will be:

Byte	Count
0x61 ('a')	7
0x62 ('b')	1
0x63 ('c')	5
0x6d ('m')	3
0x74 ('t')	6

Take these counts and generate a linked list of them in ascending order of frequency. If there is a tie in frequency, do a secondary sort and place the characters also in ascending order. For the histogram above, this list will look like:

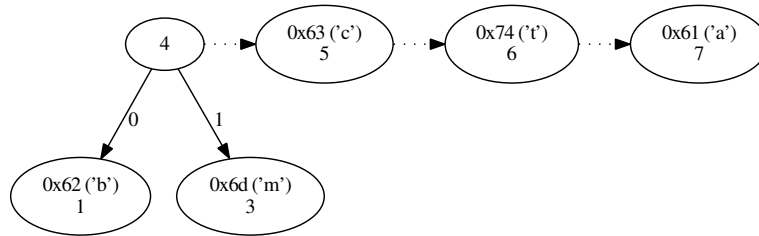


Now, remove the first two nodes on the list and construct a new node whose frequency is the sum of the two removed nodes. Make the first removed node the left child of the new node and the second removed node the right subchild:



²David A. Huffman. A method for the construction of minimum-redundancy codes. *Proceedings of the Institute of Radio Engineers*, 40(9):1098–1101, September 1952

Next, reinsert the new node into the remaining list of nodes in order. If there is a frequency tie, insert the new node before the other node in the list. After reinsertion of the new node, our list will look like:



Repeat the process until there is only one node left in the list. The rest of the process for this file is shown in Figure 1.

Now that you have the tree, to encode the file all that is necessary is to do a pass over the tree to extract the encodings into a table, then re-read the input file and translate each input character into the appropriate sequence of bits. If the file ends with a partially filled byte, pad the final byte with zeros.

The codes extracted from the tree above will be:

0x61 ('a')	: 11
0x62 ('b')	: 000
0x63 ('c')	: 01
0x6d ('m')	: 001
0x74 ('t')	: 10

File Format

The output format for the encoded file consists of two parts. First, comes a header that contains the frequency information necessary to re-create the tree, then the bits of the encoded file.

The first four bytes of the header consists of an integer containing the number of characters in the frequency table. (Remember, not all files will contain all 256 possible bytes.) Following this comes the frequency table, in alphabetical (numerical) order³. Each element of the frequency table consists of a single byte that is the byte itself, followed by a four-byte integer that is the frequency.

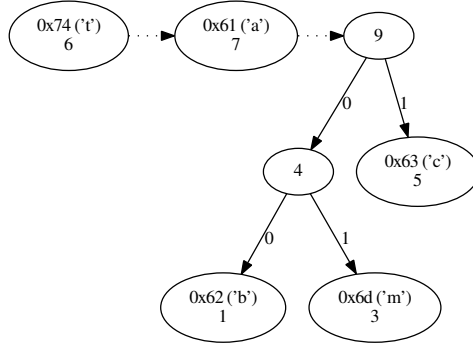
Number of chars	c1	count of c1	c2	count of c2	...	cn	count of cn
(uint32_t)	(uint8_t)	(uint32_t)	(uint8_t)	(uint32_t)	...	(uint8_t)	(uint32_t)
4 bytes	1 byte	4 bytes	1 byte	4 bytes	...	1 byte	4 bytes

Figure 2 shows the resulting encoded file for the input above (for a little-endian machine) as hexadecimal bytes. Boxes are drawn around multi-byte data for clarity. Note that there is no padding in this particular example. Also, there is no separation between the header and the body, nor is there an end of file marker. These are all determined by counting.

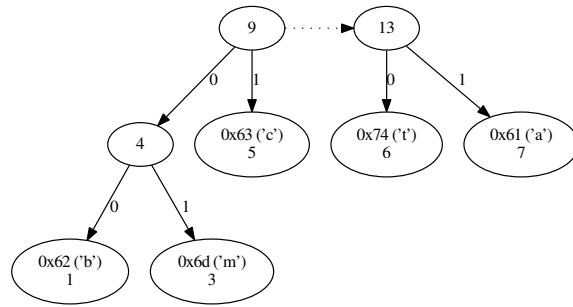
Encoding

To encode the file, build the tree, extract the encodings into a code table, then re-read the input generating the output as you go.

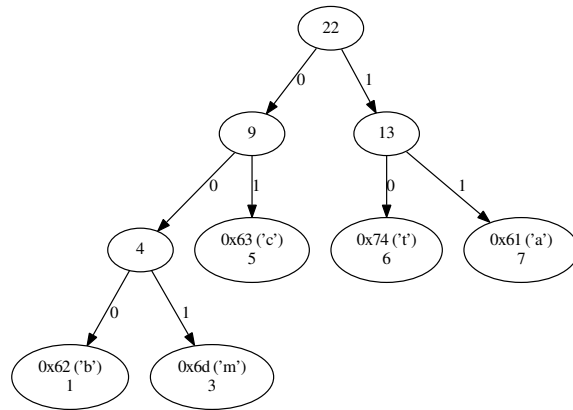
³This is not necessary for file compression, but it will make your output match the reference program's output.



(a) After a second pass



(b) After a third pass



(c) The final tree

Figure 1: Finishing tree generation

Header	Header size				'a'	Count				'b'	Count			
	05	00	00	00	61	07	00	00	00	62	01	00	00	00
	'c' Count				'm' Count				't' Count					
Body	63	05	00	00	00	6d	03	00	00	00	74	06	00	00
	File Contents													
	04	95	56	aa	bf	ff								

Figure 2: “bmmmmcccccttttttaaaaaa” as encoded by hencode.

Decoding

To decode, read the header of the encoded file, regenerate the tree, then use the bits of the file to walk the tree to regenerate the file. Start at the root. For each zero bit, go left, for each one bit, go right. Because this is a prefixless code, you will know you have decoded a character when you reach a leaf. Output this character and start again at the root.

You know how many characters are in the output from the counts encoded in the frequency table.

Tricks and Tools

od(1)	Useful for reading binary files to see what's what.
open(2) close(2) read(2) write(2) lseek(2)	The UNIX basic IO functions. Check out the <code>man</code> pages.

Figure 3: Some potentially useful library functions

Some potentially useful library functions listed in Figure 3. Also, it might be helpful to know that `<stdint.h>` defines a set of fixed-width data types which are more portable than `ints` when you really need to know how big something is. These exist in signed and unsigned versions named `intXX_t` and `uintXX_t`, where `XX` is the number of bits. For example, `uint32_t size` makes `size` a 32-bit unsigned integer.

Coding Standards and Make

See the pages on coding standards and make on the cpe 357 class web page.

Pesky Details

Endianness

Because you are reading and writing binary data, the endianness of integers will be apparent. For this assignment you do not need to be concerned with this, but know that files encoded on a big-endian machine will not decode properly on a little-endian one (and vice-versa). Note that the x86 is little-endian but if you're developing on your PowerPC, it's big-endian.

Bits and their order

- If it is necessary to pad the final byte of the file, pad it with zero bits.
- When filling bits, fill from the high order bits. That is, if the final four bits of an encoded file are 1010, the final byte of the file will be 0xA0.

What to turn in

Submit via `handin` in the CSL to the `asgn3` directory of the `pn-cs357` account:

- your well-documented source files.
- A makefile (called `Makefile`) that will build your programs when given the command “`make all`”.
- A README file that contains:
 - Your name.
 - Any special instructions for running your program.
 - Any other thing you want me to know while I am grading it.

The README file should be **plain text**, i.e, **not a Word document**, and should be named “README”, all capitals with no extension.

Sample runs

Below are some sample runs of `hencode` and `hdecode`. I have placed a runnable versions of `hencode` and `hdecode` in `~pn-cs357/demos` as `hencode`, and `hdecode`.

Here is an example of encoding the file used in the example above:

```
% htable testfile
0x61: 11
0x62: 000
0x63: 01
0x6d: 001
0x74: 10
% hencode testfile testfile.huff
% hdecode testfile.huff testfile.out
% diff testfile testfile.out
% ls -l
total 12
-rw----- 1 pnico pnico 22 Oct 18 15:00 testfile
-rw----- 1 pnico pnico 35 Oct 18 15:02 testfile.huff
-rw----- 1 pnico pnico 22 Oct 18 15:02 testfile.out
```

Note that, because `testfile` is so small, this made the “compressed” file bigger. Consider this example with a larger file, the class notes for `cpe357` so far:

```
% ls -l class.ps
-rw----- 1 pnico pnico 339520 Oct 18 15:04 class.ps
% hencode class.ps class.ps.huff
% ls -l class.*
-rw----- 1 pnico pnico 339520 Oct 18 15:04 class.ps
-rw----- 1 pnico pnico 217162 Oct 18 15:05 class.ps.huff
%
```

A reduction of 36% isn't so bad.