

Huffman Coding Project Technical Design

Victor W. Frye

Davenport University

Huffman Coding Project Technical Design

Overview

This program will compress and decompress an input ASCII text file using the Huffman coding algorithm.

Scope

This program will compress or decompress an input ASCII text file using the Huffman coding algorithm. Execution of the program will be done from the command line and determine whether to compress or decompress based on the input arguments.

The input file for compression will be read as-is and not converted in any way. The input file is assumed to be an ASCII text file and is not translated in any way. The output file for compression will be plain binary text. Additionally, a key file will be generated to decompress the file to its original state.

The input file for decompression will be assumed to be a compressed file and will not be translated in any way. The input file is assumed to be a encrypted ASCII text file with only binary representation of the original data. A key file will be required to decompress the input file. If no key file is found, the input file will be read as corrupted.

Usage

To use the program, two additional command line arguments will be needed compared to the standard executable calling. The first will either be “compress” or “decompress”, designating what action to perform upon the file. The second argument is the file path for the input file.

Note that the output file and key file are defaulted within the program. The output file will output as “output.txt” and the key file will output as “key.txt”.

Processing Outline

1. If the number of command line arguments is not two
 - a. Display usage error message
 - b. Exit
2. If the action route argument is invalid (not compress or decompress)
 - a. Display usage error message
 - b. Exit
3. If the input file does not exist
 - a. Display input error message
 - b. Exit
4. Initialize an array of character frequency objects
 - a. For each element, create a character frequency object with the character set to the corresponding ASCII value based upon the index
 - b. Set the initial character frequency count to zero for each element in the array
5. If action route argument is to compress
 - a. Open the input file
 - i. Read a character from the input file
 - ii. While not at the end of the input file
 1. Increment the character frequency count using the ASCII value of the character as the index
 2. Read the next character
 - b. Close the input file

- c. Sort the array of character frequency objects using the built-in array sort method
- d. Initialize a linked list of nodes of type character frequency
 - i. For each element in the array of character frequency objects
 - 1. If the frequency of the character frequency object is not zero (0)
 - a. Initialize a node with the character frequency object inserted as the root of the tree
 - b. Add the newly created node to the end of the linked list
- e. Construct a binary tree of the character frequency objects
 - i. While the count of linked list nodes is less than one (1)
 - 1. If only two nodes remain in the linked list and the first node is larger than the last
 - a. Swap the first and last nodes
 - 2. Remove the first two (2) nodes from the linked list
 - 3. Construct a new node with a frequency equal to the sum of the frequencies of the two (2) removed nodes
 - 4. Connect the three (3) nodes, with the newly constructed node being the parent and the two (2) removed nodes being the children (the smallest should be on the left)
 - 5. If the new node is smaller than all other nodes or no nodes exist in the linked list

- a. Add the new node to the beginning of the linked list
6. Else find the first node in the linked list that is larger than the new node and add the new node before it
- f. Construct an encoding table
- g. Create compressed file
 - i. Open the input file
 - ii. If the output file exists
 1. Delete the output file
 - iii. Open the output file
 - iv. Read a character from the input file
 - v. While not at the end of the input file
 1. Find assigned Huffman code for the read character
 2. For each bit in the code
 - a. If a full byte has been constructed
 - i. Write byte to output file
 - ii. Reset byte and counter
 - b. If the current bit in the code is a one (1)
 - i. Modify the bit corresponding to counter-designated place in the byte to be a one (1)
 - ii. Increment the counter
 - c. Else increment the counter
 3. Read next character from the input file
 - vi. Write final byte to the output file

- vii. Close the output file
 - h. If key file exists
 - i. Delete the key file
 - i. Open the key file
 - i. For each character frequency object in the character frequency array
 - 1. If the frequency of the object is not zero (0)
 - a. Write a string representation of the character frequency object to a new line in the key file
 - j. Close the key file
- 6. Else (assumed to be decompression action route)
 - a. If the key file does not exist
 - i. Display missing key file error message
 - ii. Exit
 - b. Open the key file
 - i. Instantiate counter for total characters
 - ii. Read the first line from the key file
 - iii. While the read line is not null
 - 1. Check read line for matching string representation of a character frequency object
 - 2. Convert matching string representation to character frequency objects if found
 - 3. Add the found frequency to the total characters counter

4. Read the next line
- c. Close the key file
- d. Initialize a linked list of nodes of type character frequency
 - i. For each element in the array of character frequency objects
 1. If the frequency of the character frequency object is not zero (0)
 - a. Initialize a node with the character frequency object inserted as the root of the tree
 - b. Add the newly created node to the end of the linked list
- e. Construct a binary tree of the character frequency objects
 - i. While the count of linked list nodes is less than one (1)
 1. If only two nodes remain in the linked list and the first node is larger than the last
 - a. Swap the first and last nodes
 2. Remove the first two (2) nodes from the linked list
 3. Construct a new node with a frequency equal to the sum of the frequencies of the two (2) removed nodes
 4. Connect the three (3) nodes, with the newly constructed node being the parent and the two (2) removed nodes being the children (the smallest should be on the left)
 5. If the new node is smaller than all other nodes or no nodes exist in the linked list

- a. Add the new node to the beginning of the linked list
6. Else find the first node in the linked list that is larger than
the new node and add the new node before it
- f. If the output file already exists
 - i. Delete the output file
- g. Open the input and output file
 - i. Read a character from the input file
 - ii. Instantiate counter and array to measure out full bytes
 - iii. While input file data is valid
 1. If a full byte has been measured out
 - a. Search the binary tree for leaf nodes corresponding
to path designated by bits
 - i. If bit is a one, traverse the tree to the left
 - ii. If bit is a zero, traverse the tree to the right
 - iii. If a leaf node is found and the total
remaining character to print are greater than
zero
 1. Write the character corresponding to
the leaf node to the output file
 2. Return to the root of the tree
 3. Decrease the remaining amount of
characters in file
 - b. Reset the counter and byte measuring array

2. Else

- a. Add bit to the byte measuring array
- b. Increment counter
- c. Read the next character from the input file

iv. Close the input and output file

7. Exit

Data

Input File

The input file will either be a standard ASCII text file or a compressed binary bit text file depending on whether compression or decompression is the current action.

For compression, it will be assumed that the file is an ASCII text file and will be read a single character at a time in sequential order. If a file other than an ASCII text file is used, the results will be undefined.

For decompression, it will be assumed that the file was compressed using this program and a key file already exists. If no key file is found, the program will reject the process. The input file will be read bit by bit looking for ones (1) and zeros (0). No other characters or information will be recognized. If a file other than one compressed by this program is used, the results will be undefined.

Output File

The output file for compression will be a series of ones (1) and zeros (0) as the output is a binary representation of the original data. The output of the decompression process will be the original file assuming that the key file was not corrupted or tampered with.

Data Structures

An array containing 256 character frequency objects will be used in both compression and decompression to determine the frequency of the characters found in the original file. The decompression process will rely on a key file to gather the frequencies of the characters in the original file. Regular expression will be used to decipher the key file in this process. Access to the array is handled by perfectly hashing the ASCII value of the character to the index of the

array. After all frequencies have been determined, the array is sorted using the built-in sort method and is no longer able to perform flawless $O(1)$ operations.

The sorted array is instantiated into a linked list of all character frequency objects with greater than zero occurrences registered. The linked list is used allow for quick removal of low frequency nodes to create a binary tree. A linked list is used rather than a stack so that composite nodes (nodes with a summed frequency of two smaller nodes) are still able to be re-inserted into the linked list in ascending order. This reinsertion is an $O(n)$ operation that later allows for quick access through the binary tree.

A binary search tree is the primary data structure used to perform the Huffman coding process. The binary search tree branches with smaller nodes on the left and larger nodes on the right. This allows for the Huffman code designation by assigning a one or zero respectively to the traversal.

Footnotes

¹Note: For a visible representation of the ASCII table, see Figure 1. Since some ASCII characters are invisible, they are identified as a string in the table, e.g. “NUL”, “SOH”, “BEL”, and “CR”.

Tables and Figures

Ctrl	Dec	Hex	Char	Code	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char
^@	0	00		NUL	32	20	!	64	40	@	96	60	'	128	80	Ç	160	A0	ā	192	C0	Ł	224	E0	α						
^A	1	01		SOH	33	21	!"	65	41	Å	97	61	a	129	81	Ć	161	A1	ă	193	C1	ł	225	E1	Β						
^B	2	02		STX	34	22	#\$	66	42	Β	98	62	b	130	82	Ĉ	162	A2	â	194	C2	Ł	226	E2	Γ						
^C	3	03		ETX	35	23	%	67	43	Ɔ	99	63	c	131	83	Ċ	163	A3	ă	195	C3	ł	227	E3	Π						
^D	4	04		EOT	36	24	&	68	44	D	100	64	d	132	84	Ċ	164	A4	â	196	C4	ł	228	E4	Σ						
^E	5	05		ENQ	37	25	'	69	45	E	101	65	e	133	85	Ċ	165	A5	ă	197	C5	ł	229	E5	σ						
^F	6	06		ACK	38	26	&	70	46	F	102	66	f	134	86	Ċ	166	A6	ă	198	C6	ł	230	E6	μ						
^G	7	07		BEL	39	27	,	71	47	G	103	67	g	135	87	Ċ	167	A7	ă	199	C7	ł	231	E7	Υ						
^H	8	08		BS	40	28	(72	48	H	104	68	h	136	88	Ċ	168	A8	ă	200	C8	ł	232	E8	Ϻ						
^I	9	09		HT	41	29)	73	49	I	105	69	i	137	89	Ċ	169	A9	ă	201	C9	ł	233	E9	ϻ						
^J	10	0A		LF	42	2A	*	74	4A	J	106	6A	j	138	8A	Ċ	170	AA	ă	202	CA	ł	234	EA	Ͽ						
^K	11	0B		VT	43	2B	+	75	4B	K	107	6B	k	139	8B	Ċ	171	AB	ă	203	CB	ł	235	EB	δ						
^L	12	0C		FF	44	2C	,	76	4C	L	108	6C	l	140	8C	Ċ	172	AC	ă	204	CC	ł	236	EC	Ϸ						
^M	13	0D		CR	45	2D	-	77	4D	M	109	6D	m	141	8D	Ċ	173	AD	ă	205	CD	ł	237	ED	ϸ						
^N	14	0E		SO	46	2E	.	78	4E	N	110	6E	n	142	8E	Ċ	174	AE	ă	206	CE	ł	238	EE	Ϲ						
^O	15	0F		SI	47	2F	/	79	4F	O	111	6F	o	143	8F	Ċ	175	AF	ă	207	CF	ł	239	EF	Ϻ						
^P	16	10		DLE	48	30	0	80	50	P	112	70	p	144	90	Ċ	176	B0	ă	208	D0	ł	240	F0	≡						
^Q	17	11		DC1	49	31	1	81	51	Q	113	71	q	145	91	Ċ	177	B1	ă	209	D1	ł	241	F1	±						
^R	18	12		DC2	50	32	2	82	52	R	114	72	r	146	92	Ċ	178	B2	ă	210	D2	ł	242	F2	℥						
^S	19	13		DC3	51	33	3	83	53	S	115	73	s	147	93	Ċ	179	B3	ă	211	D3	ł	243	F3	ℳ						
^T	20	14		DC4	52	34	4	84	54	T	116	74	t	148	94	Ċ	180	B4	ă	212	D4	ł	244	F4	ℴ						
^U	21	15		NAK	53	35	5	85	55	U	117	75	u	149	95	Ċ	181	B5	ă	213	D5	ł	245	F5	ℵ						
^V	22	16		SYN	54	36	6	86	56	V	118	76	v	150	96	Ċ	182	B6	ă	214	D6	ł	246	F6	+						
^W	23	17		ETB	55	37	7	87	57	W	119	77	w	151	97	Ċ	183	B7	ă	215	D7	ł	247	F7	≈						
^X	24	18		CAN	56	38	8	88	58	X	120	78	x	152	98	Ċ	184	B8	ă	216	D8	ł	248	F8	ο						
^Y	25	19		EM	57	39	9	89	59	Y	121	79	y	153	99	Ċ	185	B9	ă	217	D9	ł	249	F9	•						
^Z	26	1A		SUB	58	3A	:	90	5A	Z	122	7A	z	154	9A	Ċ	186	BA	ă	218	DA	ł	250	FA	•						
^[27	1B		ESC	59	3B	;	91	5B	[123	7B	{	155	9B	Ċ	187	BB	ă	219	DB	ł	251	FB	√						
^\	28	1C		FS	60	3C	<	92	5C	\	124	7C		156	9C	Ċ	188	BC	ă	220	DC	ł	252	FC	n						
^]	29	1D		GS	61	3D	=	93	5D]	125	7D	}	157	9D	Ċ	189	BD	ă	221	DD	ł	253	FD	2						
^^	30	1E	▲	RS	62	3E	>	94	5E	^	126	7E	~	158	9E	Ċ	190	BE	ă	222	DE	ł	254	FE	■						
^-	31	1F	▼	US	63	3F	?	95	5F	-	127	7F	À	159	9F	Ċ	191	BF	ă	223	DF	ł	255	FF							

* ASCII code 127 has the code DEL. Under MS-DOS, this code has the same effect as ASCII 8 (BS).
The DEL code can be generated by the CTRL + BKSP key.

Figure 1. The American Standard Code for Information Interchange (ASCII) table. The above images contain the decimal and hexadecimal values and their corresponding character. The characters above 128 in decimal are known as the extended character set.