

DELHI TECHNOLOGICAL UNIVERSITY



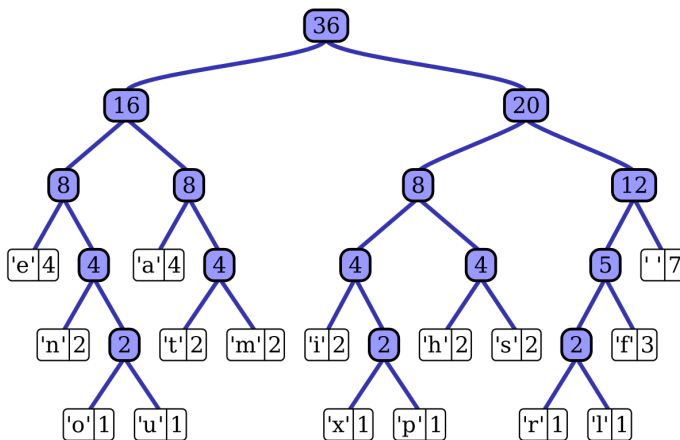
CO-201 Data Structures Innovative Project Report

COMPRESSION USING HUFFMAN CODING

GUI built on Qt5 using C++

By:

ISHAN KHANNA
2K19/EE/120



Symbol	Bits	Probability List				Bits Assignment Procedure			
S1	1	0.40	1	0.40	1	0.40	1	0.50	0
S2	01	0.20	01	0.20	01	0.20	01	0.30	00
S3	0000	0.10	0000	0.12	001	0.18	000	0.20	01
S4	0001	0.08	0001	0.10	0000	0.12	001		
S5	0011	0.07	0010	0.08	0001				
S6	00100	0.05	0011						
S7	00101								

Acknowledgement

I would like to express gratitude to Delhi Technological University for their valuable support and for providing me a platform to showcase my research and analysis skills through working on this project on Compression using Huffman Coding. Also, we would like to show gratitude to our professor, who gave us the opportunity to work on this project and provided us with required support and guidance. The success and final outcome of this project required guidance and contribution from many people and we are extremely privileged to have received it all along. All that we have accomplished is due to such supervision and assistance; we thank all who have been a part of this effort.

Contents

Acknowledgement	2
EXECUTIVE SUMMARY	4
COMPRESSION USING HUFFMAN CODING.....	5
Pseudo-Code	7
Time and space complexity analysis for Huffman Coding Algorithms.....	8
Background	9
Design and Implementation.....	10
Results.....	11
Conclusions and Reflection on Learning.....	13
References and bibliography	14

EXECUTIVE SUMMARY

Our aim is to take any text or image file and compress it without any loss in its data, that is, no change in the quality of the file being compressed.

It will enable users to share data over network with low bandwidth and thus reduced costs without any compromise in the quality of the data being transmitted. We did this using Qt5 and C++.

COMPRESSION USING HUFFMAN CODING

Huffman code is a technique that is commonly used for lossless data compression. The process of such compression proceeds through **Huffman coding**. This algorithm was developed by David A. Huffman who at the time was a Sc.D. student at MIT, USA. It was published in the 1952 paper which was titled "A Method for the Construction of Minimum-Redundancy Codes".

The output obtained from Huffman's algorithm can be seen as a variable-length code table for encoding a source symbol (such as a character in a file). The algorithm created this table from the frequency of occurrence (also known as *weight*) for each possible value of the source symbol. However, even though it is among the most optimal methods for encoding symbols separately, it is not always optimal among all compression methods.

Huffman Encoder utilises Huffman coding trees to reduce the count of bits that are needed to represent any symbol. It is a variable length encoding.

For Example :

- Message “CAB” is compressed to

011110.

- Otherwise the message would have been encoded using its ASCII representation as follows:

001000110010000100100010.

Algorithm

1. Find the frequency of each distinct symbol in the input.
2. Make a priority queue of single-node trees. Every node in the initial queue represents a distinct symbol from the set of possible symbols, and contains the frequency of that symbol in the message to be coded. Symbols with a frequency of zero are ignored.
3. Loop while there is more than 1 tree in the queue:
 - 3a. Remove the two trees from the queue that have the lowest frequency contained in their roots.
 - 3b. Create a new node that will be the root of a new tree. This new tree will have those two trees just removed in step 3a as left and right sub-trees. The frequency in the root of this new tree will be the sum of the frequency in the roots of its sub-trees.
 - 3d. Label the edge from this new root to its left sub-tree "0", and label the edge to its right sub-tree "1".
 - 3e. Insert this new tree in the queue, and go to 3.
4. Return the one tree in the queue as the Huffman code tree.

Time and space complexity analysis for Huffman Coding Algorithms

Time Complexity:

Since efficient priority queue data structures require $O(\log n)$ time complexity, and a tree that has n leaves must have $2n-1$ nodes. So, this algorithm operates in $O(n \log n)$ time, where n is the number of symbols.

Space Complexity:

If n is the count of distinct symbols in the input file, then we require $O(n)$ space for building the tree. Hence, space complexity of Huffman coding is $O(n)$.

Background

The motivation behind the project was to build something simplistic yet useful. With the large amounts of data that are available nowadays, storage space seems to run out eventually for everyone regardless of the size of the storage device.

Huffman coding is one algorithmic technique which can be used to develop a solution to this problem. Hence, it has been deployed in the project.

Design and Implementation

I have built a GUI using Qt5. It includes a text view and an image view. The text view provides basic text editing features along with the options to compress and decompress the currently opened file. It also has provisions for sharing by e-mail option printing.

The Image View includes an in-built image-display area and options to compress the file to JPEG. It also provides an option to print the image and to open an image for viewing.

Features:

Some of the features in the current version include:

Text editor – Cut, Copy, Paste, Undo, Redo

Text Compression

Text Decompression

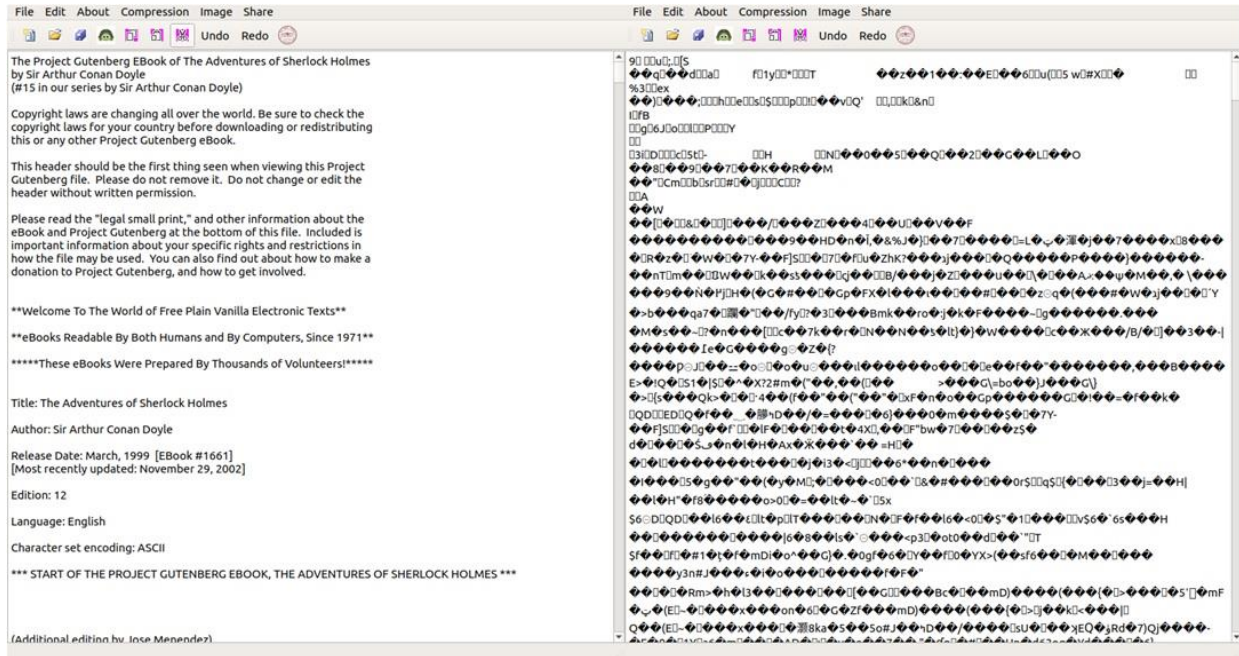
Image Viewer

Image Compression

Sharing by Mail

Print

Results





Conclusions and Reflection on Learning

Through this project, I learned about the wonderful compression technique using Huffman coding. I also learned a lot about Object Oriented Programming, Graphical User Interfaces, software organization, project management and algorithmic thinking. I have understood their value and appreciate the same. I learned a lot about the topic and explored it through the various examples, codes, derivations and applications that are contained in the topic thereby gaining immense knowledge about them.

References and bibliography

References

1. Huffman, D. (1952). "A Method for the Construction of Minimum-Redundancy Codes" (PDF). *Proceedings of the IRE*. **40** (9): 1098–1101. doi:10.1109/JRPROC.1952.273898.
2. Van Leeuwen, Jan (1976). "On the construction of Huffman trees" (PDF). *ICALP*: 382–410. Retrieved 2014-02-20.
3. Ze-Nian Li; Mark S. Drew; Jiangchuan Liu (2014-04-09). *Fundamentals of Multimedia*. Springer Science & Business Media. ISBN 978-3-319-05290-8.
4. J. Duda, K. Tahboub, N. J. Gadil, E. J. Delp, *The use of asymmetric numeral systems as an accurate replacement for Huffman coding*, Picture Coding Symposium, 2015.
5. Huffman, Ken (1991). "Profile: David A. Huffman: Encoding the "Neatness" of Ones and Zeroes". *Scientific American*: 54–58.
6. Gribov, Alexander (2017-04-10). "Optimal Compression of a Polyline with Segments and Arcs". arXiv:1604.07476 [cs.CG].
7. Gallager, R.G.; van Voorhis, D.C. (1975). "Optimal source codes for geometrically distributed integer alphabets". *IEEE Transactions on Information Theory*. **21** (2): 228–230. doi:10.1109/TIT.1975.1055357.
8. Abrahams, J. (1997-06-11). *Written at Arlington, VA, USA. Division of Mathematics, Computer & Information Sciences, Office of Naval Research (ONR). "Code and Parse*

Trees for Lossless Source Encoding". Compression and Complexity of Sequences 1997 Proceedings. Salerno: IEEE: 145–

171. CiteSeerX 10.1.1.589.4726. doi:10.1109/SEQUEN.1997.666911. ISBN 0-8186-8132-2. S2CID 124587565.

9. Hu, T. C.; Tucker, A. C. (1971). "Optimal Computer Search Trees and Variable-Length Alphabetical Codes". *SIAM Journal on Applied Mathematics*. **21** (4): 514. doi:10.1137/0121057. JSTOR 2099603.

10. Knuth, Donald E. (1998), "Algorithm G (Garsia–Wachs algorithm for optimum binary trees)", *The Art of Computer Programming, Vol. 3: Sorting and Searching* (2nd ed.), Addison–Wesley, pp. 451–453. See also History and bibliography, pp. 453–454.

Bibliography

- Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms*, Second Edition. MIT Press and McGraw-Hill, 2001. ISBN 0-262-03293-7. Section 16.3, pp. 385–392.