

The Implementation and Comparison of the BCCBT Data Compression Algorithm

Aiden Taylor - B.Sc. in Computer Science

Noah Pinel - B.Sc. in Computer Science

Ty Irving - B.Sc. in Computer Science

Abstract—When storing, transferring, or receiving files it is important to have data compression algorithms to make these processes easier and more efficient. In our project we dive into the actual implementation of a proposed data compression algorithm, and use certain factors of comparison to gauge our implementation against other open source data compression algorithms.

I. INTRODUCTION

The purpose of this paper is to take a deeper dive into the implementation and comparison of data compression algorithms. In the modern age, data compression algorithms have become an integral part of our daily lives. These algorithms help us save space on our hard drives, use up less bandwidth, speed up communications, and so much more.

Unfortunately, all these great benefits do not come without costs. In reality, data compression algorithms can be incredibly difficult to devise, and even more difficult to implement. Devising a data compression algorithm requires rigorous mathematical justifications for all the processes taken, and then implementing a data compression algorithm requires intimate knowledge of optimization and memory management techniques.

Once a data compression algorithm is devised and implemented, it is important to see how it stacks up against the other readily available options. No one algorithm excels at compressing every type of data [1], so comparing a data compression algorithm against other available options can highlight the areas where an algorithm exceeds, or highlight the areas where it falls behind. The most important factors of comparisons to note are the compression and decompression times, and

how much of the original file the algorithm is able to save without losing any data from the source file.

In this study, we decided to focus on the already devised Bit Code Complete Binary Tree (BCCBT) data compression algorithm which was proposed by Mattias Håkansson Sjöstrand in his 2005 Master's Thesis [1]. Using the pseudocode provided in Sjöstrand's thesis, we first implemented the algorithm, which was written mainly in C with some accompanying Python scripts, then we used four factors of comparison to gauge our implementation against two other open source data compression algorithms. Finally, from these comparisons, we made conclusions about the effectiveness of our implementation.

II. PROPOSED WORK

Our proposed work is intended to be a re-implementation of the BCCBT algorithm proposed by Sjöstrand for data compression. To achieve this goal, we have proposed a three-part effort that involves implementing the BCCBT algorithm, testing its performance on various sized text files, and evaluating its effectiveness compared to state-of-the-art compression algorithms. By undertaking this project, we aim to contribute to the field of data compression and provide insights into the effectiveness of the BCCBT algorithm.

To begin our work, we first configured the necessary software to emulate the BCCBT algorithm. This involved developing a collection of Python programs and a main C program to enable us to accurately simulate the algorithm's behavior. Next, we evaluated the performance of the BCCBT algorithm by conducting tests to

measure its compression effectiveness against two other popular compression algorithms, Huffman and Gzip. To perform these tests, we used a diverse set of six .txt files of varying sizes, ranging from 1MB to 80MB. By comparing the compression ratios and processing speeds of each algorithm, we were able to determine the effectiveness of the BCCBT algorithm.

Finally, we assessed the effectiveness of our implementation by replicating the evaluation process described in the original paper [1]. Using the same metrics as the original work, we measured various compression thresholds to compare our results with his. Overall, this framework allowed for us to accurately develop and test our results seamlessly while still being precise with our scientific approach.

III. EVALUATION

A. Implementation

Our implementation of the BCCBT data compression algorithm closely follows the pseudocode provided in Sjöstrand's thesis, which can be observed in Figure 1. From a high level

PSEUDOCODE OF THE BCCBT ALGORITHM

```

ENCODING
Get the frequency of each symbol from the input stream
Set the frequency table to the frequency of each symbol
Create a complete binary tree using the frequency table
Set the bit codes according to where the symbols are in the tree
While more symbols to read from the input stream
    Read one symbol from the input stream
    Write the symbol's bit code to the bit code stream
    Write the length of the bit code to the level stream
Compress the level stream with a lossless algorithm
Write the frequency table to the output stream
Write the compressed level stream and the bit code stream to the output stream

DECODING
Read the frequency table from the input stream
Create a complete binary tree using the frequency table
Read the compressed level stream from the input stream
Uncompress the compressed level stream
Read the bit code stream from the input stream
While more levels to read from the level stream
    Read one level from the level stream
    Read level bits from bit code stream
    Find the symbol in the complete binary tree using the level and the bit code
    Write the symbol to the output stream

```

Fig. 1: Pseudocode from [1]

perspective, our implementation of the BCCBT data compression algorithm can essentially be split up into two different sections, being the implementation of the Complete Binary Tree section, and the implementation of the Encoding and Decoding of files section. To help guide the understanding of both these sections we have decided to use example 3-11 from the thesis as

Symbol	Frequency
<i>a</i>	32
<i>b</i>	55
<i>c</i>	4
<i>d</i>	19
<i>e</i>	37
<i>f</i>	26
<i>g</i>	9
<i>h</i>	7

Fig. 2: Frequency table from example 3-11 [1]

this example does an excellent job of explaining how the BCCBT algorithm works at a high level [1], and can give the reader a clearer understanding of how we actually implemented the algorithm without the needing to see any code. Finally, at the end of this section, we will briefly discuss the developed software and give reference on how it can be evaluated.

1) *Complete Binary Tree*: To start example 3-11, suppose that we have an alphabet

$$\Sigma = \{a, b, c, d, e, f, g, h\}$$

taken from a source file where each symbol has a frequency given from the table in Figure 2. Looking at the pseudocode, we now want to construct a complete binary tree that will allow us to set the bit codes of each symbol based on its location in the binary tree.

Before we construct the complete binary tree in this example, let's first note some of the important properties of a complete binary tree:

- All levels of the binary tree are completely full except possibly the lowest level.
- The binary tree is filled in from top down and left to right at each level.
- The number of nodes in the binary tree at level n is 2^n ($n \in \mathbb{Z}_9$).

The last of these properties is really only important for Theorem 3-1 in the thesis [1], however we do not actually use this Theorem's optimization technique in our implementation so we will ignore it for now.

Now back to the example, the pseudocode says to construct the complete binary tree using the frequency table. This means that we insert the highest frequency symbol b as the root node, then at the

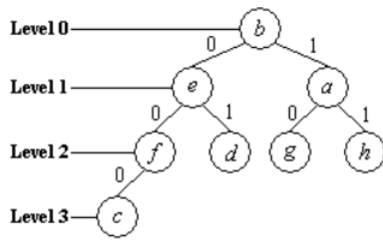


Fig. 3: Complete binary tree from example 3-11 [1]

next level of the binary tree starting from the left we fill in the next highest frequency symbol e , then moving one node to the right at the same level we fill in the next highest frequency symbol a , and so on. Continuing this process until there are no more unique symbols to add into the binary tree will result in the binary tree having the required properties of a complete binary tree, as desired. The complete binary tree generated in this example can be seen in Figure 3.

2) *Encoding and Decoding:* Below we will demonstrate the three main steps in encoding and decoding using the BCCBT algorithm, it is important to note that this is merely a toy example for the purpose of understating.

Firstly, fix the frequency table from Figure 2, where our new frequency table is ordered from the highest occurrences to lowest. This will give a better reference when we begin encoding and decoding strings of text.

Symbol	Frequency
b	55
e	37
a	32
f	26
d	19
g	9
h	7
c	4

Fig. 4: Frequency table for symbols

Note that if we were to traverse the binary tree that we constructed and record the symbols that we encounter as we move from left to right, we would obtain the same sequence of symbols as in the original frequency table, provided that the tree has been constructed correctly. Below is the corresponding complete binary tree that we will use for this demonstration.

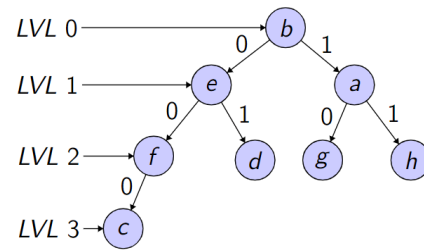


Fig. 5: Complete binary tree

We have now set the bedrock for everything we need to do with our algorithm, whether it be encoding a string, or decoding an encoded string. Before we begin encoding, we first need to establish the last piece of the puzzle, that is, the corresponding bit codes for the symbols in Σ . In order to generate the bitcodes for each node in the tree, we must follow some simple steps to retrieve them. The steps are

- 1) Position yourself at the root node of the tree
- 2) Traverse the tree, keeping note of the weights on every edge that was needed to get to your desired symbol
- 3) After reaching your symbol, you should have a string of 1's and 0's that act like coordinates for finding the specific node.

So, for example, if we wanted to find the bitcode for the symbol h , we start at the root of our tree and traverse two right edges, b to a and then a to h , concatenating the edge weights we saw in our walk we get the bitcode 11 for node h . Lastly, before we move on it is important to note that the root of the binary tree i.e. the most occurring symbol in Σ , takes no bitcode value and instead is set to null. Below is the complete table of bitcodes corresponding to all of our symbols in Σ .

BIT CODES							
a	b	c	d	e	f	g	h
1	NULL	000	01	0	00	10	11

Table of bitcodes

This is good and all, but unfortunately, there is a problem with the current bit codes that we have constructed from the complete binary tree. They are not sufficient to encode a string because they are not uniquely decodable. This means that a stranger looking at the complete binary tree alone would not be able to determine what the following string 111 corresponds to - is it three a 's or one

h followed by an a ? This lack of uniqueness means that we cannot decode the encoded string without losing information, and if we cannot decode it, we cannot expect a computer to be able to do so either. Therefore, we need to modify our bit codes to ensure that they are uniquely decodable.

The way the BCCBT algorithm handles this is quite clever, as it provides a simple yet elegant tactic that will allow us to uniquely decode. What we do is simply take the corresponding level where the symbol node lays in the tree, and prepend that level value to the bitcode, so for example, h had a bitcode 11, but after making this change the bitcode is now [2]11, So with this, there can't be any confusion, as you will have to note that the symbol is on the second level of our tree, and then take the walk given, 1, then 1 again, and then arrive at h ! This is everything we need, and we can now generate uniquely decodable strings that will allow for lossless decoding.

Lets encode the word 'edge'. Using our binary tree we take note of the following 3 bitcodes that are generated by iterating the simple bitcode generation steps.

- 1) $e \rightarrow [1]0$
- 2) $d \rightarrow [2]01$
- 3) $g \rightarrow [2]10$

we now have our encoded string

$$\text{'edge'} \rightarrow [1]0[2]01[2]10[1]0$$

we can now proceed to simulate the decoding of this string.

To decode the encoded string [1]0[2]01[2]10[1]0, we first create a complete binary tree, which we'll need for decoding. The decoding process is almost as simple as encoding. We begin by examining the first level token, [1], and extracting all the bitcodes between this level token and the next one. In this case, we find the code 0 and stop grouping since the next character is a level bracket. This gives us our first encoding block, [1]0. We know with certainty that the desired symbol is in the first level, and we can find it by traversing one 0 edge down. In this case, we land on the symbol e . We repeat this process for the rest of the string, resulting in the following decoding mappings:

- 1) $[1]0 \rightarrow e$
- 2) $[2]01 \rightarrow d$
- 3) $[2]10 \rightarrow g$
- 4) $[1]0 \rightarrow e$

which tells us that our decoded string is

$$[1]0[2]01[2]10[1]0 \rightarrow \text{'edge'}$$

and with that we conclude this thorough example of how the BCCBT algorithm goes about encoding and decoding strings of text.

3) Developed Software:

- **bccbt.c:** This file contains the implementation for both constructing and searching the Complete Binary Tree, which in turn contains the implementation for encoding symbols to character arrays of 1's and 0's, and decoding character arrays of 1's and 0's into their corresponding symbols.
- **bitarray.py:** This file contains the implementation for converting a character array of 1's and 0's to actual bits that can be written to a binary file.
- **bitpull.py:** This file contains the implementation for converting a binary file to a character array of 1's and 0's.
- **Makefile:** This file contains tests with the required sequence of commands needed in order to properly execute our implementation.

Our software can simply be evaluated through the Makefile which contains pre-built tests for various different file sizes, where each of these tests can be edited at the discretion of the user. All of this code can be found in our public GitHub repository where the link can be found in our references [2].

B. Experiments

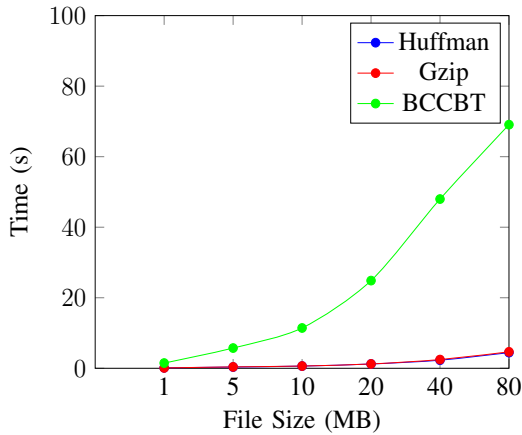
In order to effectively analyze the BCCBT algorithm we had to find two other comparable open source compression algorithms to test against. The factors of comparison that we used to actually compare the three algorithms are as follows:

- 1) Compression Time
- 2) Decompression Time
- 3) Saving % = $\frac{\text{Orig File Size} - \text{Compressed File Size}}{\text{Orig File Size}}$
- 4) Compression Ratio = $\frac{\text{Compressed File Size}}{\text{Original File Size}}$

Using these factors of comparison enabled us to test whether or not the BCCBT algorithm was effective

for compressing files of differing sizes. Note that the following results are given from our implementation of the BCCBT algorithm and may not accurately depict the results of the original implementation.

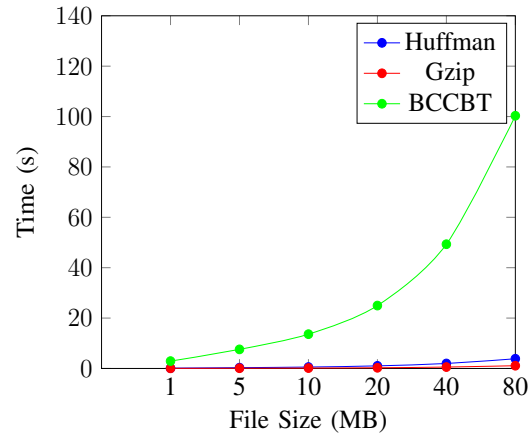
Comparison of Compression times



These first results indicate that although for low file sizes the compression time is quite similar, as the file size gets bigger the compression time exponentially grows and the gap between both Huffman and GZip gets bigger and bigger. This result indicates that it would not be practical, or in your best interest to use our implementation on file sizes above 1MB, as the other two algorithms greatly exceed our implementation's compression time.

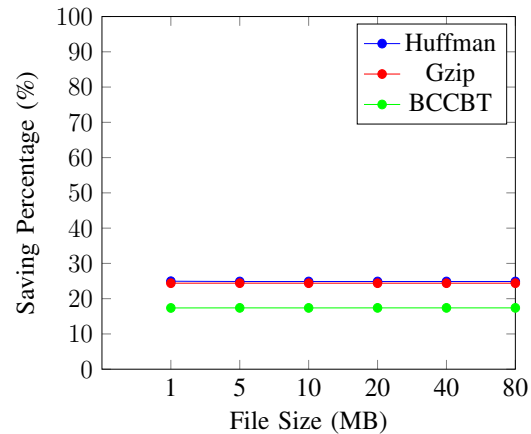
Not all hope is lost for our implementation though. Through our tests we were able to find that the compression times of our implementation of the BCCBT algorithm were slower than GZip and Huffman because of how we populate and traverse through our complete binary tree. In our implementation we are actually traversing the tree recursively rather than iteratively which can heavily impact the run time, making our implementation much slower when compared to the Huffman and Gzip compression times.

Comparison of Decompression times

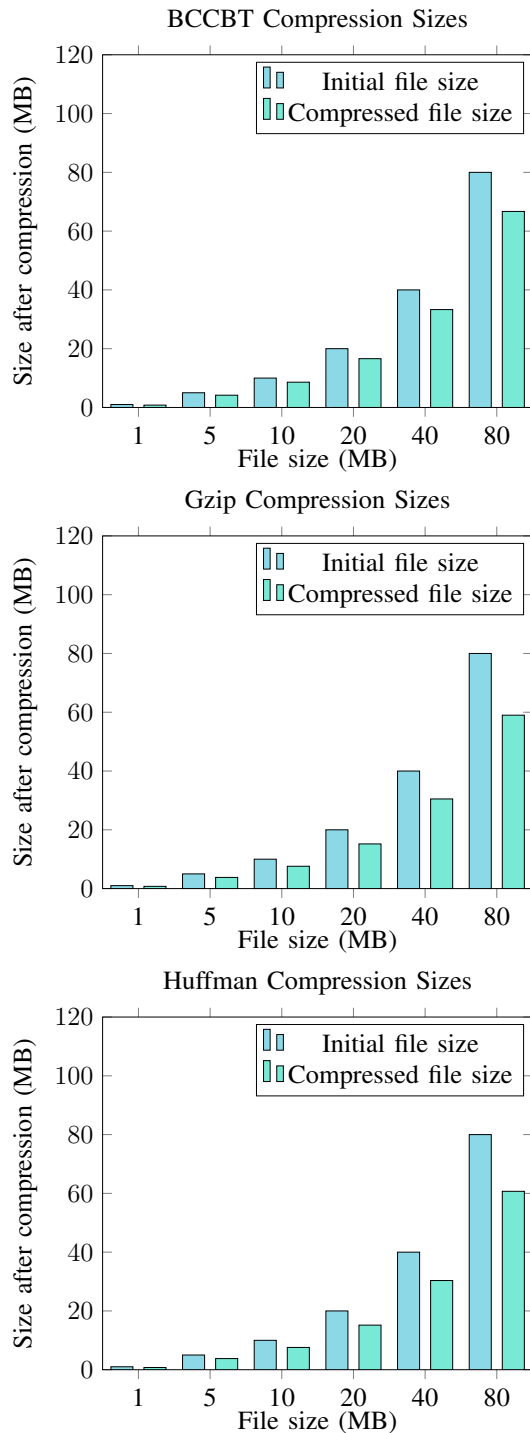


The results for the comparison of the decompression times are quite similar to the compression times mainly because decompression in this algorithm uses very similar functions to the functions used in compression. This graph again illustrates the problem we faced with the optimization of the compression and decompression functions in our implementation of the BCCBT algorithm, while at larger file sizes.

Comparison of Saving Percentage



The graph above demonstrates the saving percentage in %, which shows the comparison of the different algorithms for compression and the overall effectiveness of each algorithm depending on file size. This demonstrates that the saving percentage is linear and that the gap between the algorithms does not exponentially grow like in the compression and decompression times.



The compression size graphs for the Huffman, Gzip, and BCCBT all show that no matter what the file size is, the amount the file size is being reduced stays linear. Another important part to note about this graph is that the BCCBT algorithm does not make the compressed file size as small as the other compression algorithms, further indicating how the algorithm is not as efficient as the other open source algorithms that we compared our implementation against.

In conclusion, based on the experiment results that have been presented, the BCCBT algorithm is an efficient compression algorithm, however, our implementation of the code that compresses and decompresses the bitcodes, levels, and frequency files make the distinction of usability clear. Our implementation of the BCCBT algorithm performs much slower than the two open source algorithms, showing that it is less practical for large file sizes.

IV. FINDINGS AND CONCLUSION

Overall, the results of our experiments show that the BCCBT algorithm itself should only perform slightly worse when compared to the open source compression algorithms that were used in comparison. We must note that the results for the compression and decompression times were heavily influenced by the time constraints of this project as we had to implement the whole algorithm from scratch, and hence little focus was put into the optimization of our implementation. The results for the compression sizes tell a more convincing story as they were quite similar to the other two algorithms used in comparison. Therefore, we can say with confidence that the BCCBT data compression algorithm is actually an effective algorithm that can be used in practical scenarios, given that you optimize the compression and decompression times.

V. PLANNING AND EXECUTION

When comparing the final result of our project with our project proposal, it is clear to see that we were able to effectively achieve our main goal of implementing the BCCBT algorithm in order to compare it against other compression algorithms.

The biggest challenge we faced in our project was simply just the implementation of the BCCBT algorithm since we had to write the compression and decompression code from scratch. This took some time which resulted in code that was not as optimized as it could be. Given more time, or given the opportunity for future work, we do believe that we could yield better results by simply implementing more optimized code.

The work throughout this project was evenly split up between all group members. All group members worked on the implementation of the algorithm

by splitting it up into encoding, decoding, and the setting up of the complete binary tree and frequency table functionalities. The presentation and supporting papers were evenly split up as well, with each group member working on their assigned section representing what they focused on during the implementation and research.

REFERENCES

- [1] M. H. Sjostrand, "A study in compression algorithms," 2005. [Online]. Available: <http://bth.diva-portal.org/smash/record.jsf?pid=diva2:830266>
- [2] "Cpsc 530 project repository." [Online]. Available: <https://github.com/Aidenwjt/CPSC-530-PROJECT>
- [3] "Huffman coding implementation." [Online]. Available: <https://github.com/e-hengirmen/Huffman-Coding>
- [4] "Key value dictionary implementation." [Online]. Available: <https://gist.githubusercontent.com/kylef/86784/raw/fe97567ec9baf5c0dce3c7fcbec948e21dfcce09/dict.c>
- [5] "Binary tree implementation." [Online]. Available: <https://www.programiz.com/dsa/binary-tree>
- [6] "Copy integer array implementation." [Online]. Available: <https://stackoverflow.com/questions/8287109/how-to-copy-one-integer-array-to-another>