

# Source coding for *the Game of Thrones*

Second-order Markov Adaptive Approximation, Huffman and Fano Coding

**Name:** Jinming Ren, Yuhao Liu

**UESTC ID:** 2022190908020, 2022190908022

**UofG ID:** 2840216R, 2840218L

**Date:** November 28, 2024

**University:** UoG-UESTC Joint School

**Location:** No.6 Research Build A240

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Compression is Equivalent to Prediction</b>	<b>4</b>
<b>3</b>	<b>Mechanism of AME Coding</b>	<b>5</b>
<b>4</b>	<b>Entropy Coding Review</b>	<b>6</b>
4.1	Mechanism of Huffman coding . . . . .	6
4.2	Mechanism of Fano coding . . . . .	6
<b>5</b>	<b>Implementation</b>	<b>7</b>
5.1	Without AME . . . . .	7
5.2	With AME . . . . .	8
<b>6</b>	<b>Results and Discussion</b>	<b>9</b>
6.1	Performance comparison between Huffman and Fano coding . . . . .	9
6.2	Performance analysis of AME coding . . . . .	10
<b>7</b>	<b>Conclusion and Future Work</b>	<b>10</b>
<b>A</b>	<b>AME encode/decode algorithm</b>	<b>12</b>
<b>B</b>	<b>Python code for AME encoding</b>	<b>14</b>
<b>C</b>	<b>Python code for AME decoding</b>	<b>16</b>
<b>D</b>	<b>Huffman Tree Generated for original.txt</b>	<b>18</b>
<b>E</b>	<b>MATLAB code for huffman encode/decode procedure</b>	<b>19</b>
<b>F</b>	<b>MATLAB code for huffman encode/decode procedure</b>	<b>21</b>
<b>G</b>	<b>Part of the original.txt</b>	<b>24</b>
<b>H</b>	<b>Part of the binary_huffman.txt</b>	<b>24</b>
<b>I</b>	<b>Part of the processed.txt</b>	<b>24</b>
<b>J</b>	<b>Metrics Definitions</b>	<b>25</b>

# List of Figures

1	The AME tree construction procedure . . . . .	5
2	Flowchart of the experiment without AME . . . . .	7
3	Flowchart of the experiment with AME . . . . .	8
4	Huffman tree visualization . . . . .	18

List of Tables

1	Frequency of each character in <code>original.txt</code> . . . . .	7
3	Performance Metrics for Huffman and Fano Coding . . . . .	10
4	AME saves the length of binary file . . . . .	10

## Abstract

In this article, we successfully implemented huffman and Fano coding to encode/denode the first three chapters of *the Game of Thrones* and evaluated their raw performance in terms of the average code length, code rate, efficiency, and compression ratio in MATLAB. We also presented a new lossless coding scheme called 2nd-order Adaptive Markov Encoding (2nd-ord AME, abbreviated AME) coding and evaluated the overall performance when combined with huffman coding and fano coding.

## 1 Introduction

Source coding, a fundamental technique in data compression, plays a vital role in modern information transmission and storage. Its primary objective is to reduce the number of bits required to represent symbols from a source by leveraging their inherent statistical properties. This efficiency enables effective data storage and transmission while maintaining the integrity of the original information.

In this project, we explore the application of Huffman Coding and Fano Coding, two widely recognized entropy-based compression techniques. Huffman Coding, a lossless algorithm, is designed to assign shorter codes to frequently occurring symbols, making it highly efficient for sources with uneven symbol distributions. On the other hand, Fano Coding, while similar in principle, uses a different approach to partition symbols based on frequency.

Using MATLAB, we implemented both methods to encode and decode the first three chapters of a selected text, generating metrics such as average code length, code rate, efficiency, and compression ratio for comparison. Additionally, we introduced a important change-of-perspective that compression is essentially the same as prediction. Based on this idea, we introduced a novel pre-processing technique, Second-order Adaptive Markov Encoding (AME), designed to enhance the compression performance by uncovering structural patterns in the source text. The efficiency of AME was also evaluated.

## 2 Compression is Equivalent to Prediction

In this section, we state the motivation behind AME coding.

Both huffman coding and fano coding are based on the assumption that the source is *discrete memoryless*. They are trying to approach the compression limit (which is the entropy of the source [1]) given that the original source looks random. So they are called “entropy coding”. However, in practice, the texts inherently take some structures that cannot be explicitly captured by any relative simple models. These structures are not considered when performing entropy coding, so we wasted some potential compression ratio.

No matter what coding method we choose, such as LZ77 [6, 7], the last step them is always entropy coding. But entropy coding is well-understood and mature (Huffman being the “optimal” in some sense). Therefore, the only way to improve compression ratio is to uncover the structure of the source. We need to design a better model first to capture and hiding those structures while exposing the true “memoryless components” of the source, then the utility of entropy coding can be maximized.

The basic idea behind AME coding is that compression is equivalent to *prediction*. We are essentially building a same text predictor in both sides of the transmitter and the receiver. As long as the predicted next character matches the true one, that character is not considered as the “memoryless component” of the source (it depends on the history text). But if the predicted one doesn’t match the true one, this means there are somewhat “random” factors come into play. If we can proposed some encoding strategy that makes the predictor works exactly the same

on both the transmitter and receiver, we only need to transmit the “memoryless components” and the number of correct predictions.

For example, if we want to extract the “memoryless components” of the following content:

1 Information theory is interesting.

Feed it into the predictor, suppose we can correctly predict the character in position 2, 5-11, 15-19, 21-22, 26-34, we only need to transmit the initial letter that cannot be predicted together with the number of correct predictions:

1 I1fo7 th5i2int9

The rest letters are somehow losed some internal dependence and can be approximately considered as memoryless.

We can see in this example the efficiency of compression directly relies on how well we extract the “memoryless components” of the text, in other words, the performance of the predictor. Since human language is the product of human mind, which definitely cannot be modelled using just a few parameters. An accurate predictor would need millions of parameters, which is essentially a neural network [5]. However, due to time and space complexity, using a neural network to predict the next character is not feasible in practice. However, we can use a simpler model called markov chain to capture some of the structures in the text.

### 3 Mechanism of AME Coding

We assume no prior knowledge of the source. So the predictor is built simultaneously with the encode/decode process. The algorithm for AME encoding and decoding are shown in Algorithm 1 and 2 in Appendix A. The codes are written in Python and can be found in the appendix B and C.

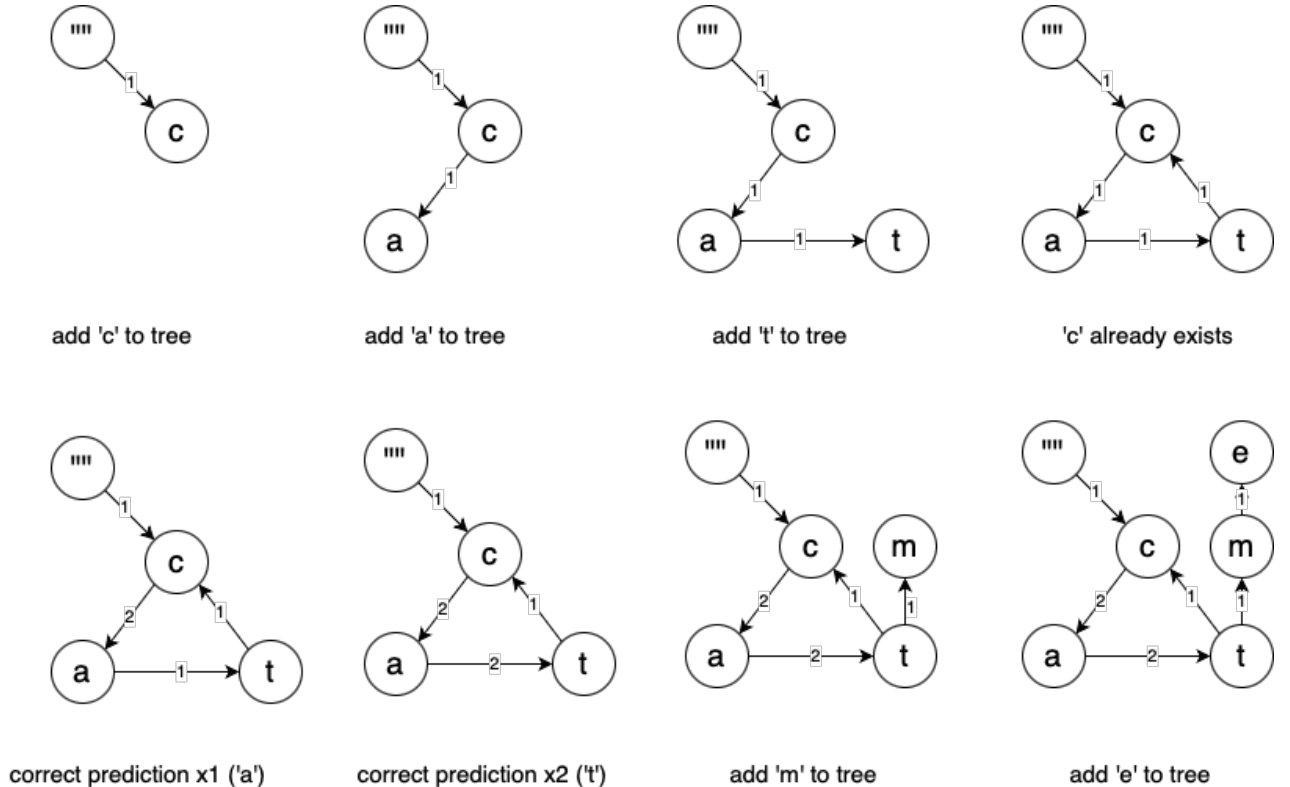


Figure 1: The AME tree construction procedure

We will take an example to show how AME works. Let's consider the following text:

```
1 catcatme
```

As shown in Figure 1. The first character is 'c', so we add 'c' to the AME tree, mark the weight 1 and add it to the output buffer. The next character is 'a', so we add 'a' to the AME tree, mark the weight also 1 and also add it to the output buffer. Repeat this process until we could possibly "predict" the next character. The first correct prediction is the 5th character 'a'. The next prediction ('t') is also right. But the character after that is 'm' and 'e', which are not in the tree, thus we then add them to the tree. Therefore, the output buffer is:

```
1 catc2me
```

It saves one character in this case, but it can be more significant in a larger text. Part of the AME-encoded `original.txt` is shown in Appendix I.

## 4 Entropy Coding Review

### 4.1 Mechanism of Huffman coding

Huffman coding is a lossless data compression algorithm proposed by David A. Huffman in 1952 [1]. It reduces the size of data by assigning variable-length codes to characters based on their frequencies of occurrence [2, 3]. The core principle behind Huffman coding is to assign shorter codes to more frequent characters and longer codes to less frequent ones. The process begins with the construction of a frequency table, which lists each character and its corresponding frequency in the data. From there, a binary tree is built by repeatedly combining the two nodes with the lowest frequencies into a new node. This is the clever part, it builds from the ground up (instead of from the node to the leaves). This process continues until all nodes are merged into a single tree, with the root node representing the entire data set. The structure of this tree enables the generation of optimal, prefix-free binary codes, where each character's code is determined by its position in the tree. In the final step, the original characters in the data are replaced with their corresponding Huffman codes, achieving the desired compression. We perform this process in MATLAB and we visualize the tree in the figure in Appendix D.

### 4.2 Mechanism of Fano coding

Fano coding is another lossless data compression algorithm that assigns variable-length codes to characters based on their frequencies, similar to Huffman coding. The key principle behind Fano coding is to recursively divide the data set into two parts, ensuring that the frequencies of the characters in each part are as balanced as possible, and then assign binary codes accordingly. The process begins with constructing a frequency table, listing each character alongside its frequency. The next step is to split the set of characters into two subsets, aiming to balance the total frequencies of both subsets. Each subset is then assigned a binary digit (0 or 1), and the process is repeated for each subset. As the binary tree is built, characters are assigned codes based on the path from the root to the leaf node representing that character. Finally, the original characters in the data are replaced by their corresponding Fano codes, resulting in compression. Like Huffman coding, Fano coding ensures that the encoded data is more compact without any loss of information. However, its encoding efficiency is quite low when dealing with source symbols that have relatively uniformly distributed probabilities.

All MATLAB codes of the encode/decode process of Huffman and Fano coding can be found in the appendix E and F.

## 5 Implementation

We performed two separated experiments for different purposes. The first is to realize Huffman and Fano coding on the original text and perform evaluations. Given that Huffman coding is the optimal method for generating codes, especially when there are significant differences in the probabilities of occurrence for different source symbols, it should significantly outperform the Fano coding technique in terms of coding effectiveness [3, 4]. The second is to combine AME coding with Huffman and Fano coding and evaluate the performance of the new scheme.

### 5.1 Without AME

The encode/decode process is shown in the form of a flowchart shown in Figure 2.

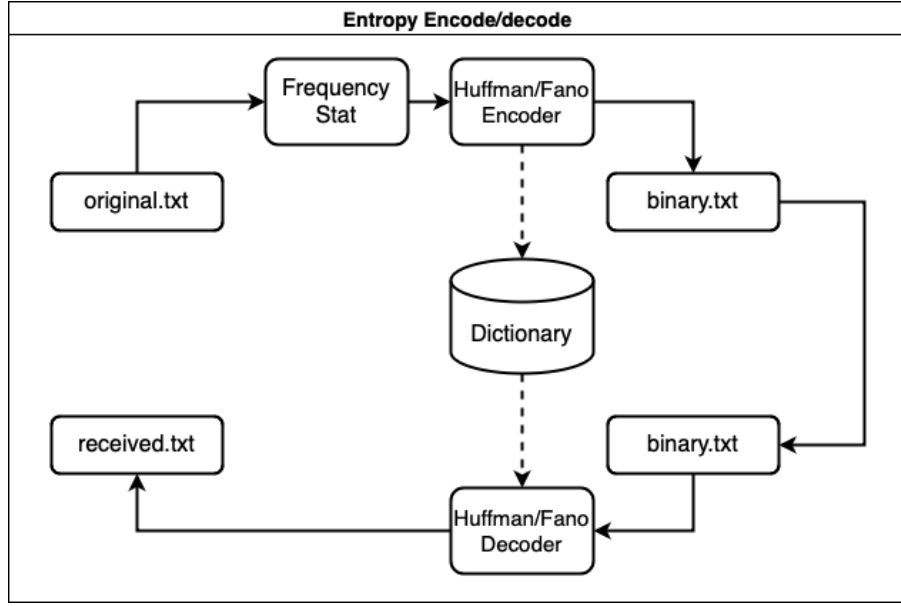


Figure 2: Flowchart of the experiment without AME

The content in `original.txt` contains the first three chapters of *the Game of Thrones*. Part of the `original.txt` is shown in Appendix G.

In order to perform huffman/fano coding on the text, we first need to calculate the frequency of each character in the text. The frequency of each character is shown in Table 1.

Table 1: Frequency of each character in `original.txt`

Index	Symbol	Probability	Index	Symbol	Probability	Index	Symbol	Probability
1	(LF)	0.0052	2	(CR)	0.0052	3	(Space)	0.1915
4	!	0.00020423	5	,	0.0127	6	-	0.001
7	.	0.0165	8	1	2.0423e-05	9	2	2.0423e-05
10	:	4.0846e-05	11	;	0.00024507	12	?	0.002
13	A	0.0013	14	B	0.0017	15	C	0.00055141
16	D	0.00036761	17	E	0.00038803	18	F	0.00079649
19	G	0.0011	20	H	0.0024	21	I	0.0026
22	J	0.001	23	K	0.00020423	24	L	0.00055141
25	M	0.00051057	26	N	0.0011	27	O	0.00046972
28	P	0.00022465	29	R	0.0017	30	S	0.0014

Index	Symbol	Probability	Index	Symbol	Probability	Index	Symbol	Probability
31	T	0.0037	32	U	0.00014296	33	V	8.1691e-05
34	W	0.0034	35	Y	0.00061268	36	a	0.0585
37	b	0.0109	38	c	0.0135	39	d	0.0414
40	e	0.0959	41	f	0.0156	42	g	0.016
43	h	0.052	44	i	0.0435	45	j	0.00061268
46	k	0.0079	47	l	0.035	48	m	0.0154
49	n	0.0476	50	o	0.0546	51	p	0.0079
52	q	0.0004493	53	r	0.0451	54	s	0.0475
55	t	0.0606	56	u	0.0174	57	v	0.0051
58	w	0.0192	59	x	0.00028592	60	y	0.0139
61	z	0.00036761	62	'	0.0026	63	"	0.0051
64	"	0.0051	65	.	4.0846e-05			

According to the frequency of each character, we can create the Huffman/fano tree based on the methods in Section 4.1 and 4.2 and generate the huffman/fano dictionaries<sup>1</sup>. We then use this dictionary to encode the text and calculate the four main metrics including average code length  $\bar{L}$ , code rate  $R$ , efficiency  $\eta$ , and compression ratio  $\xi$ . The results are shown in Section sec:result.

Then we constructed the encoded text in binary form (`binary.txt`). Part of the content of (`binary_huffman.txt`) is shown in Appendix H. Then we send the encoded text (`binary.txt`) to the receiver and decode it using the same dictionary to recover the original text.

## 5.2 With AME

The experiment procedure of the encode/decode process is shown in Figure 3.

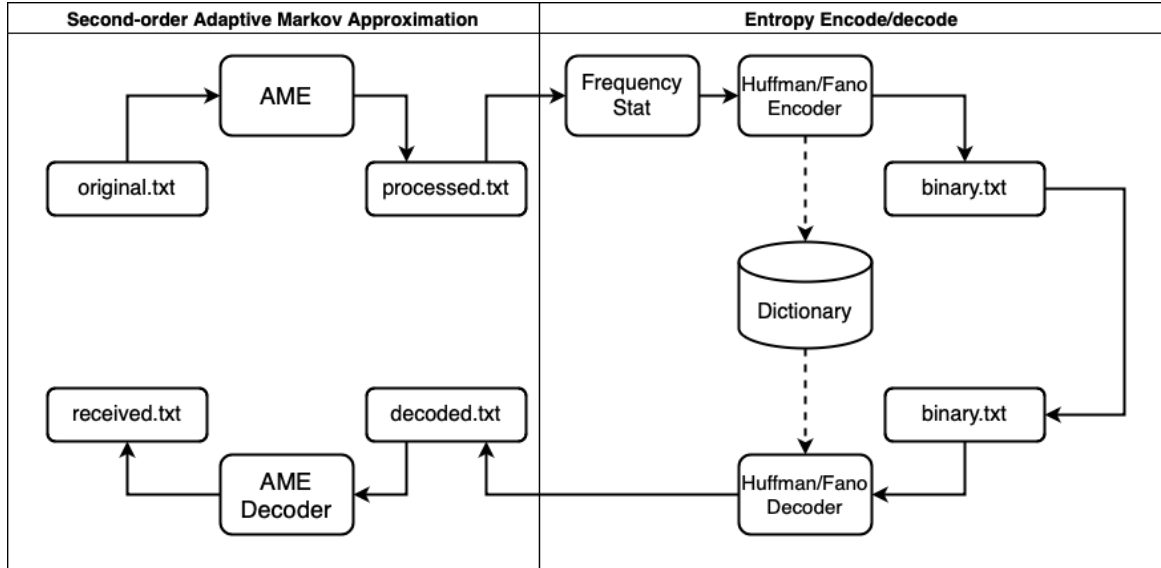


Figure 3: Flowchart of the experiment with AME

The only difference of this part is that we added a block called “Second-order Adaptive Markov Approximation” before applying the entropy encode/decode. Using the method in Sec-

<sup>1</sup>The dictionaries together with all codes can be found on <https://github.com/Marcobisky/ame-entropy-source-coding>



tion 3, we get the AME-encoded text (`processed.txt`). Part of the content of (`processed.txt`) is shown in Appendix I.

Then we apply the same procedure for `processed.txt` as we did for `original.txt`. Note this time the frequency of each character is different from the original text as shown in Table

Index	Symbol	Probability	Index	Symbol	Probability	Index	Symbol	Probability
1	(Space)	2.2477e-05	2	(CR)	0.0057	3	(LF)	0.0652
4	!	0.00022477	5	,	0.014	6	-	0.0011
7	.	0.0182	8	1	0.1721	9	2	0.0344
10	3	0.006	11	4	0.0057	12	5	0.0069
13	6	0.0017	14	7	0.00013486	15	8	2.2477e-05
16	:	4.4954e-05	17	;	0.00026972	18	?	0.0022
19	A	0.0015	20	B	0.0019	21	C	0.00060688
22	D	0.00038211	23	E	0.00042706	24	F	0.0008766
25	G	0.0012	26	H	0.0027	27	I	0.0028
28	J	0.0011	29	K	0.00022477	30	L	0.00060688
31	M	0.00056192	32	N	0.0013	33	O	0.00051697
34	P	0.00024725	35	R	0.0019	36	S	0.0016
37	T	0.0034	38	U	0.00015734	39	V	8.9908e-05
40	W	0.0036	41	Y	0.00067431	42	a	0.0571
43	b	0.0119	44	c	0.0148	45	d	0.0454
46	e	0.0484	47	f	0.0171	48	g	0.0176
49	h	0.0284	50	i	0.0459	51	j	0.00067431
52	k	0.0085	53	l	0.0319	54	m	0.0169
55	n	0.0259	56	o	0.0555	57	p	0.0087
58	q	0.00049449	59	r	0.0477	60	s	0.0501
61	t	0.039	62	u	0.0103	63	v	0.0055
64	w	0.0212	65	x	0.00031468	66	y	0.0153
67	z	0.00040459	68	'	0.0028	69	"	0.0056
70	"	0.0043	71	.	4.4954e-05			

After encoded into binary form, we send it to the receiver and decode it using the same dictionary to get the `decoded.txt`. However this is not the final result, we need to use the AME decoder to recover the original text. The AME decoder decodes the `decoded.txt` with no prior knowledge of the source. Since the markov chain is constructed in the same way, we are guaranteed to get the original text back.

## 6 Results and Discussion

### 6.1 Performance comparison between Huffman and Fano coding

The four main metrics of Huffman and Fano coding are shown in Table 3.

The four metrics are defined in Appendix J. Table 3 shown that Huffman coding outperforms Fano coding in all four metrics, which matches the fact that Huffman coding is the “optimal” coding strategy in the usual sense.

However, we noticed that the compression ratio achieved by Fano coding is greater than 1, which is unreasonable. With the following analysis, we claim it’s because this text is *unsuitable*

Table 3: Performance Metrics for Huffman and Fano Coding

Parameter	Symbol	Huffman	Fano
Average Code Length	$\bar{L}$	4.4769	8.6370
Code Rate	$R$	0.9897	0.5117
Efficiency	$\eta$	0.9897	0.5117
Compression Rate	$\xi$	0.5596	1.0796

for Fano Coding. Compared to other coding methods, the effectiveness of Fano Coding relies heavily on the uneven frequency distribution of symbols. Specifically, Fano Coding performs best when certain symbols appear very frequently while others are relatively rare. This frequency disparity allows shorter codes to be assigned to frequently occurring symbols, enabling compression.

However, in this text, the symbol frequency distribution is relatively uniform. The probability difference between frequently occurring symbols (such as spaces and letters) and rare symbols (such as punctuation marks) is not significant. In such cases of uniform symbol distribution, Fano Coding may fail to significantly reduce the file size. In fact, due to redundancy and long codewords, it may even increase the file size, leading to a compression ratio greater than 1. Moreover, we calculated the variance of the symbol frequencies in this text as follows:

```

1 >> var(probabilities)
2 ans =
3     8.9864e-04

```

It is evident that the variance of the frequencies corresponding to these characters is very small, which further validates the point that the frequency distribution of symbols in this text is relatively uniform.

All above further highlights the superiority of Huffman Coding.

## 6.2 Performance analysis of AME coding

Since we pre-processed the `original.txt` before fed into the entropy coding, we cannot directly compute the four metrics for it. So we compared the length of `binary.txt` file generated in the two methods.

Let the length of the `binary_huffman.txt` and `binary_fano.txt` be  $L_h$  and  $L_f$ , respectively. Table 4 shows how using AME could save the length of the binary file.

Table 4: AME saves the length of binary file

	$L_h$	$L_f$
With AME	211066	212752
Without AME	219213	223058
Saved	96.28%	95.38%

## 7 Conclusion and Future Work

In this project, we successfully implemented the Huffman Coding technique using MATLAB, creating compressed files optimized for transmission and storage. Additionally, by comparing the results with Fano Coding, we explored the efficiency and optimization of Huffman Coding.

Furthermore, we introduced a new lossless coding scheme called 2nd-order Adaptive Markov Encoding (AME) coding, which serves as a pre-processor before applying entropy coding. It aims to uncover the structure of the source text and improve the compression ratio by extracting the “memoryless components” of the text.

Future work could involve using higher order Markov models or new mathematical models to capture more complex structures in the text, potentially furthermore improving the compression ratio.

## References

- [1] C. E. Shannon, “A Mathematical Theory of Communication”.
- [2] Advances in Communication and Computing Technologies (ICACACT 2014), Mumbai, India, 2014, pp. 1-6, doi: 10.1109/EIC.2015.7230711.
- [3] N. Dhawale, “Implementation of Huffman algorithm and study for optimization,” 2014 International Conference on Advances in Communication and Computing Technologies (ICACACT 2014), Mumbai, India, 2014, pp. 1-6, doi: 10.1109/EIC.2015.7230711.
- [4] S. Congero and K. Zeger, “Competitive Advantage of Huffman and Shannon-Fano Codes,” in IEEE Transactions on Information Theory, vol. 70, no. 11, pp. 7581-7598, Nov. 2024, doi: 10.1109/TIT.2024.3417010.
- [5] T. Sharma et al., “A Survey on Machine Learning Techniques for Source Code Analysis,” Sep. 13, 2022, arXiv: arXiv:2110.09610. doi: 10.48550/arXiv.2110.09610.
- [6] J. Ziv and A. Lempel, A universal algorithm for sequential data compression, IEEE Trans. Inform. Theory, vol. 23, no. 3, pp. 337-343, May 1977, doi: 10.1109/TIT.1977.1055714.
- [7] A. D. Wyner and J. Ziv, The sliding-window Lempel-Ziv algorithm is asymptotically optimal, Proc. IEEE, vol. 82, no. 6, pp. 872-877, Jun. 1994, doi: 10.1109/5.286191.

## A AME encode/decode algorithm

---

**Algorithm 1** Adaptive Markov Encoding

---

**Require:** Input file `source.txt`, Output file `markov_encoded.txt`

**Ensure:** Encoded text stored in `markov_encoded.txt`

```
1: Initialize an empty tree tree, list encoded_output, and counter correct_predictions =  
  0  
2: Read source_text from source.txt  
3: for i = 1 to length(source_text) do  
4:   current_char = source_text[i]  
5:   if i == 1 then  
6:     Append current_char to encoded_output  
7:     Add root-to-current_char transition to tree  
8:   else  
9:     Predict next character using tree: prediction = predict_next(tree, prev_char)  
  
10:    if prediction == current_char then  
11:      Increment correct_predictions  
12:    else  
13:      if correct_predictions > 0 then  
14:        Append correct_predictions to encoded_output  
15:        Reset correct_predictions = 0  
16:      end if  
17:      Append current_char to encoded_output  
18:    end if  
19:    Add transition prev_char → current_char to tree  
20:  end if  
21:  prev_char = current_char  
22: end for  
23: if correct_predictions > 0 then  
24:   Append correct_predictions to encoded_output  
25: end if  
26: Write encoded_output to markov_encoded.txt
```

---

---

**Algorithm 2** Adaptive Markov Decoding

---

**Require:** Input file `markov_encoded.txt`, Output file `markov_decoded.txt`

**Ensure:** Decoded text stored in `markov_decoded.txt`

```
1: Initialize an empty tree tree, list decoded_output
2: Read encoded_text from markov_encoded.txt
3: i = 1
4: while i ≤ length(encoded_text) do
5:   current_char = encoded_text[i]
6:   if i == 1 then
7:     Append current_char to decoded_output
8:     Add root-to-current_char transition to tree
9:   else
10:    prev_char = decoded_output[last]
11:    if current_char is a digit then
12:      Extract full number as repeat_count
13:      Predict next character using tree: prediction = predict_next(tree,  

       prev_char)
14:      for j = 1 to repeat_count do
15:        Append prediction to decoded_output
16:        Add transition prev_char → prediction to tree
17:        prev_char = prediction
18:      end for
19:    else
20:      Append current_char to decoded_output
21:      Add transition prev_char → current_char to tree
22:    end if
23:  end if
24:  i = i + 1
25: end while
26: Write decoded_output to markov_decoded.txt
```

---

## B Python code for AME encoding

Listing 1: AME encoder

```
1 # adaptive_markov_encode.py
2 import json
3
4 def add_to_tree(tree, from_char, to_char):
5     """Add a transition to the tree."""
6     if from_char not in tree:
7         tree[from_char] = {"transitions": {}, "highestFrequency": 0, "
probableNextChar": ""}
8
9     transitions = tree[from_char]["transitions"]
10    if to_char in transitions:
11        transitions[to_char] += 1
12    else:
13        transitions[to_char] = 1
14
15    if transitions[to_char] >= tree[from_char]["highestFrequency"]:
16        tree[from_char]["highestFrequency"] = transitions[to_char]
17        tree[from_char]["probableNextChar"] = to_char
18
19 def predict_next(tree, current_char):
20     """Predict the next character based on tree."""
21     if current_char not in tree:
22         return ""
23     return tree[current_char]["probableNextChar"]
24
25 def encode(source_file, output_file):
26     """Encodes a file using second-order Markov approximation."""
27     with open(source_file, "rb") as f:
28         source_text = f.read().decode()
29
30     tree = {}
31     encoded_output = []
32     correct_predictions = 0
33
34     for i, current_char in enumerate(source_text):
35         if i == 0:
36             add_to_tree(tree, "", current_char)
37             encoded_output.append(current_char)
38         else:
39             prev_char = source_text[i - 1]
40             prediction = predict_next(tree, prev_char)
41
42             if prediction == current_char:
43                 correct_predictions += 1
44             else:
45                 if correct_predictions > 0:
46                     encoded_output.append(str(correct_predictions))
```

```
47         correct_predictions = 0
48         encoded_output.append(current_char)
49         add_to_tree(tree, prev_char, current_char)
50
51     if correct_predictions > 0:
52         encoded_output.append(str(correct_predictions))
53
54     with open(output_file, "wb") as f:
55         f.write("".join(encoded_output).encode())
56
57 if __name__ == "__main__":
58     encode("original.txt", "processed.txt")
```

## C Python code for AME decoding

Listing 2: AME decoder

```
1 # adaptive_markov_decode.py
2 import json
3
4 def add_to_tree(tree, from_char, to_char):
5     """Add a transition to the tree."""
6     if from_char not in tree:
7         tree[from_char] = {"transitions": {}, "highestFrequency": 0, "
probableNextChar": ""}
8
9     transitions = tree[from_char]["transitions"]
10    if to_char in transitions:
11        transitions[to_char] += 1
12    else:
13        transitions[to_char] = 1
14
15    if transitions[to_char] >= tree[from_char]["highestFrequency"]:
16        tree[from_char]["highestFrequency"] = transitions[to_char]
17        tree[from_char]["probableNextChar"] = to_char
18
19 def predict_next(tree, current_char):
20     """Predict the next character based on tree."""
21     if current_char not in tree:
22         return ""
23     return tree[current_char]["probableNextChar"]
24
25 def decode(encoded_file, decoded_file):
26     """Decodes a file using second-order Markov approximation."""
27     with open(encoded_file, "rb") as f:
28         encoded_text = f.read().decode()
29
30     tree = {}
31     decoded_output = []
32     i = 0
33
34     while i < len(encoded_text):
35         current_char = encoded_text[i]
36
37         if i == 0:
38             # First character
39             decoded_output.append(current_char)
40             add_to_tree(tree, "", current_char)
41         else:
42             prev_char = decoded_output[-1]
43             if current_char.isdigit():
44                 num_str = ""
45                 while i < len(encoded_text) and encoded_text[i].isdigit():
46                     num_str += encoded_text[i]
```



```

47         i += 1
48     i -= 1
49     repeat_count = int(num_str)
50     for _ in range(repeat_count):
51         prediction = predict_next(tree, prev_char)
52         if not prediction:
53             raise ValueError(f"Decoding error: No valid
prediction for {prev_char}.")
54         add_to_tree(tree, prev_char, prediction)
55         decoded_output.append(prediction)
56         prev_char = prediction
57     else:
58         decoded_output.append(current_char)
59         add_to_tree(tree, prev_char, current_char)
60
61     i += 1
62
63     with open(decoded_file, "wb") as f:
64         f.write("".join(decoded_output).encode())
65
66 if __name__ == "__main__":
67     decode("processed.txt", "received.txt")

```

## D Huffman Tree Generated for original.txt

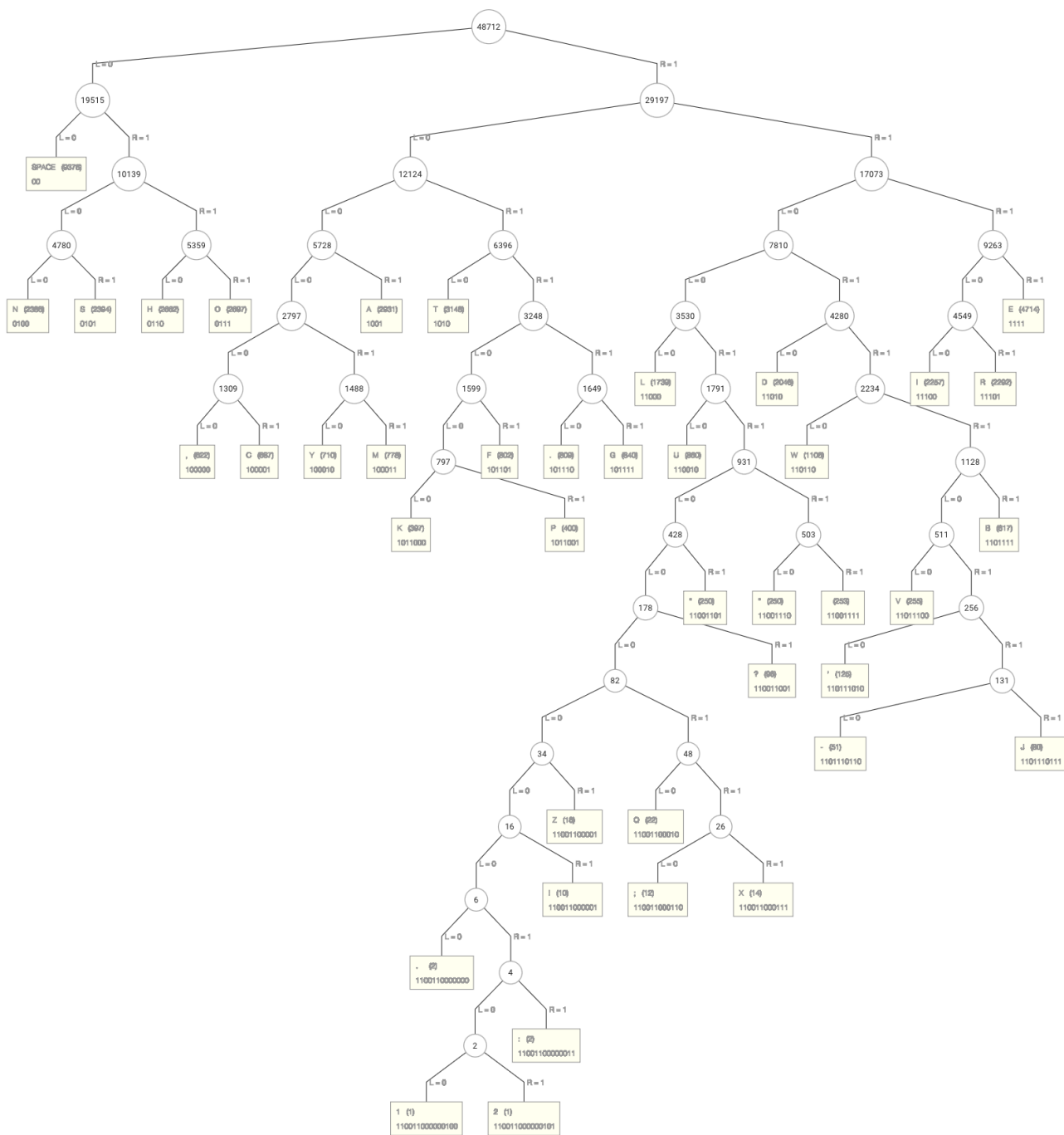


Figure 4: Huffman tree visualization

## E MATLAB code for huffman encode/decode procedure

Listing 3: Huffman encode procedure (without AME example)

```
1 function [encodedMessage, dict] = huffman_encode(text)
2     % compute frequency
3     symbols = unique(text);
4     counts = histc(text, symbols);
5
6     %change symbols into cell
7     symbols = num2cell(symbols);
8
9     % calculate probabilyies
10    probabilities = counts / sum(counts);
11
12
13    dict = huffmandict(symbols, probabilities);
14
15    % encode
16    encodedMessage = huffmanenco(num2cell(text), dict);
17 end
```

Listing 4: Huffman decode procedure (without AME example)

```
1 function decodedMessage = huffman_decode(encodedMessage, dict)
2     % Decode encoded information using Huffman dictionary
3     decodedMessage = huffmandeco(encodedMessage, dict);
4
5     % Converts the decoded cell array back to a character array
6     decodedMessage = cell2mat(decodedMessage);
7 end
```

Listing 5: Main function to perform Huffman encode/decode

```
1 % 1.initialise four different paths
2 inputFilePath = 'original.txt'; % the original text
3 encodedFilePath = 'binary_huffman.txt'; % binary-encoded text file
4 dictFilePath = 'Huffman_dictionary.txt'; % dictionary: Huffman code for each
   character
5 decodeFilePath = 'received_huffman.txt'; % decoded text file
6
7 % 2.read input text
8 fileID = fopen(inputFilePath, 'r');
9 text = fscanf(fileID, '%c');
10 fclose(fileID);
11
12 % 3.encode that original text using Huffman decoding and put that in a .txt
   file
13 [encoded_text, dict] = huffman_encode(text);
14 fileID = fopen(encodedFilePath, 'w');
15 fprintf(fileID, '%d', encoded_text);
16 fclose(fileID);
```

```

17
18 % 4.put the dictionary in a .txt file and display that
19 fileID = fopen(dictFilePath, 'w');
20 for i = 1:length(dict)
21     fprintf(fileID, '%s: %s\n', num2str(cell2mat(dict(i, 1))), num2str(
        cell2mat(dict(i, 2))));
22 end
23 fclose(fileID);
24
25 % 5.decode the encoded text and put that in a .txt file using Huffman decoding
26 decodedMessage = huffman_decode(encoded_text, dict);
27 fileID = fopen(decodeFilePath, 'w');
28 fprintf(fileID, '%s', decodedMessage);
29 fclose(fileID);
30
31 % 6.verify if decoding was successful
32 if isequal(text, decodedMessage)
33     disp('Decoding is successful!');
34 else
35     disp('Decoding failed.');
```

## F MATLAB code for huffman encode/decode procedure

Listing 6: Fano encode procedure (without AME example)

```
1 function [encodedMessage, dict] = fano_encode(text)
2     % Compute symbol frequencies
3     symbols = unique(text); % Unique characters, including spaces
4     counts = histc(text, symbols); % Count occurrences
5
6     % Convert symbols into cell format for compatibility
7     symbols = num2cell(symbols);
8
9     % Calculate probabilities
10    probabilities = counts / sum(counts);
11
12    % Sort symbols by probabilities in descending order
13    [probabilities, idx] = sort(probabilities, 'descend');
14    symbols = symbols(idx);
15
16    % Generate Fano code dictionary
17    dict = fano_create_dict(symbols, probabilities);
18
19    % Encode the text
20    encodedMessage = '';
21    for i = 1:length(text)
22        % Match each character to its code in the dictionary
23        symbol = text(i);
24        code = dict{ismember(dict(:,1), {symbol}), 2};
25        encodedMessage = strcat(encodedMessage, code);
26    end
27 end
28
29 function dict = fano_create_dict(symbols, probabilities)
30     % Initialize dictionary
31     dict = cell(length(symbols), 2);
32     dict(:, 1) = symbols; % Add symbols to the dictionary
33
34     % Recursive function to generate Fano codes
35     function generate_code(subset, code)
36         if numel(subset) == 1
37             % Assign code to the only symbol left
38             dict{ismember(dict(:,1), subset{1}), 2} = code;
39             return;
40         end
41
42         % Find the partition point
43         cumulative = cumsum(probabilities(ismember(symbols, subset)));
44         total = cumulative(end);
45         partition = find(cumulative >= total / 2, 1);
46
47         % Split symbols into two groups and assign 0/1
```

```

48     generate_code(subset(1:partition), [code, '0']);
49     generate_code(subset(partition+1:end), [code, '1']);
50 end
51
52 % Start the recursive encoding
53 generate_code(symbols, '');
54 end

```

Listing 7: Fano decode procedure (without AME example)

```

1 function decodedMessage = fano_decode(encodedMessage, dict)
2     % Decode the encoded binary string using the Fano dictionary
3
4     % Convert the dictionary into a map for fast lookup
5     codeMap = containers.Map(dict(:,2), dict(:,1));
6
7     % Initialize decoding variables
8     decodedMessage = '';
9     buffer = '';
10
11    % Iterate through the encoded message to decode
12    for i = 1:length(encodedMessage)
13        buffer = strcat(buffer, encodedMessage(i)); % Append bit to buffer
14
15        if isKey(codeMap, buffer)
16            % If the buffer matches a code in the dictionary
17            decodedChar = codeMap(buffer);
18            decodedMessage = strcat(decodedMessage, decodedChar);
19            buffer = ''; % Reset the buffer
20        end
21    end
22
23    % Ensure the decoded message matches the original text
24    if ~isempty(buffer)
25        error('Decoding error: incomplete or invalid encoding.');

```

Listing 8: Main function to perform Fano encode/decode

```

1 % 1.initialise four different paths
2 inputFilePath = 'original.txt';% the original text
3 encodedFilePath = 'binary_fano.txt';% binary-encoded text file
4 dictFilePath = 'Fano_dictionary.txt';% dictionary: Huffman code for each
   character
5 decodedFilePath = 'received_fano.txt';% decoded text file
6
7 % 2.read input text
8 fileID = fopen(inputFilePath, 'r');
9 text = fread(fileID, '*char')'; % Read all characters, preserving spaces and
   line endings
10 fclose(fileID);

```

```

11 % replace CRLF with a unique marker to ensure consistent handling
12 crlfMarker = '\r\n'; % Marker for CRLF during encoding/decoding
13 text = strrep(text, char([13 10]), crlfMarker); % Replace \r\n with marker
14
15 % 3.encode that original text using Fano decoding and put that in a .txt file
16 [encodedMessage, dict] = fano_encode(text);
17 % Ensure spaces are correctly encoded and stored in the dictionary
18 dict{strcmp(dict(:, 1), ' '), 1} = '_'; % Map space to underscore
19 fileID = fopen(encodedFilePath, 'w');
20 fprintf(fileID, '%s', encodedMessage); % Write as binary string
21 fclose(fileID);
22
23 % 4.put the Fano dictionary in a .txt file with CRLF format
24 fileID = fopen(dictFilePath, 'w');
25 for i = 1:size(dict, 1)
26     line = sprintf('%s: %s\r\n', char(dict{i, 1}), dict{i, 2}); % Use \r\n for
        CRLF
27     fwrite(fileID, line, 'char'); % Write each line explicitly with CRLF
28 end
29 fclose(fileID);
30
31 % 5.decode the encoded text and put that in a .txt file using Fano decoding
32 decodedMessage = fano_decode(encodedMessage, dict);
33 % ensure underscores are converted back to spaces
34 decodedMessage = strrep(decodedMessage, '_', ' ');
35 % restore original CRLF line endings (without markers in the output file)
36 decodedMessage = strrep(decodedMessage, crlfMarker, char([13 10])); % Replace
        marker with \r\n
37 % write the decoded message to a file
38 fileID = fopen(decodedFilePath, 'w');
39 fwrite(fileID, decodedMessage, 'char'); % Use fwrite to ensure exact output
        including spaces and line endings
40 fclose(fileID);
41
42 % 6.verify if decoding was successful
43 if isequal(strrep(text, crlfMarker, char([13 10])), decodedMessage)
44     disp('Decoding is successful!');
45 else
46     disp('Decoding failed. ');
47 end
48
49 % 7.calculate some relevant parameters and display them
50 [avgCodeLength, rate, efficiency, zip_rate] = calculate_encoding_metrics(text,
    dict);
51 disp('average code length:'); disp(avgCodeLength);
52 disp('code rate:'); disp(rate);
53 disp('efficiency:'); disp(efficiency);
54 disp('zip_rate:'); disp(zip_rate);

```

## G Part of the original.txt

```
1 PROLOGUE
2 We should start back, Gared urged as the woods began to grow dark around
  them. The wildlings are dead.
3 Do the dead frighten you? Ser Waymar Royce asked with just the hint of a
  smile.
4 Gared did not rise to the bait. He was an old man, past fifty, and he had
  seen the lordlings come and go. Dead is dead, he said. We have no business
  with the dead.
5 Are they dead? Royce asked softly. What proof have we? ...
```

## H Part of the binary\_huffman.txt

```
1 110101101011011001110010101101010101001001010010110101010100100100100111110101
2 100111100110010101100101100111111100111010000110001010101000001111010000011011
3 000100110011100011011011001010001111001010101100100110001011101101001001000111
4 000110000001101100111100011010100000000110110111000101101101010000011001110100
5 010000011000010110010100000001010001111011110110100011010100000111000001110110
6 011001110001101011001101110001110000011111011001101101101010000001011101001111
7 001011110101101110100000110011100010001000001100100000100101101010000010110111
8 000110000110011000000111001100100110010111000101011001011001111110010111110011
9 111001100011011010100000110011000000111001101101010100011001000101001010011000
10 001011111001001000001111010110100001011101110011110100000001111100111010111100
11 100010111011100011111001110011000100100100101000011011100010010110000000011011
```

## I Part of the processed.txt

```
1 PROLOGUE
2 1We should start back, Gared1urge2as the woods1began to grow dalk
  a1ound1th1m. Th2wil1lings1a1e1dead.
3 2 1Do t3dead1frighten you? Ser Wayma1 Royce1asked1w1t1 just t3hint of a smile.
4 2 1G2ed1did1not rise1to 4bait. He1was1an1ol2man, past fifty,1
  and1h2had1seen1t3lordlings1come1and1go.11ead1is1dead,1h2said.1
  We1have1no1buli1less1w1t1 t3dead.
5 2 Ar5y dead?2R1yce1ask1d1sof1ly.1What proof hav2we?1
```



## J Metrics Definitions

The four main metrics of entropy coding include average code length  $\bar{L}$ , code rate  $R$ , efficiency  $\eta$ , and compression ratio  $\xi$ . They are defined as follows:

Let  $L_i$  and  $p_i$  be the length and probability of the  $i$ -th symbol, respectively. The average code length  $\bar{L}$  is defined as:

$$\bar{L} = \sum_{i=1}^n L_i \cdot p_i. \quad (1)$$

Let the source entropy be  $H(S)$ . The code rate  $R$  is defined as:

$$R = \frac{H(S)}{\bar{L}}. \quad (2)$$

Let the code be  $n$ -ary. The efficiency  $\eta$  is defined as:

$$\eta = \frac{R}{\log_2(n)} = \frac{H(S)}{\bar{L} \cdot \log_2(n)}. \quad (3)$$

The compression ratio  $\xi$  is defined as (w.r.t. ACSII encoding):

$$\xi = \frac{\bar{L}}{8}. \quad (4)$$