

# TECHNICAL RESEARCH REPORT

## *Enhancing LZW Coding Using a Variable-Length Binary Encoding*

*by T. Acharya, J.F. Já Já*

**T.R. 95-70**



*Sponsored by  
the National Science Foundation  
Engineering Research Center Program,  
the University of Maryland,  
Harvard University,  
and Industry*

# Enhancing LZW coding using a variable-length binary encoding

Tinku Acharya

Joseph F. Já Já

Institute for Systems Research and

Institute for Advanced Computer Studies

University of Maryland, College Park, MD 20742

{*acharya, joseph*}@umiacs.umd.edu

## Abstract

We present here a methodology to enhance the LZW coding for text compression using a variable-length binary encoding scheme. The basic principle of this encoding is based on allocating a set of prefix codes to a set of integers growing dynamically. The prefix property enables unique decoding of a string of elements from this set. We presented the experimental results to show the effectiveness of this variable-length binary encoding scheme.

## 1 Introduction

The two general categories of text compression techniques are *statistical coding* and *dictionary coding*. The statistical coding is based on the statistical probability of occurrence of the characters in the text, *e.g.* Huffman coding [1], Arithmetic coding [2] etc. In dictionary coding, a dictionary of common words is generated such that the common words appearing in the text are replaced by their addresses in the dictionary [3]. Most of the adaptive dictionary based text compression algorithms belong to a family of algorithms originated by Ziv and Lempel [4, 5], popularly known as LZ coding. The basic concept of all the LZ coding algorithms is to replace the substrings with a pointer to where they have occurred earlier in the text [6]. The most popular variation is the LZW algorithm [7]. In this paper, we will describe here a methodology to enhance the compression ratio obtained by LZW algorithm using a novel technique of variable-length binary encoding based on a special binary tree data structure called the **phase in binary tree**.

We describe the LZW algorithm and redundancy in binary encoding of its output in section 2. In section 3, we describe the on-line variable length binary encoding scheme. We use this binary encoding scheme to encode the LZW codes in section 4. We describe the decoding method in section 5. The experimental results are presented in section 6.

## 2 The LZW Algorithm

Let  $S = s_1s_2\cdots s_m$  be a string (or text) over the alphabet  $\Sigma = \{a_1, a_2, \dots, a_q\}$ . The LZW algorithm maps  $S$  into the compressed string  $c(S) = p_1p_2\cdots p_n$ , where  $p_i$  is a positive integer and  $p_i \leq n + q - 1$ , for  $i = 1, \dots, n$ . This mapping can be achieved with the aid of a ‘dictionary trie’. This *dictionary trie* ( $T$ ) is constructed on-line during compression of the text as shown in Figure 1. Each node in  $T$  represents a substring found by concatenating the characters in the label of each node on the path from the root node. Each node is numbered by an integer which is used as a pointer (code value) to replace a matching substring into the text to form the output code.  $T$  is initialized as a  $q + 1$  rooted tree where the root is labeled  $(0, \lambda)$  to represent the null string  $(\lambda)$ . The root has  $q$  children nodes labeled  $(1, a_1), (2, a_2), \dots, (q, a_q)$  respectively to represent  $q$  single character strings. This is shown as  $T_1$  in Figure 1. The LZW algorithm is described below. The input text is examined character by character and the longest substring in the text which already exists in the trie, is replaced by the integer number associated with the node representing the substring in the trie. This matching substring is called a **prefix string**. This prefix string  $(\omega)$  is then extended by the next character  $(K)$  to form a new prefix string  $(\omega K)$ . A child node is created at the node representing the matching substring in the trie to represent the new prefix string.

### begin

```

Initialize the trie with single-character strings;
Initialize  $\omega$  with the first input character;
Loop : Read next input character  $K$ ;
    if no such  $K$  exists (input exhausted) then
        Output the code value of  $\omega$ ;
        EXIT from the Loop;
    end if;
    if  $\omega K$  exists in the trie then
         $\omega \leftarrow \omega K$ ;
    else /* The phrase  $\omega K$  doesn't exist in the trie */
        Output the code value of  $\omega$ ;
        Insert the phrase  $\omega K$  into the trie;
         $\omega \leftarrow K$ ;
    end if;
end Loop;
```

end.

**Example:**  $\Sigma = \{a, b, c\}$ ,  $S = abcabbcbabaaaaa$ .  $c(S) = 1, 2, 3, 4, 5, 7, 1, 10, 11$ .

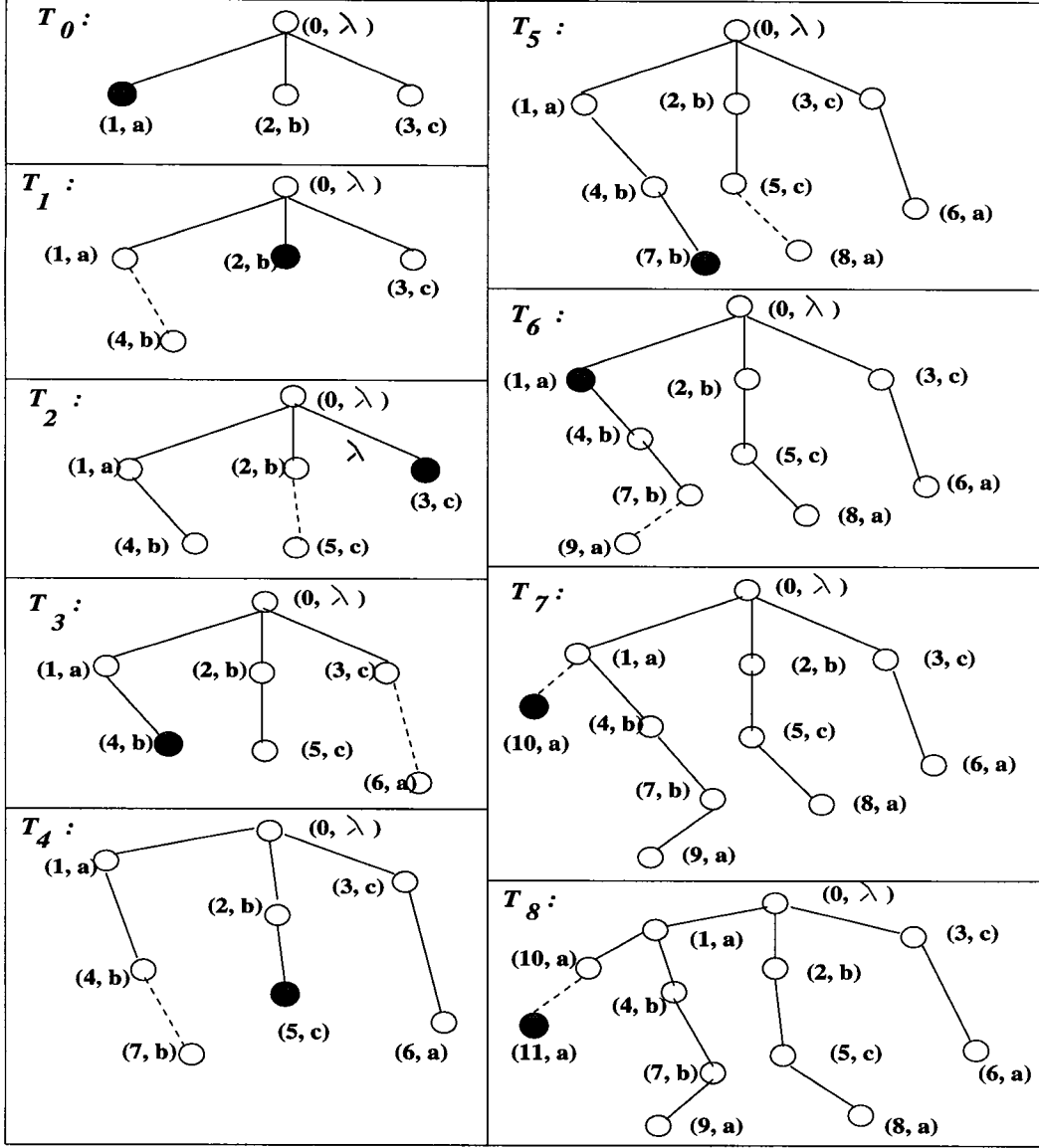


Figure 1: Example of the *dictionary Trie* using LZW algorithm.

The steps and the corresponding trie data structure to encode the above string  $S$  are shown in Figure 1. In each step  $i \geq 0$ , we output the pointer  $p_i$  which we find from the label associated with the dark node in trie  $T_i$ .  $T_i$  is then modified to  $T_{i+1}$  by inserting a node to represent the new prefix string which is shown by the node lead by the dotted line in  $T_{i+1}$ . Also note that the number of nodes in  $T_i$  is always  $i + q$  and the pointer  $p_i$  to output at step  $i$  is  $1 \leq p_i \leq i + q$ . Since the final dictionary trie  $T_8$  contains 11 prefix strings (represented by 11 nodes), each pointer will be encoded by  $\lceil 11 \rceil = 4$  bits each. As a result, the size of the compressed string  $c(S)$  is 36 bits.

The size of the pointer in LZW algorithm is predefined and hence the same fixed number of bits are output in each LZW code irrespective of the number of entries in the dictionary. As a result a large number of bits are used unnecessarily when the number of phrases into the dictionary is less than half of its maximum size. In the LZC algorithm (which is a variant of the LZW algorithm used in UNIX-Z compression), string numbers are output in binary [8]. The number of bits used to represent the string number in any step varies according to the number of phrases (say  $M$ ) currently contained in the dictionary. For example, when the value of  $M$  is in the range 256 to 511 each string number can be represented using a 9 bits binary number and when  $M$  becomes in the range 512 to 1023, each string number will be represented as a 10 bits binary number and so on. The scheme is also sub-optimal because a fractional number of bits is still wasted unless  $M$  is a power of 2. For example, when  $M$  is 513, it is possible to encode the first 512 string numbers in the dictionary (0 through 511) using 10 bits while 512 and 513 can still be represented as 2-bit binary numbers ‘10’ and ‘11’ respectively, without violating the prefix property because the first 512 binary codes will start with the bit 0 and the last two codes starts with the binary bit 1. We will describe this methodology in the following sections.

### 3 A Variable Length Binary Encoding

The variable length binary encoding is based on a binary tree called the **phase in binary tree** as defined below.

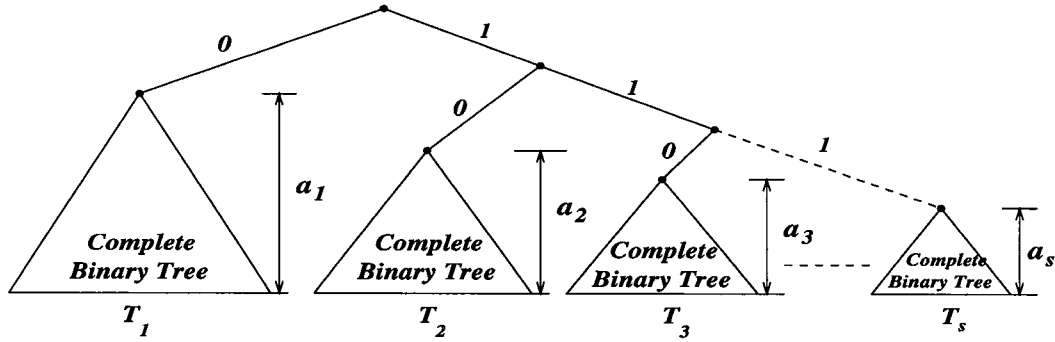


Figure 2: A phase in binary tree.

**Definition 1:** Given a positive integer  $N$  expressed as  $N = \sum_{m=1}^s 2^{a_m}$ , where  $a_1 > a_2 > \dots > a_s \geq 0$  and  $s \geq 1$ , a **phase in binary tree**  $T$  with  $N$  leaf nodes is a rooted binary tree that consists of  $s$  **complete binary trees**<sup>1</sup>  $T_1, T_2, \dots, T_s$  of heights  $a_1, a_2, \dots, a_s$

<sup>1</sup>A **complete binary tree** of height  $h$  is a binary tree with  $2^h$  leaf nodes and  $2^h - 1$  internal nodes. A complete binary tree is represented by a solid triangle in this paper.

respectively, arranged as shown in Figure 2.

**Lemma 1:** A phase in binary tree  $T$  with  $N$  leaf nodes is unique and it contains  $N - 1$  internal nodes.

**Proof :** The structure of the phase in binary tree with  $N$  leaf nodes is based on the representation of the positive integer  $N$  as

$$N = \sum_{m=1}^s 2^{a_m}$$

where  $a_1 > a_2 > \dots > a_s \geq 0$  and  $s > 0$ . Expression of  $N$  suggests that  $a_1, a_2, \dots, a_s$  are the positions of bit 1 when  $N$  is represented as a normalized unsigned binary number using  $a_1 + 1$  bits. Since the representation of an integer in the form of a normalized unsigned binary number is always unique, the structure of the phase in binary tree with  $N$  leaf nodes is also unique.

According to the definition of a phase in binary tree,  $T$  consists of  $s$  complete binary subtrees  $T_1, T_2, \dots, T_s$  as shown in Figure 2 and the number of nodes in the complete binary tree  $T_m$  is  $2^{a_m+1} - 1$ , where  $a_m$  is the height of  $T_m$  for every  $1 \leq m \leq s$ . Hence sum of the number of nodes of all the above complete binary subtrees is

$$\sum_{m=1}^s (2^{a_m+1} - 1) = 2 \sum_{m=1}^s 2^{a_m} - s = 2N - s$$

The remaining nodes in  $T$  form the path starting from the root node of  $T$  to the root node of  $T_{s-1}$ , consisting of  $s - 1$  nodes including the root node of  $T$  as shown in Figure 2. As a result, total number of nodes in the phase in binary tree  $T$  is  $2N - s + (s - 1) = 2N - 1$ . Since  $T$  has  $N$  leaf nodes, the number of internal nodes in  $T$  including the root node is  $2N - 1 - N = N - 1$ .  $\square$

**Definition 2:** Given a set of  $n$  binary codes  $V = \{V_1, V_2, \dots, V_n\}$ , such that  $|V_i| \leq \lceil \log_2 n \rceil$  for every  $1 \leq i \leq n$ , the *value*( $V_i$ ) of the binary code  $V_i \in V$  is defined by the decimal value of  $V_i$  formed by appending  $\lceil \log_2 n \rceil - |V_i|$  number of 0's at the end of  $V_i$ .

**Definition 3:** Given a phase in binary tree  $T$  with  $n$  leaf nodes, we can associate a unique set of binary codes  $V = \{V_1, V_2, \dots, V_n\}$  with the leaves by labeling every left edge of  $T$  by 0 and every right edge by 1 such that  $V_i$  is the sequence of 0's and 1's in the unique path from the root node to the  $i$ th leaf node of  $T$ , where the leaves are indexed consecutively from left to right. This set of binary codes is called a **phase in binary encoding** of size  $n$ .

**Corollary 1:** For a phase in binary encoding  $V = \{V_1, V_2, \dots, V_n\}$  of size  $n$ , the following relations always holds :  $value(V_i) < value(V_{i+1})$  and  $|V_i| \geq |V_{i+1}|$  for every  $1 \leq i < n$ .

**Corollary 2:** If  $V^a$  and  $V^b$  are two phase in binary encodings of sizes  $n_1$  and  $n_2$  respectively such that  $n_1 < n_2$ ,  $|V_i^a| \leq |V_i^b|$  and  $|V_i^a| = |V_i^b| \Rightarrow V_i^a = V_i^b$  for every  $1 \leq i \leq n_1$ .

Since each code in the phase in binary encoding is represented by a leaf node of the corresponding phase in binary tree, the phase in binary encoding forms a set of uniquely decodable prefix codes.

## 4 Allocation of the Phase in Binary Codes in LZW output

In each step of the LZW algorithm, only one phrase may be inserted into the dictionary, *i.e.* only one pointer is appended into the pointer list. We can express the set of pointers in step  $i$  as a sequence of  $q + i$  integers  $P^i = \{1, 2, \dots, q + i\}$  in increasing order. The sequence  $P^i$  of length  $q + i$  can be mapped into the leaf nodes of a phase in binary tree with  $q + i$  leaf nodes. Hence each pointer in step  $i$  can be encoded using a uniquely decodable prefix code. The interesting property of this encoding is that no code of the first  $q + i$  integers in step  $i + 1$  is a prefix of the codes in step  $i$ . As a result, the LZW codes (*i.e.* the pointers) in every step can be encoded by the phase in binary codes to form a variable-length binary encoding of text. We first show with the same example how to assign the phase in binary codes to the pointers in each step of the compression. We later describe a decoding algorithm that uniquely generates the original LZW codes. The steps are shown in Figure 3. The single character strings are represented by the *dictionary trie*  $T_0$  with the set of pointers  $P^0 = \{1, 2, 3\}$  as shown in Figure 3. These pointers can be mapped into the variable length binary encoding  $C^0 = \{00, 01, 1\}$  represented by the phase in binary tree  $B_0$ . In the next step the set of pointers is  $P^1 = \{1, 2, 3, 4\}$ , which can be mapped into the set of binary encodings  $C^1 = \{00, 01, 10, 11\}$  represented by the binary tree  $B_1$  as shown in Figure 3. Hence the output pointer 2 can be encoded using two bits ‘01’. Following the same procedure the sets of pointers  $P^2, P^3, \dots, P^8$  can be mapped into the sets of binary encodings  $C^2, C^3, \dots, C^8$  obtained from the phase in binary trees  $B_2, B_3, \dots, B_8$  respectively as shown in Figure 3. Accordingly, the output pointers 3, 4, 5, 7, 1, 10 and 11 in the next seven consecutive steps can be encoded as ‘010’, ‘011’, ‘100’, ‘110’, ‘0000’, ‘11’ and ‘11’ respectively. Hence the compressed string  $c(S)$  can be encoded as 00 01 010 011 100 110 0000 11 11 using 24 bits, instead of 36 bits using the fixed-length encoding.

## 5 Decoding of the Binary Codes

During the decompression process using the LZW algorithm, logically the same trie is created as in the compression process. But generation and interpretation of the binary tree is a little tricky although the same binary tree will be essentially generated as in the compression process. To decompress the binary codes, we start with the same phase in binary tree  $B_0$  and the corresponding trie  $T_0$  as in the compression process. In each step  $i$ , a string of characters (say  $S_i$ ) is regenerated. In step  $i$ , the phase in binary tree  $B_i$  is traversed starting from the root node. The left child (or the right child) is traversed if the next input bit is 0 (or 1) until a leaf node in  $B_i$  is reached. The label of this leaf node in  $B_i$  is the pointer to a node in  $T_i$  which represents the regenerated substring  $S_i$  to output. In step  $i$ , the phrase formed by concatenating the previous output substring  $S_{i-1}$  (only exception during the first step,  $i = 0$  and  $S_{-1} = \lambda$ ) and the first character (say  $K$ ) of the present output substring  $S_i$ , *i.e.*  $S_{i-1}K$  is inserted into the trie  $T_i$  to create the new trie  $T_{i+1}$ . If the substring  $S_{i-1}K$  already exists in  $T_i$ , the trie  $T_{i+1}$  will be identical to  $T_i$ . The new node (if any) in  $T_{i+1}$  will be represented by the pointer  $q + i$  in  $T_{i+1}$ . The phase in binary tree  $B_i$  is then incremented to  $B_{i+1}$  and the new leaf node  $B_{i+1}$  is labeled by the integer  $q + i + 1$ . The same procedure is repeated until all the bits in binary code  $c(S)$  have been exhausted.

It should be mentioned here that the decoding operation in LZW is handled differently in a special case. This special case occurs whenever the input string contains a substring of the form  $K\omega K\omega K$ , where  $K\omega$  already appears in the trie. Here  $K$  is a single character from the alphabet and  $\omega$  is a prefix string. This is explained in detail in the original paper of LZW [7]. If the original input string ( $S$ ) contains a substring of this form, during the decoding operation of our proposed binary encoding, we will find that the decoded label from the binary tree is not yet created in the trie in the corresponding step (say step  $j$ ). In this case, we know that the output string  $S_{j-1}$  (obtained in the immediate previous step  $j - 1$ ) was of the form  $K\omega$ . Hence we need to output a substring  $S_j$  of the form  $K\omega K$  in the present step and insert the prefix string  $S_j = K\omega K$  into the trie  $T_j$  to form  $T_{j+1}$  and the corresponding node representing  $S_j$  in  $T_{j+1}$  will be marked by  $q + j$ .

We describe the above decoding process using the same example to decode the binary code  $c(S) = 00\ 01\ 010\ 011\ 100\ 110\ 0000\ 11\ 11$ . This binary code  $c(S)$  was obtained by compressing  $S = abcabbcbbaaaaaa$  using the proposed phase in binary encoding scheme. We start with the phase in binary tree  $B_0$  with 3 leaf nodes and the trie  $T_0$  consisting of the single character strings as shown in Figure 4. Here the alphabet is  $\Sigma = \{a, b, c\}$  and hence size of  $\Sigma$  is  $q = 3$ . In the first step ( $i = 0$ ), we traverse the phase in binary tree  $B_0$  starting



from the root node and reach to the leaf node of label 1 after reading the first two input bits 00. This label 1 indicates a pointer in trie  $T_0$  which represents the output substring  $S_0='a'$ . The binary tree  $B_0$  is modified to  $B_1$  to represent a phase in binary tree with 4 leaf nodes and the new leaf node is now labeled by  $q + i + 1 = 4$  as shown in the phase in binary tree  $B_1$  in Figure 4. In the next step ( $i = 1$ ), the  $B_1$  is traversed and the next two input bits 01 lead to the leaf node 2 in  $B_1$ . The node of label 2 in trie  $T_1$  represents the output string  $S_1='b'$ . This character concatenated with the previous output substring 'a' forms the new prefix string 'ab' which is inserted into the trie to form  $T_2$ . This is shown by the broken edge in  $T_2$  in Figure 4. The new pointer to the prefix string 'ab' is  $q + i = 4$ . The phase in binary tree  $B_1$  is now incremented to form the phase in binary tree  $B_2$  with 5 leaf nodes. The new leaf node is labeled by  $q + i + 1 = 5$  as shown in  $B_2$  in Figure 4. In the following step ( $i = 2$ ),  $B_2$  is traversed and the next three input bits 010 lead to the leaf node of label 3 which represents the output substring  $S_2='c'$  in  $T_2$ . This is concatenated with the previous output substring  $S_1='b'$  to form the new prefix string 'bc' and inserted into the dictionary trie to form  $T_3$ . The new pointer to the prefix string 'bc' in  $T_3$  is labeled by  $q + i = 5$ . The phase in binary tree  $B_2$  is now incremented to form the phase in binary tree  $B_3$  with 6 leaf nodes. The new leaf node is labeled by  $q + i + 1 = 6$  in  $B_3$ .  $B_3$  is then traversed and the next three input bits 011 lead to the leaf node of label 4 which represents the output substring  $S_3='ab'$  as shown in  $T_3$ . The first character( $a$ ) of this output substring is concatenated with the previous output string  $S_2='c'$  to form the new prefix string 'ca' which is inserted into the trie to form  $T_4$ . The pointer of this prefix string 'ca' in  $T_4$  is now  $q + i + 1 = 6$ . The phase in binary tree  $B_3$  is now incremented to the phase in binary tree  $B_4$  with 7 leaf nodes and the new leaf node is labeled by  $q + i + 1 = 7$ . Following the same procedure in next three steps (*i.e.*  $i = 4, 5, 6$ ), we can regenerate the decoded output substrings  $S_4='bc'$ ,  $S_5='abb'$  and  $S_6='a'$  for the binary sequences 100, 110 and 0000 from the phase in binary trees  $B_4$ ,  $B_5$ ,  $B_6$  and the corresponding tries  $T_4$ ,  $T_5$  and  $T_6$  respectively. In step 7 ( $i = 7$ ), we will traverse the binary tree  $B_7$  and the next two input bits 11 will lead to the leaf node marked by the integer 10 in  $B_7$ . But there is no node in the trie  $T_7$  labeled by the pointer 10. This arises due to the special case of appearance of a substring  $K\omega K\omega K$ , as we described above. Hence the output substring should be of the form  $K\omega K$ . In this case,  $K\omega$  is the output substring  $S_6='a'$  in the previous step. As a result,  $K = a$  and  $\omega = \lambda$  here. So we output the substring  $S_7=K\omega K = aa$  and insert the prefix string 'aa' into the trie represented by the node labeled by the pointer  $q + i = 10$  to form  $T_8$ . The same situation arises, for the next two input bits 11. By traversing the tree  $B_8$  using the bits 11, we reach to the leaf node of label 11. But there is no pointer 11 in trie  $T_8$ . Hence the special case arises, where  $K = a$

and  $\omega = a$ . Hence output substring will be  $S_8 = K\omega K = aaa$ . Since all the input bits are exhausted the decoding process stops. Now concatenating all the output substrings above, we regenerate the original string  $S = abcabbcbabbaaaaaa$ .

*(Sizes are expressed in nearest Kilo-bytes)*

TEXT	Original	LZW	LZWAJ
<b>bib</b>	<b>111</b>	<b>64</b>	<b>45</b>
<b>book1</b>	<b>768</b>	<b>446</b>	<b>346</b>
<b>book2</b>	<b>610</b>	<b>346</b>	<b>259</b>
<b>geo</b>	<b>102</b>	<b>84</b>	<b>77</b>
<b>news</b>	<b>377</b>	<b>246</b>	<b>188</b>
<b>obj1</b>	<b>21</b>	<b>14</b>	<b>13</b>
<b>obj2</b>	<b>246</b>	<b>143</b>	<b>123</b>
<b>paper1</b>	<b>53</b>	<b>31</b>	<b>24</b>
<b>paper2</b>	<b>82</b>	<b>45</b>	<b>35</b>
<b>prog</b>	<b>39</b>	<b>22</b>	<b>18</b>
<b>progl</b>	<b>71</b>	<b>32</b>	<b>25</b>
<b>progp</b>	<b>49</b>	<b>19</b>	<b>18</b>
<b>trans</b>	<b>93</b>	<b>50</b>	<b>36</b>
<b>Average</b>	<b>202</b>	<b>119</b>	<b>93</b>

## 6 Experimental Results

We have implemented our scheme and tested it with texts of different sizes and characteristics. In all the cases, the phase in binary encoding method significantly improves the performance of the raw LZW technique with fixed-length pointer size of 12 bits and starting the dictionary all over again after it is full. We have performed our experiment with different maximum allowable height of the binary tree. The best performance is achieved when the binary tree is allowed to grow to a maximum height of 15. In our implementation, we allowed the binary tree to grow until it becomes a complete binary tree of height 15 and start all over again after that. We presented the experimental results in the above table to compare the compression performance based on the LZW method and our proposed variable-length encoding scheme which we call the LZWAJ scheme. The experiment was performed on the files obtained from the University of Calgary text corpus in addition to several other text files. The experimental results of our scheme illustrate the significant improvements obtained

by using our scheme over the LZW algorithm.

## 7 Conclusion

In this paper, we have presented a novel methodology for enhancing the performance of the LZW text compression algorithm. We have implemented the scheme and presented the experimental results. The experimental results show that we can achieve much better compression than these obtained by using the standard LZW algorithm.

## References

- [1] D. Huffman, 'A Method for the Construction of Minimum Redundancy Codes,' *Proc. IRE*, Vol. 40, 1952, pp. 1098-1101.
- [2] I.H. Witten, R. Neal and J.G. Cleary, 'Arithmetic coding for data compression,' *Communication of the ACM*, 30(6), 520-540, June 1987.
- [3] J.A. Storer, 'Data Compression: Methods and theory,' *Computer Science Press*, Rockville, MD, 1988.
- [4] J. Ziv and A. Lempel, 'A Universal Algorithm for Sequential Data Compression,' *IEEE Trans. on Info. Theory*, IT-23, 3, May 1977, pp. 337-343.
- [5] J. Ziv and A. Lempel, 'Compression of Individual Sequences Via Variable-rate Coding,' *IEEE Trans. Info. Theory*, IT-24,5, Sept. 1978, pp. 530-536.
- [6] Bell, T. C., Cleary, J. G. and Witten, I. H., 'Text Compression,' *Prentice Hall*, NJ, 1990.
- [7] T. Welch, 'A Technique for High-Performance Data Compression,' *IEEE Computer*, 17 (6), 8-19, June 1984, pp. 8-19.
- [8] R.N. Horspool, 'Improving LZW,' *Data Compression Conference*, 1991, pp. 332-341.

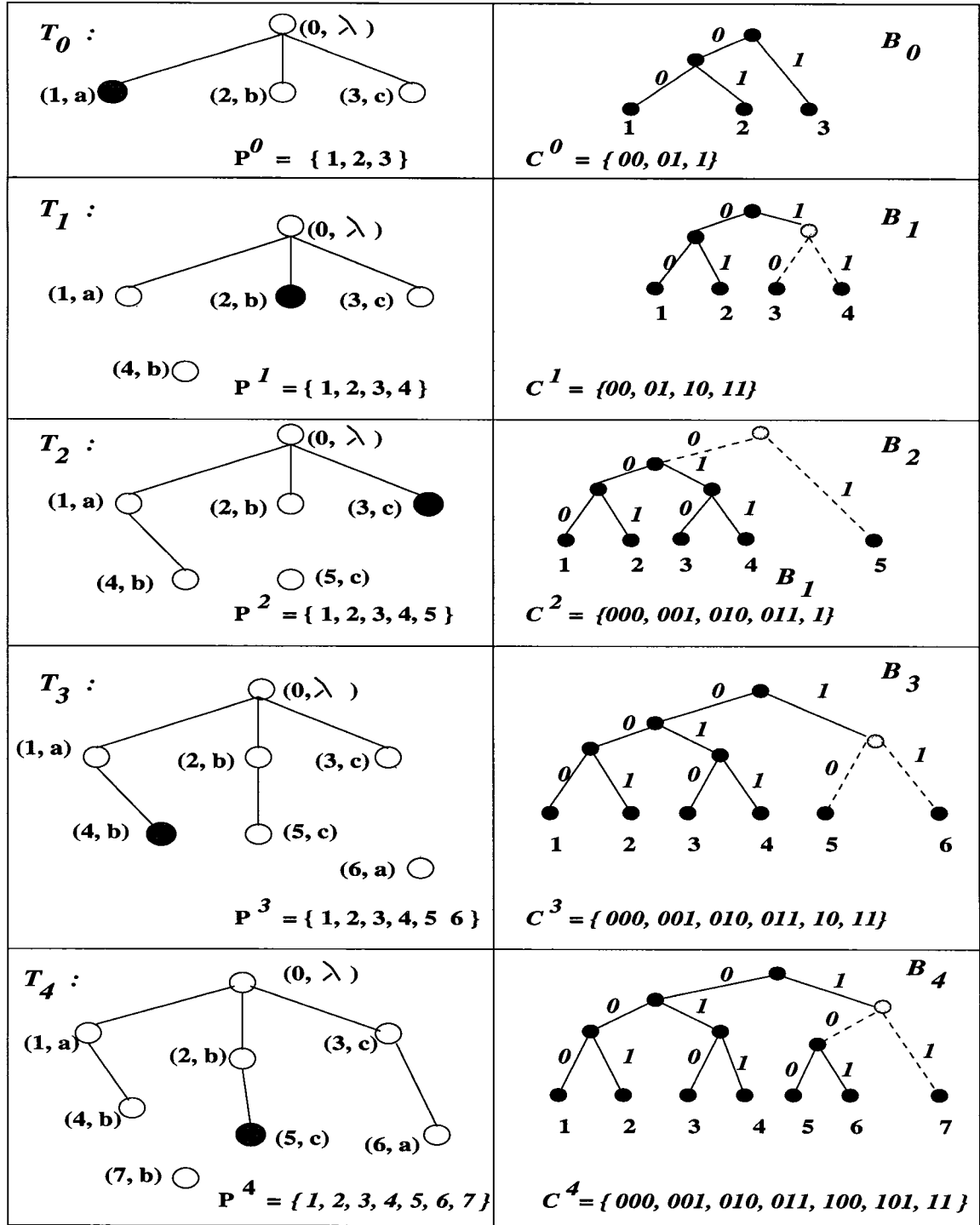


Figure 3: Example of the LZWAJ coding (Continued in the next page also).



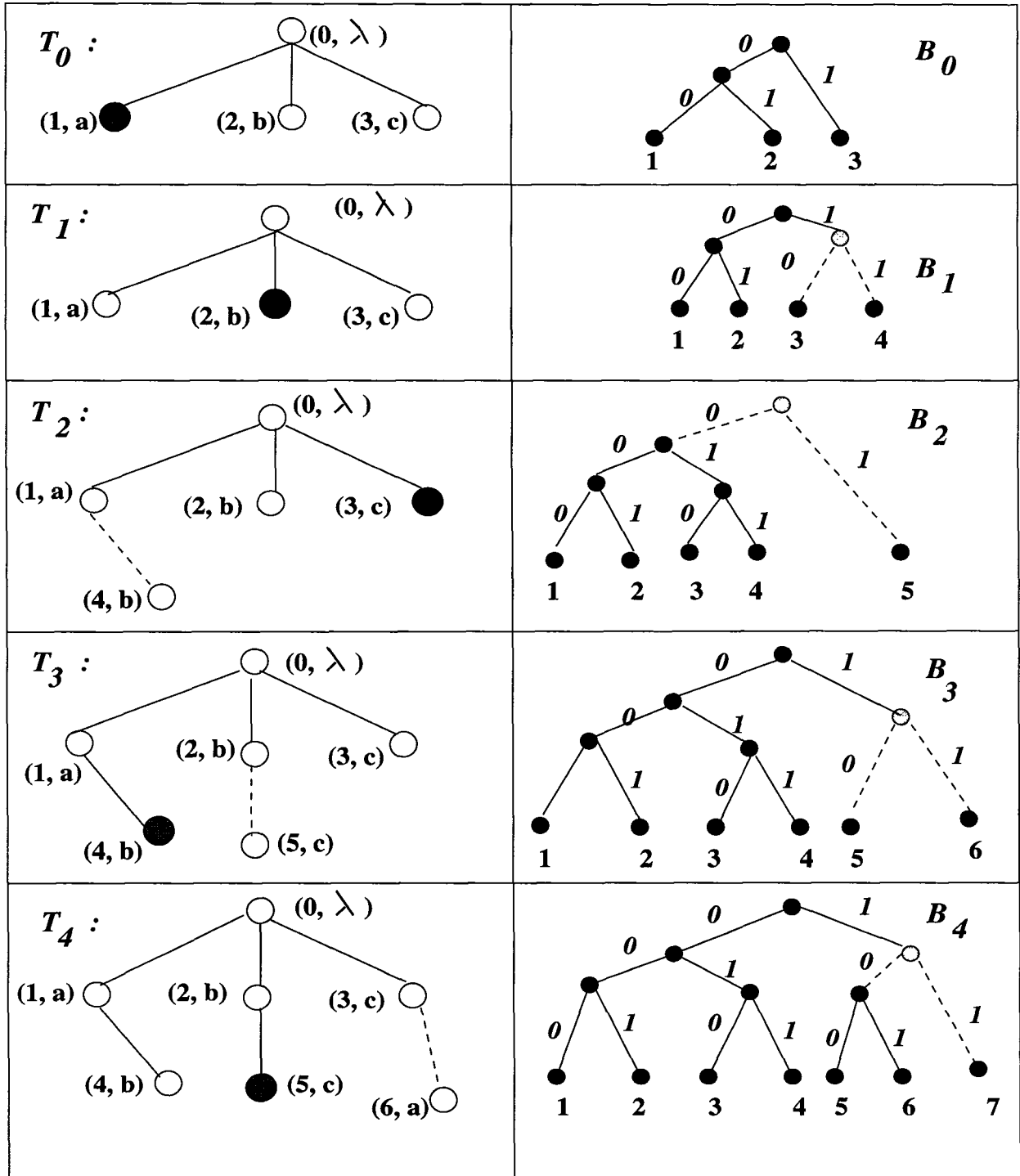


Figure 4: Example of the LZWAJ decoding of binary codes.

