

# Analysis and Practice of Uniquely Decodable One-to-One Code

Chin-Fu Liu, Hsiao-feng (Francis) Lu, and Po-Ning Chen

Dept. of Electrical & Computer Engineering

National Chiao Tung University, Taiwan 300, ROC

Email: hubert.liu.1031@gmail.com, francis@mail.nctu.edu.tw, poning@faculty.nctu.edu.tw

**Abstract**—In this paper, we consider the uniquely decodable one-to-one code (UDOOC) that is obtained by inserting a comma indicator, termed the *unique word* (UW), between consecutive one-to-one codewords for separation. As such, we analyze a class of UDOOCs and present practical algorithms for encoding and decoding such codes. Specifically, for various cases of UWs, we investigate the number of length- $n$  codewords of UDOOCs and their asymptotic growth rates in  $n$ . The proposed encoding and decoding algorithms of UDOOCs can be implemented in parallel at low computational complexity without storing the codebook. Simulation results show that for proper choices of UWs, UDOOCs can achieve better compression efficiency than Lempel-Ziv codes even when the source is not statistically independent.

## I. INTRODUCTION AND MOTIVATION

In the literature, there are basically two categories of source coding, one for the compression of a sequence of source letters and the other, for a single “one shot” source symbol. A famous representative for the former is the Huffman code, while the so-called one-to-one code (OOC) prevails the latter.

Huffman code yields the minimum average codeword length under the premise of unique decodability, a property required for the compression of a sequence of source letters. Although optimal in principle, it may encounter several obstacles in implementation. For example, the rare codewords are exceedingly long in length, thereby hampering the efficiency of decoding. Other practical obstacles include 1) the codebook may need to be priori stored for encoding and decoding, which is impractical for source with moderately large alphabet size, 2) the decoding must be sequential, not in parallel, and 3) error propagation may occur.

In contrast to the Huffman code, the OOC, when being exempted from the constraint of unique decodability, can achieve a smaller average codeword length and hence requires a codebook comprised of codewords of much shorter length.

In this work, we intend to retain the use of OOCs and propose a scheme that can be processed in parallel for the compression of a sequence of source letters. Specifically, we introduce a *unique word* (UW), behaving like a “comma” indicator, to separate the successive OOC codewords. We further prohibit the UW from appearing as an internal subword<sup>1</sup> of the concatenation of a UW, an OOC codeword and another repeated UW in sequence. Any word that does not violate the

<sup>1</sup>We say  $\mathbf{a} = a_1 \dots a_m$  is not an internal subword of  $\mathbf{b} = b_1 \dots b_n$  if there does not exist  $i$  such that  $b_i \dots b_{i+m-1} = \mathbf{a}$  for all  $1 < i < n-m+1$ .

above prohibition is a codeword. This results in a uniquely decodable OOC (abbreviated as UDOOC). The UDOOC offers many advantages in its practice and requires only a small overhead for UW. For example, the decoder can first locate all the UWs in a bit stream, parsing it into separate codewords, which allows for a subsequent parallel decoding with much shorter latency. The use of UW also limits the extent of error propagation.

We would like to stress that while the use of UW might not be new, UDOOC remains as a novel concept. It has a distinct requirement that the UW cannot appear internally in a stream consisting of a codeword prefixed and suffixed by the same UW. This is the key to the fast identification of codewords for parallel decoding. It is also possible to place two UWs side by side with nothing in between to produce a *null* codeword. We now give a formal definition of UDOOC.

*Definition 1:* Given UW  $\mathbf{k} = k_1 k_2 \dots k_L \in \mathbb{F}^L$  of length  $L$ , where  $\mathbb{F} = \{0, 1\}$  is the alphabet of code symbol, we say  $\mathcal{C}_{\mathbf{k}}(n)$  is a UDOOC of length  $n$  associated with  $\mathbf{k}$  if it contains all length- $n$  tuples  $\mathbf{b} = b_1 \dots b_n$  such that  $\mathbf{k}$  is not an internal subword of the concatenated word  $\mathbf{k}\mathbf{b}\mathbf{k}$ . In particular, we set  $\mathcal{C}_{\mathbf{k}}(0) := \emptyset$ , the empty set, corresponding to the null word. The overall UDOOC associated with  $\mathbf{k}$ , denoted by  $\mathcal{C}_{\mathbf{k}}$ , is given by

$$\mathcal{C}_{\mathbf{k}} = \bigcup_{n \geq 0} \mathcal{C}_{\mathbf{k}}(n).$$

Before giving more details, below we provide an example to illustrate how to generate a UDOOC and how it is used in encoding and decoding.

*Example 1:* Let  $\mathbf{k} = 00$  be the UW. To prohibit codewords as well as the concatenation of UW and codeword from containing  $00$  as an internal subword, the following necessary and sufficient constraints are required for codeword  $\mathbf{c} = c_1 \dots c_n$ .

- $00$  cannot be a subword of  $\mathbf{c}$ . This means that within  $\mathbf{c}$ ,
  - (C1): “0” can only be followed by “1”.
  - (C2): “1” can be followed by either “0” or “1”.
- $00$  cannot be an internal subword of  $\mathbf{k}\mathbf{c}$  or  $\mathbf{c}\mathbf{k}$ . This implies
  - (C3):  $c_1 = c_n = 1$ .

Following constraints (C1)-(C3), we can place the possible codewords on a code tree as exemplified in Fig. 1. In the code tree, each path from root node “null” to a gray-shaded node represents a codeword. Thus, the legitimate codewords for UW “00” include “null”, 1, 11, 101, 111, 1011, 1101,

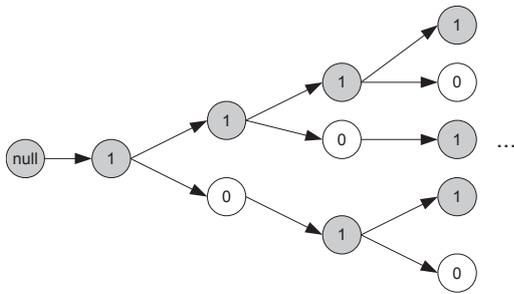


Fig. 1. The code tree for UW “00”.

and 1111, etc. Notably, we only show the codewords of length up to 4, while the tree can grow indefinitely in depth. At the decoding stage, assume the received bit-stream is 0010000110010100111100. Knowing that 00 is the UW, the decoder first locates all 00s in the bit stream and interprets the stream as “,1,null,11,101,1111,” where UWs are replaced by commas. Thus this stream is separated into five codewords and parallel decoding of them can then be conducted.

A natural question that follows is how many codewords are there for a specific UW. This quantity may be of research and practical interest for UDOOC to be useful. To calculate the number of codewords of length  $n$ , let  $a_n$  and  $b_n$  be respectively the number of “1”-nodes and the number of “0”-nodes at the  $n$ th level of the code tree. From conditions (C1) and (C2), we have the following recursions,

$$a_n = a_{n-1} + b_{n-1}, \quad b_n = a_{n-1}.$$

This yields the renowned Fibonacci sequence for  $a_n$ , i.e.,  $a_n = a_{n-1} + a_{n-2}$ . Constraint (C3) shows that only the “1”-nodes can be gray-shaded and result in a codeword. Thus, the number of codewords of length  $n$  is  $a_n$ . ■

The rest of the paper is organized as follows. Section II investigates the number of length- $n$  codewords of UDOOC and its asymptotic growth rate. Section III gives the efficient encoding and decoding algorithms for UDOOCs with two specific UWs. They can be implemented in parallel and without the need of pre-storing the codebook. Section IV presents simulations of the compression performance of UDOOCs.

## II. ASYMPTOTIC GROWTH RATE OF UDOOCs

To construct a UDOOC, the UW must be first identified. Let  $c_{\mathbf{k},n} = |\mathcal{C}_{\mathbf{k}}(n)|$  be the number of codewords in  $\mathcal{C}_{\mathbf{k}}(n)$  for UW  $\mathbf{k}$ , where  $\mathcal{C}_{\mathbf{k}}(n)$  is the UDOOC consisting of length- $n$  codewords (cf. Definition 1). To quantify how fast  $c_{\mathbf{k},n}$  grows as  $n$  increases, we introduce the following asymptotic measure.

*Definition 2:* Given UW  $\mathbf{k}$ , the asymptotic growth rate of the corresponding UDOOC is defined as

$$g_{\mathbf{k}} = \lim_{n \rightarrow \infty} \frac{c_{\mathbf{k},n+1}}{c_{\mathbf{k},n}}. \quad (1)$$

It is obvious that  $g_{\mathbf{k}} \leq 2$  for all UW  $\mathbf{k}$  as the unconstrained binary sequences have growth rate equal to 2. Notably, the limit in (1) must exist since  $g_{\mathbf{k}}$  is actually the *spectral radius* of the digraph associated with  $\mathbf{k}$  that will be discussed shortly.

Clearly, UDOOCs having larger growth rate mean that asymptotically they contain more codewords and achieve smaller average codeword length. To determine  $g_{\mathbf{k}}$  for any UW  $\mathbf{k}$ , we shall investigate the enumeration of  $c_{\mathbf{k},n}$ , and for this we introduce a digraph for UDOOC.

### A. Digraph for Enumerating $c_{\mathbf{k},n}$

For convenience of defining a digraph, we abbreviate the binary sequence  $i_s i_{s+1} \dots i_t$  as  $i_s^t$  for  $s \leq t$ , and regard  $i_s^t$  as an empty stream when  $s > t$ . In addition, we use the convenient shorthands  $(\mathbf{A})_{i,j}$  and  $(\underline{x})_j$  to denote the  $i, j$ -th entry and  $j$ -th entry of matrix  $\mathbf{A}$  and vector  $\underline{x}$ , respectively. Denote by  $G_{\mathbf{k}} = (V, E_{\mathbf{k}})$  a digraph corresponding to UW  $\mathbf{k}$ , in which  $V = \mathbb{F}^{2^{L-1}}$  is the set of all binary length- $(L-1)$  tuples, and the edge set  $E_{\mathbf{k}}$  is given by

$$E_{\mathbf{k}} = \{(i, j) : i_2^{L-1} = j_1^{L-2}, i_1 j \neq \mathbf{k}\}.$$

Define the  $2^{L-1}$ -by- $2^{L-1}$  adjacency matrix  $\mathbf{A}_{\mathbf{k}}$  for digraph  $G_{\mathbf{k}}$  as

$$(\mathbf{A}_{\mathbf{k}})_{i,j} = \begin{cases} 1, & \text{if } (i, j) \in E_{\mathbf{k}}, \\ 0, & \text{otherwise,} \end{cases}$$

where we have assumed that  $i$  (resp.  $j$ ) is the binary representation of index  $i$  (resp.  $j$ ) for  $i, j = 0, 1, \dots, 2^{L-1} - 1$  with the leftmost bit being the most significant bit. We next define a length- $2^{L-1}$  initial vector  $\underline{x}_{\mathbf{k}}$  with  $(\underline{x}_{\mathbf{k}})_j = 1$  if index  $j$  has the binary representation  $k_2^L$ , and  $(\underline{x}_{\mathbf{k}})_j = 0$ , otherwise. This vector  $\underline{x}_{\mathbf{k}}$  marks the starting node in digraph  $G_{\mathbf{k}}$ . With the above notations, it follows that the  $j$ -th entry of  $\underline{x}_{\mathbf{k}}^{\top} (\mathbf{A}_{\mathbf{k}})^n$  gives the number of length- $n$  walks on digraph  $G_{\mathbf{k}}$  that end at node  $j$  (i.e., the binary representation of  $j$ ), where superscript “ $\top$ ” stands for transpose operation.

Recall that every codeword  $\mathbf{c}$  must be followed by UW  $\mathbf{k}$  to separate consecutive codewords in a compressed stream. With  $\mathbf{k} = k_1 k_2 \dots k_L$ , it implies that  $\mathbf{c}$  is a codeword if, and only if, the sequence  $\mathbf{c} k_1^{L-1}$  is a valid walk on digraph  $G_{\mathbf{k}}$  (which ensures that  $\mathbf{k}$  is not an internal subword of  $\mathbf{c} \mathbf{k} \mathbf{c}$ ). Thus,  $c_{\mathbf{k},n}$  is the number of distinct length- $(n+L-1)$  walks from node  $k_2^L$  to node  $k_1^{L-1}$  on digraph  $G_{\mathbf{k}}$ .

From the above discussion, let  $\underline{y}_{\mathbf{k}}$  be the ending vector given by  $(\underline{y}_{\mathbf{k}})_j = 1$  if index  $j$  has the binary representation  $k_1^{L-1}$ , and  $(\underline{y}_{\mathbf{k}})_j = 0$ , otherwise. We conclude

$$c_{\mathbf{k},n} = \underline{x}_{\mathbf{k}}^{\top} (\mathbf{A}_{\mathbf{k}})^{n+L-1} \underline{y}_{\mathbf{k}}. \quad (2)$$

Using the standard technique, we can enumerate  $c_{\mathbf{k},n}$  using an indeterminate  $z$  as follows:

$$\sum_{n=0}^{\infty} c_{\mathbf{k},n} z^n = \frac{\underline{x}_{\mathbf{k}}^{\top} \text{adj}(\mathbf{I} - \mathbf{A}_{\mathbf{k}} z) \mathbf{A}_{\mathbf{k}}^{L-1} \underline{y}_{\mathbf{k}}}{\det(\mathbf{I} - \mathbf{A}_{\mathbf{k}} z)}. \quad (3)$$

where  $\mathbf{I}$  is the  $2^{L-1}$ -by- $2^{L-1}$  identity matrix. As a result, the enumeration of  $c_{\mathbf{k},n}$  depends largely upon polynomial  $\det(\mathbf{I} - \mathbf{A}_{\mathbf{k}} z)$ . In particular,  $\det(\mathbf{I} - \mathbf{A}_{\mathbf{k}} z)$  gives a linear recursion of  $c_{\mathbf{k},n}$ . Further, it can be shown that the asymptotic growth rate  $g_{\mathbf{k}}$  equals exactly the largest eigenvalue of  $\mathbf{A}_{\mathbf{k}}$ .

### B. Determining $g_k$ using a Super-Code

So far, we have seen that given any UW  $k$ , we can determine the asymptotic growth rate  $g_k$  by computing the largest eigenvalue of  $A_k$ . However, it is often difficult to find a closed-form expression for  $g_k$  due to the special structure of  $A_k$ . On the other hand, by noting that  $g_k$  is completely determined by digraph  $G_k$ , it might be possible to determine  $g_k$  using a different set of words generated by  $G_k$  that has the same growth rate as  $c_{k,n}$  and that can be more easily analyzed in combinatorics. To this end, we consider the following set for a given UW  $k$ :

$$S_k(n) = \{b \in \mathbb{F}^n : k \text{ is not an internal subword of } b\}.$$

In words,  $S_k(n)$  is the set of distinct length- $n$  walks on  $G_k$ . Let  $s_{k,n} = |S_k(n)|$  be the number of elements in  $S_k(n)$ . By an argument similar to (3), it can be shown that

$$\sum_{n=0}^{\infty} s_{k,n} z^n = \sum_{n=0}^{L-1} 2^n z^n + z^L \frac{\mathbf{1}^\top \text{adj}(\mathbf{I} - A_k z) \mathbf{1}}{\det(\mathbf{I} - A_k z)}, \quad (4)$$

where  $\mathbf{1}$  is the all-one vector of appropriate length. It is seen that the enumeration of  $s_{k,n}$  also depends upon the same polynomial  $\det(\mathbf{I} - A_k z)$  as  $c_{k,n}$  does. Furthermore,  $s_{k,n}$  has the same asymptotic growth rate as  $c_{k,n}$ .

*Theorem 1:* For any given UW  $k$ , sequences  $\{c_{k,n}\}$  and  $\{s_{k,n}\}$  have the same asymptotic growth rate  $g_k$ .

*Proof:* Clearly,  $S_k(n)$  is a super-code for  $C_k(n)$ ; hence  $s_{k,n} \geq c_{k,n}$ . On the other hand, we can show  $c_{k,n+2} \geq s_{k,n}$  for  $n \geq L$ . The details are omitted due to lack of space. ■

Determining the enumeration of  $s_{k,n}$  turns out to be much easier in comparison with  $c_{k,n}$  due to lesser constraints on sequences in  $S_k(n)$ . One approach is to employ the *Goulden-Jackson clustering method* [5], [6]. To this end, we first define the *auto-overlap function* (AOF) of UW  $k$  as follows.

*Definition 3:* Given UW  $k$  of length  $L$ , the AOF of  $k$  is defined as

$$r_k(i) = \begin{cases} 1 & \text{if } k_{i+1}^L = k_1^{L-i} \text{ for } 1 \leq i \leq L-1 \\ 0 & \text{otherwise.} \end{cases}$$

With the AOF  $r_k(i)$ , the enumeration of  $s_{k,n}$  is given completely by the following theorem.

*Theorem 2 ([5]):* Let  $r_k(i)$  be the AOF for the length- $L$  UW  $k$ ; then

$$\sum_{n \geq 0} s_{k,n} z^n = \frac{1 + \sum_{i=1}^{L-1} r_k(i) z^i}{(1-2z) \left( 1 + \sum_{i=1}^{L-1} r_k(i) z^i \right) + z^L}. \quad (5)$$

Theorem 2 provides a compact description of the determinant  $\det(\mathbf{I} - A_k z)$  for the adjacency matrix  $A_k$  using AOF  $r_k(i)$ . Furthermore, it can be used to derive an explicit linear recursion for  $c_{k,n}$ .

*Theorem 3:* Let  $A_k$  be the adjacency matrix defined as before. Then

$$\det(\mathbf{I} - A_k z) = h_k(z) = (1-2z) \left( 1 + \sum_{i=1}^{L-1} r_k(i) z^i \right) + z^L \quad (6)$$

Moreover, for  $n \geq L$ , we have

$$c_{k,n} = \left[ \sum_{i=1}^{L-1} r_k(i) (2c_{k,n-i-1} - c_{k,n-i}) \right] + 2c_{k,n-1} - c_{k,n-L}. \quad (7)$$

*Proof:* The right-hand-side (RHS) of (5) is an irreducible rational function with denominator  $h_k(z)$ ; hence  $h_k(z)$  divides  $\det(\mathbf{I} - A_k z)$  in ring  $\mathbb{Z}[z]$ . Further, we can show the degree of polynomial  $\det(\mathbf{I} - A_k z)$  is  $L$ . Thus,  $\det(\mathbf{I} - A_k z) = h_k(z)$ . The proof for (7) follows from an application of Cayley-Hamilton theorem. The details are omitted due to lack of space. ■

We have seen that  $c_{k,n}$  and  $s_{k,n}$  have enumerations both depending upon  $\det(\mathbf{I} - A_k z)$  and have the same asymptotic growth rate. The enumeration of the latter is determined completely via Theorem 2. Below we will consider the asymptotic growth rates of two specific UWs,  $a = 0 \dots 00$  and  $b = 0 \dots 01$ . It will be seen that UW  $a$  gives the largest growth rate among all UWs of length  $L$ , while interestingly UW  $b$  gives the smallest growth rate.

*Theorem 4:* Let  $a = 0 \dots 0$  and  $b = 0 \dots 01$  be two UWs of length  $L$ . Then for any UW  $k$  of length  $L$ ,

$$g_b \leq g_k \leq g_a. \quad (8)$$

Furthermore, for sufficiently large  $L$ ,

$$g_a \approx 2 - 2^{-L} \quad \text{and} \quad g_b \approx 2 - 2^{1-L}.$$

Theorem 4 indicates that compressing a source with large number of alphabets using UDOOC based on UW  $a = 0 \dots 0$  yields the best performance. For source with small number of alphabets, however, it will be seen in Section IV that the UW  $b = 0 \dots 01$  is actually a better choice. This is because the source is too small in size for asymptotics. Table I shows the values of  $g_a$  and  $g_b$  for various  $L$ . It can be verified that the approximations given in Theorem 4 are fairly accurate.

TABLE I  
THE ASYMPTOTIC GROWTH RATES FOR  $a$  AND  $b$  WITH VARIOUS  $L$

L	2	3	4	5	6	7	8
$g_a$	1.618	1.839	1.928	1.966	1.984	1.992	1.996
$g_b$	1	1.618	1.839	1.928	1.966	1.984	1.992

### III. ENCODING AND DECODING ALGORITHMS

A good source code with minimum average codeword-length should assign codewords of shorter lengths to messages with higher probabilities and reserve longer codewords only for the unlikely messages. Following this principle, let  $\mathcal{U} = \{u_1, u_2, \dots\}$  be the set of message alphabets, and let  $p_i$  be the probability that alphabet  $u_i$  occurs. Without loss of generality, we assume  $p_1 \geq p_2 \geq \dots$ . Let  $\mathcal{C}$  be a good lossless source code with encoding function  $\phi: \mathcal{U} \rightarrow \mathcal{C}$ . The above principle suggests  $\ell(\phi(u_i)) \leq \ell(\phi(u_j))$  whenever  $i \leq j$ , where  $\ell(c)$  denote the length of codeword  $c$ .

To apply the above principle to UDOOC codes, given the UW  $k$ , let  $C_k$  be the corresponding UDOOC. Then,

the corresponding encoding function  $\phi_{\mathbf{k}} : \mathcal{U} \rightarrow \mathcal{C}_{\mathbf{k}}$  can be designed as follows. Define

$$F_{\mathbf{k},n} = \begin{cases} \sum_{t=0}^n c_{\mathbf{k},t}, & \text{if } n \geq 0, \\ 0, & \text{otherwise;} \end{cases} \quad (9)$$

then  $\phi_{\mathbf{k}}$  should be a bijection between  $\{u_i : 1 \leq i \leq F_{\mathbf{k},n}\}$  and  $\bigcup_{i=0}^n \mathcal{C}_{\mathbf{k}}(i)$  for all  $n$ , provided that the sequence  $p_i = \Pr\{u_i\}$ ,  $i = 1, 2, \dots$ , is non-increasing. Any such assignment results in the same smallest average codeword-length

$$L_{\mathbf{k}} = \ell(\mathbf{k}) + \sum_{u \in \mathcal{U}} \Pr\{u\} \ell(\phi_{\mathbf{k}}(u)), \quad (10)$$

where the first term  $\ell(\mathbf{k})$  is to account for the insertion of  $\mathbf{k}$  to separate adjacent codewords.

To provide an explicit design of  $\phi_{\mathbf{k}}$ , we need to give a bijection between the set of length- $n$  codewords  $\mathcal{C}_{\mathbf{k}}(n)$  and the set of message alphabets  $\{u_m : F_{\mathbf{k},n-1} < m \leq F_{\mathbf{k},n}\}$ . For  $n \geq 1$  the proposed encoding algorithm is as follows. Given any prefix word  $\mathbf{d}$  with  $\ell(\mathbf{d}) \leq n$ , we define

$$\mathcal{C}_{\mathbf{k}}(\mathbf{d}, n) = \{c \in \mathcal{C}_{\mathbf{k}}(n) : \mathbf{d} \text{ is a prefix of } c\}.$$

Now, given message  $u_m$  with  $F_{\mathbf{k},n-1} < m \leq F_{\mathbf{k},n}$ , to determine the codeword  $\phi_{\mathbf{k}}(u_m) = \mathbf{c} = c_1 \dots c_n$ , we introduce the following branching-metric variable

$$b_0 = m - F_{\mathbf{k},n-1}$$

to keep track of the branching of each coded bit; then the first bit  $c_1$  is given by

$$c_1 = \begin{cases} 0, & \text{if } b_0 \leq |\mathcal{C}_{\mathbf{k}}(0, n)| \\ 1, & \text{if } b_0 > |\mathcal{C}_{\mathbf{k}}(0, n)|. \end{cases}$$

Similarly, having determined  $\mathbf{d} = c_1 \dots c_{i-1}$  and metric  $b_{i-1}$ , the assignment for the next bit  $c_i$  is

$$c_i = \begin{cases} 0, & \text{if } b_{i-1} \leq |\mathcal{C}_{\mathbf{k}}(\mathbf{d}0, n)|, \\ 1, & \text{if } b_{i-1} > |\mathcal{C}_{\mathbf{k}}(\mathbf{d}0, n)|. \end{cases}$$

The update for  $b_{i-1}$  is

$$b_i = \begin{cases} b_{i-1}, & \text{if } c_i = 0, \\ b_{i-1} - |\mathcal{C}_{\mathbf{k}}(\mathbf{d}0, n)|, & \text{if } c_i = 1. \end{cases}$$

It can be shown that at the end of the  $n$ th iteration, say  $\phi_{\mathbf{k}}(u_m) = \mathbf{c} = c_1 \dots c_n$ , we must have

$$m = F_{\mathbf{k},n-1} + \sum_{\substack{i:1 \leq i < n \\ c_i=1}} |\mathcal{C}_{\mathbf{k}}(c_1^{i-1}0, n)| + 1 \quad (11)$$

hence,  $\phi_{\mathbf{k}}$  is indeed a bijection between  $\{u_m : F_{\mathbf{k},n-1} < m \leq F_{\mathbf{k},n}\}$  and  $\mathcal{C}_{\mathbf{k}}(n)$ . We shall emphasize that (11) also gives the corresponding decoding function  $\psi_{\mathbf{k}} : \mathcal{C}_{\mathbf{k}} \rightarrow \mathcal{U}$  for  $\mathbf{c}$ . Specifically, given codeword  $\mathbf{c} \in \mathcal{C}_{\mathbf{k}}$  with  $\ell(\mathbf{c}) = n$ , the decoding function is exactly  $\psi_{\mathbf{k}}(\mathbf{c}) = u_m$ , where  $m$  is given by (11).

In the above encoding and decoding algorithms, although  $c_{\mathbf{k},i}$  and  $F_{\mathbf{k},i}$  can be determined by (2) and (9), we still need to determine  $|\mathcal{C}_{\mathbf{k}}(\mathbf{d}0, n)|$  for any prefix  $\mathbf{d}$ . For some specific UWs, such as  $\mathbf{a} = 0 \dots 0$ ,  $\mathbf{b} = 0 \dots 01$ , and their binary

complements, the values of  $|\mathcal{C}_{\mathbf{k}}(\mathbf{d}0, n)|$  and  $c_{\mathbf{k},n}$  can be easily determined from the following theorems. Notably, it is clear that the algorithms for encoding or decoding UDOOCs with UW  $\mathbf{a} = 0 \dots 0$  or  $\bar{\mathbf{a}} = 1 \dots 1$  are equivalent, as one can be obtained from the other by binary complementing. The same also holds for the case of  $\mathbf{b} = 00 \dots 01$  or  $\bar{\mathbf{b}} = 1 \dots 10$ . Thus, here we focus only on the cases of  $\mathbf{k} = \bar{\mathbf{a}} = 1 \dots 1$  and  $\mathbf{k} = \bar{\mathbf{b}} = 1 \dots 10$  with  $\ell(\mathbf{k}) = L$ .

*Theorem 5:* Let  $\mathbf{k} = \bar{\mathbf{a}} = 1 \dots 1$  with  $\ell(\mathbf{k}) = L$ . Then given any prefix  $\mathbf{d}$  with  $|\mathcal{C}_{\mathbf{k}}(\mathbf{d}0, n)| > 0$ , we have  $|\mathcal{C}_{\mathbf{k}}(\mathbf{d}0, n)| = c_{\mathbf{k},n-\ell(\mathbf{d})}$ . Furthermore, the values of  $c_{\mathbf{k},n}$  are given by

$$c_{\mathbf{k},n} = \begin{cases} 1, & n = 0, 1, 2, \\ 2^{n-2}, & n = 3, \dots, L \\ \sum_{i=1}^L c_{\mathbf{k},n-i} & n > L. \end{cases} \quad (12)$$

*Theorem 6:* Let  $\mathbf{k} = \bar{\mathbf{b}} = 1 \dots 10$  with  $\ell(\mathbf{k}) = L$ . Then given any prefix  $\mathbf{d}$  with  $|\mathcal{C}_{\mathbf{k}}(\mathbf{d}0, n)| > 0$ , we have  $|\mathcal{C}_{\mathbf{k}}(\mathbf{d}0, n)| = c_{\mathbf{k},n-\ell(\mathbf{d})-1}$ . Furthermore, the values of  $c_{\mathbf{k},n}$  are given by

$$c_{\mathbf{k},n} = \begin{cases} 2^n, & n = 0, 1, \dots, L-1 \\ 2c_{\mathbf{k},n-1} - c_{\mathbf{k},n-L}, & n \geq L. \end{cases} \quad (13)$$

To sum up, in Algorithms 1 and 2 we present the respective encoding and decoding functions for the UDOOCs with  $\mathbf{k} = \bar{\mathbf{a}} = 1 \dots 1$  or  $\mathbf{k} = \bar{\mathbf{b}} = 1 \dots 10$ . We assume that the set of message alphabets  $\mathcal{U} = \{u_1, u_2, \dots\}$  are ordered according to their probability masses, that is,  $\Pr\{u_i\} \geq \Pr\{u_j\}$  if  $i < j$ .

---

#### Algorithm 1 Encoding of $\mathcal{C}_{\mathbf{k}}$ with $\mathbf{k} = \bar{\mathbf{a}}$ or $\mathbf{k} = \bar{\mathbf{b}}$

---

**Input:** Index  $m$  for message  $u_m$

**Output:** Codeword  $\phi_{\mathbf{k}}(u_m) = c_1 \dots c_n$

- 1: Repeatedly compute  $c_{\mathbf{k},0}, c_{\mathbf{k},1}, \dots, c_{\mathbf{k},n}$  using (12) for  $\mathbf{k} = \bar{\mathbf{a}}$  (or (13) for  $\mathbf{k} = \bar{\mathbf{b}}$ ) to find the smallest  $n$  such that  $F_{\mathbf{k},n} \geq m$
  - 2:  $b_0 \leftarrow m - F_{\mathbf{k},n-1}$
  - 3: **for**  $i = 1$  to  $n$  **do**
  - 4:    $B \leftarrow c_{\mathbf{k},n-i+1}$  for  $\mathbf{k} = \bar{\mathbf{a}}$  or  $B \leftarrow c_{\mathbf{k},n-i}$  for  $\mathbf{k} = \bar{\mathbf{b}}$
  - 5:   **if**  $b_{i-1} \leq B$  **then**
  - 6:      $c_i \leftarrow 0$  and  $b_i \leftarrow b_{i-1}$
  - 7:   **else**
  - 8:      $c_i \leftarrow 1$  and  $b_i \leftarrow b_{i-1} - B$
  - 9:   **end if**
  - 10: **end for**
- 

#### IV. PERFORMANCE OF UDOOC

In Fig. 2, we plot the values of  $c_{\mathbf{k},n}$  for various  $\mathbf{k}$  of length  $L = 3$ . The values are presented in log-scale and are normalized against the case of  $\mathbf{k} = 0 \dots 0$ . It should be noted that we only show the curves for inequivalent UWs, by which we mean that  $\mathbf{k}$  is equivalent to  $\mathbf{k}'$  if  $c_{\mathbf{k},n} = c_{\mathbf{k}',n}$  for all  $n$ . We observe that  $c_{\mathbf{k},n}$  for  $\mathbf{k} = 0 \dots 01$  is the largest among all when  $n$  is small, but it becomes the smallest when  $n$  is large.  $c_{\mathbf{k},n}$  for  $\mathbf{k} = 0 \dots 0$  is the smallest among all when  $n$  is small, but it is the largest when  $n$  is large. This shows a tradeoff among the

**Algorithm 2** Decoding of  $\mathcal{C}_k$  with  $k = \bar{a}$  or  $k = \bar{b}$

**Input:** Codeword  $c = c_1 \dots c_n$

**Output:** Index  $m$  for message  $u_m = \psi_k(c)$

- 1: Compute  $c_{k,0}, c_{k,1}, \dots, c_{k,n}$  using (12) for  $k = \bar{a}$  (or (13) for  $k = \bar{b}$ )
- 2:  $m \leftarrow F_{k,n-1} + 1$
- 3: **for**  $i = 1$  to  $n$  **do**
- 4:  $B \leftarrow c_{k,n-i+1}$  for  $k = \bar{a}$  or  $B \leftarrow c_{k,n-i}$  for  $k = \bar{b}$
- 5: **if**  $c_i = 1$  **then**
- 6:  $m \leftarrow m + B$
- 7: **end if**
- 8: **end for**

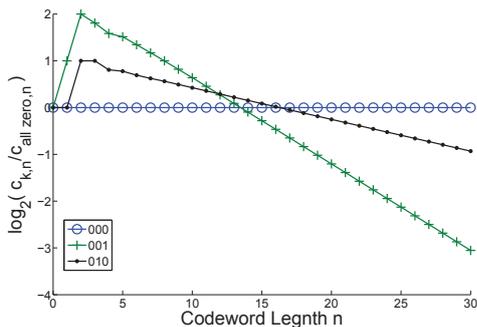


Fig. 2. Normalized number of length- $n$  codewords for  $L = 3$

choices of UWs. If the source is uniform and is extremely large in size, UDOOC using  $k = 0 \dots 0$  achieves better compression performance and has smaller average codeword length. On the other hand, if the source is nonuniform and is smaller in size,  $k = 0 \dots 01$  or its equivalence is a better choice.

We also simulate the compression performance of UDOOCs when source is uniform and independent distributed among 26 English letters [7]. The results are given in Table II, where the average codeword length has already taken into account the length of UWs. We observe that when the parser size  $M = 1$ , meaning the alphabet size is small, UDOOC with  $k = 01$  gives the best performance. However, as  $M$  increases, the alphabet size grows exponentially, and compression performance using  $k = 0 \dots 0$  eventually becomes the best. It should also be noted that for large  $M$  UDOOC with UW 01 performs extremely poor due to its asymptotic growth rate equal to 1.

TABLE II  
AVERAGE CODEWORD LENGTH IN BITS PER SOURCE SYMBOL FOR UNIFORM INDEPENDENT 26 ENGLISH LETTERS WITH PARSER SIZE  $M$

	UW = 00...0				UW = 00...01			
	M=1	M=2	M=4	M=8	M=1	M=2	M=4	M=8
L=2	6.961	6.820	6.794	6.781	5.846	12.76	159	53855
L=4	8.576	6.830	5.901	5.434	7.000	5.898	5.619	5.474
L=8	12.576	8.735	6.718	5.719	11.000	7.751	6.222	5.477

We next consider two different sources. The first is an independent English source with usual frequency [7]. The second source is the book "Alice's Adventures in Wonderland" [8],

which represents a dependent source with memory. Results of the compression performance using UDOOCs are presented in Table III and are compared with those using the conventional Huffman and Lempel-Ziv (LZ) codes [2]. It is seen that with parser size  $M = 3$ , all UDOOCs with  $k = 0 \dots 0$  perform better than LZ code.

TABLE III  
AVERAGE CODEWORD LENGTH IN BITS PER SOURCE SYMBOL FOR THE COMPRESSION OF TWO DIFFERENT SOURCES

(a)							
Type	Entropy			LZ	Huffman		
Usual Independent English Letters	M=1	M=2	M=3	6.626	M=1	M=2	M=3
	4.246	4.246	4.246		4.274	4.261	4.253
Alice's Adventures in Wonderland	M=1	M=2	M=3	6.027	M=1	M=2	M=3
	3.914	3.570	3.215		3.940	3.585	3.226

(b)							
Type		UW = 00...0			UW = 00...01		
Usual Independent English Letters	L=2	M=1	M=2	M=3	M=1	M=2	M=3
	L=4	5.591	5.550	5.637	4.557	7.771	20.907
	L=6	7.411	5.872	5.351	6.185	4.970	4.795
Alice's Adventures in Wonderland	L=2	M=1	M=2	M=3	M=1	M=2	M=3
	L=4	4.887	4.340	3.958	4.068	4.975	7.573
	L=6	6.757	4.920	4.089	5.774	4.133	3.531

V. CONCLUSION

In this paper, we proposed a class of UDOOC that allows for a fast identification of all codewords in a compressed stream, and therefore parallel decoding is possible. Many combinatorial properties of UDOOC, such as the number of length- $n$  codes, enumerations, and asymptotic growth rate, are presented and are helpful in choosing appropriate UWs to yield better compression performance. Encoding and decoding algorithms for UDOOC associated with UW  $0 \dots 0$  and  $0 \dots 01$  are given and can be implemented without storing the codebook. We simulated the compression performance using UDOOC and demonstrated that in many cases the UDOOC performs better than the LZ code.

REFERENCES

- [1] C. E. Shannon, "A mathematical theory of communication," *Bell System Technical Journal*, vol. 27, July, October 1948.
- [2] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, New York, NY: John Wiley & Sons, 1991.
- [3] W. Szpankowski, "One-to-one codes and its anti-redundancy," in *Proceedings IEEE ISIT*, September 4-9, 2005.
- [4] D. B. West, *Introduction to Graph Theory (2E)*, Prentice Hall, 2001.
- [5] I. Goulden and D. M. Jackson, "An inversion theorem for cluster decompositions of sequences with distinguished subsequences," *J. London Math. Soc.*, vol. 20, no. 2, pp. 567-576, 1979.
- [6] J. Noonan and D. Zeilberger, "The Goulden-Jackson cluster method: extensions, applications, and implementations," *J. Difference Eq. Appl.*, vol. 5, pp. 355-377, 1999.
- [7] <http://oxforddictionaries.com/words/what-is-the-frequency-of-the-letters-of-the-alphabet-in-english>
- [8] C. L. Dodgson (L. Carroll), "Alice's Adventures in Wonderland," 1865.