

# The Gapped $k$ -Deck Problem

Rebecca Golm $\ddagger$   
ECE Department  
University of Illinois  
rgolm2@illinois.edu

Mina Nahvi $\ddagger$   
Department of Mathematics  
University of Illinois  
mnavhi2@illinois.edu

Ryan Gabrys  
Calit2  
University of California, San Diego  
ryan.gabrys@gmail.com

Olgica Milenkovic  
ECE Department  
University of Illinois  
milenkov@illinois.edu

**Abstract**—The  $k$ -deck problem is concerned with finding the smallest positive integer  $S(k)$  such that there exist at least two strings of length  $S(k)$  that share the same  $k$ -deck, i.e., the multiset of subsequences of length  $k$ . We introduce the new problem of gapped  $k$ -deck reconstruction: For a given gap parameter  $s$ , we seek the smallest positive integer  $G_s(k)$  such that there exist at least two distinct strings of length  $G_s(k)$  that cannot be distinguished based on a “gapped” set of  $k$ -subsequences. The gap constraint requires the elements in the subsequences to be at least  $s$  positions apart within the original string. Our results are as follows. First, we show how to construct sequences sharing the same 2-gapped  $k$ -deck using a nontrivial modification of the recursive Morse-Thue string construction procedure. This establishes the first known constructive upper bound on  $G_2(k)$ . Second, we further improve this bound using the approach by Dudik and Schulman [6].

**Index Terms**—Gapped subsequences,  $k$ -deck, Morse-Thue sequences, String reconstruction

## I. INTRODUCTION

The problem of reconstructing strings based on evidence sets of the form of subsequences, substrings or weights of substrings has received significant attention from the theoretical computer science, bioinformatics, and information theory communities alike [1], [3], [4], [8], [10], [11], [13], [15]. One special instance of this class of problems is the  $k$ -deck problem [4], [6], [7], [9], [10], [14], of interest due to its connection to trace reconstruction [3], [5] and its applications in DNA-based data storage [16].

For a string  $x$  of length  $n$ , the multiset of the  $\binom{n}{k}$  subsequences (i.e., ordered collections of not necessarily adjacent entries) of  $x$  of length  $k$  is called the  $k$ -deck of  $x$ . We say that  $x$  is  $k$ -reconstructible if it is uniquely determined by its  $k$ -deck, meaning that there exists no other string that has the same  $k$ -deck as  $x$ . For example,  $(1, 0, 0, 1)$  and  $(0, 1, 1, 0)$  have the same 2-deck, and are hence not 2-reconstructible. A simple counting argument shows that if two sequences  $x$  and  $y$  have the same  $k$ -deck, they also have the same  $l$ -deck for all  $1 \leq l \leq k$ .

Let  $S(k)$  be the smallest positive integer  $n$  such that there exist two distinct strings of length  $n$  with the same  $k$ -deck. Kalashnik [10] raised the question of determining  $S(k)$ . Manvel, Meyerowitz, Schwenk, Smith and Stockmeyer [12] showed that  $2k \leq S(k) \leq 2^k$ . They proved the upper bound as follows. For two strings  $x$  and  $y$  of length  $n$ , let  $xy$  be

the string obtained by concatenating  $x$  and  $y$  (note that when concatenating a single bit, say 0, and a string  $x$ , we also use the notation  $(0, x)$ ). If  $x$  and  $y$  have the same  $k$ -deck, then  $xy$  and  $yx$  have the same  $(k+1)$ -deck. The upper bound follows immediately when coupled with the fact that  $(0, 1)$  and  $(1, 0)$  have the same 1-deck. The construction is often referred to as the Morse-Thue construction and the resulting strings are the well-known Morse-Thue strings [2]. Furthermore, the authors of [12] also showed that in order to prove that every string of length  $n$  is  $k$ -reconstructible, it is enough to prove that every binary string of length  $n$  is  $k$ -reconstructible. Dudik and Schulman [6] improved the above upper bound on  $S(k)$  to  $\exp\left(\frac{3+o(1)}{2 \log 3} \log^2 k\right)$ . In the literature, both bounds on the smallest  $k$  and  $n$  (for a given  $n$  and  $k$ , respectively) for unique and nonunique  $k$ -deck reconstruction have been reported.

We define the *gapped  $k$ -deck* of a binary string  $x$  as the multiset of all subsequences of length  $\leq k$  that do not include two consecutive entries in  $x$ . This definition can be extended to larger gaps between entries in  $x$ : The  *$s$ -gapped  $k$ -deck* of a binary string  $x$  is the multiset of all subsequences  $(x_{i_1}, \dots, x_{i_\ell})$ ,  $1 \leq \ell \leq k$ , such that for all  $1 \leq j \leq \ell - 1$ , we have  $i_{j+1} \geq i_j + s$ . With this definition, the gapped  $k$ -deck reduces to the 2-gapped  $k$ -deck. The problem of interest is to bound  $G_s(k)$ , the smallest positive integer  $n$  for which there exist two binary strings that share the same  $s$ -gapped  $k$ -deck. For simplicity, when  $s = 2$ , we write  $G(s)$  and refer to the corresponding setting as the *gapped  $k$ -deck*. Note that unlike the case without gaps, two strings  $x$  and  $y$  having the same multiset of gapped subsequences of length  $k$  does not imply that they also have the same multiset of gapped sequences of length  $l$  for some  $l < k$ . For example, the strings  $(0, 1, 1, 1, 0)$  and  $(1, 0, 0, 0, 1)$  have the same multiset of gapped subsequences of length 2, but they clearly have different multisets of gapped subsequences of length 1 (which by definition, is the multiset of bits (composition) of the strings). The gapped  $k$ -deck problem is of interest in molecular storage systems for which readouts are based on nanopore technologies, in which “gaps” in readouts arise due to skipping effects [16].

We initiate the study of reconstruction limits of strings given their  $s$ -gapped  $k$ -decks and present the first upper bounds on  $G_s(k)$  and  $G(k)$  in particular. In Section III we provide necessary preliminaries, while in Section IIII we describe a nontrivial extension of a Morse-Thue type construction for

The work was funded by NSF grant 2008125, Coded String Reconstruction Problems in Molecular Storage. In the author list,  $\ddagger$  denotes equal contribution.

2-gapped  $k$ -decks. In Section III we state the result for general values of  $s$  but omit the proof. Section IV presents an improvement of the upper bound for  $G(k)$  from Section III based on an adaptation of the method described in [6].

## II. PRELIMINARIES

For a string  $x = (x_1, \dots, x_n) \in \{0, 1\}^n$ , let

$$\mathcal{B}^{(k)}(x) = \{(x_{i_1}, x_{i_2}, \dots, x_{i_\ell}) : i_j \geq i_{j-1} + 2, 0 \leq \ell \leq k\} \quad (1)$$

denote the multiset of all subsequences of  $x$  of length  $\leq k$  such that the index of every entry used in a subsequence is nonadjacent in the original string. Also, let

$$\mathcal{D}^{(k)}(x) = \{(x_{i_1}, x_{i_2}, \dots, x_{i_k}) : i_j \geq i_{j-1} + 2\},$$

be the *exact* gapped  $k$ -deck of  $x$ . Here, we assume that  $\mathcal{B}^{(0)}(x) = \mathcal{D}^{(0)}(x) = \emptyset$ . Clearly,  $\mathcal{B}^{(k)}(x) = \bigcup_{i=0}^k \mathcal{D}^{(i)}(x)$ . As mentioned in the introduction, unlike the classical (un-gapped) case, the problem of reconstructing  $x$  from  $\mathcal{B}^{(k)}(x)$  differs from that of reconstructing  $x$  from  $\mathcal{D}^{(k)}(x)$ . Our focus is on finding  $G(k)$ , the smallest integer  $n$  such that there exist two distinct binary strings of length  $n$  with the same gapped  $i$ -deck for all  $1 \leq i \leq k$ . Alternatively,  $G(k)$  is the smallest integer  $n$  such that there exist two distinct binary strings of length  $n$ ,  $x$  and  $y$ , satisfying  $\mathcal{B}^{(k)}(x) = \mathcal{B}^{(k)}(y)$ . It is worth pointing out that if  $n$  is the smallest integer such that there exist two strings  $x$  and  $y$  of length  $n$  with  $\mathcal{D}^{(k)}(x) = \mathcal{D}^{(k)}(y)$ , then  $n = 2k - 1$ . We have  $n \geq 2k - 1$  because a string of length less than  $2k - 1$  has no gapped subsequence of length  $k$ . On the other hand, for any string  $z = (z_1 \dots z_k)$  of length  $k$ , all the strings of length  $2k - 1$  of the form  $(z_1 x_1 z_2 x_2 \dots z_{k-1} x_{k-1} z_k)$  have the same gapped  $k$ -deck because the only gapped  $k$ -subsequence of  $x$  is  $z$ . This observation generalizes for  $s$ -gapped  $k$ -decks and  $n = sk - 1$ .

In Section III we prove that  $G(k) \leq 4(2^k - 1) - 2$ . We also provide an upper bound on  $G_s(k)$ , the smallest integer  $n$  such that there exist two distinct strings of length  $n$  with the same  $s$ -gapped  $i$ -deck for all  $1 \leq i \leq k$ , where  $s \geq 2$ . The bound reads as  $G_s(k) \leq (5s - 2)2^{k-1} - 5s + 4$ , but the accompanying proof is omitted due to space limitations. The proof of our first bound on  $G(k)$  builds upon the next lemma.

**Lemma 1.** [I2] *If  $x = (x_1, x_2, \dots, x_m)$  and  $y = (y_1, y_2, \dots, y_m)$  have the same  $k$ -deck, then the two concatenation strings  $xy = (x_1, \dots, x_m, y_1, \dots, y_m)$  and  $yx = (y_1, \dots, y_m, x_1, \dots, x_m)$  have the same  $(k + 1)$ -deck.*

*Proof:* The following correspondence proves the claim: Pick any subsequence  $z$  of  $xy$  of length at most  $k + 1$ . If  $z$  is fully contained within the  $x$  (or  $y$ ) substring, let  $\phi(z)$  be the same subsequence in the  $x$  (or  $y$ ) substring of  $yx$ . Now, assume  $z = z_1 z_2$ , where  $z_1$  is a subsequence of  $x$  and  $z_2$  is a subsequence of  $y$ . Note that  $z_1$  and  $z_2$  have length at most  $k$ , therefore there exists a subsequence  $w_1$  of  $y$  that equals  $z_1$ , due to the fact that  $x$  and  $y$  have the same  $i$ -deck for all  $1 \leq i \leq k$ . Similarly,  $x$  contains a subsequence  $w_2$  that equals  $z_2$ . Now, let  $\phi(z) = w_1 w_2$ . Therefore,  $xy$  and  $yx$  have the same  $(k + 1)$ -deck. ■

Using the strings  $x = (0, 1)$ ,  $y = (1, 0)$  and  $k = 1$  to initialize the recursion, we can see that  $(0, 1, 1, 0)$  and  $(1, 0, 0, 1)$  have the same 2-deck. Repeating the process, we find that  $(0, 1, 1, 0, 1, 0, 0, 1)$  and  $(1, 0, 0, 1, 0, 1, 1, 0)$  have the same 3-deck and so on. However, this construction does not work for the gapped case. For example,  $(1, 0)$  and  $(0, 1)$  have the same gapped 1-deck (i.e., composition), but  $(1, 0, 0, 1)$  and  $(0, 1, 1, 0)$  do not have the same gapped 2-deck. The reason why the construction fails is that we cannot pick both  $x_m$  and  $y_1$  (as defined in Lemma I) when choosing a gapped subsequence of  $xy$ . Hence, we need to “pad” the boundary between the two concatenated strings in an adequate manner.

## III. THE PADDED MORSE-THUE SEQUENCE APPROACH

We prove the existence of two strings  $x, y \in \{0, 1\}^n$ , where  $n = 4(2^k - 1) - 2$ , that satisfy  $\mathcal{B}^{(k)}(x) = \mathcal{B}^{(k)}(y)$ , using induction. We start with a few definitions. For a binary string  $x = (x_1, \dots, x_n) \in \{0, 1\}^n$  let

$$\mathcal{B}_L^{(k)}(x) := \mathcal{B}^{(k)}(x_2, x_3, \dots, x_n), \quad (2)$$

$$\mathcal{B}_R^{(k)}(x) := \mathcal{B}^{(k)}(x_1, x_2, \dots, x_{n-1}), \quad (3)$$

$$\mathcal{B}_{LR}^{(k)}(x) := \mathcal{B}^{(k)}(x_2, \dots, x_{n-1}). \quad (4)$$

Note that (2) represents the multiset of all gapped subsequences formed by puncturing  $x$  on the left, (3) represents the multiset of all gapped subsequences formed by puncturing  $x$  on the right, while (4) represents the multiset of all gapped subsequences formed by puncturing  $x$  on both ends. We define the sets  $\mathcal{D}_L^{(k)}(x)$ ,  $\mathcal{D}_R^{(k)}(x)$ , and  $\mathcal{D}_{LR}^{(k)}(x)$  analogously.

We initialize two strings for the “degenerate” case of  $k = 1$ , corresponding to equal compositions, as follows:

$$x^{(1)} = (0, 0, 1, 0), \quad y^{(1)} = (0, 1, 0, 0). \quad (5)$$

Puncturing the first bit from both  $x^{(1)}$  and  $y^{(1)}$  produces strings that still share the same gapped 1-deck. The same claim holds for the case when one punctures the last bit from both  $x^{(1)}$  and  $y^{(1)}$ . Finally, the claim is true when one punctures both the first and the last bit from both strings. Hence, for  $i = 1$ ,

$$\mathcal{B}^{(i)}(x^{(i)}) = \mathcal{B}^{(i)}(y^{(i)}), \quad \mathcal{B}_{LR}^{(i)}(x^{(i)}) = \mathcal{B}_{LR}^{(i)}(y^{(i)}) \quad (6)$$

$$\mathcal{B}_L^{(i)}(x^{(i)}) = \mathcal{B}_L^{(i)}(y^{(i)}), \quad \mathcal{B}_R^{(i)}(x^{(i)}) = \mathcal{B}_R^{(i)}(y^{(i)}).$$

Let  $G^*(k)$  be the smallest integer  $n$  such that there exist two distinct binary strings of length  $n$ ,  $x^{(k)}$  and  $y^{(k)}$ , for which (6) holds for the case  $i = k$ .

**Theorem 2.** *With  $x^{(1)}$  and  $y^{(1)}$  defined as in (5) and*

$$x^{(k)} = (0, x^{(k-1)}, 0, 0, y^{(k-1)}, 0), \quad (7)$$

$$y^{(k)} = (0, y^{(k-1)}, 0, 0, x^{(k-1)}, 0),$$

*defined recursively, we have that (6) holds for all  $i$ . As a result,  $G^*(k) \leq 4(2^k - 1)$  and  $G(k) \leq 4(2^k - 1) - 2$ .*

*Proof:* We split the proof into four subproofs, in order to show that each of the four conditions in (6) hold for

$i = k$  if they hold for  $i = k - 1$ . We do this by partitioning each deck in (6) with respect to whether each padded 0 is included in a subsequence or not, and by showing that there exists a correspondence between each pair of decks. The bound follows since the length of  $\mathbf{x}^{(k)}$  equals  $4(2^k - 1)$  and  $G(k) \leq G^*(k) - 2$ , given that one can remove the padded 0s.

*Part 1: Proof that  $\mathcal{B}_{LR}^{(k)}(\mathbf{x}^{(k)}) = \mathcal{B}_{LR}^{(k)}(\mathbf{y}^{(k)})$ .*

By definitions (4) and (7), this is equivalent to showing that

$$\mathcal{B}^{(k)}(\mathbf{x}^{(k-1)}, 0, 0, \mathbf{y}^{(k-1)}) = \mathcal{B}^{(k)}(\mathbf{y}^{(k-1)}, 0, 0, \mathbf{x}^{(k-1)}).$$

We can partition  $\mathcal{B}^{(k)}(\mathbf{x}^{(k-1)}, 0, 0, \mathbf{y}^{(k-1)})$  depending on which of the two 0s, if any, is included in the subsequence:

- 1)  $\{\left(\mathcal{D}^{(k_1)}(\mathbf{x}^{(k-1)}), \mathcal{D}^{(K-k_1)}(\mathbf{y}^{(k-1)})\right)\}, K \leq k$ ;
- 2)  $\{\left(\mathcal{D}_R^{(k_1)}(\mathbf{x}^{(k-1)}), 0, \mathcal{D}^{(K-k_1)}(\mathbf{y}^{(k-1)})\right)\}, K \leq k - 1$ ;
- 3)  $\{\left(\mathcal{D}^{(k_1)}(\mathbf{x}^{(k-1)}), 0, \mathcal{D}_L^{(K-k_1)}(\mathbf{y}^{(k-1)})\right)\}, K \leq k - 1$ ,

where  $k_1$  varies from 0 to  $K$ , for all  $K$ . First, we consider the case where neither of the two 0s is used and show that

$$\begin{aligned} &\{\left(\mathcal{D}^{(k_1)}(\mathbf{y}^{(k-1)}), \mathcal{D}^{(K-k_1)}(\mathbf{x}^{(k-1)})\right)\} = \\ &\{\left(\mathcal{D}^{(k_1)}(\mathbf{x}^{(k-1)}), \mathcal{D}^{(K-k_1)}(\mathbf{y}^{(k-1)})\right)\}, \end{aligned} \quad (8)$$

for any  $K \leq k$ . In this case, each string comprises  $k_1 \leq K$  symbols from  $\mathbf{x}^{(k-1)}$  and  $K - k_1$  symbols from  $\mathbf{y}^{(k-1)}$ , where  $K \leq k$  denotes the length of the resulting string. When  $k_1 = K$  or  $k_1 = 0$ , (8) holds, since the subsequences  $\mathbf{x}^{(k-1)}, \mathbf{y}^{(k-1)}$  appear in both  $\mathbf{x}^{(k)}$  and  $\mathbf{y}^{(k)}$ . Otherwise, when  $0 < k_1 < K$ , since the subsequences  $\mathbf{x}^{(k-1)}, \mathbf{y}^{(k-1)}$  appear (and are “nonadjacent”) in both  $\mathbf{x}^{(k)}$  and  $\mathbf{y}^{(k)}$ ,  $\mathcal{D}^{(k_1)}(\mathbf{x}^{(k-1)}) = \mathcal{D}^{(k_1)}(\mathbf{y}^{(k-1)})$ , and  $\mathcal{D}^{(K-k_1)}(\mathbf{y}^{(k-1)}) = \mathcal{D}^{(K-k_1)}(\mathbf{x}^{(k-1)})$  (since both  $k_1 < k$  and  $K - k_1 < k$ ), it follows that (8) also holds for  $0 < k_1 < K$ .

The multiset of subsequences covered by case 2 contains strings that are formed by concatenating  $k_1$  bits from  $\mathbf{x}^{(k-1)}$ , the first 0 between the subsequences  $\mathbf{x}^{(k-1)}$  and  $\mathbf{y}^{(k-1)}$  and  $K - k_1$  bits from  $\mathbf{y}^{(k-1)}$ . Next, we show that

$$\begin{aligned} &\{\left(\mathcal{D}_R^{(k_1)}(\mathbf{x}^{(k-1)}), 0, \mathcal{D}^{(K-k_1)}(\mathbf{y}^{(k-1)})\right)\} = \\ &\{\left(\mathcal{D}_R^{(k_1)}(\mathbf{y}^{(k-1)}), 0, \mathcal{D}^{(K-k_1)}(\mathbf{x}^{(k-1)})\right)\} \end{aligned} \quad (9)$$

for  $K \leq k - 1$ . Since  $K \leq k - 1$ , we have  $\mathcal{D}_R^{(k_1)}(\mathbf{x}^{(k-1)}) = \mathcal{D}_R^{(k_1)}(\mathbf{y}^{(k-1)})$  and  $\mathcal{D}^{(K-k_1)}(\mathbf{y}^{(k-1)}) = \mathcal{D}^{(K-k_1)}(\mathbf{x}^{(k-1)})$ , which implies that we can form strings by concatenating  $k_1$  bits from  $\mathbf{y}^{(k-1)}$ , the first 0 between the substrings  $\mathbf{y}^{(k-1)}$  and  $\mathbf{x}^{(k-1)}$ , and  $K - k_1$  bits from  $\mathbf{x}^{(k-1)}$ . Thus, (9) holds.

Using the same approach, it can be shown that

$$\begin{aligned} &\{\left(\mathcal{D}^{(k_1)}(\mathbf{x}^{(k-1)}), 0, \mathcal{D}_L^{(K-k_1)}(\mathbf{y}^{(k-1)})\right)\} = \\ &\{\left(\mathcal{D}^{(k_1)}(\mathbf{y}^{(k-1)}), 0, \mathcal{D}_L^{(K-k_1)}(\mathbf{x}^{(k-1)})\right)\}, \end{aligned} \quad (10)$$

for any  $K \leq k - 1$ . From (8), (9), and (10), it then follows that  $\mathcal{B}_{LR}^{(k)}(\mathbf{x}^{(k)}) = \mathcal{B}_{LR}^{(k)}(\mathbf{y}^{(k)})$ .

*Part 2: Proof that  $\mathcal{B}_L^{(k)}(\mathbf{x}^{(k)}) = \mathcal{B}_L^{(k)}(\mathbf{y}^{(k)})$ .*

By definitions (2) and (7), this is equivalent to showing that

$$\mathcal{B}^{(k)}(\mathbf{x}^{(k-1)}, 0, 0, \mathbf{y}^{(k-1)}, 0) = \mathcal{B}^{(k)}(\mathbf{y}^{(k-1)}, 0, 0, \mathbf{x}^{(k-1)}, 0).$$

We first partition  $\mathcal{B}_L^{(k)}(\mathbf{x}^{(k)})$  into two multisets: The first contains subsequences that include the last (trailing) 0 while the second contains those which do not (equivalent to  $\mathcal{B}_{LR}^{(k)}(\mathbf{x}^{(k)})$ ). The first multiset can be partitioned into three classes:

- 1)  $\{\left(\mathcal{D}^{(k_1)}(\mathbf{x}^{(k-1)}), \mathcal{D}_R^{(K-k_1)}(\mathbf{y}^{(k-1)}), 0\right)\}, K \leq k - 1$ ;
- 2)  $\{\left(\mathcal{D}_R^{(k_1)}(\mathbf{x}^{(k-1)}), 0, \mathcal{D}_R^{(K-k_1)}(\mathbf{y}^{(k-1)}), 0\right)\}, K \leq k - 2$ ;
- 3)  $\{\left(\mathcal{D}^{(k_1)}(\mathbf{x}^{(k-1)}), 0, \mathcal{D}_{LR}^{(K-k_1)}(\mathbf{y}^{(k-1)}), 0\right)\}, K \leq k - 2$ .

Using an almost identical argument as the one described in *Part 1*, one can show that  $\mathcal{B}_L^{(k)}(\mathbf{x}^{(k)}) = \mathcal{B}_L^{(k)}(\mathbf{y}^{(k)})$ .

*Part 3: Proof that  $\mathcal{B}_R^{(k)}(\mathbf{x}^{(k)}) = \mathcal{B}_R^{(k)}(\mathbf{y}^{(k)})$ .*

The proof of this case follows by symmetry from *Part 2*.

*Part 4: Proof that  $\mathcal{B}^{(k)}(\mathbf{x}^{(k)}) = \mathcal{B}^{(k)}(\mathbf{y}^{(k)})$ .*

The final step in the proof is to show that

$$\begin{aligned} &\mathcal{B}^{(k)}(0, \mathbf{x}^{(k-1)}, 0, 0, \mathbf{y}^{(k-1)}, 0) \\ &= \mathcal{B}^{(k)}(0, \mathbf{y}^{(k-1)}, 0, 0, \mathbf{x}^{(k-1)}, 0) \end{aligned}$$

Using a similar approach as before, we now partition the subsequences in  $\mathcal{B}^{(k)}(\mathbf{x}^{(k)})$  according to whether they

- 1) contain the leading 0, but not the trailing 0;
- 2) contain the trailing 0, but not the leading 0;
- 3) contain neither the trailing nor the leading 0;
- 4) contain both the leading and the trailing 0.

This is equivalent to:

- 1)  $\mathcal{B}_R^{(k)}(\mathbf{x}^{(k)}) \setminus \mathcal{B}_{LR}^{(k)}(\mathbf{x}^{(k)})$ ;
- 2)  $\mathcal{B}_L^{(k)}(\mathbf{x}^{(k)}) \setminus \mathcal{B}_{LR}^{(k)}(\mathbf{x}^{(k)})$ ;
- 3)  $\mathcal{B}_{LR}^{(k)}(\mathbf{x}^{(k)})$ ;
- 4)  $\mathcal{B}^{(k)}(\mathbf{x}^{(k)}) \setminus (\mathcal{B}_R^{(k)}(\mathbf{x}^{(k)}) \cup \mathcal{B}_L^{(k)}(\mathbf{x}^{(k)}))$ .

From the first three parts of the proof, we know that the first three multisets are the same for  $\mathbf{x}^{(k)}$  and  $\mathbf{y}^{(k)}$ . We only need to prove that the fourth multiset is the same as well. Again we partition the multiset of interest into three classes:

- 1)  $\{\left(0, \mathcal{D}_L^{(k_1)}(\mathbf{x}^{(k-1)}), \mathcal{D}_R^{(K-k_1)}(\mathbf{y}^{(k-1)}), 0\right)\}$ ;
- 2)  $\{\left(0, \mathcal{D}_{LR}^{(k_1)}(\mathbf{x}^{(k-1)}), 0, \mathcal{D}_R^{(K-k_1)}(\mathbf{y}^{(k-1)}), 0\right)\}$ ;
- 3)  $\{\left(0, \mathcal{D}_L^{(k_1)}(\mathbf{x}^{(k-1)}), 0, \mathcal{D}_{LR}^{(K-k_1)}(\mathbf{y}^{(k-1)}), 0\right)\}$ ,

where for case 1,  $K \leq k - 2$ , and for cases 2 and 3,  $K \leq k - 3$ . Using similar arguments as before completes the proof. ■

Using a similar approach, we can extend the bound to the s-gapped case to get  $G_s(k) \leq (5s - 2)2^{k-1} - 5s + 4$ . This is

done by adding  $s - 1$  0s on the outside and  $s$  0s between  $xy$  and  $yx$ . We remove  $2s - 2$  0s for the bound since the  $s$ -gapped  $k$ -deck does not need to satisfy the extra conditions required by the recursive construction.

We numerically computed  $G(k)$  for  $k = 2, 3, 4$ . The results are displayed below, which clearly indicate that the upper

$k$	$G(k)$	Confusable pairs (examples)
2	6	(0,1,0,0,1,1), (0,0,1,1,0,1)
3	13	(1,1,0,1,1,1,1,0,1,0,1,1,1), (1,1,1,0,1,0,1,1,1,1,0,1,1)
4	24	(1,1,0,0,1,1,0,1,0,1,0,1,0,0,1,1,0,0,1,1,0,1,0,0), (1,1,0,1,0,0,1,1,0,0,1,1,0,1,0,1,0,1,0,0,1,1,0,0)

bound  $4(2^k - 1) - 2$  is loose for larger values of  $k$ : For  $k = 4$ , the bound equals 58 while the correct value is only 24. Also, the exact values of  $G(k)$  are significantly larger than those for the ungapped case, for which we know that  $S(k) = 4, 7, 12$  (compared to  $G(k) = 6, 13, 24$ ) for  $k = 2, 3, 4$ , respectively. We therefore turn our attention to improving the bound on  $G(k)$  using more sophisticated counting arguments.

#### IV. IMPROVED UPPER BOUNDS FOR GAPPED $k$ -DECKS

We find the following definitions and notation from [6] useful for our subsequent derivations. Let  $\Gamma = \{X, Y\}$  and let  $J$  denote a “wildcard”. For integers  $0 \leq r \leq k$  let

$$U_r(k) = \{w \in \bigcup_{j=r}^k (\Gamma \cup \{J\})^j : w \text{ has exactly } r \text{ non-}J \text{ symbols}\}.$$

For  $t \geq 1$  and  $k_1 \geq \dots \geq k_t \geq t$ , let

$$U(k_1, \dots, k_t) = U_1(k_1) \cup U_2(k_2) \cup \dots \cup U_t(k_t). \quad (11)$$

We restrict our attention to  $U(k_1, k_2) = U_1(k_1) \cup U_2(k_2)$ , the set of all strings of length at most  $k_1$  that have exactly one non- $J$  character and the set of strings of length at most  $k_2$  that have exactly two non- $J$  characters.

When we refer to the multiplicity with which a string  $w$  that contains wildcards ( $J$ 's) occurs as a subsequence of a string  $p$  that contains no wildcards (denoted by  $N(w, p)$ ), we map each wildcard to either  $X$  or  $Y$ . For example, if  $w = (J, X)$  and  $p = (Y, X, Y, X)$ , we have  $N(w, p) = 4$  because  $(X, X)$  and  $(Y, X)$  occur as subsequences of  $p$  with multiplicity 1 and 3, respectively. Let  $p$  and  $q$  be two binary strings. We write  $p \sim_{U_r(k)} q$  if  $N(w, p) = N(w, q)$  for all  $w \in U_r(k)$ . In addition, we write  $p \sim_{U(k_1, k_2)} q$  if  $N(w, p) = N(w, q)$  for all  $w \in U(k_1, k_2)$ .

Next, let  $S_U(k_1)$  be the smallest integer  $m$  for which there exist distinct strings  $p$  and  $q$  of length  $m$  such that  $p \sim_{U(k_1)} q$ . Similarly, let  $S_U(k_1, k_2)$  be the smallest integer  $m$  for which there exist distinct strings  $p$  and  $q$  of length  $m$  such that  $p \sim_{U(k_1, k_2)} q$ . The following lemma is used in our subsequent derivations.

**Lemma 3.** [6] *Let  $k_1 \geq k_2 \geq 2$  and  $\kappa = k_1^2 + k_2^2(k_2 - 1)/2$ . Then  $S_U(k_1, k_2) \leq \kappa(\lg \kappa + \lg \lg \kappa + 1) = (1 + o(1))\kappa \lg \kappa$ .*

Let  $N_g(w, x)$  be the number of times a string  $w$  appears as a gapped subsequence of  $x$  (i.e., so that all indices in  $x$  are nonadjacent). When  $N_g(w, x) = N_g(w, y)$  for all strings  $w$  of length  $\leq k$ , then we write  $x \sim^{k(s)} y$ , i.e.  $\mathcal{B}^{(k)}(x) = \mathcal{B}^{(k)}(y)$ .

Let  $\Gamma = \{X, Y\}$  and let  $\Sigma$  be an arbitrary alphabet. For a finite-length string  $x$  over  $\Sigma$ , define  $x_0$  to be the string obtained by padding  $x$  with one 0 at both ends. For a finite-length string  $p$  over  $\Gamma$  and two finite-length strings  $x$  and  $y$  over  $\Sigma$ , let  $h_{x,y}(p)$  be the string obtained from  $p$  by replacing each  $X$  by the string  $x_0$  and each  $Y$  by the string  $y_0$ . For example, if  $p = (X, Y)$  and  $x = (0, 1, 0, 1)$  and  $y = (1, 1, 0, 0)$ , then  $x_0 = (0, 0, 1, 0, 1, 0)$ ,  $y_0 = (0, 1, 1, 0, 0, 0)$  and  $h_{x,y}(p) = (0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 0)$ . We are now ready to prove an analogue of Lemma 9 from [6] for the case of gapped  $k$ -decks.

**Lemma 4.** *Let  $x$  and  $y$  be two distinct strings in  $\Sigma^n$  such that  $x \sim^{k(s)} y$ ,  $\mathcal{B}_L^{(k)}(x) = \mathcal{B}_L^{(k)}(y)$ ,  $\mathcal{B}_R^{(k)}(x) = \mathcal{B}_R^{(k)}(y)$  and  $\mathcal{B}_{LR}^{(k)}(x) = \mathcal{B}_{LR}^{(k)}(y)$ . Let  $p$  and  $q$  be two distinct binary strings in  $\Gamma^m$ , such that for some  $\sigma \in \{0, 1, 2\}$ , we have  $p \sim_{U(2k+\sigma, k+\sigma)} q$ . Then,  $h_{x,y}(p)$  and  $h_{x,y}(q)$  are distinct and we have  $h_{x,y}(p) \sim^{3k+\sigma(s)} h_{x,y}(q)$ . The same result holds when puncturing  $h_{x,y}(p)$  and  $h_{x,y}(q)$  on the left, right, and on both sides by one bit.*

*Proof:* Due to space limitations, we only provide a sketch of the proof. Let  $w$  be a string of length at most  $3k + \sigma$  in  $\Sigma$ , and  $p = (p_1, \dots, p_m)$ . The idea of the original proof [6] for the ungapped case is as follows: Each mapping that takes  $w$  to  $h_{x,y}(p)$  (as a subsequence) defines a splicing of  $w$  of the form  $w = (w_1, w_2, \dots, w_m)$ , where  $w_i$  is the preimage of  $h_{x,y}(p_i)$ . Note that some  $w_i$  strings may be empty. Hence, we can write the set of all mappings which take  $w$  to  $h_{x,y}(p)$  as the union of direct products (see [6] for the specific notation)  $\bigcup_{t \geq 1} \bigcup_l \bigcup_r \prod_{i=1}^t \mathcal{N}(w_{l,i}, h_{x,y}(p_{r(i)}))$ , where  $t$  is the number of nonempty segments, the second union is taken over all functions  $l$  which partition  $w$  into  $t$  nonempty segments for a fixed  $t$ , the third union is taken over all functions  $r$  mapping the  $t$  chosen nonempty segments in  $w$  to  $t$  of the segments in  $h_{x,y}(p)$  (equivalently, each  $r$  corresponds to a way in which we pick  $t$  out of  $m$  segments of  $h_{x,y}(p)$ ), and finally,  $\mathcal{N}(w_{l,i}, h_{x,y}(p_{r(i)}))$  denotes the set of mappings taking the  $i$ -th nonempty segment of  $w$  into the corresponding chosen segment in  $h_{x,y}(p)$ . Note that for every mapping  $f$  which takes  $w$  to  $h_{x,y}(p)$ , there is a specific  $t$ ,  $l$  and  $r$  that correspond to  $f$ . Furthermore,  $f$  is the direct product of  $t$  mappings, each taking one of the  $t$  nonempty segments of  $w$  to the corresponding segment of  $h_{x,y}(p)$  (which is all uniquely determined by fixing  $t$ ,  $l$  and  $r$ ). Now, by converting this expression into a corresponding sum, we get  $N(w, h_{x,y}(p)) = \sum_{t \geq 1} \sum_l \sum_r \prod_{i=1}^t N(w_{l,i}, h_{x,y}(p_{r(i)}))$ . Since  $w$  has length at most  $3k + \sigma \leq 3k + 2$ , it has at most two segments of length  $> k$ . Therefore, we have three types of  $l$ 's (i.e., ways of partitioning  $w$  into  $t$  nonempty segments for a fixed  $t$ ): The ones with no segments of length  $> k$ , the ones with one such segment and the ones with two such segments. This means that  $N(w, h_{x,y}(p))$  is a triple sum. Since  $p \sim_{U(2k+\sigma, k+\sigma)} q$ , after some calculations we get  $N(w, h_{x,y}(p)) = N(w, h_{x,y}(q))$ .

To adapt this procedure for the gapped case, we need to show that  $N_g(w, h_{x,y}(p)) = N_g(w, h_{x,y}(q))$ . The first difference is that each gapped mapping that takes  $w$  to  $h_{x,y}(p)$



(as a gapped subsequence) defines a splicing of  $w$  of the form  $w = w_1 z_1 w_2 z_2 \dots z_{m-1} w_m$ , where  $w_i$  is the preimage of  $h_{x,y}(p_i)$  and  $z_i$  is the preimage of the  $i$ -th pair of 0s that we added between the  $x$  and between the  $y$  strings and between the  $x$  and  $y$  strings when constructing  $h_{x,y}(p)$ . Again, note that some  $w_i$  and  $z_i$  strings may be empty.

However, an important difference between the gapped and ungapped case is that we need to consider different cases based on whether  $z_i$  is empty or not, for all indices  $i$ . This is because if  $z_1$  is nonempty, for example, then we need to make sure that we do not use the rightmost bit in  $h_{x,y}(p_1)$  (or the leftmost bit in  $h_{x,y}(p_2)$ , depending on whether the gapped mapping takes  $z_1$  to the first or second 0 from the pair of 0s in  $h_{x,y}(p)$  that are positioned between  $h_{x,y}(p_1)$  and  $h_{x,y}(p_2)$ ). This case ( $z_1$  nonempty) gives rise to several additional cases that need to be considered, depending on which of the strings  $z_2, z_3, \dots, z_{m-1}$  are empty. In other words, we can write out  $h_{x,y}(p) = (0, s_1, 0, 0, s_2, 0, 0, \dots, 0, 0, s_m, 0)$  where each  $s_i \in \{x, y\}$ . Any gapped subsequence will then be of the form  $(J_0, \mathcal{D}_{\alpha_1}^{(i_1)}(s_1), J_1, \mathcal{D}_{\alpha_2}^{(i_2)}(s_2), J_2, \dots, J_{m-1}, \mathcal{D}_{\alpha_m}^{(i_m)}(s_m), J_m)$  where  $J_j \in \{0, \emptyset\}$  and  $\alpha_j \in \{\emptyset, L, R, LR\}$ . Here, each  $J_j$  represents a 0 in the padding and whether it is a part of the subsequence or not. Then, depending on whether or not we use the padding, we puncture  $s_j$  on the left, right, both, or neither. This is captured by the indices  $\alpha_j$ 's. We also have that  $\sum_{j=1}^m i_j + \sum_{j=1}^m |J_j| \leq 3k + \sigma$ . Since by our assumptions  $\mathcal{D}_{\alpha}^{(i)}(p_j) = \mathcal{D}_{\alpha}^{(i)}(q_j)$  for all  $i \leq k$ ,  $1 \leq j \leq m$ ,  $\alpha_j \in \{\emptyset, L, R, LR\}$ , we have an equivalence between  $h_{x,y}(p)$  and  $h_{x,y}(q)$  for each  $i_j \leq k$ . By the summation constraint, there are at most two indices  $j$  such that  $i_j \geq k + 1$ . Let us consider the case when there is exactly one such  $j$ , denoted by  $j^*$ . In this case we have to pick fewer than  $2k + \sigma$  of the remaining characters to obtain the final subsequence. We can also divide  $w$  into a collection of  $w_j$ 's, where  $w_j$  is the string mapped to one of the blocks and  $|w_j| < k$  for  $j \neq j^*$ . The multiplicity of  $w$  can be seen to be  $N_g(w_{j^*}, (s_{j^*})_{\alpha_{j^*}}) \prod_{j \neq j^*} N_g(w_j, (x)_{\alpha_j})$ , where  $\alpha$  once again depends on whether the bit used for padding is included in the subsequence. By our assumption we have  $N(J^{a_1} A J^{a_2}, p) = N(J^{a_1} A J^{a_2}, q)$ , where  $J^{a_1}$  is a sequence of  $a_1$  concatenated wildcard characters,  $A \in \{x, y\}$ ,  $a_1 + a_2 < 2k + \sigma$ . Hence, there are equally many  $s_{j^*}$ 's in  $h_{x,y}(p)$  and  $h_{x,y}(q)$ . Using similar arguments and the fact that  $N(J^{a_1} A J^{a_2} B J^{a_3}, p) = N(J^{a_1} A J^{a_2} B J^{a_3}, q)$  we can also prove the equivalence for the case of two indices  $j$  for which  $i_j \geq k + 1$ . This leads to  $N_g(w, h_{x,y}(p)) = N_g(w, h_{x,y}(q))$ . ■

The lemma gives rise to the following important Corollary.

**Corollary 5.** For every  $\sigma \in \{0, 1, 2\}$ , one has  $G^*(3k + \sigma) \leq (G^*(k) + 2)(S_U(2k + \sigma, k + \sigma))$ .

Combining the above corollary with Lemma 3 and Lemma 4 leads to an upper bound for  $G(k)$  as follows. First, we set

$$\kappa = (2k + \sigma)^2 + (k + \sigma)^2(k + \sigma - 1)/2,$$

which equals

$$\begin{aligned} & \frac{1}{2}(\sigma^3 + (3k + 1)\sigma^2 + 3\sigma(k^2 + 2k) + (1 + 7/k)k^3) \\ & = \left(\frac{1}{2} + o(1)\right)k^3, \end{aligned}$$

in Lemma 3 to obtain

$$S_U(2k + \sigma, k + \sigma) \leq C(k)k^3 \log_3 k, \quad (12)$$

where  $C(k) = \frac{3 \log 3}{2} + o(1)$ . We also have that  $C(k) \leq 10$  for  $k \geq 9$ ,  $C(k) \leq 3$  for  $k \geq 3^5$  [6]. Using the inequality from Corollary 5 and (12), we set  $k_0 = k$  and for  $i > 0$ ,  $k_i = \lfloor k_{i-1}/3 \rfloor \leq k/3^i$ . We stop the recursion with  $i = i_0$ , where  $k_{i_0} \leq 4$  (so  $G^*(k_{i_0}) \leq 4(2^4 - 1) = 60$ ) and get

$$\begin{aligned} G^*(k) & \leq G^*(k_{i_0}) \prod_{i=1}^{i_0} S_U(2k_i + \sigma_i, k_i + \sigma_i) \\ & \quad + 2 \sum_{i=1}^{i_0} \prod_{j=1}^i S_U(2k_i + \sigma_i, k_i + \sigma_i). \end{aligned} \quad (13)$$

Combining the above bound with that on  $S_U$  we obtain

$$\begin{aligned} G^*(k) & \leq G^*(k_{i_0}) \prod_{i=1}^{i_0} C(k_i)(k_i)^3 \log_3(k_i) \\ & \quad + 2 \sum_{i=1}^{i_0} \prod_{j=1}^i C(k_i)(k_i)^3 \log_3(k_i) \\ & \leq 3^{\log_3(60) + \sum_{i=1}^{\lfloor \log_3(k/4) \rfloor} [O(1) + 3(\log_3 k - i) + \log_3(\log_3 k - i)]} \\ & \quad + \sum_{i=1}^{\lfloor \log_3(k/4) \rfloor} 3^{\log_3 2 + \sum_{j=1}^i [O(1) + 3(\log_3 k - j) + \log_3(\log_3 k - j)]} \\ & = 3^{O(1) + O(\log_3 k) + O(\log_3^2 k) + O(\log_3 k \log_3 \log_3 k)} \\ & \quad + O(\log_3(k)) 3^{O(1) + O(\log_3 k) + O(\log_3^2 k) + O(\log_3 k \log_3 \log_3 k)} \\ & = O(\log_3(k)) 3^{O(\log_3^2 k)}. \end{aligned}$$

Since we have  $G(k) \leq G^*(k) - 2$ , we also have  $G(k) \leq O(\log_3(k)) 3^{O(\log_3^2 k)}$ . By bounding  $G^*(k)$  we also obtain

$$\begin{aligned} G(k) & \leq (4^{2^{k/3}} - 1) + 2)C(k/3)(k/3)^3 \log_3(k/3) - 2 \\ & \leq 4/27 * 2^{k/3} * k^3 \log_3(k/3) * C(k/3) - 2 \end{aligned}$$

for  $k \geq 28$ . Since in this case  $C(k/3) \leq 10$ , we arrive at

$$G(k) \leq 1.482 * 1.26^k * k^3 \log_3(k/3) - 2, \quad (14)$$

In comparison, the general bound for the ungapped case, derived in [6], reads as

$$S(k) \leq 1.2\Gamma(\log_3 k) 3^{(3/2) \log_3^2 k - (1/2) \log_3 k}, \quad k \geq 85.$$

The bounds are summarized in the tables below.

$k$	Bound
2-4	Exact values: 6,13,24
5-27	$4(2^k - 1)$
$> 27$	$1.482 * 1.26^k * k^3 \log_3(k/3) - 2$

$k$	28	29	30	31	32	33
$G(k) \leq$	42742211	60773950	86039831	121319982	170424514	238563374

## REFERENCES

- [1] J. Acharya, H. Das, O. Milenkovic, A. Orlitsky, and S. Pan, String reconstruction from substring compositions, *SIAM Journal on Discrete Mathematics* 29, no. 3 (2015): 1340-1371.
- [2] J-P. Allouche and J. Shallit, "The ubiquitous Prouhet-Thue-Morse Sequence," *In Sequences and their Applications*, pp. 1-16. Springer, London, 1999.
- [3] T. Batu, S. Kannan, S. Khanna, and A. McGregor, Reconstructing strings from random traces, *Departmental Papers (CIS)* (2004): 173.
- [4] Z. Chase, Separating words and trace reconstruction, *In Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, (2021) 21–31.
- [5] M. Cheraghchi, R. Gabrys, O. Milenkovic, and J. Ribeiro, "Coded trace reconstruction," *IEEE Transactions on Information Theory*, 66, no. 10, pp. 6084-6103, 2020.
- [6] M. Dudik and L.J. Schulman, Reconstruction from subsequences, *Journal of Combinatorial Theory, Series A* 103(2) (2003), 337–348.
- [7] J. Chrisnata, H. M. Kiah, S. Rao, A. Vardy, E. Yaakobi, A. Yao, "On the number of distinct k-decks: Enumeration and bounds," *19th International Symposium on Communications and Information Technologies (ISCIT)* (2019) 519–524.
- [8] R. Gabrys and O. Milenkovic, "Unique reconstruction of coded strings from multiset substring spectra," *IEEE Transactions on Information Theory* 65, no. 12 (2019): 7682–7696.
- [9] R. Gabrys and O. Milenkovic, "The hybrid k-deck problem: Reconstructing sequences from short and long traces," *IEEE International Symposium on Information Theory (ISIT)* (2017) 1306–1310.
- [10] L.O. Kalashnik, The reconstruction of a word from fragments, *Numerical mathematics and computer technology* (1973), 56–57.
- [11] H. M. Kiah, G. J. Puleo, and O. Milenkovic, Codes for DNA sequence profiles, *IEEE Transactions on Information Theory* 62, no. 6 (2016): 3125-3146.
- [12] B. Manvel, A. Meyerowitz, A. Schwenk, K. Smith and P. Stockmeyer, Reconstruction of sequences, *Discrete Mathematics* 94(3) (1991), 209–219.
- [13] D. Margaritis and S. S. Skiena, Reconstructing strings from substrings in rounds, *In Proceedings of IEEE 36th Annual Foundations of Computer Science* (1995) 613–620.
- [14] A. D. Scott, Reconstructing sequences, *Discrete Mathematics* (1997) 175 1–3 231-238.
- [15] E. Ukkonen, Finding approximate patterns in strings, *Journal of Algorithms*, 6, no. 1 (1985): 132-137.
- [16] S.M.H. Yazdi, R. Gabrys and O. Milenkovic, "Portable and error-free DNA-based data storage," *Scientific Reports* 7, no. 1 pp. 1-6, 2016 (online) / 2017 (print).