

On the Number of Distinct k -Decks: Enumeration and Bounds

Johan Chrisnata*, Han Mao Kiah*, Sankeerth Rao†, Alexander Vardy†, Eitan Yaakobi‡, and Hanwen Yao†

*School of Physical and Mathematical Sciences, Nanyang Technological University, Singapore

†Department of Electrical & Computer Engineering, University of California San Diego, LA Jolla, CA, USA

‡Department of Computer Science, Technion Israel Institute of Technology, Haifa, Israel

Emails: johanchr001@ntu.edu.sg, hmkih@ntu.edu.sg, sankeerth1729@gmail.com, avardy@ucsd.edu,

yaakobi@cs.technion.ac.il, hay125@eng.ucsd.edu

Abstract—The k -deck of a sequence is defined to be the multiset of all its subsequences of length k and let $D_k(n)$ denote the number of distinct k -decks for binary sequences of length n . In this paper, we determine the exact value of $D_k(n)$ for small values of k and n and provide asymptotic estimates of $D_k(n)$ when k is fixed.

Specifically, for fixed k , we provide a trellis-based method to compute $D_k(n)$ in time polynomial in n . We then compute $D_k(n)$ for $k \in \{3, 4, 5, 6\}$ and $k \leq n \leq 30$. We also improve the asymptotic upper bound on $D_k(n)$ and in particular, show $D_k(n) = O(n^{(k-1)2^{k-1}+1})$. For the specific case when $k = 3$, we show $D_3(n) = \Omega(n^6)$ while the upper bound states that $D_3(n) = O(n^9)$.

I. INTRODUCTION

A protein macromolecule is a long string of amino acids. However, current sequencing technology either is unable to determine the *long* sequence directly or reads the sequence at a high error rate. Therefore, most sequencing methods obtain information about its *short* substrings or subsequences and attempt to infer or *reconstruct* the original string from this partial information. This gives rise to a myriad of combinatorial problems, known as string reconstruction problems [1], [2], [5], [6], [8], [10].

In this paper, we study the reconstruction problem involving k -decks. First described by Kalashnik [7], the k -deck of a sequence is defined to be the multiset of all its subsequences of length k . Traditionally, the k -deck problem is to determine $S(k)$, the smallest value of n such that all sequences of length n have unique k -decks. The exact values of $S(k)$ are known for $k \leq 5$ and both upper and lower bounds for $S(k)$ have been extensively studied [3], [4], [7], [9], [12], [15].

Motivated by applications in DNA-based data storage (see [16] for a broad overview), we study the *coded* version of the k -deck problem. Consider sequences or words of length n . Instead of requiring *all* words to have different k -decks, we choose a subset of these words, or a *codebook*, such that every codeword in this codebook can be uniquely identified by its k -deck. In this setting, we consider the following fundamental problem: *how large can this codebook be?* Equivalently, this problem can be restated as an enumeration problem:

Consider all words of length n . How many distinct k -decks are there?

Let $D_k(n)$ denote this quantity and $D_k(n)$ is the object of study for this paper. In another context, Rigo and Salimov used the term *k -binomial equivalence* to describe words with the same k -deck and provided rudimentary upper bounds on $D_k(n)$ [14]. In particular, Rigo and Salimov determined $D_2(n)$ and showed that $D_k(n) = O(n^{\Delta(k)})$, where $\Delta(k) = 2((k-1)2^k + 1)$. The values

of $D_3(n)$ for $n \leq 16$ are listed on the On-Line Encyclopedia of Integer Sequences [13].

In this paper, we provide a trellis-based method (see [11] for the definition of a trellis) to compute $D_k(n)$ and determine the exact values of $D_k(n)$ for $k \in \{3, 4, 5, 6\}$ and $n \leq 30$. We also provide asymptotic estimates for the case where k is fixed. In particular, we improve the asymptotic upper bound to $O(n^{\Delta(k)/4+1/2})$ for $k \geq 2$.

The paper is organized as follows. Section II formally defines the problem, summarizes previous results and states our contributions. Section III details a polynomial-time algorithm that computes $D_k(n)$. Section IV then provides an upper bound on $D_k(n)$ for general k , while Section V provides a lower bound for the case $k = 3$.

II. PROBLEM STATEMENT AND CONTRIBUTIONS

Let $X = x_1x_2 \cdots x_n$ be a binary word of length n . For $A \subseteq \{1, 2, \dots, n\}$, we use X_A to denote the subsequence with indices in A . In other words, $X_A = x_{a_1}x_{a_2} \cdots x_{a_k}$ where $a_1 < a_2 < \cdots < a_k$ and $A = \{a_1, a_2, \dots, a_k\}$. For $k \leq n$, the k -deck of X , denoted by $D_k(X)$, refers to the multiset of all $\binom{n}{k}$ subsequences of length k . We represent the k -deck of a word X by an integer-valued vector of length 2^k . Specifically, $D_k(X) \triangleq (X_\alpha)_{\alpha \in \{0,1\}^k}$, where X_α denotes the number of occurrences of α as a subsequence of X and the indices in $\{0, 1\}^k$ are presented in an increasing lexicographic order.

Example 1. Let $X = 110011$. Then $X_{\{3,5,6\}} = X_{\{4,5,6\}} = 011$ and we check that $X_{011} = 2$. Furthermore, $D_1(X) = (2, 4)$, $D_2(X) = (1, 4, 4, 6)$ and $D_3(X) = (0, 2, 0, 2, 2, 8, 2, 4)$. \square

Two words X and Y are said to be k -equivalent, or $X \sim_k Y$ if their k -decks are the same, i.e. $D_k(X) = D_k(Y)$.

Example 2. Let $Y = 101101$. Then $D_1(Y) = (2, 4)$, $D_2(Y) = (1, 4, 4, 6)$ and $D_3(Y) = (0, 1, 2, 3, 1, 6, 3, 4)$. Hence, $X \sim_k Y$ for $k \in \{1, 2\}$, but $X \not\sim_3 Y$. \square

It can be shown that \sim_k defines an equivalence relationship on all binary words. Furthermore, if X and Y have the same k -deck, then the lengths of X and Y are necessarily the same. Hence, we fix n and partition the binary words of length n using the relation \sim_k . Then we set $D_k(n)$ to be the resulting number of equivalence classes. In this paper, we determine the exact value of $D_k(n)$ for $k \in \{3, 4, 5, 6\}$ and $k \leq n \leq 30$ and provide asymptotic estimates of $D_k(n)$ when k is fixed.

For $k \in \{1, 2\}$, the exact values on $D_k(n)$ and characterization of $D_k(\mathbf{X})$ have been determined [14, Lemma 4].

Proposition 1. Suppose that \mathbf{X} is a binary word of length n .

- (i) Then $D_1(\mathbf{X}) = (n - w, w)$, where $X_1 = w$. Therefore, $D_1(n) = n + 1$.
- (ii) Then $D_2(\mathbf{X}) = \left(\binom{n-w}{2}, t, w(n-w) - t, \binom{n-w}{2}\right)$, where $X_1 = w$ and $X_{01} = t$. Therefore, $D_2(n) = (n^3 + 5n + 6)/6$.

For $k \geq 3$, the best known upper bound on $D_k(n)$ is below.

Theorem 2 (Rigo and Salimov [14, Proposition 5]). For all $n \geq k$, we have that $D_k(n) \leq \prod_{\ell=1}^k \left(\binom{n}{\ell} + 1\right)^{2^{\ell-1}}$. When k is fixed, we have that $D_k(n) = O\left(n^{2((k-1)2^k+1)}\right)$.

In addition to $D_k(n)$, we define $S(k) \triangleq \min\{n : D_k(n) < 2^n\}$. The exact values of $S(k)$ have been determined for $k \in \{3, 4, 5\}$ (see [4] for a survey). The first open case is $S(6)$ and the best known upper bound is given by Manvel et al. who constructed a pair of words of length thirty with the same 6-deck [12, Example 4].

Theorem 3. $S(6) \leq 30$.

For completeness, we present the best known upper and lower bounds for $S(k)$ that were summarized in [4].

Theorem 4. We have that $S(k) = \Omega(k^2)$ and that

$$S(k) \leq \begin{cases} 1.75 \cdot 1.62^k, & \text{for } 7 \leq k \leq 28, \\ 0.25 \cdot 1.17^k k^3 \log k, & \text{for } 29 \leq k \leq 84, \\ 3^{(3/2+o(1)) \log_3^2 k}, & \text{for } k \geq 85. \end{cases}$$

A. Our Contributions

(A) In Section III, we use a trellis structure to describe the k -decks and using this insight, we then provide an algorithm that enumerates all k -decks efficiently. When k is a constant, the technique computes $D_k(n)$ in polynomial time. We then compute $D_k(n)$ for $k \in \{3, 4, 5, 6\}$ and $k \leq n \leq 30$ and establish that $S(6) = 30$.

(B) In Section IV, we improve the asymptotic upper bound on $D_k(n)$. In particular, we show that

$$D_k(n) = O\left(n^{(k-1)2^{k-1}+1}\right). \quad (1)$$

(C) In Section V, we look at the specific case for $k = 3$ and show that $D_3(n) = \Omega(n^6)$. On the other hand, we note that (1) shows that $D_3(n) = O(n^9)$.

III. POLYNOMIAL-TIME ENUMERATION

In this section, we introduce a trellis-based algorithm that calculates $D_k(n)$ for a fixed value of k . To compute $D_k(n)$, we construct a trellis with levels $i \in \{k, k+1, \dots, n\}$ such that we are able to compute $D_k(i)$ at level i . At level $i+1$, instead of naively enumerating all k -decks for all possible $\mathbf{X} \in \{0, 1\}^{i+1}$, our algorithm runs recursively and calculates the set $\{D_k(\mathbf{X}) : \mathbf{X} \in \{0, 1\}^{i+1}\}$ from the set $\{D_k(\mathbf{X}) : \mathbf{X} \in \{0, 1\}^i\}$, which reduces the complexity from 2^{i+1} down to $D_k(i)$.

To do so, we make two important combinatorial observations. In [14], it was observed that $\mathbf{X} \underset{k}{\sim} \mathbf{Y}$ implies $\mathbf{X} \underset{s}{\sim} \mathbf{Y}$ for all $1 \leq s < k$. The next proposition gives an explicit method to compute the s -deck from a k -deck.

Proposition 5. Let $\mathbf{X} \in \{0, 1\}^n$, $\alpha \in \{0, 1\}^s$ with $1 \leq s < k$, then

$$\binom{n-s}{k-s} \mathbf{X}_\alpha = \sum_{\beta \in \{0, 1\}^k} \beta_\alpha \mathbf{X}_\beta \quad (2)$$

Proof. For $A \subset [n]$, recall that \mathbf{X}_A denotes the subsequence of \mathbf{X} with indices in A . To demonstrate (2), we consider the two collections \mathcal{A} and \mathcal{B} of tuples. First, we set

$$\mathcal{A} \triangleq \{(A; S) : A, S \subset [n], \mathbf{X}_A = \alpha, |S| = k - s, A \cap S = \emptyset\}.$$

For each occurrence of α in \mathbf{X} , we fix A and have $\binom{n-s}{k-s}$ choices for S . Therefore, the left hand side of (2) counts the number of tuples in \mathcal{A} .

Next, we set

$$\mathcal{B} \triangleq \{(\beta; B; T) : B, T \subset [n], \mathbf{X}_B = \beta, T \subseteq B, \mathbf{X}_{B \setminus T} = \alpha\}.$$

For each $\beta \in \{0, 1\}^k$ and each occurrence of β in \mathbf{X} , we have β_α choices for T . Therefore, the right hand side of (2) counts the number of tuples in \mathcal{B} .

To establish (2), it remains to exhibit a bijection between \mathcal{A} and \mathcal{B} . Consider the map $\phi : \mathcal{A} \rightarrow \mathcal{B}$ such that $\phi(A; S) = (\beta; B; T)$, where $B = A \cup S$ and $\beta = \mathbf{X}_B$.

For the inverse, we consider the $\psi : \mathcal{B} \rightarrow \mathcal{A}$ such that $\psi(\beta; B; T) = (A; S)$, where $A = B \setminus T$. It is not difficult to verify that both $\phi \circ \psi$ and $\psi \circ \phi$ are identity maps on their respective domains. Therefore, we establish (2). \square

For $\mathbf{X} \in \{0, 1\}^n$ and $a \in \{0, 1\}$, let $(\mathbf{X}|a)$ denote the concatenation of \mathbf{X} and a . Our second observation states that we can compute $D_k(\mathbf{X}|a)$ from $D_{k-1}(\mathbf{X})$ and $D_k(\mathbf{X})$.

Proposition 6. Let $\mathbf{X} \in \{0, 1\}^n$, $\alpha \in \{0, 1\}^k$, then

$$(\mathbf{X}|0)_\alpha = \mathbf{X}_\alpha + \mathbf{X}_\beta, \text{ where } (\beta|0) = \alpha, \quad (3)$$

$$(\mathbf{X}|1)_\alpha = \mathbf{X}_\alpha + \mathbf{X}_\beta, \text{ where } (\beta|1) = \alpha. \quad (4)$$

Proof. Consider the collection of index subsets:

$$S = \{A \subset [n] : (\mathbf{X}|0)_A = \alpha\}.$$

Then S can be written as a disjoint union of S_1 and S_2 where

$$S_1 = \{A \in S : n+1 \notin A\}, \quad S_2 = \{A \in S : n+1 \in A\}.$$

Since $(\mathbf{X}|0)_\alpha = |S|$, $\mathbf{X}_\alpha = |S_1|$ and $\mathbf{X}_\beta = |S_2|$, we have (3). Equation (4) can be proved in the same manner. \square

A more general version of Proposition 6 was given in [14]. As the authors did not furnish a proof, we provide one here for completeness.

We are ready to present our trellis. As mentioned earlier, the trellis has levels $i \in \{k, k+1, \dots, n\}$. Each vertex at level i represents a k -deck of some word of length i and we denote the vertices at level i with $\mathbb{D}_k(i)$. Therefore,

$$\mathbb{D}_k(i) = \{D_k(\mathbf{X}) : \mathbf{X} \in \{0, 1\}^i\}.$$

Using (2), (3), and (4), each vertex in $\mathbb{D}_k(i)$ is extended by two edges labeled '0' and '1' to two vertices in $\mathbb{D}_k(i+1)$. The resulting trellis is *biproper*, in other words, every vertex has exactly two outgoing arcs with distinct labels and at most two incoming arcs with distinct labels. Furthermore, $\mathbb{D}_k(i)$ is the set of all k -decks of words of length i and the set of paths to a vertex,

Algorithm 1: Compute $D_k(n)$

```

1 initialize  $\mathbb{D}_k(k)$  as follow:
    $\mathbb{D}_k(k) = \{(1, 0, \dots, 0), (0, 1, \dots, 0), \dots, (0, \dots, 0, 1)\}$ 
2 for  $i = k$  to  $n - 1$  do
3   initialize  $\mathbb{D}_k(i + 1)$  as an empty set
4   for every  $k$ -deck  $D \in \mathbb{D}_k(i)$  do
5     for  $a \in \{0, 1\}$  do
6       Let  $X$  be a word such that  $D_k(X) = D$ .
7       Using equation (2) with (3) or (4),
8       compute  $D' = D_k(X|a)$ 
9       if  $D' \notin \mathbb{D}_k(i + 1)$  then
10        insert  $D'$  to the set  $\mathbb{D}_k(i + 1)$ 
11 return  $D_k(n) = |\mathbb{D}_k(n)|$ 
    
```

or a k -deck, is the set of all binary sequences having this k -deck. See Figure 1 for a trellis section for $k = 2$ and levels $i \in \{4, 5\}$. A formal description of the enumeration method is detailed in Algorithm 1.

We discuss the computational complexity of Algorithm 1. First, we note that X need not be explicitly found in Line 6. To compute $D_k(X|a)$, it suffices to apply (2), (3) and (4) to $D = D_k(X)$. Also, since Equations (2), (3) and (4) involve sums with at most 2^k terms, Lines 6 and 7 take constant time.

The time complexity for Line 7 depends on the data structure we used for $\mathbb{D}_k(n)$. If we use a binary search tree, we can insert each “new” k -deck in $O(D_k(n) \log D_k(n))$ time using $O(D_k(n))$ space. Therefore, Algorithm 1 runs in $O(nD_k(n) \log D_k(n))$ time using $O(nD_k(n))$ space. For fixed k , since $D_k(n)$ is polynomial in n by Theorem 2, the algorithm has space and time complexity polynomial in n .

To conclude this section, we compute the values of $D_k(n)$ for $k \in \{3, 4, 5, 6\}$ and $k \leq n \leq 30$ and present them in Table I. In particular, we computed that $D_6(n) = 2^n$ for $n \leq 29$. Together with Theorem 3, we established the following.

Theorem 7. $S(6) = 30$.

IV. UPPER BOUNDS ON $D_k(n)$

Fix $k \geq 3$. In this section, we derive an upper bound on $D_k(n)$. To this end, we fix some $(k - 1)$ -deck D' and consider the collection \mathcal{F} of all words of length n whose $(k - 1)$ -deck is given by D' . Suppose that the number of k -equivalence classes in \mathcal{F} is at most U for all choice of $(k - 1)$ -decks. Then an upper bound for $D_k(n)$ is simply given by $D_{k-1}(n)U$.

To find U , we consider an additional parameter $1 \leq m \leq k - 1$ and define $J(k, m)$ to be all binary words of length k and weight m . Similar to [12], we relax the notion of k -equivalence and define the (k, m) -equivalence relation: $X \underset{(k, m)}{\sim} Y$ if and only if $(X_\beta)_{\beta \in J(k, m)} = (Y_\beta)_{\beta \in J(k, m)}$. Suppose that the number of (k, m) -equivalence classes in \mathcal{F} is at most $U(m)$ for all choice of $(k - 1)$ -decks. Then we can obtain $U = \prod_{m=1}^{k-1} U(m)$.

We now proceed to estimate $U(m)$. Now, suppose $X, Y \in \mathcal{F}$. Since $X \underset{k-1}{\sim} Y$, we have that $X \underset{1}{\sim} Y$. In other words, X and Y have the same weight. Hence, we let w denote the weight of any word in \mathcal{F} .

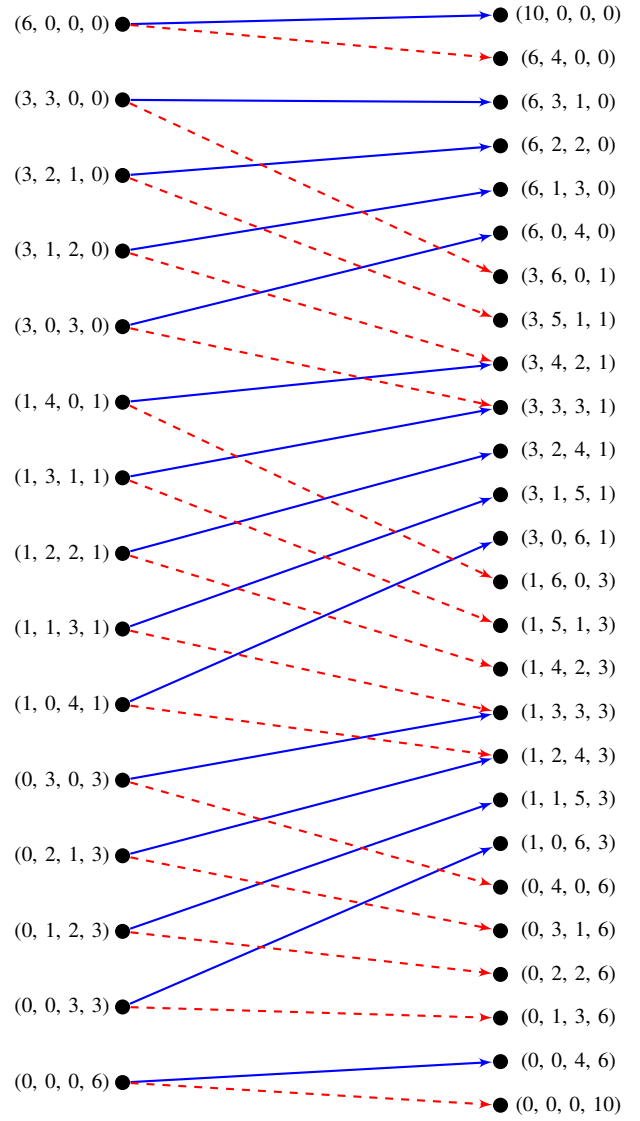


Fig. 1. Trellis section for $k = 2$ and levels $i \in \{4, 5\}$. Left vertices belong to $\mathbb{D}_2(4)$, while right vertices belong to $\mathbb{D}_2(5)$. Blue solid edges corresponds to label ‘0’, while red dashed edges corresponds to label ‘1’.

Let $X \in J(n, w)$, the set of all words of length n and weight w . Suppose that $\beta \in J(k, m)$ and let α be a subsequence of length $k - 1$ of β . Then α necessarily belongs to $J(k - 1, m - 1) \cup J(k - 1, m)$. Recall that β_α is the number of occurrences of α in β . Below we obtain a combinatorial relationship between X_α , X_β and β_α , which is a refinement of Proposition 5.

Proposition 8. Let $X \in J(n, w)$ and $1 \leq m \leq k - 1$.

If $\alpha \in J(k - 1, m - 1)$, then

$$(w - m + 1)X_\alpha = \sum_{\beta \in J(k, m)} \beta_\alpha X_\beta. \quad (5)$$

If $\alpha \in J(k - 1, m)$, then

$$(n - k + m - w + 1)X_\alpha = \sum_{\beta \in J(k, m)} \beta_\alpha X_\beta. \quad (6)$$

Proof. Let $X = x_1 x_2 \dots x_n$. Suppose that $\alpha \in J(k - 1, m - 1)$. To demonstrate (5), we proceed in a similar manner as in the

$n \setminus k$	3	4	5	6
3	8	—	—	—
4	16	16	—	—
5	32	32	32	—
6	64	64	64	64
7	126	128	128	128
8	247	256	256	256
9	480	512	512	512
10	926	1024	1024	1024
11	1764	2048	2048	2048
12	3337	4092	4096	4096
13	6208	8176	8192	8192
14	11408	16328	16384	16384
15	20608	32604	32768	32768
16	36649	65075	65534	65536
17	63976	129824	131064	131072
18	109866	258906	262120	262144
19	185012	516168	524212	524288
20	306285	1028448	1048360	1048576
21	497190	2048272	2096586	2097152
22	792920	4077316	4192896	4194304
23	1241936	8111400	8385216	8388608
24	1913566	16124458	16769254	16777216
25	2898574	32034016	33536094	33554432
26	4323980	63579386	67067294	67108864
27	6353060	126076522	134124596	134217728
28	9206137	249736704	268228914	268435456
29	13158574	494124382	536416730	536870912
30	18576644	976302888	1072750464	1073741820

TABLE I

VALUES OF $D_k(n)$ FOR $3 \leq k \leq 6$ AND $k \leq n \leq 30$. VALUES HIGHLIGHTED IN BOLD CORRESPOND TO $D_k(S(k))$.

proof of Proposition 5 and consider the following two collections \mathcal{A} and \mathcal{B} of tuples. Set

$$\mathcal{A}^* \triangleq \{(A; s) : A \subset [n], \mathbf{X}_A = \boldsymbol{\alpha}, x_s = 1, s \notin A\}.$$

Since $\boldsymbol{\alpha}$ has weight $m - 1$ and \mathbf{X} has weight w , we have $(w - m + 1)$ choices for s for each occurrence of $\boldsymbol{\alpha}$ in \mathbf{X} . Therefore, the left hand side of (5) counts the number of tuples in \mathcal{A} .

Next, set

$$\mathcal{B}^* \triangleq \{(B; t) : B \subset [n], \mathbf{X}_B = \boldsymbol{\beta}, t \in B, \mathbf{X}_{B \setminus \{t\}} = \boldsymbol{\alpha}\}.$$

For each $\boldsymbol{\beta} \in J(k, m)$ and each occurrence of $\boldsymbol{\beta}$ in \mathbf{X} , we have $\boldsymbol{\beta}_\alpha$ choices for t . Therefore, the right hand side of (5) counts the number of tuples in \mathcal{B} .

To establish (5), it remains to exhibit a bijection between \mathcal{A}^* and \mathcal{B}^* . Consider the maps ϕ and ψ defined in the proof of Proposition 5. When we restrict the domains of ϕ and ψ to \mathcal{A}^* and \mathcal{B}^* , respectively, the maps are well-defined bijections from \mathcal{A}^* to \mathcal{B}^* and vice versa. Hence, we obtain (5). When $\boldsymbol{\alpha} \in J(k - 1, m)$, Equation (6) can be similarly established by requiring $x_s = 0$ in the definition of \mathcal{A}^* . \square

Define $\mathbf{H}^{(k, m)}$ to be the $\binom{k}{m} \times \binom{k}{m}$ matrix whose rows and columns are indexed by $J(k - 1, m - 1) \cup J(k - 1, m)$ and $J(k, m)$, respectively. The entries of $\mathbf{H}^{(k, m)}$ are given by $\mathbf{H}_{\boldsymbol{\alpha}, \boldsymbol{\beta}}^{(k, m)} \triangleq \boldsymbol{\beta}_\alpha$.

Further define a column vector \mathbf{Z} of length $\binom{k}{m}$ such that the first $\binom{k-1}{m-1}$ entries are given by $((w - m + 1)\mathbf{X}_\alpha)_{\alpha \in J(k-1, m-1)}$ and the next $\binom{k-1}{m}$ entries are given by $((n - k + m - w + 1)\mathbf{X}_\alpha)_{\alpha \in J(k-1, m)}$. Then (5)

and (6) imply that

$$\mathbf{H}^{(k, m)}(\mathbf{X}_\beta)_{\beta \in J(k, m)} = \mathbf{Z}. \quad (7)$$

Example 3. Let $k = 3$ and $m = 2$. Then

$$\mathbf{H}^{(3, 2)} = \begin{pmatrix} 2 & 1 & 0 \\ 0 & 1 & 2 \\ 1 & 1 & 1 \end{pmatrix}.$$

Consider $\mathbf{X} = 110011$ and so, $n = 6$, $w = 4$. Also,

$$\begin{aligned} (\mathbf{X}_\beta)_{\beta \in J(3, 2)}^T &= (2, 8, 2), \\ \mathbf{Z}^T &= (3 \cdot 4, 3 \cdot 4, 2 \cdot 6) = (12, 12, 12). \end{aligned}$$

We verify that (7) holds. \square

The following lemma then characterizes the (k, m) -equivalence of two words when they share the same $(k - 1)$ -deck.

Lemma 9. Let $\mathbf{H}^{(k, m)}$ be as defined above. If $\mathbf{X} \sim_{k-1} \mathbf{Y}$, then $(\mathbf{X}_\beta - \mathbf{Y}_\beta)_{\beta \in J(k, m)}$ belongs to the nullspace of $\mathbf{H}^{(k, m)}$.

Proof. Since $\mathbf{X} \sim_{k-1} \mathbf{Y}$, we have that $\mathbf{H}^{(k, m)}(\mathbf{X}_\beta)_{\beta \in J(k, m)} = \mathbf{H}^{(k, m)}(\mathbf{Y}_\beta)_{\beta \in J(k, m)}$. Hence, $\mathbf{H}^{(k, m)}(\mathbf{X}_\beta - \mathbf{Y}_\beta)_{\beta \in J(k, m)} = \mathbf{0}$ and the lemma follows. \square

Therefore, it remains to provide an upper bound on the nullity of $\mathbf{H}^{(k, m)}$.

Proposition 10. The nullity of $\mathbf{H}^{(k, m)}$ is at most $\binom{k-2}{m-1}$.

Proof. We write $\mathbf{H} = \mathbf{H}^{(k, m)}$ for short. Recall that the columns of \mathbf{H} are indexed by $J(k, m)$ and we arrange the columns in an increasing lexicographic order as in $J(k, m)$. We demonstrate that the nullity of \mathbf{H} is at most $\binom{k-2}{m-1}$ by exhibiting $\binom{k}{m} - \binom{k-2}{m-1}$ columns with leading coefficients.

We have the following cases.

- Let $\boldsymbol{\beta} = \beta_1 \beta_2 \cdots \beta_k \in J(k, m)$ with $\beta_1 = 0$. Then consider the row $\boldsymbol{\alpha} \triangleq \beta_2 \cdots \beta_k \in J(k - 1, m)$ and clearly, $\mathbf{H}_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \geq 1$. Suppose $\boldsymbol{\beta}' \in J(k, m)$ and $\boldsymbol{\beta}'_\alpha \geq 1$. Then $\boldsymbol{\beta}$ is necessarily lexicographically smaller than or equal to $\boldsymbol{\beta}'$. In other words, $\mathbf{H}_{\boldsymbol{\alpha}, \boldsymbol{\beta}''} = 0$ for all words $\boldsymbol{\beta}''$ that are lexicographically smaller than $\boldsymbol{\beta}$. Therefore, $\mathbf{H}_{\boldsymbol{\alpha}, \boldsymbol{\beta}}$ is the leading coefficient of row $\boldsymbol{\alpha}$.
- Let $\boldsymbol{\beta} = \beta_1 \beta_2 \cdots \beta_k \in J(k, m)$ with $\beta_k = 1$. Then consider the row $\boldsymbol{\alpha} \triangleq \beta_1 \beta_2 \cdots \beta_{k-1} \in J(k - 1, m - 1)$ and as before, $\mathbf{H}_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \geq 1$. Proceeding as before, we observe that $\mathbf{H}_{\boldsymbol{\alpha}, \boldsymbol{\beta}''} = 0$ for all words $\boldsymbol{\beta}'' \in J(k, m)$ that are lexicographically smaller than $\boldsymbol{\beta}$. Therefore, $\mathbf{H}_{\boldsymbol{\alpha}, \boldsymbol{\beta}}$ is the leading coefficient of row $\boldsymbol{\alpha}$.

Hence, the columns with possibly no leading coefficients start with a one and end with a zero. Therefore, there are $\binom{k-2}{m-1}$ such columns and the proposition follows. \square

Finally, we state the main theorem for this section and provide an upper bound on $D_k(n)$. Recall that $\mathbf{X} \sim_{(k, m)} \mathbf{Y}$ if and only if

$$(\mathbf{X}_\beta)_{\beta \in J(k, m)} = (\mathbf{Y}_\beta)_{\beta \in J(k, m)}.$$

Theorem 11. The number of (k, m) -equivalence classes for words of length n with the same $(k - 1)$ -deck is $O\left(n^{\binom{k-2}{m-1}}\right)$.

Therefore, the number of distinct k -decks with the same $(k - 1)$ -deck is $O\left(n^{k2^{k-2}}\right)$ and hence, $D_k(n) = O\left(n^{(k-1)2^{k-1}+1}\right)$.

Proof. Fix X to be of length n . Suppose that $Y \stackrel{k-1}{\sim} X$. Then Lemma 9 states that $(Y_\beta - X_\beta)_{\beta \in J(k,m)}$ belongs to the nullspace of $H^{(k,m)}$. Since the nullity of $H^{(k,m)}$ is at most $\binom{k-2}{m-1}$ and every entry of $(Y_\beta)_{\beta \in J(k,m)}$ is at most $\binom{n}{k} = O(n^k)$, the number of choices for $(Y_\beta)_{\beta \in J(k,m)}$ is $O\left(n^{k\binom{k-2}{m-1}}\right)$.

Therefore, the number of distinct k -decks with the same $(k-1)$ -deck is

$$O\left(\prod_{m=1}^{k-1} n^{k\binom{k-2}{m-1}}\right) = O\left(n^{k\left(\sum_{m=1}^{k-1} \binom{k-2}{m-1}\right)}\right) = O\left(n^{k2^{k-2}}\right).$$

Finally, it follows from simple induction that

$$D_k(n) = D_{k-1}(n) \cdot O\left(n^{k2^{k-2}}\right) = O\left(n^{(k-1)2^{k-1}+1}\right). \quad \square$$

V. LOWER BOUNDS ON $D_3(n)$

In this section, we focus on the case $k = 3$ and prove that $D_3(n) = \Omega(n^6)$. As with Section IV, we consider the words with the same $(k-1)$ -deck, or 2-deck, and determine the number of 3-equivalence classes amongst these words.

Let $X \in J(n, w)$. Following [12], we consider the notion of zero-vectors. The zero-vector of X , denoted by $u_X = (u_0, u_1, \dots, u_w)$, is the vector of length $w+1$, where u_0 is the number of zeroes in front of the first one, u_w is the number of zeroes after the last one, and u_j is the number of zeroes between the j th one and the $(j+1)$ th one for any $1 \leq j \leq w-1$. In other words, if $u_X = (u_0, u_1, \dots, u_w)$, then $X = 0^{u_0}10^{u_1}1 \dots 10^{u_w}$.

Recall that X_α is the number of occurrences of α as a subsequence of X . Set $X_{01} = t$ and $X_1 = w$. Our objective is to estimate the possible values of X_{011} . To this end, we have the following lemma from [12].

Lemma 12 ([12, Lemma 13]). *For $k \geq 1$, define the following $k \times (w+1)$ -integer-valued matrix:*

$$M_k \triangleq \begin{pmatrix} \binom{w-1}{k-1} & \binom{w-1}{k-1} & \binom{w-2}{k-1} & \dots & 1 & 0 & \dots & 0 \\ 0 & \binom{w-1}{k-2} & 2\binom{w-2}{k-2} & \dots & (w-k+1)\binom{k-1}{k-2} & (w-k+2) & \dots & 0 \\ 0 & 0 & \binom{w-2}{k-3} & \dots & \binom{w-k+1}{2}\binom{k-1}{k-3} & \binom{w-k+2}{2}\binom{k-2}{k-3} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \binom{w-k+1}{k-1} & \binom{w-k+2}{k-1} & \dots & \binom{w}{k-1} \end{pmatrix}.$$

Then $M_k u_X = (X_\alpha)_{\alpha \in J(k, k-1)}$.

Example 4. Let $X = 110011$ and $k = 3$. So, $u_X = (0, 0, 2, 0, 0)$ and

$$M_k = \begin{pmatrix} 6 & 3 & 1 & 0 & 0 \\ 0 & 3 & 4 & 3 & 0 \\ 0 & 0 & 1 & 3 & 6 \end{pmatrix}.$$

We verify that $(X_\alpha)_{\alpha \in J(3,2)}^T$ is indeed given by $M_3 u_X = (2, 8, 2)^T$.

Lemma 13. *Set Γ to be the following set of vectors of length $w+1$:*

$$\Gamma \triangleq \{(1, -2, 1, 0, 0, \dots, 0), (0, 1, -2, 1, 0, \dots, 0), (0, 0, 1, -2, 1, \dots, 0), \dots, (0, \dots, 1, -2, 1)\}.$$

If $u_{X'} = u_X + c$ for some $c \in \Gamma$, then $X_1 = X'_1$, $X_{01} = X'_{01}$ and $X'_{011} = X_{011} + 1$.

Proof. It follows from the definition of zero-vectors that $X_1 = X'_1$. For the other two equalities, we have that $M_k u_{X'} = M_k u_X + M_k c$. Applying Lemma 12 for $k \in \{2, 3\}$ and considering the first coordinate, we have

$$X'_{01} = X_{01} + M_2 c \text{ and } X'_{011} = X_{011} + M_3 c.$$

Then $M_2 c = 0$ because $a - 2(a+1) + (a+2) = 0$ for all a . On the other hand, $M_3 c = 1$ because $\binom{a}{2} - 2\binom{a+1}{2} + \binom{a+2}{2} = 1$ for all a . The lemma is then immediate. \square

Example 5. As before, let $X = 110011$ and $k = 3$. Further set $X' = 101101$ and so, $u_{X'} = (0, 1, 0, 1, 0)$ and $X' = X + (0, 1, -2, 1, 0)$. We verify that $X_1 = X'_1 = 4$, $X_{01} = X'_{01} = 4$, and $X'_{011} = X_{011} + 1 = 3$.

Next, we have the following proposition.

Proposition 14. *Fix n , $0 \leq w \leq n$ and $0 \leq t \leq w(n-w)$. Then there exists X and Y of length n such that the following hold:*

- (i) $X_1 = Y_1 = w$ and $X_{01} = Y_{01} = t$;
- (ii) $X_{011} \leq (n-w)\binom{q+1}{2}$ where $q = \lfloor t/(n-w) \rfloor$;
- (iii) $Y_{011} \geq \frac{w-1}{2} \left(t - \binom{w}{2}\right)$;
- (iv) for any $X_{011} \leq s \leq Y_{011}$, there exists Z such that $Z_1 = w$, $Z_{01} = t$ and $Z_{011} = s$.

Proof. Write that $t = q(n-w) + r$. Set X to be the word whose zero vector is (A_0, A_1, \dots, A_w) , where $A_{w-q-1} = r$, $A_{w-q} = n-w-r$ and all others are 0. Recall that $X_{01} = \sum_{i=0}^w u_i(w-i)$ and indeed, $X_{01} = r(q+1) + (n-w-r)q = t$. Furthermore, $X_{011} = r\binom{q+1}{2} + (n-w-r)\binom{q}{2} \leq (n-w)\binom{q+1}{2}$.

To construct Y , we iteratively add some c from Γ to u_X . Specifically, we start with $X^{(0)} = X$ and $u^{(0)} = u_X$ and suppose that we have $u^{(0)}, u^{(1)}, \dots, u^{(i)}$. If $u^{(i)}$ has a component with value at least two at index j with $1 \leq j \leq w-1$, we choose $c^{(i)} = (0, 0, \dots, 0, 1, -2, 1, 0, \dots, 0)$ with the minus two at index j . Then we set $u^{(i+1)} = u^{(i)} + c^{(i)}$ and $X^{(i+1)}$ to be the corresponding word. It follows from Lemma 13 that the two-deck of $X^{(i+1)}$ is the same as $X^{(i)}$ and the number of 011 in $X^{(i)}$ is given by $X_{011} + i$.

Hence, we terminate the process when the components of $u^{(i)}$ are at most one on the indices from 1 to $w-1$. Let (B_0, B_1, \dots, B_w) be the final zero-vector and Y be the corresponding word. It then remains to show Proposition 14(iii).

Again, we have $0 \leq B_i \leq 1$ for all $1 \leq i \leq w-1$. Furthermore, we have $Y_{01} = X_{01} = t$, which means

$$t = wB_0 + \sum_{i=1}^{w-1} (w-i)B_i \leq wB_0 + \sum_{i=1}^{w-1} w-i \leq wB_0 + \binom{w}{2}.$$

Therefore, $B_0 \geq \frac{t}{w} - \frac{w-1}{2}$ and so, $Y_{011} \geq B_0 \binom{w}{2} \geq \frac{w-1}{2} \left(t - \binom{w}{2}\right)$. \square

Following this proposition, for purposes of brevity, we write

$$s' \triangleq (n-w)\binom{q+1}{2},$$

$$s'' \triangleq \frac{w-1}{2} \left(t - \binom{w}{2}\right).$$

Therefore, a lower bound for $D_3(n)$ is

$$\sum_{w=0}^n \sum_{t=0}^{w(n-w)} \max\{0, s'' - s'\}. \quad (8)$$

We estimate the terms of (8). Note that

$$\begin{aligned} \sum_{w=0}^n \sum_{t=0}^{w(n-w)} s'' &= \sum_{w=0}^n \sum_{t=0}^{w(n-w)} \frac{w-1}{2} \left(t - \binom{w}{2} \right) \\ &= \sum_{w=0}^n \frac{w-1}{2} \sum_{t=1}^{w(n-w)} t - \binom{w}{2} \\ &\geq \sum_{w=0}^n \frac{w-1}{2} \left(\frac{1}{2} w^2 (n-w)^2 - \binom{w}{2} w(n-w) \right) \\ &\geq \sum_{w=0}^n \frac{w-1}{4} (w^2 (n-w)^2 - w^3 (n-w)) \\ &= \sum_{w=0}^n \frac{w-1}{4} (w^2 (n-w)(n-2w)) \\ &= \sum_{w=0}^n \frac{1}{4} w^3 (n-w)(n-2w) + O(n^4). \end{aligned} \quad (9)$$

On the other hand,

$$\begin{aligned} \sum_{w=0}^n \sum_{t=0}^{w(n-w)} s' &= \sum_{w=0}^n \sum_{t=0}^{w(n-w)} (n-w) \binom{q+1}{2} \\ &= \sum_{w=0}^n \sum_{t=0}^{w(n-w)} \frac{1}{2} (n-w) q^2 + O(n^2) \\ &= \sum_{w=0}^n \sum_{q=0}^w \frac{1}{2} (n-w)^2 q^2 + O(n^3) \\ &\quad \text{since for each } q, r \text{ can go from } 0 \text{ to } n-w-1 \\ &= \sum_{w=0}^n \frac{1}{6} (n-w)^2 w^3 + O(n^4). \end{aligned} \quad (10)$$

Combining (9) and (10) into (8), we have that the number of 3-decks is at least

$$\begin{aligned} \sum_{w=0}^{n/4} \frac{1}{4} w^3 (n-w)(n-2w) - \frac{1}{6} (n-w)^2 w^3 + O(n^4) \\ &= \sum_{w=0}^{n/4} w^3 (n-w) \left(\frac{1}{4} n - \frac{1}{2} w - \frac{1}{6} n + \frac{1}{6} w \right) + O(n^4) \\ &= \sum_{w=0}^{n/4} \frac{1}{12} w^3 (n-w)(n-4w) + O(n^4) = \Omega(n^6). \end{aligned} \quad (11)$$

We summarize our discussion in the following theorem.

Theorem 15. $D_3(n) = \Omega(n^6)$.

Remark 6. The statement in Theorem 15 can be made more precise. We have demonstrated that the number of $(3, 2)$ -equivalence classes amongst all words of length n is $\Omega(n^6)$. It then follows from Theorem 11 that this estimate is tight.

Remark 7. Implicit in the proof of Proposition 14 is an efficient method that encodes messages into words with distinct 3-decks. Specifically, let the message set be

$$\mathcal{M} \triangleq \{(w, t, s) : 0 \leq w \leq n, 0 \leq t \leq w(n-w), s' \leq s \leq s''\}.$$

Given any triple $(w, t, s) \in \mathcal{M}$, we can construct X in linear time such that $X_1 = w$, $X_{01} = t$ and $X_{011} = s'$.

Following the procedure described in the proof of Proposition 14, we choose a sequence of $c^{(0)}, c^{(1)}, \dots \in \Gamma$ to add to u_X . Since $s' \leq s \leq s''$, there is a sequence such that the resulting zero-vector corresponds to Y and $Y_{011} = s$. Therefore, Y is the codeword encoding the message (w, t, s) and Y can be computed in $O(n^3)$ time.

VI. CONCLUSION

We provide an efficient trellis-based method to compute the number of distinct k -decks and determined the exact value of $D_k(n)$ for $k \in \{3, 4, 5, 6\}$ and $k \leq n \leq 30$. An interesting consequence is that we established the fact that $S(6) = 30$.

We also establish an asymptotic upper bound on $D_k(n)$ for general k and an asymptotic lower bound for $D_3(n)$. In summary, we have $D_3(n) = O(n^9)$ and $D_3(n) = \Omega(n^6)$. It remains open to determine tight bound on the asymptotic growth rate of $D_3(n)$.

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, 1340-1371, 2015.
- [2] Z. Chang, J. Chrisnata, M. F. Ezerman and H. M. Kiah, "Rates of DNA sequence profiles for practical values of read lengths," *IEEE Transactions on Information Theory*, vol. 63, no. 11, pp. 7166-7177, Nov. 2017.
- [3] C. Choffrut, J. Karhumäki, *Combinatorics of words*, in: G. Rozenberg, A. Salomaa (Eds.), *Handbook of Formal Languages*, vol. I, Springer, Berlin, 1997, pp. 329-438.
- [4] M. Dudik and L. J. Schulman, "Reconstruction from subsequences," *Journal of Combinatorial Theory*, vol. 103, no. 2, pp. 337-348, 2003.
- [5] R. Gabrys and O. Milenkovic, "The hybrid k -deck problem: Reconstructing sequences from short and long traces," in 2017 IEEE International Symposium on Information Theory (ISIT), Jun. 2017, pp. 1306-1310.
- [6] R. Gabrys and O. Milenkovic, "Unique reconstruction of coded sequences from multiset substring spectra," in 2018 IEEE International Symposium on Information Theory (ISIT), Jun. 2018, pp. 2540-2544.
- [7] Kalashnik, L. O. "The reconstruction of a word from fragments," *Numerical Mathematics and Computer Technology*, Akad. Nauk. Ukrain. SSR Inst. Mat., Preprint IV (1973): 56-57.
- [8] H. M. Kiah, G. J. Puleo, and O. Milenkovic, "Codes for DNA sequence profiles," *IEEE Transactions on Information Theory*, vol. 62, no. 6, pp. 3125-3146, 2016.
- [9] I. Krasikov and Y. Roditty, "On a reconstruction problem for sequences," *Journal of Combinatorial Theory*, vol. 77, no. 2, pp. 344-348, 1997.
- [10] P. Ligeti and P. Sziklai, "Reconstruction from subwords," in 6th International Conference on Applied Informatics, Jan. 2004, pp. 1-7.
- [11] J. L. Massey, "Foundation and methods of channel encoding," in *Proc. Int. Conf. Information Theory and Systems*, vol. 65, NTG-Fachberichte, 1978.
- [12] B. Manvel, A. Meyerowitz, A. Schwenk, K. Smith, and P. Stockmeyer, "Reconstruction of sequences," *Discrete Math*, vol. 94, no. 3, pp. 209-219, 1991.
- [13] OEIS Foundation Inc. (2019), *The On-Line Encyclopedia of Integer Sequences*, <http://oeis.org/A258585>
- [14] M. Rigo, P. Salimov, "Another generalization of abelian equivalence: Binomial complexity of infinite words," *Theoretical Computer Science*, vol. 601, pp. 47-57, 2015.
- [15] A.D. Scott, "Reconstructing sequences," *Discrete Mathematics*, vol. 175, no. 1-3, pp. 231-238, 1997.
- [16] S. M. H. T. Yazdi, H. M. Kiah, E. Garcia-Ruiz, J. Ma, H. Zhao, and O. Milenkovic, "DNA-based storage: Trends and methods," *IEEE Transactions on Molecular, Biological and Multi-Scale Communications*, vol. 1, no. 3, pp. 230-248, Sep. 2015.