

Permutation Recovery against Deletion Errors for DNA Data Storage

Charchit Gupta[†], Shubhransh Singhvi^{*}, Avital Boruchovsky[‡],
Yuval Goldberg[‡], Han Mao Kiah[§] and Eitan Yaakobi[‡]

^{*}Signal Processing & Communications Research Center, IIIT Hyderabad, India

[†]IIIT Hyderabad, India

[‡]Department of Computer Science, Technion—Israel Institute of Technology, Haifa 3200003, Israel

[§]School of Physical and Mathematical Sciences, Nanyang Technological University, Singapore

Abstract

Owing to its immense storage density and durability, DNA has emerged as a promising storage medium. However, due to technological constraints, data can only be written onto many short DNA molecules called *data blocks* that are stored in an unordered way. To handle the unordered nature of DNA data storage systems, a unique *address* is typically prepended to each data block to form a *DNA strand*. However, DNA storage systems are prone to errors and generate multiple noisy copies of each strand called *DNA reads*. Thus, we study the problem of *permutation recovery* for DNA data storage.

The permutation recovery problem for DNA data storage requires one to reconstruct the addresses or in other words to uniquely identify the noisy reads. By successfully reconstructing the addresses, one can essentially determine the correct order of the data blocks, effectively solving the clustering problem.

We first show that we can almost surely identify all the noisy reads under certain mild assumptions. We then propose a permutation recovery procedure and demonstrate that on average the procedure uses only a fraction of \mathcal{M}^2 data comparisons (when \mathcal{M} is the number of reads).

I. INTRODUCTION

The need for a more durable and compact storage system has become increasingly evident with the explosion of data in modern times. While magnetic and optical disks have been the primary solutions for storing large amounts of data, they still face limitations in terms of storage density and physical space requirements. Storing a zettabyte of data using these traditional technologies would necessitate a vast number of units and considerable physical space.

The idea of using macromolecules for ultra-dense storage systems was recognized as early as the 1960s when physicist Richard Feynman outlined his vision for nanotechnology in his talk ‘There is plenty of room at the bottom’. Using DNA is an attractive possibility because it is extremely dense (up to about 1 exabyte per cubic millimeter) and durable (half-life of over 500 years). Since the first experiments conducted by Church et al. in 2012 [?] and Goldman et al. in 2013 [?], there have been a flurry of experimental demonstrations (see [?], [?] for a survey). Amongst the various coding design considerations, in this work, we study the unsorted nature of the DNA storage system [?], [?].

A DNA storage system consists of three important components. The first is the DNA synthesis which produces the oligonucleotides, also called *strands*, that encode the data. The second part is a storage container with compartments which stores the DNA strands, however without order. Finally, to retrieve the data, the DNA is accessed using next-generation sequencing, which results in several noisy copies, called *reads*. The processes of synthesizing, storing, sequencing, and handling strands are all error prone. Due to this unordered nature of DNA-based storage systems, when the user retrieves the information, in addition to decoding the data, the user has to determine the identity of the data stored in each strand. A typical solution is to simply have a set of addresses and store this address information as a prefix to each DNA strand. As the addresses are also known to the user, the user can identify the information after the decoding process. As these addresses along with the stored data are prone to errors, this solution needs further refinements.

In [?], the strands (strand = address + data) are first clustered with respect to the edit distance. Then the authors determine a consensus output amongst the strands in each cluster and finally, decode these consensus outputs using a **classic** concatenation scheme. For this approach, the clustering step is computationally expensive. When there are \mathcal{M} reads, the usual clustering method involves \mathcal{M}^2 pairwise comparisons to compute distances. This is costly when the data strands are long, and the problem is further exacerbated if the metric is the edit distance. Therefore, in [?], a distributed approximate clustering algorithm was proposed and the authors clustered 5 billion strands in 46 minutes on 24 processors.

In [?], the authors proposed and investigated an approach that avoids clustering, by studying a generalisation of the *bee identification problem*. Informally, the bee identification problem requires the receiver to identify M ‘bees’ using a set of M unordered noisy measurements [?]. Later, in [?], the authors generalized the setup to multi-draw channels where every bee (addresses) results in N noisy outputs (noisy addresses). The task then is to identify each of the M bees from the MN noisy outputs and it turns out that this task can be reduced to a minimum-cost network flow problem. In contrast to previous

works, the approach in [?] utilizes only the address information, which is of significantly shorter length, and the method does not take into account the associated noisy data. Hence, this approach involves no data comparisons.

However, as evident, the clustering and bee identification based approaches do not completely take into account the nature of the DNA storage system. In particular, the clustering approaches do not utilise the uncorrupted set of addresses which can be accessed by the receiver and the bee identification approach uses solely the information stored in address and neglects the noisy data strands.

In this work, we overcome these shortcomings by utilising both the address and data information to identify the noisy reads. Specifically for the binary deletion channel, we first show that we can almost surely correctly identify all the reads under certain mild assumptions. Then we propose our permutation recovery procedure and demonstrate that on average the procedure uses only a fraction of \mathcal{M}^2 data comparisons (when there are \mathcal{M} reads).

II. PROBLEM FORMULATION

Let N and M be positive integers. Let $[M]$ denote the set $\{1, 2, \dots, M\}$. An N -permutation ψ over $[M]$ is an NM -tuple $(\psi(j))_{j \in [MN]}$ where every symbol in $[M]$ appears exactly N times, and we denote the set of all N -permutations over $[M]$ by $\mathbb{S}_N(M)$. Let the set of addresses be denoted by $\mathcal{A} \subseteq \{0, 1\}^n$ and $M \triangleq |\mathcal{A}|$. We will use the terms addresses and codewords interchangeably. We assume that every codeword $\mathbf{x}_i \in \mathcal{A}$ is attached to a length- L data part $\mathbf{d}_i \in \{0, 1\}^L$ to form a strand, which is the tuple, $(\mathbf{x}_i, \mathbf{d}_i)$. Let the multiset of data be denoted by $D = \{\{\mathbf{d}_i : i \in [M]\}\}$ and the set of strands by $R = \{(\mathbf{x}_i, \mathbf{d}_i) : i \in [M]\}$. Throughout this paper, we assume that D is drawn uniformly at random over $\{0, 1\}^L$. Let $\mathcal{S}_N((\mathbf{x}, \mathbf{d}))$ denote the multiset of channel outputs when (\mathbf{x}, \mathbf{d}) is transmitted N times through the channel \mathcal{S} . Assume that the entire set R is transmitted through the channel \mathcal{S} , hence an unordered multiset, $R' = \{(\mathbf{x}'_1, \mathbf{d}'_1), (\mathbf{x}'_2, \mathbf{d}'_2), \dots, (\mathbf{x}'_{MN}, \mathbf{d}'_{MN})\}$, of MN noisy strands (reads) is obtained, where for every $j \in [MN]$, $(\mathbf{x}'_j, \mathbf{d}'_j) \in \mathcal{S}_N((\mathbf{x}_{\pi(j)}, \mathbf{d}_{\pi(j)}))$ for some N -permutation π over $[M]$, which will be referred to as the *true N -permutation*. Note that the receiver, apart from the set of reads R' , has access to the set of addresses \mathcal{A} but does not know the set of data D .

Let $\mathbf{x}, \mathbf{y} \in \{0, 1\}^*$. We denote the shortest common supersequence between \mathbf{x}, \mathbf{y} by $\text{SCS}(\mathbf{x}, \mathbf{y})$ and the longest common subsequence by $\text{LCS}(\mathbf{x}, \mathbf{y})$.

In [], the authors introduced the following problems:

Problem 1. Given \mathcal{S}, ϵ and \mathcal{A} , find the region $\mathcal{R} \in \mathbb{Z}_+^2$, such that for $(N, L) \in \mathcal{R}$, it is possible to identify the true permutation with probability at least $1 - \epsilon$ when the data D is drawn uniformly at random.

For $(N, L) \in \mathcal{R}$, we can find the true permutation by making at least $(NM)^2$ data comparisons. This may be expensive when the data parts are long, i.e., when L is large. Therefore, our second objective is to reduce the number of data comparisons.

Problem 2. Let $\kappa < 1$. Given $\mathcal{S}, \epsilon, \mathcal{A}$ and $(N, L) \in \mathcal{R}$, design an algorithm to identify the true permutation with probability at least $1 - \epsilon$ using $\kappa(NM)^2$ data comparisons. As before, D is drawn uniformly at random.

In Section ??, we demonstrate that the algorithm in [] identifies the true permutation with a vanishing probability as n grows. In Section ??, we address Problem 1 and identify the region \mathcal{R} for which there exists only one valid permutation, viz. the true permutation. In Section ??, we describe our algorithm that identifies the true-permutation with probability at least $1 - \epsilon$ when $(N, L) \in \mathcal{R}$. In Section ??, we analyse the expected number of data comparisons performed by the algorithm.

III. BEE-IDENTIFICATION OVER MULTI-DRAW DELETION CHANNEL

The algorithm in [] uses solely the information stored in the addresses to identify the true-permutation, and does not take into consideration the noisy data that is also available to the receiver. The first step in their algorithm is to construct a bipartite graph $\mathcal{G} = (\mathcal{X} \cup \mathcal{Y}, E)$.

- (i) Nodes. The left nodes are the addresses ($\mathcal{X} = \mathcal{A}$) and the right nodes are the noisy reads ($\mathcal{Y} = R'$).
- (ii) Demands. For each left node / codeword $\mathbf{x} \in \mathcal{X}$, we assign the demand $\delta(\mathbf{x}) \triangleq -N$, while for each right node / output $(\mathbf{y}, \mathbf{d}') \in \mathcal{Y}$, we assign the demand $\delta(\mathbf{y}) \triangleq 1$.
- (iii) Edges. There exists an edge between $\mathbf{x} \in \mathcal{X}$ and $(\mathbf{y}, \mathbf{d}') \in \mathcal{Y}$ if and only if $P(\mathbf{y}|\mathbf{x}) > 0$, where $P(\mathbf{y}|\mathbf{x})$ is the likelihood probability of observing \mathbf{y} given that \mathbf{x} was transmitted. For $\mathbf{x} \in \mathcal{X}$ and $(\mathbf{y}, \mathbf{d}') \in \mathcal{Y}$, let $E_{\mathbf{x}}$ and $E_{(\mathbf{y}, \mathbf{d}')}$ denote the multiset of neighbours of \mathbf{x} and the set of neighbours of $(\mathbf{y}, \mathbf{d}')$ in \mathcal{G} , respectively, i.e., $E_{\mathbf{x}} = \{(\mathbf{y}, \mathbf{d}') | (\mathbf{x}, (\mathbf{y}, \mathbf{d}')) \in E\}$, $E_{(\mathbf{y}, \mathbf{d}')} = \{\mathbf{x} | (\mathbf{x}, (\mathbf{y}, \mathbf{d}')) \in E\}$. Note that the degree of every left node is at least N as $\mathcal{S}_N((\mathbf{x}, \mathbf{d})) \subseteq E_{\mathbf{x}}$.
- (iv) Costs. For an edge $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$, we assign the cost $\gamma(\mathbf{x}, \mathbf{y}) = -\log P(\mathbf{y} | \mathbf{x})$. Note that the cost is well-defined as the value $P(\mathbf{y} | \mathbf{x})$ is necessarily positive.

Let us now show that the above algorithm (Chrisnata et al.[]) identifies the true permutation with a vanishing probability as n grows, when $M = 2^n$.

Let $a \in \mathcal{Y}$ and $x \in \mathcal{X}$ and let $|a| = n - k, k \in [0, n]$. Consider the cost of an edge between a and x

$$\text{cost}(x, a) = -\log(P(a \mid x)),$$

For a deletion channel with deletion probability p , we have that

$$P(a \mid x) = \text{Emb}(x, a) \cdot p^k (1 - p)^{n-k}.$$

A sufficient condition for the algorithm to fail is the existence of $x, y \in \{0, 1\}^n$, $a \in \{0, 1\}^k$, $b \in \{0, 1\}^{k'}$ where a, b are noisy reads of x, y , respectively and $0 \leq k, k' \leq n$. Now consider the following event:

$$\begin{aligned} \text{cost}(y, a) &< \text{cost}(x, a) \wedge \text{cost}(x, b) < \text{cost}(y, b) \\ P(a \mid y) &> P(a \mid x) \wedge P(b \mid x) > P(b \mid y) \end{aligned}$$

If the above event is true then the minimum-cost flow algorithm will necessarily assign a, b wrongfully. However, note that it might not assign b to x and a to y . Since,

$$\begin{aligned} P(a \mid x) &= \text{Emb}(x, a) \cdot p^k (1 - p)^{n-k} \\ P(a \mid y) &= \text{Emb}(y, a) \cdot p^k (1 - p)^{n-k} \end{aligned}$$

Therefore, $P(a \mid y) > P(a \mid x)$ if and only if $\text{Emb}(y, a) > \text{Emb}(x, a)$.

We now prove that there exists a set of pairs $x, y \in \{0, 1\}^n$ for a large enough $n \in \mathbb{N}$ that have noisy reads a, b (respectively) such that

$$\text{Emb}(y, a) > \text{Emb}(x, a) \wedge \text{Emb}(x, b) > \text{Emb}(y, b).$$

Let us denote $x_{m,n} = 0^m 10^{n-m-3} 11$ and $y_{m,n} = 0^m 10^{n-m-4} 100$ for each $m, n \in \mathbb{N}$ such that $m \leq n - 4$. Also, denote for each $0 \leq k \leq n - 3$: $a_k = 0^k, b_k = 0^{k-1} 1$. For each m, n, k , it holds that

$$\begin{aligned} \text{Emb}(x_{m,n}, a_k) &= \binom{n-3}{k} < \binom{n-2}{k} = \text{Emb}(y_{m,n}, a_k) \\ \text{Emb}(x_{m,n}, b_k) &= \begin{cases} 2 \binom{n-3}{k-1} & k-1 > m \\ 2 \binom{n-3}{k-1} + 3 \binom{m}{k-1} & k-1 \leq m \end{cases} \\ \text{Emb}(y_{m,n}, b_k) &= \begin{cases} \binom{n-4}{k-1} & k-1 > m \\ \binom{n-4}{k-1} + 2 \binom{m}{k-1} & k-1 \leq m \end{cases} \end{aligned}$$

Therefore,

$$\text{Emb}(x_{m,n}, b_k) > \text{Emb}(y_{m,n}, b_k).$$

Therefore, when a_k is a noisy version of $x_{m,n}$ and $b_{k'}$ is a noisy version of $y_{m,n}$ (respectively, for some fixed m, n), the min cost algorithm wrongfully assigns a_k and b_k .

Let P_0 denote the probability that the min cost algorithm fails to identify the true permutation. As shown above, P_0 is lower bounded by the probability that the noisy versions of $x_{m,n}$ and $y_{m,n}$ ($m \leq n - 4$) are $a_k, b_{k'}$ ($k, k' \leq n - 3$). It can be verified that

$$\begin{aligned} P(a_k \mid x_{m,n}) &= p^3, \\ P(b_{k'} \mid y_{m,n}) &\geq p^3(1 - p). \end{aligned}$$

Therefore, we have that

$$\begin{aligned} P_0 &\geq 1 - \prod_{m=0}^{n-4} (1 - P(a_k \mid x_{m,n}) P(b_{k'} \mid y_{m,n})) \\ &\geq 1 - (1 - p^6(1 - p))^{n-3} \\ &\xrightarrow{n \rightarrow \infty} 1. \end{aligned}$$

IV. UNIQUENESS OF N -PERMUTATION

The task of identifying the true permutation π , can be split into two steps. We can first identify the *partitioning* $\{\mathcal{S}_N((x_i, d_i)) : i \in [M]\}$ and then for each *partition* $(\mathcal{S}_N((x_i, d_i)))$ identify the *label*, viz. the channel input (x_i) , where $i \in [M]$. Hence, given R' and \mathcal{A} , we are able to find the true permutation if and only if there exists only one valid partitioning and one valid labelling.

In this section, we study Problem 1 when $\mathcal{S} = \text{BDL}(p)$. Specifically, in Lemmas ?? and ??, we determine the values L_{Th} and N_{Th} , respectively, such that for all $L \geq L_{\text{Th}}$ and $N \geq N_{\text{Th}}$, we are able to find the true permutation with high probability. The result is formally stated in Theorem ??.

Before formally defining partitioning and labelling, we introduce some notations. For $\mathbf{a}, \mathbf{b} \in \{0, 1\}^n$, let \mathbf{a}', \mathbf{b}' be the channel outputs through $\text{BDL}(p)$ of \mathbf{a}, \mathbf{b} , respectively. We say that \mathbf{a}' and \mathbf{b}' are *congruent*, denoted by $\mathbf{a}' \cong \mathbf{b}'$, if $|\text{SCS}(\mathbf{a}', \mathbf{b}')| \leq n$. Further, by abuse of notation, we would denote the event of all sequences in $A \subseteq \mathcal{S}_N(\mathbf{a})$ being congruent with all sequences in $B \subseteq \mathcal{S}_N(\mathbf{b})$ by $A \cong B$, where $\mathcal{S}_N(x_i)$ denotes the multiset of channel outputs when x_i is transmitted N times through the channel \mathcal{S} , $i \in \{1, 2\}$.

A read $(\mathbf{y}, \mathbf{d}') \in R'$ is said to be **faulty** if there exists some other read $(\tilde{\mathbf{y}}, \tilde{\mathbf{d}}') \in R' \setminus \{(\mathbf{y}, \mathbf{d}')\}$ with $(\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((x, \mathbf{d}))$ and $(\tilde{\mathbf{y}}, \tilde{\mathbf{d}}') \in \mathcal{S}_N((\tilde{x}, \tilde{\mathbf{d}}))$ such that $(\mathbf{y}, \mathbf{d}') \cong (\tilde{\mathbf{y}}, \tilde{\mathbf{d}}')$. Let R_{faulty} denote the multiset of such faulty reads. In the next lemma, we give a bound on the length of the longest common subsequence of two uniformly chosen binary strings, which we would be crucial in analyzing the probability of a read being faulty.

Lemma 1. [] Let $X_{k,\ell}$ denote the length of the longest common subsequence of two uniformly chosen binary strings of length k and ℓ , respectively and $\lambda > 0$, then

$$\mathbb{E}[X_{k,\ell}] = \gamma_2 \max\{k, \ell\}, \quad (1)$$

where γ_2 [bounds –Shubhransh]. Further,

$$P(|X_{k,\ell} - \mathbb{E}(X_{k,\ell})| \geq \lambda) \leq 2e^{-\frac{\lambda^2}{8 \max\{k, \ell\}}}. \quad (2)$$

Let A^* denote the event that all noisy reads have lengths between $[(1-p)L - c_L\sqrt{L}, (1+p)L + c_L\sqrt{L}]$.

Lemma 2. The probability of the event A^* is at least

$$P(A^*) \geq \left(1 - 2e^{(-2c_L^2)}\right)^{N2^n}.$$

Proof. For $c_L > 0$, Let $d'_1 \in \{0, 1\}^*$, then it can be verified that $\mathbb{E}(|d'_1|) = (1-p)L$. Further, using Hoeffding's inequality, we have that

$$\begin{aligned} P\left(|d'_1| - (1-p)L \geq c_L\sqrt{L}\right) &\leq 2e^{\left(-\frac{2(c_L\sqrt{L})^2}{L}\right)}, \\ P\left(|d'_1| - (1-p)L \leq -c_L\sqrt{L}\right) &\leq 2e^{(-2c_L^2)}. \end{aligned}$$

Therefore, the probability of the event A^* is at least

$$P(A^*) \geq \left(1 - 2e^{(-2c_L^2)}\right)^{N2^n}.$$

□

In the next lemma, we calculate the probability of a read being faulty.

Lemma 3. For $(\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((x, \mathbf{d}))$,

$$P(\mathbb{I}_{(\mathbf{y}, \mathbf{d}') \in R_{\text{faulty}}} | A^*) =$$

Proof. Let $|d'| \triangleq \ell \in [(1-p)L - c_L\sqrt{L}, (1+p)L + c_L\sqrt{L}]$. Consider another noisy strand d'_1 obtained from $d_1 \in \{0, 1\}^L$ and let $|d'_1| \triangleq \ell_1 \in [(1-p)L - c_L\sqrt{L}, (1+p)L + c_L\sqrt{L}]$. From Lemma??, we have that

$$\mathbb{E}(|\text{LCS}(d', d'_1)| \mid |d'| = \ell, |d'_1| = \ell_1) = \gamma_2 \max\{\ell, \ell_1\},$$

where $\gamma_2 \approx 0.8$. Further, it is known that $|\text{SCS}(d'_1, d'_2)| = |d'_1| + |d'_2| - |\text{LCS}(d'_1, d'_2)|$. Therefore,

$$\mathbb{E}(|\text{SCS}(d', d'_1)| \mid |d'| = \ell, |d'_1| = \ell_1) = \ell + \ell_1 - \gamma_2 \max\{\ell, \ell_1\}.$$

Given A^* and the lengths of the noisy strands, we have that

$$\gamma_2 \left((1-p)L + c_L\sqrt{L} \right) \geq \mathbb{E}(|\text{LCS}(d', d'_1)| \mid |d'| = \ell, |d'_1| = \ell_1, A^*), \quad (3)$$

and

$$|\text{SCS}(d', d'_1)| \geq 2 \left((1-p)L - c_L \sqrt{L} \right) - |\text{LCS}(d'_1, d'_2)|.$$

We now analyze the probability that $|\text{SCS}(d', d'_1)| > L$ given A^* and the lengths of the noisy strands.

$$\begin{aligned} & P\left(|\text{SCS}(d', d'_1)| > L \mid |d'| = \ell, |d'_1| = \ell_1, A^*\right) \\ & \leq P\left(2 \left((1-p)L - c_L \sqrt{L} \right) - |\text{LCS}(d', d'_1)| > L \mid |d'| = \ell, |d'_1| = \ell_1, A^*\right) \\ & = P\left(\mathbb{E}[|\text{LCS}(d', d'_1)|] - |\text{LCS}(d', d'_1)| > \mathbb{E}[|\text{LCS}(d', d'_1)|] - 2 \left((1-p)L - c_L \sqrt{L} \right) + L \mid |d'| = \ell, |d'_1| = \ell_1, A^*\right) \\ & \leq P\left(\mathbb{E}[|\text{LCS}(d', d'_1)|] - |\text{LCS}(d', d'_1)| > L - 2 \left((1-p)L - c_L \sqrt{L} \right) - \gamma_2 \left((1-p)L + c_L \sqrt{L} \right) \mid |d'| = \ell, |d'_1| = \ell_1, A^*\right). \end{aligned}$$

For $c_L > \max\left(0, \frac{(2 + \gamma_2)(1-p) - 1}{2 - \gamma_2}\right) \sqrt{L}$, from Lemma ??, it follows that

$$\begin{aligned} P\left(|\text{SCS}(d', d'_1)| > L \mid |d'| = \ell, |d'_1| = \ell_1, A^*\right) & \leq 2e^{-\frac{\left(L - 2 \left((1-p)L - c_L \sqrt{L} \right) - \gamma_2 \left((1-p)L + c_L \sqrt{L} \right)\right)^2}{8 \max\{\ell, \ell_1\}}} \\ & \leq 2e^{-\frac{\left(L(1 - (2 + \gamma_2)(1-p)) + c_L \sqrt{L}(2 - \gamma_2)\right)^2}{8 \gamma_2 \left((1-p)L + c_L \sqrt{L} \right)}} \end{aligned}$$

□

Definition 1. A *partitioning* $\mathcal{P} = \{P_1, P_2, \dots, P_M\}$ of \mathcal{Y} is defined as the collection of disjoint submultisets of \mathcal{Y} , each of size N , such that for $i \in [M]$, for $(j, k) \in \binom{[N]}{2}$, $|\text{SCS}(\mathbf{y}_j, \mathbf{y}_k)| \leq n$ & $|\text{SCS}(\mathbf{d}'_j, \mathbf{d}'_k)| \leq L$, where $(\mathbf{y}_j, \mathbf{d}'_j), (\mathbf{y}_k, \mathbf{d}'_k) \in P_i$.

V. PERMUTATION RECOVERY ALGORITHM

Algorithm 1 Permutation Recovery Algorithm

```

1: procedure PRUNE( $\mathcal{G}, (\tilde{\mathbf{y}}, \tilde{\mathbf{d}}')$ )
2:    $(\tilde{\mathbf{y}}, \tilde{\mathbf{d}}') \rightarrow \text{Pruned}, \mathcal{T} = \{\}$ 
3:   for  $(\mathbf{y}, \mathbf{d}') \in \mathcal{N}_{(\tilde{\mathbf{y}}, \tilde{\mathbf{d}}')}$  do
4:     if  $(\mathbf{y}, \mathbf{d}') \cong (\tilde{\mathbf{y}}, \tilde{\mathbf{d}}')$  then
5:        $(\mathbf{y}, \mathbf{d}') \rightarrow \mathcal{T}$ 
6:   if  $|\mathcal{T}| = N - 1$  then
7:     Let  $\mathcal{X}^* = \bigcap_{(\mathbf{y}, \mathbf{d}') \in \mathcal{T}} E_{(\mathbf{y}, \mathbf{d}')}$ 
8:     for  $(\mathbf{y}, \mathbf{d}') \in \mathcal{T}$  do
9:       Remove  $\{(\mathbf{x}, (\mathbf{y}, \mathbf{d}')) : \mathbf{x} \notin \mathcal{X}^*\}$  from  $E$ 
10:     $(\mathbf{y}, \mathbf{d}') \rightarrow \text{Pruned}$ 

11: procedure PRUNING ALGORITHM( $\mathcal{P}_{\mathcal{G}}, \mathcal{G}$ )
12:   Pruned =  $\{\}$ 
13:   while  $|\text{Pruned}| < N2^n$  do
14:      $(\tilde{\mathbf{y}}, \tilde{\mathbf{d}}') = \arg \min\{|\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}| : (\mathbf{y}, \mathbf{d}') \in \mathcal{Y}\}$ 
15:     PRUNE( $\mathcal{G}, (\tilde{\mathbf{y}}, \tilde{\mathbf{d}}')$ )
16:   return PMA( $\mathcal{P}_{\mathcal{G}}, \mathcal{G}$ )

```

Lemma 4. Let $(\mathbf{x}, \mathbf{d}) \in R, (\mathbf{y}, \mathbf{d}') \in \mathcal{S}((\mathbf{x}, \mathbf{d}))$. Then, we have that

$$\mathbb{E}[\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}] = NW_{\mathcal{A}, \mathbf{x}}(2p - p^2) - 1$$

Proof. Let $(\tilde{x}, \tilde{d}) \in R$, such that $d_L(x, \tilde{x}) = r$, i.e. $\text{LCS}(x, \tilde{x}) = n - r$. Let $(\tilde{y}, \tilde{d}') \in \mathcal{S}_N((\tilde{x}, \tilde{d}))$. Then for $(y, d') \in \mathcal{N}_{(\tilde{y}, \tilde{d})}$, there must be at least one erasure in x or \tilde{x} at all the positions where they differ, which happens with probability $(1 - (1 - p)^2)^r = (2p - p^2)^r$. Hence, $\mathbb{E}[\mathcal{N}_{(y, d')}] = \sum_{r=1}^n W_{r,x} (2p - p^2)^r N + (N - 1) = NW_{\mathcal{A},x} (2p - p^2) - 1$. \square

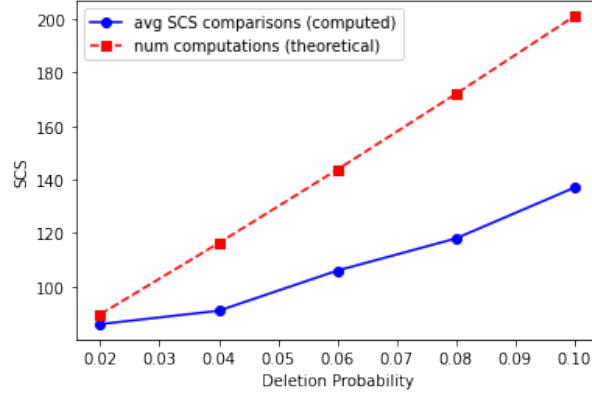


Fig. 1: Computed vs Theoretical SCS comparisons.

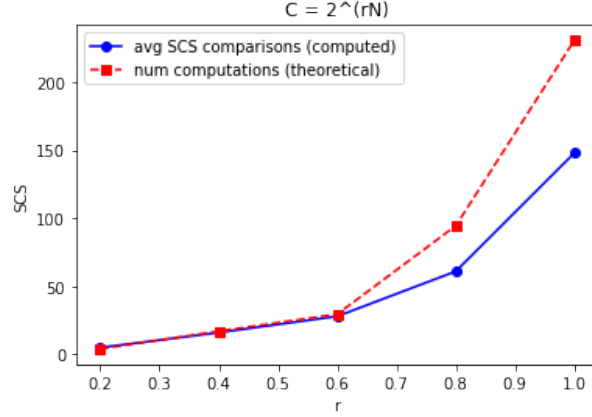


Fig. 2: Computed vs Theoretical SCS comparisons.

VI. REFERENCES

- [1] S. Singhvi, A. Boruchovsky, H. M. Kiah and E. Yaakobi, Data-Driven Bee Identification for DNA Strands, 2023.
- [2] J. Chrisnata, H. M. Kiah, A. Vardy and E. Yaakobi, "Bee Identification Problem for DNA Strands," IEEE Journal on Selected Areas in Information Theory, vol. 4, pp. 190-204, 2023.
- [3] C. Rashtchian, K. Makarychev, M. Racz, S. Ang, D. Jevdjic, S. Yekhanin, L. Ceze and K. Strauss, "Clustering Billions of Reads for DNA Data Storage," in Advances in Neural Information Processing Systems, 2017.
- [4] G. S. Lueker, "Improved bounds on the average length of longest common subsequences," Journal of the ACM (JACM), vol. 56, p. 1-38, 2009.