

Permutation Recovery for DNA Data Storage

Shubhansh Singhvi*, Charchit Gupta[†], Avital Boruchovsky[‡],
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). Specifically, ...

Index Terms

DNA Data Storage, Permutation Recovery, Clustering, Bee Identification.

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 [4] and Goldman et al. in 2013 [5], there have been a flurry of experimental demonstrations (see [11], [13] for a survey). Amongst the various coding design considerations, in this work, we study the unsorted nature of the DNA storage system [8], [11].

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 [9], 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 [10], a distributed approximate clustering algorithm was proposed and the authors clustered 5 billion strands in 46 minutes on 24 processors.

In [3], 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 [12]. Later, in [3], 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 [3] 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 erasure 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 M^2 data comparisons (when there are M reads).

II. RELATED WORK

A brief history of the Bee Identification Problem: Consider M bees flying in a beehive, each tagged with a unique barcode. We take a snapshot of the bees and obtain an unordered set of noisy barcodes. The *bee-identification problem* - proposed and formally defined by Tandon et al. [1] - requires one to uniquely identify each bee from the multiset of the noisy barcodes. The naive way is to separately look at each barcode and decode them *independently*. However, due to noise, certain bees may be assigned to the same barcode, and in this case, we fail to identify all the bees. In contrast, one can *jointly* look at all the barcodes and determine the best way to assign the barcodes so that the accuracy/likelihood of correct identification is maximized. Tandon et al. studied error exponents and showed that decoding the barcodes jointly results in a significantly smaller error exponent compared to decoding the barcodes independently. The authors further posited that joint decoding entails a computationally prohibitive exhaustive search. However, to the contrary, Kiah et al. [2] demonstrated that efficient joint decoding is achievable. Specifically, for the binary erasure channel (BEC) and binary symmetric channel (BSC), they reduced the bee-identification problem to the problem of finding a perfect matching and minimum-cost matching, respectively. Hence, applying the well-known Hopcraft-Karp algorithm [3] and Hungarian method [4], respectively, one can identify the bees in time polynomial in M . They further reduce the running times by using classical tools like peeling decoders and list-decoders, and also show that the identifier algorithms when used with Reed-Muller codes terminates in almost linear and quadratic time for BEC and BSC, respectively. Chrisnata et al. [5], proposed and analysed a generalization of the Bee Identification problem to the multi-draw channels. For multi-draw deletion channel, similar to [2], they reduce the identification problem to the problem of finding a minimum-cost matching. Then, applying the Edmonds-Karp or Tomizawa algorithm [6], they show that the bee identification problem for multi-draw deletion channels can be solved in $\mathcal{O}(M^3)$ time.

Clustering Problem in DNA Data Storage: ...

III. 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 denoted by \mathcal{A} be a (n, k, d_{\min}) -binary linear code such that $M = 2^k$, and 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 the weight enumerator polynomial of \mathcal{A} be denoted by $W_{\mathcal{A}}(z) = \sum_{i=0}^n w_i z^i$, where w_i denotes the number of codewords in \mathcal{A} of hamming weight i .

In this work, we first consider the following problem.

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 IV, we first propose an extension of the Peeling Matching Algorithm [6] to the multi-draw erasure channel. We then demonstrate that the peeling matching algorithm identifies the true permutation with a vanishing probability as n grows. In Section V, we address Problem 1 and identify the region \mathcal{R} for which there exists only one valid permutation, viz. the true permutation. In Section VI, we describe our algorithm that identifies the true-permutation with probability at least $1 - \epsilon$ when $(N, L) \in \mathcal{R}$. In Section VII, we analyse the expected number of data comparisons performed by the algorithm.

IV. BEE-IDENTIFICATION OVER MULTI-DRAW ERASURE CHANNEL

In this section, we extend the Peeling Matching Algorithm (PMA), presented in [6] for $N = 1$ to a general value of N . The PMA-based approach 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 the peeling matching algorithm is to construct a bipartite undirected graph $\mathcal{G} = (\mathcal{X} \cup \mathcal{Y}, E)$, where the left nodes are the addresses ($\mathcal{X} = \mathcal{A}$) and the right nodes are the noisy reads ($\mathcal{Y} = R'$). There exists an edge between $x \in \mathcal{X}$ and $(y, d') \in \mathcal{Y}$ if and only if $P(y|x) > 0$, where $P(y|x)$ is the likelihood probability of observing y given that x was transmitted. For $x \in \mathcal{X}$ and $(y, d') \in \mathcal{Y}$, let E_x and $E_{(y, d')}$ denote the multiset of neighbours of x and the set of neighbours of (y, d') in \mathcal{G} , respectively, i.e., $E_x = \{(y, d') | (x, (y, d')) \in E\}$, $E_{(y, d')} = \{x | (x, (y, d')) \in E\}$. Note that the degree of every left node is at least N as $\mathcal{S}_N((x, d)) \subseteq E_x$. For left nodes and right nodes with degrees N and 1 respectively, the corresponding neighbor(s) can be matched with certainty. For ease of exposition, we refer to such nodes as *good* nodes.

Definition 1. A node $(y, d') \in \mathcal{Y}$ is said to be a *good right node* if $|E_{(y, d')}| = 1$. A node $x \in \mathcal{X}$ is said to be a *Type-A good left node* if $|E_x| = N$ or a *Type-B good left node* if $|\{(y, d') \in E_x : |E_{(y, d')}| = 1\}| = N$.

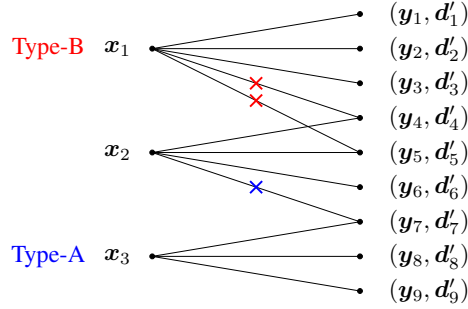


Fig. 1: Let $N = 3$. If x_1 is peeled, then x_2 becomes a Type-B good left node, and if x_3 is peeled then x_2 becomes a Type-A good left node.

Let $\mathcal{Y}_g, \mathcal{X}_{g_A}$ and \mathcal{X}_{g_B} denote the set of good right nodes, Type-A good left nodes and Type-B good left nodes, respectively. The peeling matching algorithm when executed over \mathcal{G} , finds good left nodes and identifies the corresponding N channel outputs until there are no good left nodes. Let $\mathcal{P}_G = (\mathcal{X} \cup \mathcal{Y}, \mathcal{P}_E)$ denote the bipartite matching identified by the Peeling Matching Algorithm.

Algorithm 1 Peeling Matching Algorithm

```

1: procedure PEEL( $\mathcal{P}_G, \mathcal{G}, x$ )
2:   if  $x \in \mathcal{X}_{g_A}$  then
3:     for  $(y, d') \in E_x$  do
4:       Remove  $\{(\tilde{x}, (y, d')) : \tilde{x} \in E_{(y, d')}\}$  from  $E$ 
5:       Add  $(x, (y, d'))$  to  $\mathcal{P}_E$ 
6:       Remove  $(y, d')$  from  $\mathcal{Y}$ 
7:     Remove  $x$  from  $\mathcal{X}$ 
8:   else if  $x \in \mathcal{X}_{g_B}$  then
9:     Add  $\{(x, (y, d')) : (y, d') \in E_x \cap \mathcal{Y}_g\}$  to  $\mathcal{P}_E$ 
10:    Remove  $\{(x, (y, d')) : (y, d') \in E_x\}$  from  $E$  and remove  $E_x \cap \mathcal{Y}_g$  from  $\mathcal{Y}$ 
11:    Remove  $x$  from  $\mathcal{X}$ 
12: procedure PMA( $\mathcal{P}_G, \mathcal{G}$ )
13: for  $x \in \mathcal{X}_{g_A} \cup \mathcal{X}_{g_B}$  do

```

```

14:   PEEL( $\mathcal{P}_G, \mathcal{G}, x$ )
15:   if  $|\mathcal{P}_E| = N2^n$  then
16:     return  $\mathcal{P}_G$ 
17:   else
18:     return FAILURE

```

Note that as shown in Fig. 1, peeling Type-A and Type-B good left nodes might generate new Type-A and Type-B good left nodes, respectively. Thus, at any instant during the course of the algorithm, we assume X_{g_A}, X_{g_B} to reflect the Type-A and Type-B good left nodes, respectively, at that instant.

Proposition 1. [7] *Algorithm 1 finds the true permutation if only if there are no cycles in \mathcal{G} .*

Let the multiset of right nodes that are in a cycle be denoted by Y_{cycle} . In the next lemma, we derive a lower bound on the probability of observing at least one cycle in \mathcal{G} .

Lemma 1. *For $\mathcal{A} = \{0, 1\}^n$, the probability of observing at least one cycle in \mathcal{G} is lower bounded by*

$$P(|Y_{\text{cycle}}| > 1) > 1 - \frac{\mathcal{U}_{\text{cycle}}}{N2^n(1 - \mathcal{U}_{\text{cycle}})},$$

where $\mathcal{U}_{\text{cycle}} \triangleq 2^{-N((1+p^2)^n - 1)}$.

Proof. Let Y_{cycle}^* denote the multiset of right nodes that are in a cycle of size four. Let $(y, d') \in \mathcal{S}_N((x, d))$ and $(\tilde{y}, \tilde{d}') \in \mathcal{S}_N((\tilde{x}, \tilde{d}))$ such that $d_H(x, \tilde{x}) = r > 0$. It can be verified that $\{x, (y, d'), \tilde{x}, (\tilde{y}, \tilde{d}')\}$ forms a cycle with probability p^{2r} . Therefore, the probability that (y, d') is not in a cycle is

$$\begin{aligned} P((y, d') \notin Y_{\text{cycle}}^*) &= \prod_{r=1}^n (1 - p^{2r})^{N \binom{n}{r}} \\ &= 2^{\log\left(\prod_{r=1}^n (1 - p^{2r})^{N \binom{n}{r}}\right)} \\ &= 2^{N \sum_{r=1}^n \left(\binom{n}{r} \log(1 - p^{2r})\right)}. \end{aligned}$$

Using Jensen's inequality, it can be verified that

$$\begin{aligned} P((y, d') \notin Y_{\text{cycle}}^*) &\leq 2^{N(2^n - 1) \log\left(\frac{\sum_{r=1}^n \binom{n}{r} (1 - p^{2r})}{2^n - 1}\right)} \\ &= 2^{N(2^n - 1) \log\left(1 - \frac{(1 + p^2)^n - 1}{2^n - 1}\right)}. \end{aligned}$$

Using Taylor series expansion, for $|a| < 1$, $\log(1 + a) \leq a$. Therefore, we have that

$$P((y, d') \notin Y_{\text{cycle}}^*) \leq \mathcal{U}_{\text{cycle}} \triangleq 2^{-N((1 + p^2)^n - 1)}.$$

From linearity of expectation, we have that

$$\begin{aligned} \mathbb{E}[|Y_{\text{cycle}}^*|] &= \sum_{(y, d') \in \mathcal{Y}} \mathbb{E}[\mathbb{I}_{\{(y, d') \in Y_{\text{cycle}}^*\}}] \\ &= N2^n \left(1 - \prod_{r=1}^n (1 - p^{2r})^{N \binom{n}{r}}\right). \end{aligned}$$

Further, from the linearity of variances of indicator random variables, we have that

$$\begin{aligned} \text{Var}[|Y_{\text{cycle}}^*|] &= \sum_{(y, d') \in \mathcal{Y}} \text{Var}[\mathbb{I}_{\{(y, d') \in Y_{\text{cycle}}^*\}}] \\ &= N2^n \left(1 - \prod_{r=1}^n (1 - p^{2r})^{N \binom{n}{r}}\right) \left(\prod_{r=1}^n (1 - p^{2r})^{N \binom{n}{r}}\right). \end{aligned}$$

Next, from Chebyshev's inequality and the upper bound on $P((y, d') \notin Y_{\text{cycle}})$, it can be verified that

$$P(|Y_{\text{cycle}}^*| < 1) < \frac{\text{Var}[|Y_{\text{cycle}}^*|]}{(\mathbb{E}[|Y_{\text{cycle}}^*|])^2} < \frac{\mathcal{U}_{\text{cycle}}}{N2^n(1 - \mathcal{U}_{\text{cycle}})}.$$

Therefore, the results follows. \square

Hence, from Proposition 1 and Lemma 1, it is highly improbable (vanishingly low probability) to find the true permutation using only the addresses and its noisy measurements (i.e., \mathbf{x} and \mathbf{y} 's). In the next section, we see that by making use of the data parts, we can find the true permutation under certain mild assumptions.

V. UNIQUENESS OF THE N -PERMUTATION

The task of identifying the true permutation π , can be split into two steps. We can first identify the *partitioning* $\{\mathcal{S}_N((\mathbf{x}_i, \mathbf{d}_i)) : i \in [M]\}$ and then for each *partition* $(\mathcal{S}_N((\mathbf{x}_i, \mathbf{d}_i)))$ identify the *label*, viz. the channel input (\mathbf{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 $\mathcal{S} = \text{BEC}(p)$. Specifically, in Lemmas 3 and 4, 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 1.

Before formally defining partitioning and labelling, we introduce some notations. Let $\mathbf{x}_1, \mathbf{x}_2 \in \{0, 1\}^\ell$. For $i \in \{1, 2\}$, let $\mathbf{b}_i \in \{0, 1, *\}^\ell$ be the output of \mathbf{x}_i through $\text{BEC}(p)$. We denote the event of \mathbf{b}_1 and \mathbf{b}_2 agreeing at the non-erased positions by $\mathbf{b}_1 \cong \mathbf{b}_2$. For example, let $\mathbf{x}_1 = 00000$, $\mathbf{x}_2 = 00011$ and let $\mathbf{b}_1 = 0000*$, $\mathbf{b}_2 = 000*1$ then $\mathbf{b}_1 \cong \mathbf{b}_2$. Further, by abuse of notation, we would denote the event of all sequences in $A \subseteq \mathcal{S}_N(\mathbf{x}_1)$ agreeing at the non-erased positions with all sequences in $B \subseteq \mathcal{S}_N(\mathbf{x}_2)$ by $A \cong B$, where $\mathcal{S}_N(\mathbf{x}_i)$ denotes the multiset of channel outputs when \mathbf{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 $(\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((\mathbf{x}, \mathbf{d}))$ and $(\tilde{\mathbf{y}}, \tilde{\mathbf{d}}') \in \mathcal{S}_N((\tilde{\mathbf{x}}, \tilde{\mathbf{d}}))$ 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 calculate the probability of a read being faulty.

Lemma 2. For $(\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((\mathbf{x}, \mathbf{d}))$,

$$P(\mathbb{I}_{(\mathbf{y}, \mathbf{d}') \in R_{\text{faulty}}}) = 1 - \prod_{r=1}^n \left(1 - (2p - p^2)^r \left(1 - \frac{1}{2}(1 - p)^2 \right)^L \right)^{Nw_r}.$$

Furthermore,

$$P(\mathbb{I}_{(\mathbf{y}, \mathbf{d}') \in R_{\text{faulty}}}) \leq N \left(1 - \frac{1}{2}(1 - p)^2 \right)^L (W_{\mathcal{A}}(2p - p^2) - 1).$$

Proof. For $(\mathbf{x}, \mathbf{d}), (\tilde{\mathbf{x}}, \tilde{\mathbf{d}}) \in R$, let $(\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((\mathbf{x}, \mathbf{d}))$ and $(\tilde{\mathbf{y}}, \tilde{\mathbf{d}}') \in \mathcal{S}_N((\tilde{\mathbf{x}}, \tilde{\mathbf{d}}))$ with $d_H(\mathbf{x}, \tilde{\mathbf{x}}) = r > 0$. For $\mathbf{y} \cong \tilde{\mathbf{y}}$, the positions where \mathbf{x} and $\tilde{\mathbf{x}}$ differ must be erased in at least one of them, which happens with probability $(1 - (1 - p)^2)^r = (2p - p^2)^r$. For index $i \in [L]$, the probability that both \mathbf{d}' and $\tilde{\mathbf{d}}'$ are not erased and disagree on i is $\frac{1}{2}(1 - p)^2$. Therefore, $P(\mathbf{d}' \cong \tilde{\mathbf{d}}') = (1 - \frac{1}{2}(1 - p)^2)^L$. Hence, $P((\mathbf{y}, \mathbf{d}') \cong (\tilde{\mathbf{y}}, \tilde{\mathbf{d}}')) = (2p - p^2)^r (1 - \frac{1}{2}(1 - p)^2)^L$. Thus,

$$\begin{aligned} P(\mathbb{I}_{(\mathbf{y}, \mathbf{d}') \in R_{\text{faulty}}}) &= 1 - \prod_{(\tilde{\mathbf{y}}, \tilde{\mathbf{d}}') \in R' / \mathcal{S}_N((\mathbf{x}, \mathbf{d}'))} \left(1 - P((\mathbf{y}, \mathbf{d}') \cong (\tilde{\mathbf{y}}, \tilde{\mathbf{d}}')) \right) \\ &= 1 - \prod_{r=1}^n \left(1 - (2p - p^2)^r \left(1 - \frac{1}{2}(1 - p)^2 \right)^L \right)^{Nw_r}. \end{aligned}$$

Using Weierstrass Inequality, we have that

$$\begin{aligned} P(\mathbb{I}_{(\mathbf{y}, \mathbf{d}') \in R_{\text{faulty}}}) &\leq \sum_{r=1}^n Nw_r (2p - p^2)^r \left(1 - \frac{1}{2}(1 - p)^2 \right)^L \\ &= N \left(1 - \frac{1}{2}(1 - p)^2 \right)^L \left(\sum_{r=1}^n w_r (2p - p^2)^r \right) \\ &= N \left(1 - \frac{1}{2}(1 - p)^2 \right)^L (W_{\mathcal{A}}(2p - p^2) - 1). \end{aligned}$$

\square

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

We will refer to $\mathcal{P}^* \triangleq \{\mathcal{S}_N((\mathbf{x}_i, \mathbf{d}_i)) : i \in [M]\}$ as the *true partitioning* of R' . Let $\mathbb{P}_{R'}$ denote the set of all possible partitionings of R' . Note that if $|\mathbb{P}_{R'}| = 1$ then $\mathbb{P}_{R'} = \{\mathcal{P}^*\}$. Let $\mathcal{G}' = (\mathcal{Y}, E')$, where $\mathcal{Y} = R'$. For $(\mathbf{y}, \mathbf{d}'), (\tilde{\mathbf{y}}, \tilde{\mathbf{d}}') \in \mathcal{Y}$, $((\mathbf{y}, \mathbf{d}'), (\tilde{\mathbf{y}}, \tilde{\mathbf{d}}')) \in E'$ if $(\mathbf{y}, \mathbf{d}') \equiv (\tilde{\mathbf{y}}, \tilde{\mathbf{d}}')$. Note that a partitioning $\mathcal{P} \in \mathbb{P}_{R'}$ corresponds to partitioning the graph \mathcal{G}' into M cliques each of size N .

Proposition 2. $|\mathbb{P}_{R'}| = 1$ if and only if there exists a unique partitioning of the graph \mathcal{G}' into M cliques each of size N .

In the next lemma, we derive a threshold on L such that for $L \geq L_{\text{Th}}$, $\mathbb{P}_{R'} = \{\mathcal{P}^*\}$ with probability at least $1 - \epsilon_1$.

Lemma 3. For $L \geq L_{\text{Th}} \triangleq \log_{(1-\frac{1}{2}(1-p)^2)} \left(\frac{\sqrt[N]{\epsilon_1/2^k}}{N(W_{\mathcal{A}}(2p-p^2)-1)} \right)$, we have that $\mathbb{P}_{R'} = \{\mathcal{P}^*\}$ with probability at least $1 - \epsilon_1$.

Proof. Note that for every $\mathbf{x} \in \mathcal{A}$, if there exists at least one $(\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((\mathbf{x}, \mathbf{d}))$ such that $(\mathbf{y}, \mathbf{d}')$ is not faulty, then the only valid partitioning is \mathcal{P}^* . Let $\mathbf{X}_{\text{faulty}}$ denote the set of left nodes with $\mathcal{S}_N((\mathbf{x}, \mathbf{d})) \subset \mathbf{R}_{\text{faulty}}$. From Markov Inequality,

$$\begin{aligned} P(\mathbf{x} \in \mathbf{X}_{\text{faulty}}) &= P(\mathbb{I}_{\mathcal{S}_N((\mathbf{x}, \mathbf{d})) \subset \mathbf{R}_{\text{faulty}}} \geq 1) \\ &\leq \mathbb{E} [\mathbb{I}_{\mathcal{S}_N((\mathbf{x}, \mathbf{d})) \subset \mathbf{R}_{\text{faulty}}}] = (P(\mathbb{I}_{(\mathbf{y}, \mathbf{d}') \subset \mathbf{R}_{\text{faulty}}} = 1))^N. \end{aligned}$$

Therefore, from Lemma 2, $P(\mathbf{x} \in \mathbf{X}_{\text{faulty}})$ is at most

$$\left(N \left(1 - \frac{1}{2}(1-p)^2 \right)^L (W_{\mathcal{A}}(2p-p^2)-1) \right)^N.$$

From linearity of expectation, $\mathbb{E} [|\mathbf{X}_{\text{faulty}}|]$ is at most

$$2^k \left(N \left(1 - \frac{1}{2}(1-p)^2 \right)^L (W_{\mathcal{A}}(2p-p^2)-1) \right)^N.$$

From Markov inequality, $P(|\mathbf{X}_{\text{faulty}}| \geq 1) \leq \mathbb{E} [|\mathbf{X}_{\text{faulty}}|]$. Hence, $P(|\mathbf{X}_{\text{faulty}}| < 1)$ is at least

$$= 1 - 2^k \left(N \left(1 - \frac{1}{2}(1-p)^2 \right)^L (W_{\mathcal{A}}(2p-p^2)-1) \right)^N.$$

Lastly, it can be verified that $P(|\mathbf{X}_{\text{faulty}}| < 1) \geq 1 - \epsilon_1$ if $L > \log_{(1-\frac{1}{2}(1-p)^2)} \left(\frac{\sqrt[N]{\epsilon_1/2^k}}{N(W_{\mathcal{A}}(2p-p^2)-1)} \right)$. □

The next corollary follows from the fact that $W_{\{0,1\}^n}(z) = \sum_{i=0}^n \binom{n}{i} z^i = (1+z)^n$.

Corollary 1. For $\mathcal{A} = \{0, 1\}^n$, $L_{\text{Th}} = \log_{(1-\frac{1}{2}(1-p)^2)} \left(\frac{\sqrt[N]{\epsilon_1/2^n}}{N((1+2p-p^2)^n-1)} \right)$.

Corollary 2. For given n, ϵ_1, p , $\mathcal{A} = \{0, 1\}^n$, L_{Th} is minimum at $N = \ln \left(\frac{2^n}{\epsilon_1} \right)$.

Proof. From Corollary 1, observe that

$$\arg \min_{N \in \mathbb{Z}_+} L_{\text{Th}} = \arg \min_{N \in \mathbb{Z}_+} \log_2 \left(\frac{N((1+2p-p^2)^n-1)}{\sqrt[N]{\epsilon_1/2^n}} \right),$$

which can be readily verified to be minimum at $N = \ln \left(\frac{2^n}{\epsilon_1} \right)$. □

Corollary 3. For given N, ϵ_1, p , $\mathcal{A} = \{0, 1\}^n$, $L_{\text{Th}} = \mathcal{O}(n)$. Specifically,

$$L_{\text{Th}} < \left\lceil \frac{\log_2(N(1+2p-p^2)) + \frac{1-\log_2(\epsilon_1)}{N}}{1 - \log_2(1+2p-p^2)} \right\rceil n.$$

Proof. From Corollary 1, $L_{\text{Th}} = \frac{\log_2 \left(\frac{N((1+2p-p^2)^n-1)}{\sqrt[N]{\epsilon_1/2^n}} \right)}{(1-\log_2(1+2p-p^2))}$. Therefore, we have that

$$\frac{L_{\text{Th}}}{n} < \frac{\log_2 \left(\frac{N(1+2p-p^2)^n}{\sqrt[N]{\epsilon_1/2^n}} \right)}{n(1 - \log_2(1+2p-p^2))}$$

$$\begin{aligned}
&= \frac{\log_2(N) + n \log(1 + 2p - p^2) - \left(\frac{\log(\epsilon_1) - n}{N}\right)}{n(1 - \log_2(1 + 2p - p^2))} \\
&= \frac{\log_2(1 + 2p - p^2) + \frac{1}{N}}{(1 - \log_2(1 + 2p - p^2))} + \frac{\log_2(N) - \frac{\log_2(\epsilon_1)}{N}}{n(1 - \log_2(1 + 2p - p^2))} \\
&< \frac{\log_2(N(1 + 2p - p^2)) + \frac{1 - \log_2(\epsilon_1)}{N}}{1 - \log_2(1 + 2p - p^2)}.
\end{aligned}$$

□

Definition 3. Given a partitioning $\mathcal{P} = \{P_1, P_2, \dots, P_M\}$, we define a **labelling**, denoted by \mathcal{L} , as a length- M vector of distinct addresses from \mathcal{A} such that $\mathcal{L}[i] \in \{\mathbf{x} : \forall (\mathbf{y}, \mathbf{d}') \in P_i, P(\mathbf{x}|\mathbf{y}) > 0\}$, where $\mathcal{L}[i]$ denotes the i -th element of \mathcal{L} , and $i \in [M]$.

We denote the set of all possible labellings for a given partitioning \mathcal{P} by $\mathbb{L}_{\mathcal{P}, R'}$. Given the true partitioning \mathcal{P}^* , we define the *true labelling*, denoted by \mathcal{L}^* , as the labelling in which for each partition $\mathcal{S}_N((\mathbf{x}_i, \mathbf{d}_i))$, the assigned label is \mathbf{x}_i , where $i \in [M]$. Note that if $\mathcal{P} \neq \mathcal{P}^*$ then $\mathcal{L}^* \notin \mathbb{L}_{\mathcal{P}, R'}$. Further, if $|\mathbb{L}_{\mathcal{P}^*, R'}| = 1$ then $\mathbb{L}_{\mathcal{P}^*, R'} = \{\mathcal{L}^*\}$. Let $\mathcal{G}'' = (\mathcal{X}, E'')$, where $\mathcal{X} = \mathcal{A}$. There is a directed edge $\mathbf{x} \rightarrow \tilde{\mathbf{x}}$ if all of the N channel outputs of \mathbf{x} are erased at the positions where \mathbf{x} and $\tilde{\mathbf{x}}$ differ, i.e., $\{\tilde{\mathbf{x}}\} \in \{\bigcap_{(\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((\mathbf{x}, \mathbf{d}))} E(\mathbf{y}, \mathbf{d}')\}$.

Proposition 3. $|\mathbb{L}_{\mathcal{P}^*, R'}| = 1$ if and only if there are no directed cycles in \mathcal{G}'' .

In the next lemma, we derive a threshold on N such that for $N \geq N_{\text{Th}}$, $\mathbb{L}_{\mathcal{P}^*, R'} = \{\mathcal{L}^*\}$ with probability at least $1 - \epsilon_2$.

Lemma 4. For $N \geq N_{\text{Th}} \triangleq \log_p(W_{\mathcal{A}}^{-1}(1 + \frac{\epsilon_2}{2^k}))$, we have that $\mathbb{L}_{\mathcal{P}^*, R'} = \{\mathcal{L}^*\}$ with probability at least $1 - \epsilon_2$.

Proof. Let X_{faulty} denote the set of nodes in \mathcal{G}'' that have at least one outgoing edge. For $\mathbf{x}, \tilde{\mathbf{x}} \in \mathcal{X}$, let $r = d_H(\mathbf{x}, \tilde{\mathbf{x}})$. Note that $\mathbf{x} \rightarrow \tilde{\mathbf{x}}$ happens with probability p^{rN} . Therefore, the probability that \mathbf{x} has no outgoing edges is $\prod_{r=1}^n (1 - p^{rN})^{w_r}$. Hence, from linearity of expectation,

$$\begin{aligned}
\mathbb{E}[|X_{\text{faulty}}|] &= \sum_{\mathbf{x} \in \mathcal{X}} \mathbb{E}[\mathbb{I}_{\{\mathbf{x} \in X_{\text{faulty}}\}}] \\
&= 2^k \left(1 - \prod_{r=1}^n (1 - p^{rN})^{w_r}\right).
\end{aligned}$$

From Weierstrass inequality, we have that $\prod_{r=1}^n (1 - p^{rN})^{w_r} \geq 1 - \sum_{r=1}^n w_r p^{rN} = 2 - W_{\mathcal{A}}(p^N)$. Therefore,

$$\mathbb{E}[|X_{\text{faulty}}|] \leq 2^k (W_{\mathcal{A}}(p^N) - 1).$$

From Markov inequality, $P(|X_{\text{faulty}}| \geq 1) \leq \mathbb{E}[|X_{\text{faulty}}|]$. Hence,

$$P(|X_{\text{faulty}}| < 1) \geq 1 - 2^k (W_{\mathcal{A}}(p^N) - 1).$$

Lastly, it can be verified that $P(|X_{\text{faulty}}| < 1) \geq 1 - \epsilon_2$ if $N \geq \log_p(W_{\mathcal{A}}^{-1}(1 + \frac{\epsilon_2}{2^k}))$. □

The next corollary follows from the fact that $W_{\{0,1\}^n}(z) = \sum_{i=0}^n \binom{n}{i} z^i = (1+z)^n$.

Corollary 4. For $\mathcal{A} = \{0, 1\}^n$, $N_{\text{Th}} \triangleq \log_p(\sqrt[n]{1 + \frac{\epsilon_2}{2^n}} - 1)$.

Corollary 5. For given p, ϵ_2 , $\mathcal{A} = \{0, 1\}^n$, $N_{\text{Th}} = \Theta(n)$. Specifically,

$$\frac{n + \log_2\left(\frac{n}{\epsilon_2 \ln(2)}\right)}{\log_2\left(\frac{1}{p}\right)} > N_{\text{Th}} > \frac{n + \log_2\left(\frac{n}{\epsilon_2}\right)}{\log_2\left(\frac{1}{p}\right)}.$$

Proof. From Lemma 4, $N_{\text{Th}} = \frac{\log_2(\sqrt[n]{1 + \frac{\epsilon_2}{2^n}} - 1)}{\log_2(p)}$. For $0 < a < 1, 0 < b < 1$, it can be verified that $1 + ab \ln(2) < 2^{ab} < (1+a)^b < 1 + ab$. Therefore, we have that

$$\begin{aligned}
&\frac{\log_2\left(\frac{\epsilon_2 \ln(2)}{n 2^n}\right)}{\log_2(p)} > N_{\text{Th}} > \frac{\log_2\left(\frac{\epsilon_2}{n 2^n}\right)}{\log_2(p)} \\
&\iff \frac{n + \log_2\left(\frac{n}{\epsilon_2 \ln(2)}\right)}{\log_2\left(\frac{1}{p}\right)} > N_{\text{Th}} > \frac{n + \log_2\left(\frac{n}{\epsilon_2}\right)}{\log_2\left(\frac{1}{p}\right)}.
\end{aligned}$$

□

Thus, we define the region \mathcal{R} as $\mathcal{R} \triangleq \{(\beta, N) : \beta \geq \beta_{\text{Th}}, N \geq N_{\text{Th}}\}$. In the next theorem, we give a sufficient condition for the existence of a unique N -permutation.

Theorem 1. For $(L, N) \in \mathcal{R}$, it is possible to identify the true permutation with probability at least $1 - \epsilon$, when $\epsilon_1, \epsilon_2 < \frac{\epsilon}{2}$.

Proof. From Lemma 3 and 4, it follows that for $L > L_{\text{Th}}$ and $N > N_{\text{Th}}$, $\mathbb{P}_{R'} = \{\mathcal{P}^*\}$ with probability $(1 - \epsilon_1)$ and $\mathbb{P}_{\mathcal{P}^*, R'} = \{\mathcal{L}^*\}$ with probability $(1 - \epsilon_2)$, respectively. Hence, for $L > L_{\text{Th}}$ and $N > N_{\text{Th}}$, there exists only one valid permutation with probability $(1 - \frac{\epsilon}{2})^2 > (1 - \epsilon)$. □

From Corollaries 3 and 5, for $\mathcal{A} = \{0, 1\}^n$, we observe that $L_{\text{Th}} < \lambda^* n$ and $N_{\text{Th}} < \nu^* n$ for some constants λ^* and ν^* . This means that we only require data parts to be of length $L = \lambda^* n$ and the number of reads to be $N = \nu^* n$ so that correct identification occurs with high probability. In the next section, we design an algorithm to find the true permutation with a small number of data comparisons.

VI. ACHIEVABILITY: PERMUTATION RECOVERY ALGORITHM

As the receiver has access to the set of addresses, we design an algorithm that reduces the number of data comparisons by comparing a pair of reads if and only if they agree at the positions that are not erased in the address part. Hence, similar to the peeling matching algorithm, we first build the bipartite graph $\mathcal{G} = (\mathcal{X} \cup \mathcal{Y}, E)$ as described in Section IV. Let $\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}$ denote the two-hop neighborhood of $(\mathbf{y}, \mathbf{d}')$ in \mathcal{G} .

Proposition 4. For $(\mathbf{y}, \mathbf{d}'), (\tilde{\mathbf{y}}, \tilde{\mathbf{d}}') \in \mathcal{Y}$, $(\mathbf{y}, \mathbf{d}') \in \mathcal{N}_{(\tilde{\mathbf{y}}, \tilde{\mathbf{d}}')}$ if and only if $\mathbf{y} \cong \tilde{\mathbf{y}}$.

In the next lemma, we calculate the expected value of $|\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}|$.

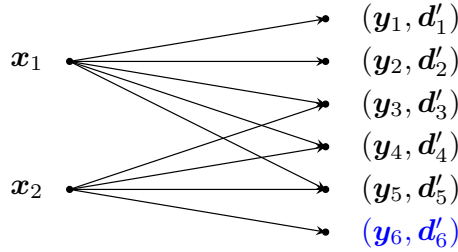


Fig. 2: Let $N = 3$. For $(\mathbf{y}_6, \mathbf{d}'_6)$, we can potentially identify the remaining 2 copies by performing only $|\mathcal{N}_{(\mathbf{y}_6, \mathbf{d}'_6)}| = 3$ data comparisons.

Lemma 5. Let $(\mathbf{x}, \mathbf{d}) \in R$, $(\mathbf{y}, \mathbf{d}') \in \mathcal{S}((\mathbf{x}, \mathbf{d}))$. Then, we have that

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

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

Lemma 6. For $\mathcal{A} = \{0, 1\}^n$ and $(\mathbf{y}, \mathbf{d}') \in \mathcal{Y}$,

$$\mathbb{E}[|\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}] \mid (\mathbf{y}, \mathbf{d}') = N2^r(1 + p)^{n-r} - 1,$$

where r denotes the number of erasures in \mathbf{y} . Further, $\mathbb{E}[|\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}] = N(1 + 2p - p^2)^n - 1$.

Proof. Let the number of erasures in \mathbf{y} be r . Let $(\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((\mathbf{x}, \mathbf{d}))$.

- 1) For $(\tilde{\mathbf{x}}, \tilde{\mathbf{d}}) \in \mathcal{X}$ that does not differ from (\mathbf{x}, \mathbf{d}) at the non-erased positions in $(\mathbf{y}, \mathbf{d}')$, $\mathcal{S}_N((\tilde{\mathbf{x}}, \tilde{\mathbf{d}})) \subset \mathcal{N}_{(\mathbf{y}, \mathbf{d}')}$. Note that there are $2^r - 1$ such $(\tilde{\mathbf{x}}, \tilde{\mathbf{d}})$.
- 2) For $(\tilde{\mathbf{x}}, \tilde{\mathbf{d}}) \in \mathcal{X}$ that differ from (\mathbf{x}, \mathbf{d}) at i out of the $n - r$ non-erased positions in $(\mathbf{y}, \mathbf{d}')$, we have that for $(\tilde{\mathbf{y}}, \tilde{\mathbf{d}}') \in \mathcal{S}_N((\tilde{\mathbf{x}}, \tilde{\mathbf{d}}))$, $P((\tilde{\mathbf{y}}, \tilde{\mathbf{d}}') \in \mathcal{N}_{(\mathbf{y}, \mathbf{d}')}) = p^i$. Note that there are $2^r \binom{n-r}{i}$ such $(\tilde{\mathbf{x}}, \tilde{\mathbf{d}})$.

Hence, from linearity of expectations,

$$\mathbb{E}[|\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}] \mid (\mathbf{y}, \mathbf{d}') = (N(2^r - 1) + N - 1) + N2^r \sum_{i=1}^{n-r} \binom{n-r}{i} p^i$$

$$= N2^r(1+p)^{n-r} - 1.$$

□

The permutation recovery algorithm as described below, iteratively selects the right node $(\mathbf{y}, \mathbf{d}')$ with the smallest two-hop neighborhood in \mathcal{Y} and then as shown in Fig. 2, performs $|\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}|$ data comparisons to identify the remaining $N - 1$ copies. Note that the algorithm finds the remaining $N - 1$ copies if and only if $(\mathbf{y}, \mathbf{d}') \notin \mathcal{R}_{\text{faulty}}$. Let $\mathcal{P}_G = (\mathcal{X} \cup \mathcal{Y}, \mathcal{P}_E)$ denote the bipartite matching identified by the permutation recovery algorithm.

Algorithm 2 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( $\mathcal{P}_G, \mathcal{G}$ )
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( $\mathcal{P}_G, \mathcal{G}$ )

```

Proposition 5. For $(L, N) \in \mathcal{R}$, Algorithm 2 finds the true permutation with probability at least $1 - \epsilon$, when $\epsilon_1, \epsilon_2 < \frac{\epsilon}{2}$.

Proof. For $L > L_{\text{Th}}$, every left node has at least one non-faulty channel output with probability at least $(1 - \epsilon_1)$. Thus, the permutation recovery algorithm identifies the true partitioning with probability at least $(1 - \epsilon_1)$. For $N > N_{\text{Th}}$, there exists only one valid labelling, viz. the true labelling with probability at least $(1 - \epsilon_2)$. Thus, the permutation recovery algorithm identifies the true permutation with probability at least $(1 - \frac{\epsilon}{2})^2 > (1 - \epsilon)$. □

VII. ANALYSIS OF PERMUTATION RECOVERY ALGORITHM FOR BEC

In this section, we analyse the expected number of data comparisons performed by Algorithm 2 for three subregions of \mathcal{R} . In the next lemma, we give an upper bound on the expected number of data comparisons performed by Algorithm 2 when $(L, N) \in \mathcal{R}$.

Lemma 7. The expected number of data comparisons performed by Algorithm 2 when $(L, N) \in \mathcal{R}$ is at most

$$\mathcal{U}_0 \triangleq N^2 2^n (1 + 2p - p^2)^n.$$

Proof. Note that the number of data comparisons performed by Algorithm 2 is at most $\sum_{(\mathbf{y}, \mathbf{d}') \in \mathcal{Y}} |\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}|$. From linearity of expectations, $\mathbb{E} \left[\sum_{(\mathbf{y}, \mathbf{d}') \in \mathcal{Y}} |\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}| \right] = \sum_{(\mathbf{y}, \mathbf{d}') \in \mathcal{Y}} \mathbb{E} [|\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}|]$. From Lemma 6, the result follows. □

Let β_0 be a threshold on β such that for $\beta \geq \beta_0$, $P(|\mathcal{R}_{\text{faulty}}| > 1) < \epsilon_1$. In the next lemma, we derive this threshold β_0 .

Lemma 8. For $\beta \geq \beta_0 \triangleq \frac{\log_2 \left(\frac{\epsilon_1}{2^n N^2 ((1+2p-p^2)^n - 1)} \right)}{n \log_2 (1 - \frac{1}{2}(1-p)^2)}$, $P(|\mathcal{R}_{\text{faulty}}| > 1) < \epsilon_1$.

Proof. If $\mathbb{E} [|\mathcal{R}_{\text{faulty}}|] < \epsilon_1$ then from Markov inequality the result follows. From Lemma 2, we have that

$$\mathbb{E} [|\mathcal{R}_{\text{faulty}}|] = N2^n \left(1 - \prod_{r=1}^n \left(1 - (2p - p^2)^r \left(1 - \frac{1}{2}(1-p)^2 \right)^L \right)^{N \binom{n}{r}} \right).$$

Hence, to show that $\mathbb{E}[|R_{\text{faulty}}|] < \epsilon_1$ it is sufficient to show that

$$1 - \frac{\epsilon_1}{2^n N} < \prod_{r=1}^n \left(1 - (2p - p^2)^r \left(1 - \frac{1}{2}(1-p)^2 \right)^L \right)^{N \binom{n}{r}}.$$

Using Weierstrass inequality we have that

$$\begin{aligned} & \prod_{r=1}^n \left(1 - (2p - p^2)^r \left(1 - \frac{1}{2}(1-p)^2 \right)^L \right)^{N \binom{n}{r}} \\ & \geq 1 - \sum_{r=1}^n N \binom{n}{r} (2p - p^2)^r \left(1 - \frac{1}{2}(1-p)^2 \right)^L \\ & = 1 - N \left(\left(1 - \frac{1}{2}(1-p)^2 \right)^L ((1 + 2p - p^2)^n - 1) \right), \end{aligned}$$

and thus it is enough to show that

$$N \left(\left(1 - \frac{1}{2}(1-p)^2 \right)^L ((1 + 2p - p^2)^n - 1) \right) < \frac{\epsilon_1}{2^n N}.$$

Lastly, it can be verified that $\mathbb{E}[|R_{\text{faulty}}|] < \epsilon_1$ if $\log_{(1-\frac{1}{2}(1-p)^2)} \left(\frac{\epsilon_1}{2^n N^2 ((1+2p-p^2)^n - 1)} \right) < L = \beta n$. \square

We define $\mathcal{R}' \subseteq \mathcal{R}$ as $\mathcal{R}' \triangleq \{(\beta, N) : \beta \geq \beta_0, N \geq N_{\text{Th}}\}$. To analyse the expected number of data comparisons performed by Algorithm 2 when $(\beta, N) \in \mathcal{R}'$, we define the notion of order of a left node.

Definition 4. A node $x \in \mathcal{X}$ has order s if $\min\{|E_{(\mathbf{y}, \mathbf{d}')}| : (\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((x, \mathbf{d}))\} = s$.

For $s \in [2^n]$, let \mathbf{X}_s denote the set of left nodes with order s . In the next lemma, we calculate the probability that a left node has order s .

Lemma 9. For $x \in \mathcal{X}$, before the initiation of Algorithm 2, $P(x \in \mathbf{X}_s)$ is

$$\begin{cases} \left(\sum_{i=\ell}^n \binom{n}{i} p^i (1-p)^{n-i} \right)^N - \left(\sum_{i=\ell+1}^n \binom{n}{i} p^i (1-p)^{n-i} \right)^N & s \in \{2^\ell, \ell \in [0 : n]\} \\ 0 & \text{otherwise.} \end{cases}$$

Proof. Note that before any edges are removed from \mathcal{G} , the degree of a right node can take values only from the set $\{2^\ell, \ell \in [0 : n]\}$. Therefore, for $s \notin \{2^\ell, \ell \in [0 : n]\}$, $P(x \in \mathbf{X}_s) = 0$. For $s \in \{2^\ell, \ell \in [0 : n]\}$, it can be verified that $P(x \in \mathbf{X}_s) = P(\bigcup_{(\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((x, \mathbf{d}))} |E_{(\mathbf{y}, \mathbf{d}')}| \geq s) - P(\bigcup_{(\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((x, \mathbf{d}))} |E_{(\mathbf{y}, \mathbf{d}')}| \geq s+1)$. Therefore, the result follows. \square

In the next lemma, we derive an upper bound on the expected number of data comparisons performed by Algorithm 2 when $(\beta, N) \in \mathcal{R}'$.

Lemma 10. For $(\beta, N) \in \mathcal{R}'$, the expected number of data comparisons performed by Algorithm 2 is at most

$$\mathcal{U}_1 \triangleq \sum_{r=0}^n \mathbb{E}[|\mathbf{X}_{2^r}|] N 2^r ((1+p)^{n-r}).$$

Proof. Note that for $\beta > \beta_0$, there are no faulty right nodes with probability at least $1 - \epsilon_1$. Hence, the expected number of data comparisons performed by Algorithm 2 to identify the N channel outputs of x is at most $\mathbb{E}[\min\{|\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}| : (\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((x, \mathbf{d}))\}]$. For $x \in \mathcal{X}$, by law of total expectation

$$\begin{aligned} & \mathbb{E}[\min\{|\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}| : (\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((x, \mathbf{d}))\}] \\ & = \mathbb{E}[\mathbb{E}[\min\{|\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}| : (\mathbf{y}, \mathbf{d}') \in \mathcal{S}_N((x, \mathbf{d}))\}] \mid x \in \mathbf{X}_s] \end{aligned}$$

From Lemma 5, $\mathbb{E}[|\mathcal{N}_{(\mathbf{y}, \mathbf{d}')}| \mid (\mathbf{y}, \mathbf{d}')] = N 2^r (1+p)^{n-r} - 1$,

$$= \sum_{r=0}^n P(x \in \mathbf{X}_{2^r}) (N 2^r (1+p)^{n-r} - 1).$$

From linearity of expectation, the result follows. \square

We now define the notion of confusability for left nodes.

Definition 5. Let $x, \tilde{x} \in \mathcal{X}$ then x is *confusable* with \tilde{x} , denoted by $x \rightarrow \tilde{x}$, if there exists at least one $(\tilde{y}, \tilde{d}') \in \mathcal{S}_N((\tilde{x}, \tilde{d}'))$ such that $E_{(\tilde{y}, \tilde{d}')} = \{x, \tilde{x}\}$.

Next, we build a graph of left nodes, $T = (\mathcal{X}, E_{\text{conf}})$. Let $x, \tilde{x}, x' \in \mathcal{X}$. Note that before the initiation of Algorithm 2, for $x \rightarrow \tilde{x}$, it must be that $d_H(x, \tilde{x}) = 1$. For ease of analysis, we do not consider the confusable edges that would be generated over the course of Algorithm 2. Thus, there is an edge $x \rightarrow \tilde{x} \in E_{\text{conf}}$ if and only if x is confusable with \tilde{x} before the initiation of the algorithm. In the next lemma, we derive the probability that x has edges to all nodes in $S \subseteq \{x' : d_H(x, x') = 1\}$.

Lemma 11. Let $x \in \mathcal{X}$ and let $S \subseteq \{x' : d_H(x, x') = 1\}$. Then,

$$P\left(\bigcup_{j=1}^{|S|} (x \rightarrow x_i)\right) = \prod_{j=1}^{|S|} \left(1 - (1 - p(1 - p)^{n-1})^{N-j+1}\right),$$

where $x_i \in S$ for $i \in [|S|]$.

Proof. From Definition 5, $x \rightarrow \tilde{x}$ there must exist a $y \in \mathcal{S}_N((x, d))$ such that $E_{(y, d')} = \{x, \tilde{x}\}$, which happens if and only if y is erased only at the position where x and \tilde{x} differ, which happens with the probability $p(1 - p)^{n-1}$. Therefore, $P\left(\bigcup_{j=1}^{|S|} (x \rightarrow x_i)\right) = \prod_{j=1}^{|S|} \left(1 - (1 - p(1 - p)^{n-1})^{N-j+1}\right)$. \square

Next, let $G_A = (\mathcal{X}, \mathcal{E})$ be a directed n -cube [2]. A vertex $x \in \mathcal{X}$ has outgoing edges to the vertices $\{x' : d_H(x, x') = 1, x' \in \mathcal{X}\}$. Let $G_A(p_e)$ denote a random sub-graph of G_A where every edge in \mathcal{E} is selected with probability p_e .

Proposition 6. The probability of the appearance of a connected component is greater in T than in $G_A(p_T)$, where $p_T \triangleq (1 - (1 - p(1 - p)^{n-1})^{N-n+1})$.

Proof. Since the probability that there is an edge $x \rightarrow x'$ is independent of the existence of the edge $\tilde{x} \rightarrow x'$, the proposition follows from Lemma 11. \square

Lemma 12. For $N > N_0 \triangleq n - \frac{1}{\log(1 - p(1 - p)^{n-1})} = \mathcal{O}_p\left(\frac{1}{p(1 - p)^{n-1}}\right)$, T is almost surely connected.

Proof. From [2], we know that $G_A(p_e)$ is almost surely connected if $p_e > \frac{1}{2}$. It can be verified that for $N > n - \frac{1}{\log(1 - p(1 - p)^{n-1})}$, $p_e = p_T > \frac{1}{2}$. Then from Proposition 6, the result follows. \square

We define region $\mathcal{R}'' \subseteq \mathcal{R}'$ as $\mathcal{R}'' \triangleq \{(\beta, N) : \beta \geq \beta_0, N \geq N_0\}$.

Lemma 13. The expected number of data comparisons performed by Algorithm 2 when $(\beta, N) \in \mathcal{R}''$ is at most

$$\mathcal{U}_2 \triangleq N2^n (1 + p)^n.$$

Proof. Since, the graph T is connected, Algorithm 2 will always prune an order 1 node. Hence, the result follows. \square

Hence, from Lemmas 7, 10 and 13, the expected number of data comparisons performed by Algorithm 2 is only a $\kappa_{\beta, N}$ -fraction of data comparisons required by clustering based approaches agnostic to the nature of DNA storage systems, $\left(\frac{1+2p-p^2}{2}\right)^n \geq \kappa_{\beta, N} \geq \left(\frac{1+p}{2}\right)^n$.

REFERENCES

- [1] S. Singhvi, A. Boruchovsky, H. M. Kiah and E. Yaakobi, "Data-Driven Bee Identification for DNA Strands", *arXiv preprint*, arXiv:2305.04597, 2023.
- [2] B. Bollobás, C. Gotsman, and E. Shamir, "Connectivity and dynamics for random subgraphs of the directed cube," *ISRAEL JOURNAL OF MATHEMATICS*, vol. 83, pp. 321–328, 1993.
- [3] J. Chrisnata, H. M. Kiah, A. Vardy, and E. Yaakobi, "Bee identification problem for DNA strands," *IEEE International Symposium on Information Theory (ISIT)*, pp. 969–974, June, 2022.
- [4] G. M. Church, Y. Gao, and S. Kosuri. "Next-generation digital information storage in DNA," *Science*, vol. 337, no. 6102, pp. 1628–1628, 2012.
- [5] N. Goldman, P. Bertone, S. Chen, C. Dessimoz, E. M. LeProust, B. Sipos, and E. Birney. "Towards practical, high-capacity, low-maintenance information storage in synthesized DNA," *Nature*, vol. 494, no. 7435, pp. 77–80, 2013.
- [6] H. M. Kiah, A. Vardy, and H. Yao, "Efficient bee identification," *IEEE International Symposium on Information Theory (ISIT)*, pp. 1943–1948, July, 2021.
- [7] H. M. Kiah, A. Vardy, and H. Yao, "Efficient algorithms for the bee-identification problem," *arXiv preprint* arXiv:2212.09952, 2022.
- [8] A. Lenz, P. H. Siegel, A. Wachter-Zeh and E. Yaakobi, "Coding over sets for DNA storage," *IEEE Transactions on Information Theory*, vol. 66, no. 4, pp. 2331–2351, April 2020.
- [9] L. Organick, S. Ang, Y.J. Chen, R. Lopez, S.Yekhanin, K. Makarychev, M. Racz, G. Kamath, P. Gopalan, B. Nguyen, C. Takahashi, S. Newman, H. Y. Parker, C. Rashtchian, K. Stewart, G. Gupta, R. Carlson, J. Mulligan, D. Carmean, G. Seelig, L. Ceze, and K. Strauss, "Random access in largescale DNA data storage," *Nature Biotechnology*, vol. 36, no. 3, pp. 242–248, 2018.
- [10] 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," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [11] I. Shomorony, and R. Heckel, "Information-theoretic foundations of DNA data storage," *Foundations and Trends® in Communications and Information Theory*, 19(1), 1–106, 2022
- [12] A. Tandon , V.Y.F. Tan, and L.R. Varshney, "The bee-identification problem: Bounds on the error exponent," *IEEE Transactions on Communications*, vol. 67, issue no.11, pp. 7405–7416, November, 2019.
- [13] S. Yazdi, H. M. Kiah, E. R. Garcia, J. Ma, H. Zhao, and O. Milenkovic, "DNA-based storage: Trends and methods," *IEEE Trans. Molecular, Biological, Multi-Scale Commun.*, vol. 1, no. 3, pp. 230–248, 2015.