Correcting multiple short duplication and substitution errors

Yuanyuan Tang*, Shuche Wang[†], Ryan Gabrys[‡], and Farzad Farnoud*

* Electrical & Computer Engineering, University of Virginia, U.S.A., {yt5tz,farzad}@virginia.edu

† Institute of Operations Research and Analytics, National University of Singapore, shuche.wang@u.nus.edu

‡ Calit2, University of California-San Diego, U.S.A., rgabrys@eng.ucsd.edu

Abstract—Due to its higher data density, longevity, energy efficiency, and ease of generating copies, DNA is considered a promising storage technology for satisfying future needs. However, a diverse set of errors including deletions, insertions, duplications, and substitutions may arise in DNA at different stages of data storage and retrieval. The current paper constructs error-correcting codes for simultaneously correcting short (tandem) duplications and at most p substitutions, where a short duplication generates a copy of a substring with length ≤ 3 and inserts the copy following the original substring. Compared to the state-of-the-art codes for duplications only, the proposed codes correct up to p substitutions (in addition to duplications) at the additional cost of roughly $8p(\log_a n)(1 + o(1))$ symbols of redundancy, thus achieving the same asymptotic rate, where q > 4 is the alphabet size. Furthermore, the time complexities of both the encoding and decoding processes are polynomial when p is a constant with respect to n.

I. INTRODUCTION

With recent advances in sequencing and biological synthesis, deoxyribonucleic acid (DNA) is considered a promising candidate for satisfying future data storage needs [1], [2]. In particular, experiments in [1], [3]–[7] demonstrate that data can be stored on and subsequently retrieved from DNA. Compared to traditional data storage media, DNA has the advantages of higher data density, longevity, energy efficiency, and ease of generating copies [1], [7]. However, a diverse set of errors may occur at different stages of the data storage and retrieval process, such as deletions, insertions, duplications, and substitutions. Many recent works, such as [7]–[24], have been devoted to protecting the data against these errors. The current paper constructs error-correcting codes for duplication and substitution errors.

A (tandem) duplication in a DNA sequence generates a copy of a substring and then inserts it directly following the original substring [8], where the duplication length is the length of the copy. For example, given ACTG, a tandem duplication may generate ACTCTG, where CTCT is a (tandem) repeat. Bounded-length duplications are those whose length is at most a given constant. In particular, we refer to duplications of length at most 3 as short duplications. Correcting fixed-length duplications [8], [10]—[12], [25] and bounded-length duplications [8], [23], [26]—[29] have been both studied recently. In particular, the code in [8], which has a polynomial-time encoder, provides the highest known asymptotic rate for

This work was supported in part by NSF grants CIF 1816409 and CIF 1755773.

correcting any number of short duplications. For an alphabet of size q, the rate is $\log r$, where r is the largest positive real root of $x^3-(q-2)x^2-(q-3)x-(q-2)=0$. For q=4 this rate is $\log 2.6590$ and as q increases, the rate is approximately $\log(q-1)$ $\boxed{23}$.

For channels with both duplication and substitution errors, restricted substitutions [12], [25], which occur only in duplicated copies, and unrestricted substitutions [12], [29]–[31] have been studied. The closest work to the current paper, [31], constructed error-correcting codes for short duplications and at most p unrestricted substitutions with an asymptotic code rate lower bounded by $\log(q-2)$ [31]. However, compared to the codes for only duplications [8], the codes in [31] incur an asymptotic rate loss in order to correct the additional $\leq p$ substitutions. The current paper focuses on constructing error-correcting codes for short duplications and at most p (unrestricted) substitutions with significantly less redundancy than the approach from [31]. We note that the short duplications and substitutions can occur in an arbitrary order.

One of the challenging aspects of correcting duplications and at most p substitutions simultaneously is that a single substitution may be duplicated many times and affect an unbounded segment. However, assuming that the input to the channel is irreducible, i.e., it has no repeats of length ≤ 3 , after removing all tandem copies with length ≤ 3 from the output, the effects of short duplications and at most p substitutions can be localized in at most p substrings, each with length ≤ 17 [29], [31]. Therefore, similar to [31], we construct our error-correcting codes as a subset of irreducible strings, but leverage the syndrome compression technique to substantially reduce the redundancy. Syndrome compression has been recently used to provide explicit constructions for correcting a wide variety of errors with the redundancy as low as roughly twice the Gilbert-Varshamov bound [32]-[35]. More specifically, we protect the data by appending a vector with length around $8p(\log_a n)(1 + o(1))$ to each input of length n, where the appended vector is used to distinguish all confusable inputs. We ensure that the appended vector is itself protected against errors and can be decoded correctly. We then use it to recover the data by eliminating incorrect confusable inputs. Compared to the explicit code for duplications only [8], the proposed code corrects $\leq p$ substitutions in addition to duplications at the extra cost of roughly $8p(\log_q n)(1+o(1))$ symbols of redundancy for $q \ge 4$, and achieves the same asymptotic code rate. This improves

upon the approach from [31], which suffers an asymptotic rate loss. Both time complexities of the encoding and decoding processes are polynomial when p is a constant.

The paper is organized as follows. Section III presents the notation and preliminaries. In Section IIII, we derive an upper bound on the size of the confusable set for an irreducible string, which is a key step of the syndrome compression technique used to construct our error-correcting codes. Finally, Section IV presents the code construction, as well as a discussion of the redundancy and an analysis of the encoding and decoding complexities. Due to lack of space, some of the proofs are omitted or only sketched.

II. NOTATION AND PRELIMINARIES

Let $\Sigma_q = \{0,1,2,\cdots,q-1\}$ represent a finite alphabet of size q, Σ_q^n the set of all strings of length n over Σ_q , and Σ_q^* the set of all finite strings over Σ_q . Given two integers a,b with $a \leq b$, the set $\{a,a+1,\cdots,b\}$ is shown as [a,b]. We simplify [1,b] as [b]. Unless otherwise stated, logarithms are to the base 2.

We use bold symbols to denote strings over Σ_q , i.e., $x, y_j \in \Sigma_q^*$. The entries of a string are represented by plain typeface, e.g., the *i*th elements of $x, y_j \in \Sigma_q^*$ are $x_i, y_{ji} \in \Sigma_q$ respectively. For two strings $x, y \in \Sigma_q^*$, let xy denote their concatenation. Given four strings $x, u, v, w \in \Sigma_q^*$, if x = uvw, then v is called a substring of x. Furthermore, we let |x| represent the length of a string $x \in \Sigma_q^n$, and let ||S|| denote the size (the number of elements) of a set S.

A (tandem) duplication (TD) of length k (k-TD) is the operation of generating a copy of a substring and inserting it directly following the substring, where k is the length of the copy. For example, for x = uvw with |v| = k, a k-TD may generate uvvw, where vv is called a (tandem) repeat with length 2k. A duplication of length at most k is denoted as a $\leq k$ -TD. We focus on $\leq k$ -TDs with k = 3, which we call short duplications. For example, given $x = 213012 \in \Sigma_4^*$, a sequence of ≤ 3 -TDs may produce

$$x = 213012 \rightarrow 213\underline{213}012 \rightarrow 2132130\underline{30}12$$

 $\rightarrow 2132\underline{2}1303012 = x',$ (1)

where the duplicated copies are marked with underlines. We call x' a *descendant* of x, i.e., a string generated from x by a sequence of ≤ 3 -TDs. Furthermore, for a string $x \in \Sigma_q^*$, let $D_{\leq k}^*(x)$ be the set of all descendants generated from x by an arbitrary number of $\leq k$ -TDs.

A deduplication of length k replaces a repeat vv by v with |v|=k. Then $\leq k$ -deduplications are deduplications with length upper bound by k. In this paper, we focus on ≤ 3 -deduplications, simply called deduplications in the rest of the paper. For example, the string x in (1) can be recovered from x' by three deduplications.

The set of $\leq k$ -irreducible strings of length n, denoted $\operatorname{Irr}_{\leq k}(n)$, consists of strings without repeats vv, where $|v| \leq k$. Furthermore, $\operatorname{Irr}_{\leq k}(*)$ represents all $\leq k$ -irreducible strings of finite length. A duplication root of x' is a $\leq k$ -irreducible string x such that x' is a descendant of x. Equivalently, x can be obtained from x' by performing all possible $\leq k$ -deduplications. The set of duplication roots of x' is denoted

 $R_{\leq k}(x')$, i.e., $R_{\leq k}(x') = \{x \in \operatorname{Irr}(*) \mid x' \in D_{\leq k}^*(x)\}$. For ≤ 3 -TDs, the work [8] showed that $R_{\leq 3}(x')$ is a singleton, and so we treat it as a string instead of a set. The uniqueness of the root for k=3 implies that if x'' is a descendant of x', we have $R_{\leq 3}(x') = R_{\leq 3}(x'')$.

Besides \leq 3-TDs, we consider substitution errors, where each substitution replaces a symbol by another one from the same alphabet. Continuing the example in (I), two substitutions and two duplications applied to x' may produce

$$\mathbf{x}' = 213221303012 \rightarrow 213211303012 \rightarrow 213213211303012$$

 $\rightarrow 213213211323012 \rightarrow 213213211323323012 = \mathbf{x}'',$

where the substituted symbols are marked in red. Let $D^{\alpha,p}_{\leq k}(\boldsymbol{x})$ represent the set of strings derived from \boldsymbol{x} by $\alpha \leq k$ -TDs and p substitutions. Furthermore, let $D^{*,\leq p}_{\leq k}(\boldsymbol{x})$ represent the set of strings generated by an arbitrary number of $\leq k$ -TDs and at most p substitutions. In the example above, we have $\boldsymbol{x}'' \in D^{*,\leq 2}_{\leq k}(\boldsymbol{x})$.

For simplicity, when k=3, we drop the ≤ 3 subscript and write $D^*(\cdot), R(\cdot), \operatorname{Irr}(\cdot), D^{\alpha,p}(\boldsymbol{x})$, and $D^{*,\leq p}(\boldsymbol{x})$. In the rest of the paper, unless otherwise stated, duplications are assumed to be ≤ 3 -TDs, and irreducible strings represent ≤ 3 -irreducible strings.

We define a *substring edit* in a string $x \in \Sigma_q^*$ as the operation of replacing a substring u with a string v, where at least one of u, v is nonempty. The length of the substring edit is $\max\{|u|, |v|\}$. An *L-substring edit* is one whose length is at most L. Furthermore, an *L-burst deletion* in $x \in \Sigma_q^*$ is defined as removing a substring v of x, where |v| = L is the length of the burst deletion.

Given a sequence $x \in \Sigma_q^n$, we define the binary representation matrix $\mathcal{U}(x)$ of x as

$$\begin{bmatrix} u_{1,1} & u_{1,2} & \cdots & u_{1,n} \\ u_{2,1} & u_{2,2} & \cdots & u_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ u_{\lceil \log q \rceil, 1} & u_{\lceil \log q \rceil, 2} & \cdots & u_{\lceil \log q \rceil, n} \end{bmatrix} \in \{0, 1\}^{\lceil \log q \rceil \times n},$$

where the *j*th column of $\mathcal{U}(x)$ is the binary representation of the *j*th symbol of x for $j \in [n]$. The *i*th row of $\mathcal{U}(x)$ is denoted as $\mathcal{U}_i(x)$ for $i \in \lceil \log q \rceil$.

In order to construct error-correcting codes by applying the syndrome compression technique [32], we first introduce some auxiliary definitions and a theorem.

Suppose $q \geq 3$ is a constant. We start with the definition of confusable sets for a given channel and a given set of strings $S \subseteq \Sigma_q^n$. In our application, S is the set of irreducible strings, upon which the proposed codes will be constructed.

Definition 1. A confusable set $B(x) \subseteq S$ of $x \in S$ consists of all $y \in S$, excluding x, such that x and y can produce the same output when passed through the channel.

Definition 2. Let $\mathcal{R}(n)$ be an integer function of n. A labeling function for the confusable sets $B(x), x \in S$, is a function

$$f:\Sigma_q^n\to \Sigma_{2^{\mathcal{R}(n)}}$$

such that, for any $x \in S$ and $y \in B(x)$, $f(x) \neq f(y)$.

Theorem 3. (c.f. [32] Theorem 5]) Let $f: \Sigma_q^n \to \Sigma_{2^{\mathcal{R}(n)}}$, where $\mathcal{R}(n) = o(\log\log n \cdot \log n)$, be a labeling function for the confusable sets $B(\boldsymbol{x}), \boldsymbol{x} \in S$. Then there exists an integer $a \leq 2^{\log \|B(\boldsymbol{x})\| + o(\log n)}$ such that for all $\boldsymbol{y} \in B(\boldsymbol{x})$, we have $f(\boldsymbol{x}) \not\equiv f(\boldsymbol{y}) \bmod a$.

The above definitions and theorem are used in our code construction based on syndrome compression, presented in Section [IV] The construction and analysis rely on the confusable sets for the channel, discussed in the next section.

III. CONFUSABLE SETS FOR CHANNELS WITH SHORT DUPLICATION AND SUBSTITUTION ERRORS

In this section, we study the confusable sets of input strings passing through channels with an arbitrary number of \leq 3-TDs and at most p substitutions.

Since \leq 3-TDs and deduplications do not alter the duplication root of the input and because the duplication root is unique, Irr(n) is a code capable of correcting \leq 3-TDs. The decoding process simply removes all tandem repeats. In other words, if we append a deduplication block, which removes all repeats, to the channel with duplication errors, any irreducible sequence passes through this concatenated channel with no errors. This approach produces codes with the same asymptotic rate as that of [8], achieving the highest known asymptotic rate.

Similar to $\boxed{31}$, we extend this strategy to correct duplication and substitution errors. First, we take the code to be a subset of irreducible strings. Second, we find the code for a new channel obtained by concatenating a deduplication block to the channel with duplication and substitution errors (recall that duplications and substitution errors can occur in any order). Denote the channel that introduces any number of duplications and $\leq p$ substitutions, followed by a deduplication block that removes all repeats, as the DSD(p) channel. It is clear that an error-correcting code for the DSD(p) channel is also an error-correcting code for the channel with duplications and $\leq p$ substitutions. We define the confusable sets over Irr(n) for the DSD(p) channel.

Definition 4. Suppose $x \in Irr(n)$ is an irreducible string of length n. Let

$$B_{\operatorname{Irr}}^{\leq p}(\boldsymbol{x}) = \{ \boldsymbol{y} \in \operatorname{Irr}(n) : \boldsymbol{y} \neq \boldsymbol{x}, \\ R(D^{*, \leq p}(\boldsymbol{x})) \cap R(D^{*, \leq p}(\boldsymbol{y})) \neq \varnothing \}$$
(3)

denote the irreducible-confusable set of x.

We now find a bound on $\|B_{\operatorname{Irr}}^{\leq p}(x)\|$, which is needed to construct the code and determine its rate. Since deduplications can be undone by duplications, instead of the $\operatorname{DSD}(p)$ channel, we can consider a concatenation of p DSD(1) channels, without reducing the size of the confusable set. The input of each DSD(1) channel suffers a number of duplications, at most one substitution, and then all possible deduplications. Fig. I shows a confusable string z, obtainable from both $x, y \in \operatorname{Irr}(n)$, after passing through p DSD(1) channels, each represented by a solid arrow. More precisely, $x_i \in R(D^{*,\leq 1}(x_{i-1}))$ and $y_i \in R(D^{*,\leq 1}(y_{i-1}))$, where $x = x_0, y = y_0, z = x_p = y_p$.

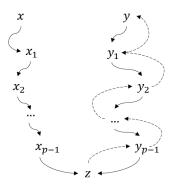


Figure 1. Confusable strings for a concatenation of channels with deduplications, at most 1 substitution, and all possible deduplications.

Lemma 5. Given an irreducible string $x \in Irr(n)$, the irreducible-confusable set of x satisfies

$$\left\| B_{\operatorname{Irr}}^{\leq p}(\boldsymbol{x}) \right\| \leq \max_{\substack{\{\boldsymbol{x}_i\}, \\ \{\boldsymbol{y}_j\}}} \prod_{i=0}^{p-1} \left\| R(D^{*,\leq 1}(\boldsymbol{x}_i)) \right\| \prod_{j=1}^{p} \left\| R(D^{*,\leq 1}(\boldsymbol{y}_i)) \right\|$$

where the maximums are over $x_i \in R(D^{*,\leq i}(x))$, $y_j \in R(D^{*,\leq j}(y))$ with $y \in B^{\leq p}_{\operatorname{Irr}}(x)$, and $x_p = y_p$.

Proof sketch: As discussed above, if $y \in B^{\leq p}_{\operatorname{Irr}}(x)$, then there exists z that can be obtained from both x and y via p concatenated DSD(1) channels, as shown in Fig. \square Critically, the DSD(1) channel is reversible. That is, if u and v are irreducible sequences such that $u \in R(D^{*,\leq 1}(v))$, we also have $v \in R(D^{*,\leq 1}(u))$. The dashed lines in the figure represent the reverse channels. So the total number of possibilities for y can be bounded by considering the number of possibilities in each step from x to y. The lemma then follows by induction on the number of possibilities for each string in the sequence $x \to x_1 \to \cdots \to z \to u$.

 $x \to x_1 \to \cdots \to z \to y_{p-1} \to \cdots \to y$. Our next step in bounding $\|B_{\operatorname{Irr}}^{\leq p}(x)\|$ is to bound $\|R(D^{*,\leq 1}(u))\|$ for an irreducible string u. We need only consider $\|R(D^{*,1}(u))\|$ as $\|R(D^{*,\leq 1}(u))\| \leq \|R(D^{*,1}(u))\| + 1$. The next lemma provides a bound on $\|R(D^{*,1}(u))\|$, which depends on the length of u. The proof relies on first reducing the problem to a special case in which |u| = 5 and then leveraging the regular language representing $D^*(u)$.

Lemma 6. Let $u \in Irr(n) \subseteq \Sigma_q^n$ be an irreducible string of length n with $q \geq 3$. Then

$$||R(D^{*,1}(\boldsymbol{u}))|| \le 968nq.$$

Proof sketch: We first show that the effect of the substitution can be isolated to a substring of length 5, i.e.,

$$||R(D^{*,1}(\boldsymbol{u}))|| \le n \max_{abcde \in \Sigma_q^5} ||R(D^{*,1}(abcde))||.$$

It can then be shown that the maximum on the right side can be obtained by assuming all symbols are distinct, i.e., replacing abcde by 01234. We then use a finite automaton to represent the set $D^*(01234)$, where paths correspond to descendants of 01234. A substitution can be represented by changing the label of an edge. With certain modifications, the automaton can also

be used to represent $R(D^{*,1}(01234))$ and to bound its size, which is determined to be 968q.

According to Lemma 6, the upper bound of $||R(D^{*,1}(u))||$ depends on the length n of u. Therefore, we can obtain upper bounds on $||R(D^{*,\leq 1}(x_i))||$ and $||R(D^{*,\leq 1}(y_j))||$ in Lemma 5 if upper bounds on $||x_i||$ and $||y_j||$ are available. To obtain these bounds, we recall a theorem from ||36||.

Theorem 7. [36] Theorem 5] Given strings $\mathbf{x} \in \Sigma_q^n$ and $\mathbf{v} \in D^{*,1}(\mathbf{x})$, $R(\mathbf{v})$ can be obtained from $R(\mathbf{x})$ by deleting a substring of length at most $\mathcal{L} = 17$ and inserting a substring of length at most \mathcal{L} in the same position.

In Lemma 5, we have $x_i \in R(D^{*,\leq i}(x))$ and $y_j \in R(D^{*,\leq j}(y))$ with $i,j\leq p$, implying that

$$|\boldsymbol{x}_i| \le n + p\mathcal{L}, \quad |\boldsymbol{y}_i| \le n + p\mathcal{L}.$$
 (4)

The next theorem follows from Lemmas 5 and 6

Theorem 8. Let $x \in Irr(n) \subseteq \Sigma_q^n$ be an irreducible string of length n, with $q \geq 3$. Then

$$||B_{\text{Irr}}^{\leq p}(\boldsymbol{x})|| \leq (968q(n+p\mathcal{L})+1)^{2p}.$$

IV. ERROR-CORRECTING CODES

As stated in Section $\boxed{\text{III}}$ our error-correcting code for correcting duplications and substitutions is a subset of irreducible strings of a given length. In this section, we construct this subset by applying the syndrome compression technique $\boxed{32}$, where we will make use of the size of the irreducible-confusable set $\|B_{\operatorname{Irr}}^{\leq p}(x)\|$ derived in Section $\boxed{\operatorname{III}}$ In this section, unless otherwise stated, we assume $q \geq 4$.

A. Code construction

We will start from a preliminary code, given in (5), and address its shortcomings, building up to the final code given in Construction 12.

Suppose p is constant with respect to n. Furthermore, suppose there exists a labeling function f and, for each $\boldsymbol{x} \in \operatorname{Irr}(n)$, an integer a such that for any $\boldsymbol{y} \in B^{\leq p}_{\operatorname{Irr}}(\boldsymbol{x})$, $f(\boldsymbol{x}) \not\equiv f(\boldsymbol{y}) \bmod a$. Let \boldsymbol{r} be a vector which encodes the information $(a, f(\boldsymbol{x}) \bmod a)$, to be precisely determined later. The set

$$\{xr: x \in Irr(n)\}\tag{5}$$

is a code capable of correcting duplications and at most p substitutions provided that, given the output $\boldsymbol{w} \in R(D^{*,\leq p}(\boldsymbol{x}\boldsymbol{r}))$, the following conditions are satisfied: 1) we can recover $(a, f(\boldsymbol{x}) \bmod a)$ and 2) we can find some $\boldsymbol{v} \in R(D^{*,\leq p}(\boldsymbol{x}))$. To see this, observe that if $\boldsymbol{y} \neq \boldsymbol{x}$ can also produce \boldsymbol{v} , then $\boldsymbol{y} \in B^{\leq p}_{\operatorname{Irr}}(\boldsymbol{x})$, and hence it can be eliminated as an input candidate since $f(\boldsymbol{y}) \not\equiv f(\boldsymbol{x}) \bmod a$.

The first condition can be addressed by adapting the code given in [31]. Construction 10], which has asymptotic rate $\geq \log(q-2)$. More precisely, a straightforward extension of [31]. Theorem 11] leads to the following lemma.

Lemma 9. Let $\sigma = 01020$. There exists an encoder \mathcal{E}_1 : $\Sigma_2^m \to \operatorname{Irr}(m')$ such that i) $\sigma \mathcal{E}_1(u) \in \operatorname{Irr}(*)$ and ii) for

any string $x \in \operatorname{Irr}(*)$ with $x\sigma \mathcal{E}_1(u) \in \operatorname{Irr}(*)$, we can recover u from any $w \in R(D^{*,\leq p}(x\sigma \mathcal{E}_1(u)))$. Asymptotically, $m' < m/\log(q-2)(1+o(1))$.

We use $\mathcal{E}_1(a, f(\boldsymbol{x}) \bmod a)$ to denote $\mathcal{E}_1(\boldsymbol{u})$, where \boldsymbol{u} is a binary sequence representing the pair $(a, f(\boldsymbol{x}) \bmod a)$. From the lemma, letting $\boldsymbol{r} = \boldsymbol{\sigma}\mathcal{E}_1(a, f(\boldsymbol{x}) \bmod a)$ will enable us to recover $(a, f(\boldsymbol{x}) \bmod a)$ from any $\boldsymbol{w} \in R(D^{*, \leq p}(\boldsymbol{xr}))$ for any $\boldsymbol{x} \in \operatorname{Irr}(*)$ provided that $\boldsymbol{xr} \in \operatorname{Irr}(*)$. We will discuss ensuring \boldsymbol{xr} is irreducible later.

The second condition requires us to find some \boldsymbol{v} that is only a function of \boldsymbol{x} rather than \boldsymbol{xr} . This is more challenging as the boundary between \boldsymbol{x} and \boldsymbol{r} becomes unclear or may not even exist after duplication and substitution errors, making it difficult to find $\boldsymbol{v} \in R(D^{*,\leq p}(\boldsymbol{x}))$ from $\boldsymbol{w} \in R(D^{*,\leq p}(\boldsymbol{xr}))$. To address this, instead of finding \boldsymbol{v} , we find some string \boldsymbol{s} that is only a function of \boldsymbol{x} as follows. Denote by $D^{*,\leq p,\leq 2p\mathcal{L}}(\boldsymbol{x})$ the set of strings that can be obtained by deleting a suffix of length at most $2p\mathcal{L}$ from some $\boldsymbol{v} \in R(D^{*,\leq p}(\boldsymbol{x}))$.

Lemma 10. Let x be an irreducible string of length n and r be any string such that xr is irreducible. Let $w \in R(D^{*, \leq p}(xr))$ and s be the prefix of w of length $n - p\mathcal{L}$. Then s can be obtained from some $v \in R(D^{*, \leq p}(x))$ by deleting a suffix of length at most $2p\mathcal{L}$. That is, $s \in D^{*, \leq p, \leq 2p\mathcal{L}}(x)$.

Hence, by choosing the first $n-p\mathcal{L}$ elements of $\boldsymbol{w} \in R(D^{*,\leq p}(\boldsymbol{xr}))$, we find $\boldsymbol{s} \in D^{*,\leq p,\leq 2p\mathcal{L}}(\boldsymbol{x})$, which is a function of only \boldsymbol{x} rather than \boldsymbol{xr} . But in doing so, we have introduced an additional error, namely deleting a suffix of length at most $2p\mathcal{L}$. As a result, we need to replace the labeling function f with a stronger labeling function g_q that, in addition to handling both substitutions and duplications, can handle deleting a suffix of \boldsymbol{x} . More precisely, g_q is a labeling function for the confusable set

$$B_{\operatorname{Irr}}^{\leq p, \leq 2p\mathcal{L}}(\boldsymbol{x}) = \{ \boldsymbol{y} \in \operatorname{Irr}(n) : \boldsymbol{y} \neq \boldsymbol{x}, \\ D^{*, \leq p, \leq 2p\mathcal{L}}(\boldsymbol{x}) \cap D^{*, \leq p, \leq 2p\mathcal{L}}(\boldsymbol{y}) \neq \varnothing \}. \quad (6)$$

The details of determining g_q will be discussed in Section [V-B]

The final piece of the construction is ensuring that the concatenation of x and r is irreducible. This can be done by adding a buffer b_x between them.

Lemma 11. For $q \geq 3$ and any irreducible string \boldsymbol{x} over Σ_q , there is a string $\boldsymbol{b_x}$ of length c_q such that $\boldsymbol{xb_x\sigma}$ is irreducible. Furthermore, it suffices to choose $c_3 = 13$, $c_4 = 7$, and $c_q = 6$ for $q \geq 5$.

The lemma implies xb_xr is irreducible, because r starts with σ , which has length 5, and because short repeats have length at most 6. So any repeat must be contained in $xb_x\sigma$ or in r, which is not possible as they are both irreducible.

We are now ready to present the code construction and then a theorem that establishes its error-correcting capability. The proof of the theorem summarizes our preceding discussion.

Construction 12. Let g_q be a labeling function for the confusable sets $B_{\operatorname{Irr}}^{\leq p, \leq 2p\mathcal{L}}(\boldsymbol{x}), \boldsymbol{x} \in \operatorname{Irr}(n)$. Furthermore, for each

x, let a_1 be an integer such that $g_q(x) \not\equiv g_q(y) \mod a_1$ for $y \in B_{\operatorname{Irr}}^{\leq p, \leq 2p\mathcal{L}}(x)$. Let

$$C_n = \{ xb_x r : x \in Irr(n), r = \sigma \mathcal{E}_1(a_1, g_q(x) \bmod a_1) \}.$$

Note that for simplicity, we index the code by the length of x rather than the length of the codewords xb_xr , i.e., n in \mathcal{C}_n refers to the length of x. The length of r is discussed in Subsection [V-C] below.

Theorem 13. The code in Construction $\boxed{12}$ can correct any number of short duplications and at most p substitutions.

Proof: Let the retrieved word be $\boldsymbol{w} \in R(D^{*,\leq p}(\boldsymbol{x}\boldsymbol{b_x}\boldsymbol{r}))$. From Lemma \boldsymbol{g} given \boldsymbol{w} , we can find a_1 and $g_q(\boldsymbol{x}) \mod a_1$. Let \boldsymbol{s} be the $(n-p\mathcal{L})$ -prefix of \boldsymbol{w} . By Lemma \boldsymbol{g} is irreducible. Then, by Lemma \boldsymbol{g} the $(n-p\mathcal{L})$ -prefix of \boldsymbol{w} , denoted \boldsymbol{s} , satisfies $\boldsymbol{s} \in D^{*,\leq p,\leq 2p\mathcal{L}}(\boldsymbol{x})$. By definition, for all $\boldsymbol{y} \neq \boldsymbol{x}$ that could produce the same \boldsymbol{s} , we have $\boldsymbol{y} \in B^{\leq p,\leq 2p\mathcal{L}}(\boldsymbol{x})$. But then, $g_q(\boldsymbol{y}) \not\equiv g_q(\boldsymbol{x}) \mod a_1$, and so we can determine \boldsymbol{x} by exhaustive search.

B. The labeling function

In this subsection, we present the labeling function g_q such that $g_q(x) \neq g_q(y)$ for $y \in B_{\operatorname{Irr}}^{\leq p, \leq 2p\mathcal{L}}(x)$. As shown in Theorem \overline{I} , each substitution is reflected in the root of the sequence as a substring edit of length at most \mathcal{L} . Considering also the suffix deletion of length at most $2p\mathcal{L}$, it follows that $s \in B_{\operatorname{Irr}}^{\leq p, \leq 2p\mathcal{L}}(x)$ can be obtained from both x and y by at most $3p\mathcal{L}$ deletions and at most $p\mathcal{L}$ insertions. Hence, it suffices to find g_q such that $g_q(x) \neq g_q(y)$ if there is a string s that can be obtained from both s and s through s indeed in the redundancy of the sufficient of the redundancy, which is still primarily controlled by $\max_{s \in \operatorname{Irr}(n)} \|B_{\operatorname{Irr}}^{\leq p, \leq 2p\mathcal{L}}(s)\|$. To find s we use the labeling function for binary sequences given in s whose properties are discussed in the following theorem.

Theorem 14. [35] There exists a labeling function $g: \{0,1\}^n \to \Sigma_{2^{\mathcal{R}(t,n)}}$ such that for any two distinct strings \mathbf{c}_1 and \mathbf{c}_2 confusable under t insertions, deletions, and substitutions, we have $g(\mathbf{c}_1) \neq g(\mathbf{c}_2)$, where $\mathcal{R}(t,n) = [(t^2+1)(2t^2+1)+2t^2(t-1)]\log n + o(\log n)$.

Since $s \in D^{*,\leq p,\leq 2p\mathcal{L}}(x)$ can be obtained from x via $4p\mathcal{L}$ indels, $\mathcal{U}_i(s)$ can be derived from $\mathcal{U}_i(x)$ by at most $4p\mathcal{L}$ indels for $i \in [\lceil \log q \rceil]$. Based on Theorem [14] by letting $t = 4p\mathcal{L}$, we can obtain a labeling function g for recovering $\mathcal{U}_i(x)$ from $\mathcal{U}_i(s)$ under at most $4p\mathcal{L}$ indels. Therefore, $g_q: \Sigma_q^n \to \Sigma_{2\lceil \log q \rceil \mathcal{R}(t,n)}$,

$$g_q(\mathbf{x}) = \sum_{i=1}^{\lceil \log q \rceil} 2^{\mathcal{R}(t,n)(i-1)} g(\mathcal{U}_i(\mathbf{x})), \tag{7}$$

where $t=4p\mathcal{L}$, is a labeling function for the confusable sets $B_{\mathrm{Irr}}^{\leq p,\leq 2p\mathcal{L}}(\boldsymbol{x}), \ \boldsymbol{x}\in\mathrm{Irr}(n).$ For each \boldsymbol{x} , a value a_1 needs to be also determined such that $g_q(\boldsymbol{x})\not\equiv g_q(\boldsymbol{y}) \bmod a_1$ for $\boldsymbol{y}\in B_{\mathrm{Irr}}^{\leq p,\leq 2p\mathcal{L}}(\boldsymbol{x}).$ The existence of such a_1 , satisfying $\log a_1\leq \log \|B_{\mathrm{Irr}}^{\leq p,\leq 2p\mathcal{L}}(\boldsymbol{x})\|+o(\log n)$, is guaranteed by Theorem 3 provided that p is a constant (ensuring that $p^4=o(\log\log n)$).

C. The redundancy of the error-correcting codes

In this section, we study the rate and the redundancy of the code proposed in Construction [12] and compare these to those of the short-duplication-correcting code given in [8], which has the highest known asymptotic rate. A simplified version of the construction of [8] is given below.

Construction 15. (c.f. [8]) For a positive integer n, let

$$C_n^d = \{ \boldsymbol{x} \in \Sigma_q^n : \boldsymbol{x} \in Irr(n) \}.$$

Given $q \geq 4$, this code has the same asymptotic rate as the original in [8] and its size differs by only a constant factor. We thus compare the proposed code with \mathcal{C}_n^d .

The codes \mathcal{C}_n and \mathcal{C}_n^d have the same size but the length of \mathcal{C}_n is larger by $|\boldsymbol{b_xr}|$ symbols. We have seen in Lemma 11 that $|\boldsymbol{b_x}| \leq 13$. The extra redundancy is then $|\boldsymbol{r}| + O(1)$, which depends on $\|B_{\mathrm{Irr}}^{\leq p, \leq 2p\mathcal{L}}(\boldsymbol{x})\|$, investigated in the next lemma.

Lemma 16. For $x \in Irr(n)$ over the alphabet Σ_q with $q \geq 3$,

$$||B_{\operatorname{Irr}}^{\leq p, \leq 2p\mathcal{L}}(\boldsymbol{x})|| \leq q^{4p\mathcal{L}}(n+p\mathcal{L})^{2p}.$$

The proof relies on Theorem 8 and takes into account the effect of the suffix deletion.

Lemma 17. Suppose both p and $q \ge 4$ are constant with respect to n. Given $\mathbf{x} \in \operatorname{Irr}(n) \subseteq \Sigma_q^n$, the length L of $\mathbf{r} = \boldsymbol{\sigma} \mathcal{E}_1(a_1, g_q(\mathbf{x}) \bmod a_1)$ satisfies $L \le 8p(\log_q n)(1 + o(1))$.

Proof: From the previous subsection, assuming p is a constant, we have that $\log a_1 \leq \log \|B_{\operatorname{Irr}}^{\leq p, \leq 2p\mathcal{L}}(\boldsymbol{x})\| + o(\log n)$. Since $(g_q(\boldsymbol{x}) \bmod a_1) \leq a_1$, we need $L_r = 2\log \|B_{\operatorname{Irr}}^{\leq p, \leq 2p\mathcal{L}}(\boldsymbol{x})\| + o(\log n)$ bits to represent the pair $(a_1, g_q(\boldsymbol{x}) \bmod a_1)$. By Lemma [16] assuming q is a constant, we have $L_r \leq 4p\log n + o(\log n) = 4p\log n(1+o(1))$. Then, by Lemma [9] $|\mathcal{E}_1((a_1, g_q(\boldsymbol{x}) \bmod a_1))| \leq 4p\log n(1+o(1))/\log(q-2)$. The lemma follows from the facts that $\log q \log q \geq 2$ for $q \geq 4$ and that $|\sigma| = 5$.

 $\frac{\log q}{\log(q-2)} \le 2$ for $q \ge 4$ and that $|\sigma| = 5$.

Using Lemma [17], the next theorem gives the extra redundancy of correcting p additional substitutions compared to [8]. It also shows that there is no asymptotic rate penalty, in contrast to prior work [31], which also corrects duplications and p substitutions.

Theorem 18. Assuming p and q are constants, compared to C_n^d , the proposed code C_n has the same size but is longer by $8p \log_q n(1 + o(1))$ symbols. The codes have the same asymptotic rate, which equals $\log 2.6590$ for q = 4.

D. The time complexities of encoding and decoding

The encoding process relies on determining a value for a_1 satisfying the condition discussed in Subsection IV-B among at most $2^{\log \|B_{\operatorname{Irr}}^{\leq p, \leq 2p\mathcal{L}}(\boldsymbol{x})\|+o(\log n)}$ possibilities. This step has complexity $O(n^{4p+1})$ in n, making the total complexity of encoding polynomial in n. Decoding requires deduplication, which is linear in the length of the retrieved sequence, and a brute-force search among all inputs that can lead to the same output $(n-p\mathcal{L})$ -prefix of the root of the retrieved sequence, which is polynomial in n. Hence, decoding is polynomial in the length of the retrieved sequence.

REFERENCES

- [1] S. 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, 2015.
- [2] S. Jain, F. Farnoud, M. Schwartz, and J. Bruck, "Noise and uncertainty in string-duplication systems," in 2017 IEEE International Symposium on Information Theory (ISIT). IEEE, 2017, pp. 3120–3124.
- [3] S. H. T. Yazdi, Y. Yuan, J. Ma, H. Zhao, and O. Milenkovic, "A rewritable, random-access DNA-based storage system," *Scientific re*ports, vol. 5, no. 1, pp. 1–10, 2015.
- [4] Y. Erlich and D. Zielinski, "DNA fountain enables a robust and efficient storage architecture," *Science*, vol. 355, no. 6328, pp. 950–954, 2017.
- [5] M. Blawat, K. Gaedke, I. Huetter, X.-M. Chen, B. Turczyk, S. Inverso, B. W. Pruitt, and G. M. Church, "Forward error correction for DNA data storage," *Procedia Computer Science*, vol. 80, pp. 1011–1022, 2016.
- [6] L. Organick, S. D. Ang, Y.-J. Chen, R. Lopez, S. Yekhanin, K. Makarychev, M. Z. Racz, G. Kamath, P. Gopalan, B. Nguyen et al., "Random access in large-scale DNA data storage," *Nature biotechnol-ogy*, vol. 36, no. 3, pp. 242–248, 2018.
- [7] H. H. Lee, R. Kalhor, N. Goela, J. Bolot, and G. M. Church, "Terminator-free template-independent enzymatic DNA synthesis for digital information storage," *Nature communications*, vol. 10, no. 1, pp. 1–12, 2019.
- [8] S. Jain, F. Farnoud, M. Schwartz, and J. Bruck, "Duplication-correcting codes for data storage in the DNA of living organisms," *IEEE Transac*tions on Information Theory, vol. 63, no. 8, pp. 4996–5010, 2017.
- [9] S. L. Shipman, J. Nivala, J. D. Macklis, and G. M. Church, "CRISPR-Cas encoding of a digital movie into the genomes of a population of living bacteria," *Nature*, vol. 547, no. 7663, pp. 345–349, Jul. 2017.
- [10] M. Kovačević and V. Y. Tan, "Asymptotically optimal codes correcting fixed-length duplication errors in DNA storage systems," *IEEE Commu*nications Letters, vol. 22, no. 11, pp. 2194–2197, 2018.
- [11] Y. Yehezkeally and M. Schwartz, "Reconstruction codes for DNA sequences with uniform tandem-duplication errors," *IEEE Transactions* on *Information Theory*, vol. 66, no. 5, pp. 2658–2668, 2020.
- [12] Y. Tang, Y. Yehezkeally, M. Schwartz, and F. Farnoud, "Single-error detection and correction for duplication and substitution channels," *IEEE Transactions on Information Theory*, vol. 66, no. 11, pp. 6908–6919, 2020.
- [13] 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, 2020.
- [14] K. Cai, Y. M. Chee, R. Gabrys, H. M. Kiah, and T. T. Nguyen, "Optimal codes correcting a single indel/edit for DNA-based data storage," arXiv preprint arXiv:1910.06501, 2019.
- [15] O. Elishco, R. Gabrys, and E. Yaakobi, "Bounds and constructions of codes over symbol-pair read channels," *IEEE Transactions on Information Theory*, vol. 66, no. 3, pp. 1385–1395, 2020.
- [16] A. Lenz, Y. Liu, C. Rashtchian, P. H. Siegel, A. Wachter-Zeh, and E. Yaakobi, "Coding for efficient DNA synthesis," in *IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2020, pp. 2885–2890.
- [17] R. Gabrys, S. Pattabiraman, and O. Milenkovic, "Mass error-correction codes for polymer-based data storage," in *IEEE International Symposium* on *Information Theory (ISIT)*, 2020, pp. 25–30.
- [18] S. Jain, F. Farnoud, M. Schwartz, and J. Bruck, "Coding for optimized writing rate in DNA storage," in *IEEE International Symposium on Information Theory (ISIT)*, 2020, pp. 711–716.
- [19] H. M. Kiah, T. Thanh Nguyen, and E. Yaakobi, "Coding for sequence reconstruction for single edits," in *IEEE International Symposium on Information Theory (ISIT)*, 2020, pp. 676–681.
- [20] Y. Yehezkeally and M. Schwartz, "Uncertainty of reconstructing multiple messages from uniform-tandem-duplication noise," in *IEEE Interna*tional Symposium on Information Theory (ISIT), 2020, pp. 126–131.
- [21] T. T. Nguyen, K. Cai, K. A. S. Immink, and H. M. Kiah, "Constrained coding with error control for DNA-based data storage," in *IEEE Inter*national Symposium on Information Theory (ISIT). IEEE, 2020, pp. 694–699.
- [22] J. Sima, N. Raviv, and J. Bruck, "Robust indexing-optimal codes for DNA storage," in *IEEE International Symposium on Information Theory* (ISIT). IEEE, 2020, pp. 717–722.
- [23] Y. M. Chee, J. Chrisnata, H. M. Kiah, and T. T. Nguyen, "Efficient encoding/decoding of GC-balanced codes correcting tandem duplications," *IEEE Transactions on Information Theory*, vol. 66, no. 8, pp. 4892–4903, 2020.

- [24] Y. Tang and F. Farnoud, "Correcting deletion errors in DNA data storage with enzymatic synthesis," in 2021 IEEE Information Theory Workshop (ITW), 2021, pp. 1–6.
- [25] ——, "Error-correcting codes for noisy duplication channels," *IEEE Transactions on Information Theory*, vol. 67, no. 6, pp. 3452–3463, 2021.
- [26] S. Jain, F. Farnoud, and J. Bruck, "Capacity and expressiveness of genomic tandem duplication," *IEEE Transactions on Information Theory*, vol. 63, no. 10, pp. 6129–6138, 2017.
- [27] M. Kovačević, "On the maximum number of non-confusable strings evolving under short tandem duplications," *Problems of Information Transmission*, vol. 58, no. 2, pp. 111–121, 2022.
- [28] Y. M. Chee, J. Chrisnata, H. M. Kiah, and T. T. Nguyen, "Deciding the confusability of words under tandem repeats in linear time," ACM Transactions on Algorithms (TALG), vol. 15, no. 3, pp. 1–22, 2019.
- [29] Y. Tang and F. Farnoud, "Error-correcting codes for short tandem duplication and substitution errors," in *IEEE International Symposium* on *Information Theory (ISIT)*. IEEE, 2020, pp. 734–739.
- [30] —, "Error-correcting codes for noisy duplication channels," in 2019 57th Annual Allerton Conference on Communication, Control, and Computing (Allerton). IEEE, 2019, pp. 140–146.
- [31] Y. Tang, H. Lou, and F. Farnoud, "Error-correcting codes for short tandem duplications and at most p substitutions," in 2021 IEEE International Symposium on Information Theory (ISIT). IEEE, 2021, pp. 1835–1840.
- [32] J. Sima, R. Gabrys, and J. Bruck, "Syndrome compression for optimal redundancy codes," in 2020 IEEE International Symposium on Information Theory (ISIT). IEEE, 2020, pp. 751–756.
- [33] J. Sima and J. Bruck, "On optimal k-deletion correcting codes," *IEEE Transactions on Information Theory*, vol. 67, no. 6, pp. 3360–3375, 2020.
- [34] J. Sima, R. Gabrys, and J. Bruck, "Optimal codes for the q-ary deletion channel," in 2020 IEEE International Symposium on Information Theory (ISIT). IEEE, 2020, pp. 740–745.
- [35] —, "Optimal systematic t-deletion correcting codes," in 2020 IEEE International Symposium on Information Theory (ISIT). IEEE, 2020, pp. 769–774.
- [36] Y. Tang and F. Farnoud, "Error-correcting codes for short tandem duplication and edit errors," *IEEE Transactions on Information Theory*, vol. 68, no. 2, pp. 871–880, 2022.