

Error-correcting Codes for Short Tandem Duplication and Edit Errors

Yuanyuan Tang and Farzad Farnoud

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

Abstract

Due to its high data density and longevity, DNA is considered a promising medium for satisfying ever-increasing data storage needs. However, the diversity of errors that occur in DNA sequences makes efficient error-correction a challenging task. This paper aims to address simultaneously correcting two types of errors, namely, short tandem duplication and edit errors, where an edit error may be a substitution, deletion, or insertion. We focus on tandem repeats of length at most 3 and design codes for correcting an arbitrary number of duplication errors and one edit error. Because an edited symbol can be duplicated many times (as part of substrings of various lengths), a single edit can affect an unbounded substring of the retrieved word. However, we show that with appropriate preprocessing, the effect may be limited to a substring of finite length, thus making efficient error-correction possible. We construct a code for correcting the aforementioned errors and provide lower bounds for its rate. Compared to optimal codes correcting only duplication errors, numerical results show that the asymptotic cost of protecting against an additional edit is only 0.003 bits/symbol when the alphabet has size 4, an important case corresponding to data storage in DNA.

I. INTRODUCTION

Recent advances in DNA synthesis and sequencing technologies [2] have made DNA a promising candidate for rising data storage needs. Compared to traditional storage media, DNA storage has several advantages, including higher data density, longevity, and ease of generating copies [2]. However, DNA is subject to a diverse set of errors that may occur during the various stages of data storage and retrieval, including substitutions, duplications, insertions, and deletions. This poses a challenge to the design of error-correcting codes and has led to many recent works studying the subject, including [2]–[16]. The current paper focuses on correcting short duplication and edit errors, where an edit is a substitution, insertion, or deletion.

A (tandem) duplication error generates a copy of a substring of the DNA sequence and inserts it after the original substring [3]. For example, from ACGT we may obtain ACGCGT. The length of the duplication is the length of the substring being copied, which is 2 in the preceding example. In the literature, both fixed-length duplication [3]–[6] and bounded-length duplication, where the duplication length is bounded from above [3], [17]–[19] have been studied. For duplications whose length is at most

This work was supported in part by NSF grants under grant nos. 1816409 and 1755773. This paper was presented in part at the 2020 IEEE Symposium of Information Theory (ISIT) in 2020 [1].

Yuanyuan Tang is with the Department of Electrical and Computer Engineering, University of Virginia, Charlottesville, VA, 22903, USA, (email: yt5tz@virginia.edu).

Farzad Farnoud (Hassanzadeh) is with the Department of Electrical and Computer Engineering and the Department of Computer Science, University of Virginia, Charlottesville, VA, 22903, USA, (email: farzad@virginia.edu).

Copyright (c) 2020 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

3, the case most relevant to this paper, Jain et al. [3] proposed error-correcting codes that were shown to have an asymptotically optimal rate by Kovačević [18].

In an edit event, a symbol in the sequence is substituted or deleted, or a new symbol is inserted. Among these, substitution errors have been studied in the literature in conjunction with *fixed-length* duplication errors, including substitution errors that are restricted to the inserted copies, reflecting the noisiness of the copying mechanism during the duplication process [20]–[22], and substitution errors that may occur anywhere in the string [6].

We focus on correcting errors that may arise from channels with many duplication errors of length at most 3, which we refer to as *short duplications*, and one edit error, which may occur in any position in the string. Considering a single edit error reveals important insights into the interactions between edit and duplication errors and will be of use for studying the general case of t edit errors. As a simple example of this channel, the input ACG may become $ACG \rightarrow AC\underline{CC}G \rightarrow AC\underline{T}CG \rightarrow AC\underline{T}ACTACTCG \rightarrow ACT\underline{CT}ACTACTCG$, where the duplication copies are marked with underlines, and the occurrences of the symbol T result from copies of the substitution $C \rightarrow T$. Given that an arbitrary number of duplications are possible, an unbounded segment of the output word may be affected by the errors, and, for example, the substituted symbol may appear many times. However, relying on the fact that short tandem duplications lead to regular languages, we show that with an appropriate construction and preprocessing of the output of the channel, the deleterious effects of the errors may be localized. We leverage constrained coding and maximum distance separable codes to design codes for correcting the resulting errors, establish a lower bound on the code rate, and provide an asymptotic analysis that shows that the code has rate at least $\log(q-2)$, where q is the size of the alphabet and the log is in base 2. We note that the rate of the code correcting only short duplications is upper bounded by $\log(q-1)$. When $q = 4$, the case corresponding to DNA storage, we provide a computational bound for the code rate, showing that asymptotically its rate is only 0.003 bits/symbol smaller than the code that corrects short duplications but no edits.

We will first consider only substitution edits and construct error-correcting codes capable of correcting many short duplications and a substitution. We will then prove that the same code can correct any number of duplications and an edit error by transforming insertion and deletion errors to substitution errors.

The paper is organized as follows. In Section II, we provide the notation and relevant background. Section III analyzes the error patterns that result from passing through duplication and substitution channels. After that, the code construction as well as the code size are presented in Section IV. In particular, in Subsection IV-A, we construct codes that can correct any number of duplications and one substitution and in Subsection IV-B, we show that the same codes can correct duplication and edit errors. Finally, Section V presents our concluding remarks.

II. NOTATION AND PRELIMINARIES

Let $\Sigma_q = \{0, 1, \dots, q-1\}$ denote a finite alphabet of size q . To avoid trivial cases, we assume $q \geq 3$, which in particular includes the case of $q = 4$, relevant to DNA data storage. The set of all strings of finite length over Σ_q is denoted by Σ_q^* , while Σ_q^n represents the strings of length n . In particular, Σ_q^* contains the empty string Λ . Let $[n]$ denote the set $\{1, \dots, n\}$.

Strings over Σ_q are denoted by bold symbols, such as \mathbf{x} and \mathbf{y}_j , or by capital letters. The elements of strings are shown with plain typeface, e.g., $\mathbf{x} = x_1x_2 \cdots x_n$ and $\mathbf{y}_j = y_{j1}y_{j2} \cdots y_{jm}$, where $x_i, y_{ji} \in \Sigma_q$. Given two strings $\mathbf{x}, \mathbf{y} \in \Sigma_q^*$, \mathbf{xy} denotes their concatenation and \mathbf{x}^m denotes the concatenation of m copies of \mathbf{x} . We use $|\mathbf{x}|$ to denote the length of a word $\mathbf{x} \in \Sigma_q^*$. For four words $\mathbf{x}, \mathbf{u}, \mathbf{v}, \mathbf{w} \in \Sigma_q^*$, if \mathbf{x} can be expressed as $\mathbf{x} = \mathbf{uvw}$, then \mathbf{v} is a *substring* of \mathbf{x} . Furthermore, \mathbf{u} is a prefix of \mathbf{x} and \mathbf{w} is a suffix of \mathbf{x} . For a string \mathbf{x} , two substrings \mathbf{u} and \mathbf{v} are said to *overlap* if we can write $\mathbf{x} = \mathbf{abcde}$, where $\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}, \mathbf{e}$ are nonempty, $\mathbf{u} = \mathbf{bc}$, and $\mathbf{v} = \mathbf{cd}$.

Given a word $\mathbf{x} \in \Sigma_q^*$, a *tandem duplication* (TD) of length k copies a substring of length k and inserts it after the original. This is referred to as a k -TD. For example, a 2-TD may generate $\mathbf{abc}b\mathbf{cde}$ from \mathbf{abcde} . Here, $\mathbf{bc}b\mathbf{c}$ is called a (*tandem*) *repeat* of length 2. Our focus in this paper is on TDs of length bounded by k , denoted $\leq k$ -TD, for $k = 3$. For example, given $\mathbf{x} = 1201210$ we may obtain via ≤ 3 -TDs

$$\begin{aligned} \mathbf{x} = 1201210 &\rightarrow 1201\underline{201}210 \rightarrow \\ &120120\underline{201}210 \rightarrow 1201202\underline{201}210 = \mathbf{x}', \end{aligned} \quad (1)$$

where the underlined substrings are the inserted copies. We say that \mathbf{x}' is a *descendant* of \mathbf{x} , i.e., a sequence resulting from \mathbf{x} through a sequence of duplications.

Let $\text{Irr}_{\leq k}(n) \subseteq \Sigma_q^n$ denote the set of *irreducible strings* (more precisely, $\leq k$ -irreducible strings) of length n , i.e., strings without repeats of length at most k . We use $\text{Irr}_{\leq k}(\cdot)$ denotes $\leq k$ -irreducible strings of arbitrary lengths. Furthermore, let $D_{\leq k}^*(\mathbf{x})$ denote the *descendant cone* of \mathbf{x} , containing all the descendants of \mathbf{x} after an arbitrary number of $\leq k$ -TDs.

Given a string \mathbf{x}' , let

$$R_{\leq k}(\mathbf{x}') = \{\mathbf{x} \in \text{Irr}_{\leq k}(\cdot) \mid \mathbf{x}' \in D_{\leq k}^*(\mathbf{x})\}$$

denote the set of *duplication roots* of \mathbf{x}' , i.e., irreducible sequences of which \mathbf{x}' is a descendant. Note that \mathbf{x} is the root of itself if $\mathbf{x} \in \text{Irr}_{\leq k}(\cdot)$. For a set S of strings, $R_{\leq k}(S)$ is the set of strings each of which is a root of at least one string in S . If $R_{\leq k}(\cdot)$ is a singleton, we may view it as a string rather than a set. A root can be obtained from \mathbf{x} by repeatedly replacing all repeats of the form \mathbf{aa} with \mathbf{a} , where $|\mathbf{a}| \leq k$ (each such operation is called a *deduplication*). For ≤ 3 -TDs, the duplication root is unique [3], i.e., $|R_{\leq 3}(\cdot)| = 1$. If \mathbf{x}' is a descendant of \mathbf{x} , we have $R_{\leq 3}(\mathbf{x}) = R_{\leq 3}(\mathbf{x}')$. For $k = 3$, we may drop the ≤ 3 subscript from the notation and write $D^*(\cdot), R(\cdot), \text{Irr}(\cdot)$.

We also consider substitution errors, although our attention is limited to at most one error of this kind. Continuing the example given in (1), a substitution occurring in the descendant \mathbf{x}' of \mathbf{x} may result in \mathbf{x}'' :

$$\mathbf{x}' = 1201202201210 \rightarrow \mathbf{x}'' = 1201202\underline{1}01210, \quad (2)$$

where the substituted symbol is underlined.

We denote by $D_{\leq k}^{t,p}(\mathbf{x})$ the set of strings that can be obtained from \mathbf{x} through t TDs of length at most k and p substitutions, *in any order*. We note that substitutions are *unrestricted* in the sense that they may occur in any position in the string, unlike the *noisy duplication* setting, where they are restricted to the inserted copies [6], [20]. Replacing t with $*$ denotes any number

of $\leq k$ -TDs and replacing p with $\leq p$ denotes at most p substitutions. We again drop $\leq k$ from the notation when $k = 3$. In the example given in (1) and (2), we have $x'' \in D^{*,1}(x)$, denoting that x'' is a descendant generated from x after an arbitrary number of ≤ 3 -TDs and a substitution error.

III. CHANNELS WITH MANY ≤ 3 -TDs AND ONE SUBSTITUTION ERROR

In this section, we study channels that alter the input string by applying an arbitrary number of duplication errors and at most one substitution error, where the substitution may occur at any time in the sequence of errors. We will first study the conditions a code must satisfy to be able to correct such errors. Then, we will investigate the effect of such channels on the duplication root of sequences, which is an important aspect of designing our error-correcting codes.

A code C is able to correct an arbitrary number of ≤ 3 -TDs and a substitution if and only if for any two distinct codewords $c_1, c_2 \in C$, we have

$$D^{*,\leq 1}(c_1) \cap D^{*,\leq 1}(c_2) = \emptyset.$$

To satisfy this condition, it is sufficient to have

$$R(D^{*,\leq 1}(c_1)) \cap R(D^{*,\leq 1}(c_2)) = \emptyset. \quad (3)$$

Condition (3) implies that for distinct codewords c_1 and c_2 , $R(c_1) \neq R(c_2)$. Since both $R(c_1)$ and $R(c_2)$ are singleton, this latter condition is in fact sufficient for correcting only ≤ 3 -TDs since this type of error does not alter the duplication root. For correcting only ≤ 3 -TDs, defining the code as the set of irreducible strings of a given length leads to asymptotically optimal codes [3], [18]. The decoding process is simply finding the root of the received word.

We take a similar approach to correct many ≤ 3 -TDs and a substitution. More specifically, the proposed code C is a subset of ≤ 3 -irreducible strings, i.e., $R(c) = c$ for $c \in C$. To recover c from the received word y , we find $R(y)$ and from that recover $R(c) = c$, as will be discussed.

We start by studying the effect of ≤ 3 -TDs and one substitution on the root of a string. Specifically, for strings x and $x'' \in D^{*,\leq 1}(x)$, it is of interest to determine how $R(x'')$ differs from $R(x)$. We either have $x'' \in D^*(x)$, i.e., x'' suffers only duplications, or $x'' \in D^{*,1}(x)$. In the former case $R(x'') = R(x)$. Hence, below we consider only $x'' \in D^{*,1}(x)$. Note that duplications that occur after the substitution do not affect the root and so in our analysis we may assume that the substitution is the last error. We start by providing a useful definition and an auxiliary lemma and then a lemma that considers a simple case.

Let s and \bar{s} be strings of length n , and let A be the set of symbols in s and \bar{A} the set of symbols in \bar{s} . We say that s *dominates* \bar{s} if there exists a function $f : A \rightarrow \bar{A}$ such that $\bar{s} = f(s)$, where $f(s) = f(s_1) \cdots f(s_n)$. For example, 0102 dominates 1212 (using the mapping $f(0) = 1, f(1) = 2, f(2) = 2$) but 0102 does not dominate 0010. The string $012 \cdots k$ dominates any string of length $k + 1$.

The following lemma helps reduce the number of cases we need to consider using the dominance relationship defined above.

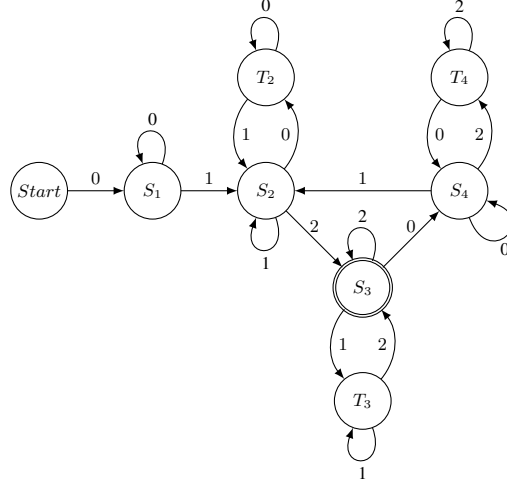


Figure 1: Finite automaton for the regular language $D^*(012)$ based on [17].

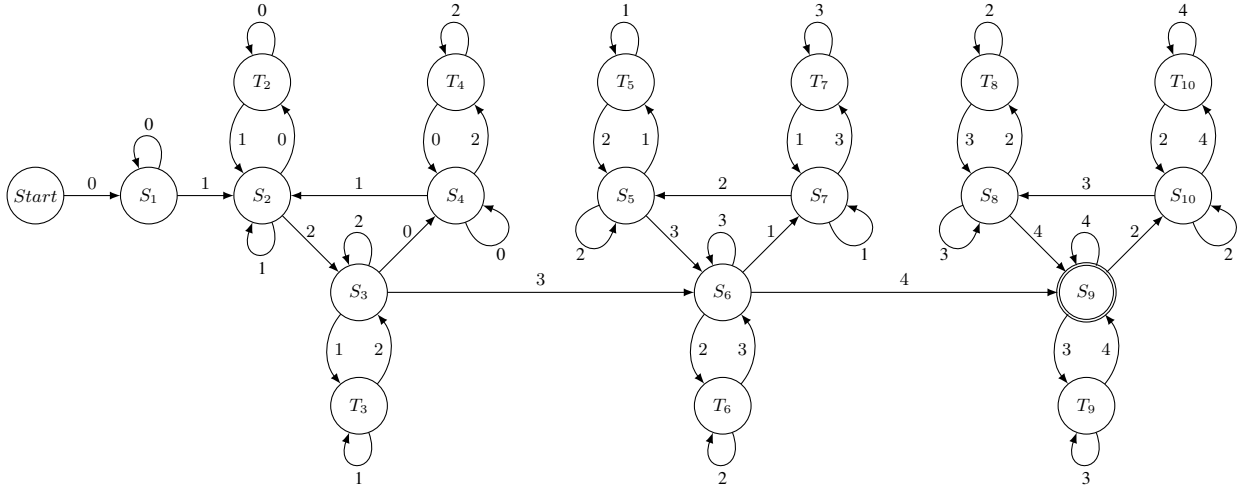


Figure 2: Finite automaton for the regular language $D^*(01234)$ based on [17].

Lemma 1. Suppose s dominates \bar{s} . The following hold:

- i) Suppose we apply the same duplication in both s and \bar{s} (that is, in the same position and with the same length). Let the resulting strings be s' and \bar{s}' , respectively. Then s' dominates \bar{s}' .
- ii) If a deduplication is possible in s , a deduplication in the same position and with the same length is possible in \bar{s} . Let the result of applying this deduplication to s and \bar{s} be denoted by s' and \bar{s}' , respectively. Then s' dominates \bar{s}' .
- iii) Let \bar{s}' be obtained from \bar{s} via a substitution in position i and let s' be obtained from s by substituting the symbol in position i with a symbol x not present in s . Then, s' dominates \bar{s}' .
- iv) We have $|R(\bar{s})| \leq |R(s)|$.

Before proving the lemma, we provide an example for each statement below, where duplicated, deduplicated, and substituted symbols are underlined. For i) consider

$$s = 0102 \rightarrow s' = 010\underline{1}02$$

$$\bar{s} = 0101 \rightarrow \bar{s}' = 010\underline{1}01,$$

for ii) consider

$$s = 0102\underline{1}021 \rightarrow s' = 01021$$

$$\bar{s} = 0101\underline{1}011 \rightarrow \bar{s}' = 01011,$$

for iii) consider

$$s = 01021021 \rightarrow s' = 0102\underline{3}021$$

$$\bar{s} = 01011011 \rightarrow \bar{s}' = 0101\underline{0}011,$$

and for iv) consider

$$s = 0102\underline{1}021 \rightarrow 01021 = R(s)$$

$$\bar{s} = 0101\underline{1}011 \rightarrow 010\underline{1}1 \rightarrow 01\underline{1} \rightarrow 01 = R(\bar{s}).$$

Proof: Let f be a function that can show s dominates \bar{s} . i) For the first statement, the same mapping f also shows that s' dominates \bar{s}' . ii) For the second statement, consider a repeat aa in s . Then the repeat $f(a)f(a)$ is present in \bar{s} in the same position. So the deduplication is possible in \bar{s} . The same mapping f proves that s' dominates \bar{s}' . iii) For the third statement, let the substitution in \bar{s} alter the symbol in position i to some symbol a . If we extend f by mapping x to a , then f proves that s' dominates \bar{s}' . iv) From ii), any sequence of deduplications applied to s can also be applied to $\bar{s} = f(s)$. In particular, the sequence of deduplications that takes s to its root $R(s)$ takes $\bar{s} = f(s)$ to $f(R(s))$. The root of \bar{s} can be obtained by removing any remaining repeats in $f(R(s))$ (recall that the root is unique so all sequences of deduplications must lead to the same sequence). Hence $|R(\bar{s})| \leq |f(R(s))|$. Noting $|f(R(s))| = |R(s)|$ completes the proof. ■

Lemma 2. For any alphabet Σ_q ,

$$\max_{x \in \Sigma_q^3} \max_{x'' \in D^{*,1}(x)} |R(x'')| = 13, \quad (4)$$

$$\max_{x \in \Sigma_q^5} \max_{x'' \in D^{*,1}(x)} |R(x'')| \leq 17. \quad (5)$$

Proof: For the first statement, it suffices to consider only $x = 012$ and assume that the substitution that leads to x'' replaces a symbol in x with some symbol other than 0, 1, and 2, e.g., 3. To see this, consider any string \bar{x} of length 3 over any alphabet. The string \bar{x} is dominated by x . Now consider any $\bar{x}'' \in D^{*,1}(\bar{x})$. There is a sequence of “errors” consisting of duplications, a substitution, and more duplications that transforms \bar{x} to \bar{x}'' . By Lemma 1, i) and iii), there is a corresponding sequence of errors, consisting of duplications, a substitution, and duplications, that when applied to x will result in x'' , where x'' dominates \bar{x}'' (the substitution in the sequence of errors for x substitutes the existing symbol with a symbol not in the set $\{0, 1, 2\}$). Then by Lemma 1, iv), we have $|R(\bar{x}'')| \leq |R(x'')|$. Since this is true for any choice of \bar{x} and any $\bar{x}'' \in D^{*,1}(\bar{x})$,

Table I: Paths representing irreducible strings starting from and ending at specific states.

State	Irreducible paths from ‘Start’ to state	Irreducible paths from State to S_3
S_1	0	012, 1012, 12, 12012,
S_2	01, 01201	012,1012, 12, 12012, 2, 2012, 212, 212012
S_3	012	012, 02012, 12, 12012, 2, 2012, 212, 212012
S_4	0120	012, 02012, 1012, 12, 12012, 2012
T_2	010, 012010	012, 1012,12, 12012
T_3	0121	12, 12012, 2, 2012, 212, 212012
T_4	01202	012, 02012, 2012

it suffices to find

$$\max_{\mathbf{x}'' \in D^{*,1}(012)} |R(\mathbf{x}'')|,$$

where the substitution resulting in \mathbf{x}'' replaces the existing symbol with a symbol not present in $\mathbf{x} = 012$. Henceforth, we assume $\mathbf{x} = 012$.

As shown in [17], $D^*(\mathbf{x})$ is a regular language whose words can be described as paths from ‘Start’ to S_3 in the finite automaton given in Figure 1, where the word associated with each path is the sequence of the edge labels. Let $\mathbf{x}' \in D^*(\mathbf{x})$ and $\mathbf{x}'' \in D^{0,1}(\mathbf{x}')$. Assume $\mathbf{x}' = \mathbf{u}\mathbf{w}\mathbf{z}$ and $\mathbf{x}'' = \mathbf{u}\hat{\mathbf{w}}\mathbf{z}$, where \mathbf{u}, \mathbf{z} are strings and w and $\hat{w} \notin \{0, 1, 2\}$ are distinct symbols. The string \mathbf{u} represents a path from ‘Start’ to some state U and the string \mathbf{z} represents a path from some state Z to S_3 in the automaton, where there is an edge with label w from U to Z .

Since $R(\mathbf{x}'') = R(R(\mathbf{u})\hat{w}R(\mathbf{z}))$, we have $|R(\mathbf{x}'')| \leq |R(\mathbf{u})| + 1 + |R(\mathbf{z})|$ (recall that $R(s)$ is a singleton for a string s). The maximum value for $|R(\mathbf{u})|$ is the length of some path from ‘Start’ to U such that the corresponding sequence does not have any repeats (henceforth, called an *irreducible path*). All such paths/sequences are listed in the second column of Table I for all choices of U . Similarly, the maximum value for $|R(\mathbf{z})|$ is the length of some irreducible path from Z to S_3 ; all such possibilities are listed in the third column of Table I. An inspection of Table I shows that choosing $U = T_2$ and $Z = S_2$ leads to the largest value of $|R(\mathbf{u})| + 1 + |R(\mathbf{z})|$, namely $6 + 1 + 6 = 13$. We note that the specific sequence achieving this length is $\mathbf{x}'' = 0120103212012$, which can be obtained via the sequence $\mathbf{x} \rightarrow 012\underline{012}\underline{012} \rightarrow 01201\underline{012}012 \rightarrow 01201012\underline{12}012 \rightarrow 012010\underline{3}212012 = \mathbf{x}''$ with a substitution $1 \rightarrow \mathbf{3}$ in the last step, where we have combined non-overlapping duplications into a single step.

Let us now prove the second statement. Again we need only consider $\mathbf{x} = 01234$, for which $D^*(\mathbf{x})$ is the regular language whose automaton is shown in Figure 2. In a similar manner to the proof of the previous part, we can show that the length of the longest irreducible path from ‘Start’ to any state in the automaton is at most 8 and the length of the longest irreducible path from any state to S_9 is also at most 8. Hence, $|R(\mathbf{x}'')| \leq 8 + 1 + 8 = 17$, completing the proof. ■

We now consider changes to the roots of arbitrary strings when passed through a channel with arbitrarily many ≤ 3 -TDs and one substitution. The next lemma is used in the main result of this section, Theorem 5, which shows that even though a substituted symbol may be duplicated many times, the effect of a substitution on the root is bounded.

Lemma 3. Let x be any string of length at least 5 and $x' \in D^*(x)$. For any decomposition of x as

$$x = r \, ab \, t \, de \, s,$$

for $a, b, d, e \in \Sigma_q$ and $r, t, s \in \Sigma_q^*$, with t nonempty, there is a decomposition of x' as

$$x' = u \, ab \, w \, de \, v$$

such that $u, w, v \in \Sigma_q^*$, $uab \in D^*(rab)$, $abwde \in D^*(abtde)$, and $dev \in D^*(des)$.

Proof: If $x = x'$, the claim is true since we may choose $u = r, w = t, v = s$. It suffices to consider the case in which x' is obtained from x via a single duplication. The case of more duplications can be proved inductively.

First suppose the length of the duplication transforming x to x' is 1. If this duplication occurs in r , we choose u to be the descendant of r and let $w = t$ and $v = s$, satisfying the claim. Duplication of a single symbol in t or s is handled similarly. If a is duplicated, we let $u = ra, w = t, v = s$. If b is duplicated, we let $u = r, w = bt, v = s$. The cases for d and e are similar.

Second, consider a duplication of length 2 or 3. Such a duplication is fully contained in rab , $abtde$, or des . A duplication of length 2 or 3 applied to a string z does not alter the first two and the last two symbols of z . So, for example, if the duplication occurs in rab , then we can choose u such that $uab \in D^1(rab)$ and let $w = t$ and $v = s$. The cases of duplications contained in the other strings are similar. ■

We now provide an example in which we illustrate how the root of a string can be altered by several duplications and one substitution.

Example 4. Fix $\Sigma_4 = \{0, 1, 2, 3\}$ as the alphabet. In the following examples, x is an irreducible string, $x' \in D^*(x)$, and $x'' \in D^{0,1}(x')$. We compare $R(x) = x$ with $R(x'')$. In particular, we will decompose $R(x)$ and $R(x'')$ as $R(x) = \alpha\beta\gamma$ and $R(x'') = \alpha\beta'\gamma$. In other words, $R(x'')$ can be obtained from $R(x)$ by deleting β and inserting β' .

- Let $x = 012302$, $x' = 011\underline{2012012}30202$, and $x'' = 011\underline{2013012}30202$, where the underlined symbols result from duplication and the bold symbol from substitution. Then $R(x'') = 012013012302$ and the change from $R(x)$ to $R(x'')$ can be viewed as

$$R(x) = \underbrace{012}_{\alpha} \underbrace{302}_{\gamma} \rightarrow R(x'') = \underbrace{012}_{\alpha} \underbrace{013012}_{\beta'} \underbrace{302}_{\gamma},$$

with $\beta = \Lambda$.

- Let $x = 13203103$, $x' = 1313213203103\underline{103}$, and $x'' = 1313213103103\underline{103}$. Then $R(x'') = 13213103$ and the change from $R(x)$ to $R(x'')$ can be viewed as

$$R(x) = \underbrace{132}_{\alpha} \underbrace{0}_{\beta} \underbrace{3103}_{\gamma} \rightarrow R(x'') = \underbrace{132}_{\alpha} \underbrace{1}_{\beta'} \underbrace{3103}_{\gamma}.$$

- Let $x = 012010321201230$, $x' = 01201201032120\underline{2012012}30$, and $x'' = 01201201012120\underline{2012012}30$. Then $R(x'') =$

01230 and the change from $R(x)$ to $R(x'')$ can be viewed as

$$R(x) = \underbrace{012}_{\alpha} \underbrace{0103212012}_{\beta} \underbrace{30}_{\gamma} \rightarrow R(x'') = \underbrace{012}_{\alpha} \underbrace{30}_{\gamma}, \quad (6)$$

with $\beta' = \Lambda$.

Let \mathcal{L} be the smallest integer, if it exists, such that for any alphabet Σ_q , any $x \in \Sigma_q^*$, and any $x'' \in D^{*,1}(x)$, we can obtain $R(x'')$ from $R(x)$ by deleting a substring of length at most \mathcal{L} and inserting a substring of length at most \mathcal{L} in the same position. The example given in (6) shows that \mathcal{L} , if it exists, satisfies $\mathcal{L} \geq 10$. We note however that the definition does not guarantee that \mathcal{L} exists as we may be able to produce examples in which the length of the deleted or the inserted substring is arbitrarily long. The next theorem shows that such examples cannot be constructed by providing an explicit upper bound on \mathcal{L} .

Theorem 5. \mathcal{L} exists (i.e., it is finite). Moreover, $\mathcal{L} \leq 17$.

Proof: We may assume x is irreducible. If it is not, let $x_0 = R(x)$ so that $x'' \in D^{*,1}(x) \subseteq D^{*,1}(x_0)$. If the statement of the theorem holds for x_0 , it also holds for x since $R(x) = R(x_0)$.

We will find $\alpha, \beta, \beta', \gamma \in \Sigma_q^*$ with $R(x) = \alpha\beta\gamma$ and $R(x'') = \alpha\beta'\gamma$ such that $|\beta'| \leq 17$. By symmetry, it suffices to prove $|\beta'| \leq 17$ for all irreducible x . To see this symmetry, note that $\alpha\beta'\gamma$ is obtained from $\alpha\beta\gamma$ by applying, in order, duplications, a single substitution, more duplications, and finally removing all repeats (performing all possible deduplications). Recall that for ≤ 3 -TDs, the root is unique and regardless of the order in which deduplications are applied, we will arrive at the same root. In other words, applying a sequence of duplications to a string s and then removing all repeats is equivalent to removing all repeats from s . Hence, we may instead assume that the process transforming $\alpha\beta\gamma$ to $\alpha\beta'\gamma$ is as follows: duplications, substitution, deduplications. Since this process is reversible, general statements that hold for β' also hold for β .

Let $x' \in D^*(x)$ be obtained from x through duplications and x'' be obtained from x' through a substitution. We assume that $x = rabcde s$, where $r, s \in \Sigma_q^*$ and $a, b, c, d, e \in \Sigma_q$, such that the substituted symbol in x' is a copy of c . Note that if $|x| < 5$ or if a copy of one of its first two symbols or its last two symbols are substituted, then we can no longer write x as described. To avoid considering these cases separately, we may append two dummy symbols to the beginning of x and two dummy symbols to the end of x , where the four dummy symbols are distinct and do not belong to Σ_q , and prove the result for this new string. Since these dummy symbols do not participate in any duplication, substitution, or deduplication events, the proof is also valid for the original x .

With the above assumption and based on Lemma 3, we can write

$$\begin{aligned} x &= r \ ab \ c \ de \ s \\ x' &= u \ ab \ w \ de \ v \in D^*(x), \\ x'' &= u \ ab \ z \ de \ v \in D^{0,1}(x'), \end{aligned} \quad (7)$$

where $uab \in D^*(rab)$, $abwde \in D^*(abcde)$, $dev \in D^*(des)$, and z is obtained from w by substituting an occurrence of c . From (7), $R(x'') = R(rR(abzde)s)$, where $R(abzde)$ starts with ab and ends with de (which may fully or partially

overlap). The outer R in $R(rR(abzde)s)$ may remove some symbols at the end of r , beginning and end of $R(abzde)$, and the beginning of s , leading to $\alpha\beta'\gamma$, where α is a prefix of r , β' is a substring of $R(abzde)$, and γ is a suffix of s . Hence, $|\beta'| \leq |R(abzde)|$. But $abzde \in D^{*,1}(abcde)$ and thus by Lemma 2, $|R(abzde)| \leq 17$, completing the proof. ■

IV. ERROR-CORRECTING CODES

Having studied how duplication roots are affected by tandem duplication and substitution errors, in Subsection IV-A, we construct codes that can correct such errors. In Subsection IV-B, we show that the same codes can correct duplication and edit errors. We will also determine the rate of these codes and compare it with the rate of codes that only correct duplications, which provides an upper bound.

A. Code constructions for correcting duplications and a substitution error

As noted in the previous section, the effect of a substitution error on the root of the stored codeword is local in the sense that a substring of bounded length may be deleted and another substring of bounded length may be inserted in its position. A natural approach to correcting such errors is to divide the codewords into blocks such that this alteration can affect a limited number of blocks. In particular, we divide the string into *message blocks* that are separated by *marker blocks* known to the decoder. We start with an auxiliary construction.

Construction 6. Let l, m, N be positive integers with $m > l$ and $\sigma \in \text{Irr}(l)$. The code C_σ of length $n = N(m + l) - l$ over Σ_q consists of irreducible strings x obtained by alternating between message blocks of length m and copies of the marker sequence σ , i.e.,

$$x = B_1\sigma B_2\sigma \cdots \sigma B_N,$$

such that $x \in \text{Irr}(N(m + l) - l)$, $B_i \in \text{Irr}(m) \subseteq \Sigma_q^m$, $i \in [N]$, and there are exactly two occurrences of σ in $\sigma B_i\sigma$, for all $i \in [N]$. (Thus, there are precisely $N - 1$ occurrences of σ in x .)

We remark that for our purposes, we can relax the condition on $\sigma B_i\sigma$ for $i = 1, N$. Specifically, it suffices to have exactly one occurrence of σ in $B_1\sigma$ and one occurrence of σ in σB_N . For simplicity however, we do not use these relaxed conditions.

Example 7. Let $m = 6$, $N = 5$, and $\sigma = 01231$ with $l = 5$. Then the code C_σ in Construction 6 will contain the codeword

$$x = \textcolor{red}{012013}01231\textcolor{red}{030121}01231\textcolor{red}{202312}01231\textcolor{red}{303203}01231\textcolor{red}{203023},$$

where the message blocks $B_1 = 012013$, $B_2 = 030121$, $B_3 = 202312$, $B_4 = 303203$, $B_5 = 203023$ are marked in red. Furthermore, $\sigma B_i\sigma \in \text{Irr}(16)$ for $i \in [5]$, which, as will be shown in Lemma 9 below, implies that the codeword $x \in \Sigma_4^{50}$ is irreducible.

Given an input x with N message blocks, let y be the root of the output after duplications and at most one substitution. We define a *block* in y as a maximal substring that does not overlap with any separator sequence σ . Note that a block in y may or may not be an error-free message block from x .

With Construction 6 in hand, in the next theorem, we show that the effect of one substitution and many tandem duplications is limited to a small number of message blocks. We note that \mathcal{L} appearing in the theorem below was defined before Theorem 5 and satisfies $\mathcal{L} \leq 17$.

Theorem 8. *Let \mathcal{C}_σ be the code defined in Construction 6. If $m > \mathcal{L}$, then there exists a decoder \mathcal{D}_σ that, for any $x \in \mathcal{C}_\sigma$ and $y \in R(D^{*, \leq 1}(x))$, outputs $z = \mathcal{D}_\sigma(y)$ such that, relative to x , either two of the message blocks B_i are substituted in z or four of them are erased.*

Proof: Let $x = \alpha\beta\gamma$ and $y = \alpha\beta'\gamma$, where by Theorem 5, $|\beta|, |\beta'| \leq \mathcal{L}$. To avoid a separate treatment for blocks B_1 and B_N , the decoder appends σ to the beginning and the end of y and assumes that the codewords are of the form $\sigma B_1 \sigma \cdots \sigma B_N \sigma$. The decoder considers two cases depending on whether the marker sequences σ are in the same positions in y as in the codewords in \mathcal{C}_σ . If the markers are in the same positions, as shown in Figure 3, then $|x| = |y|$, and consequently, $|\beta| = |\beta'| \leq \mathcal{L}$. Since $\mathcal{L} < m = |B_i|$, at most two (adjacent) blocks B_i are affected by substituting β by β' and thus $z = y$ differs from x in at most two block.

On the other hand, if the markers are in different positions in y compared to the codewords in \mathcal{C}_σ , the decoder uses the location of the markers to identify the position of the message blocks that may be affected and erases them, as described below. By the definition of blocks in y , since the markers are in different positions in x and y , there is at least one block B in y whose length differs from m . Hence, y has a substring u of length $m + 2l$ that starts with σ and contains part or all of B but does not end with σ . Two examples where such a situation may arise are shown in Figure 4. On the left, $|\beta| = |\beta'|$ and a marker is absent from y due to substituting β with β . On the right, β and β' have different lengths, and this causes the markers to move.

Let $\delta = |x| - |y| = |\beta| - |\beta'|$ and $\delta^+ = \max(0, \delta)$. Note that $|\beta'| = |\beta| - \delta \leq \mathcal{L} - \delta$. Furthermore, $|\beta'| \leq \mathcal{L}$ and so $|\beta'| \leq \min(\mathcal{L}, \mathcal{L} - \delta) = \mathcal{L} - \delta^+$. Let y' be obtained by removing u along with $\mathcal{L} - \delta^+ - 1$ elements from each of its sides from y . This removes β' from y . More formally, we claim y' can be obtained via a deletion from x . First, suppose $|\beta'| = 0$. Then $y = \alpha\gamma$. Note that u is not a substring of x since every substring of x that has length $m + 2l$ and starts with σ also ends with σ . Hence, it must overlap with both α and γ . After deleting u from y , we obtain a string $y' = \alpha'\gamma'$ such that α' is a prefix of α and γ' is a suffix of γ , proving the claim. Next, suppose that $|\beta'| > 0$ and recall that $y = \alpha\beta'\gamma$. Since u is not a substring of x , at least one of the following holds: i) u overlaps with both α and β' or ii) u overlaps with both β' and γ . Case i) is shown in Figure 4. In either case, the substring of y consisting of u and the $\mathcal{L} - \delta^+ - 1$ elements on each side of y contains β' , proving the claim. Hence, y' is a sequence that relative to x suffers a deletion of length at most $m + 2l + 2\mathcal{L} - 2\delta^+ - 2 + |\beta| - |\beta'| < 3m + 2l$ from a known position. The deletion affects at most 4 message blocks and since its location is known, the decoder can mark these message blocks as erased. ■

In Construction 6, the constraint that x must be irreducible creates interdependence between the message blocks, making the code more complex. The following lemma allows us to treat each message block independently provided that σ is sufficiently long.



Figure 3: If marker sequences, shown as gray, are in the same positions in the codeword x and the retrieved string y , then β and β' have the same length and at most two of the message blocks are affected by the errors, as discussed in the proof of Theorem 8.

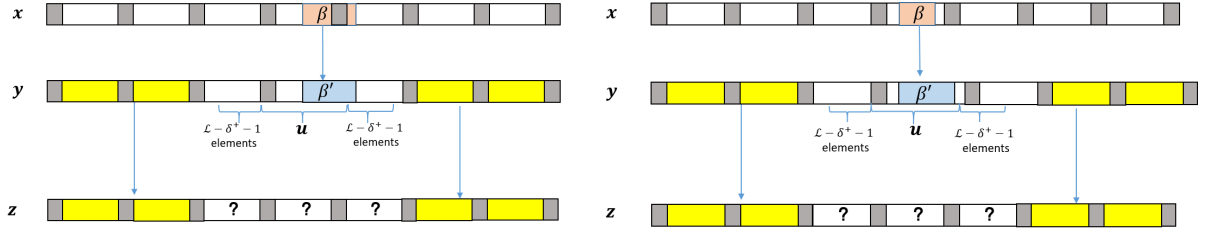


Figure 4: If marker sequences, shown as gray, are in different positions in the codeword x and the retrieved string y , then a substring u is identified and then expanded to ensure it contains β' . Those blocks in y that intersect with this expanded substring are marked as erasures while other blocks are error-free message blocks, as described in the proof of Theorem 8).

Lemma 9. Let x be as defined in Construction 6 and assume $l \geq 5$. The condition $x \in \text{Irr}(N(m+l) - l)$ is satisfied if

$$\sigma B_i \sigma \in \text{Irr}(m+2l), \quad \text{for all } i \in [N]. \quad (8)$$

Proof: Suppose that x has a repeat aa , with $|a| \leq 3$. Since $|aa| \leq 6$ and $|\sigma| \geq 5$, there is no i such that the repeat lies in $B_i \sigma B_{i+1}$ and overlaps both B_i and B_{i+1} . So it must be fully contained in $B_1 \sigma$, σB_N , or $\sigma B_i \sigma$ for some $2 \leq i \leq N-1$, contradicting assumption (8). ■

We now present a code based on Construction 6 and prove that it can correct any number of tandem duplications and one substitution error.

Construction 10. Let l, m be positive integers with $m > l \geq 5$, and $\sigma \in \text{Irr}(l)$. Furthermore, let \mathcal{B}_σ^m denote the set of sequences B such that $\sigma B \sigma \in \text{Irr}(m+2l)$ has exactly two occurrences of σ , and $M = M_\sigma^{(m)} = |\mathcal{B}_\sigma^m|$. Finally, let t be a positive integer such that $2^t \leq M$ and $\zeta : \mathbb{F}_{2^t} \rightarrow \mathcal{B}_\sigma^m$ be an injective mapping. We define \mathcal{C}_{MDS} as

$$\mathcal{C}_{\text{MDS}} = \{\zeta(c_1)\sigma\zeta(c_2)\sigma \cdots \sigma\zeta(c_N) : c \in \text{MDS}(N, N-4, 5)\},$$

where $\text{MDS}(N, N-4, 5)$ denotes an MDS code over \mathbb{F}_{2^t} of length $N = 2^t - 1$, dimension $N-4$, and Hamming distance $d_H = 5$.

Note that the mapping ζ exists because $|\mathbb{F}_{2^t}| \leq |\mathcal{B}_\sigma^m|$ by the choice of t . For example, we can sort the elements of \mathbb{F}_{2^t} and \mathcal{B}_σ^m lexicographically and map the i th element of \mathbb{F}_{2^t} to the i th element of \mathcal{B}_σ^m .

Theorem 11. If $m > \mathcal{L}$, then the error-correcting code \mathcal{C}_{MDS} in Construction 10 can correct any number of ≤ 3 -TDs and at most one substitution error.

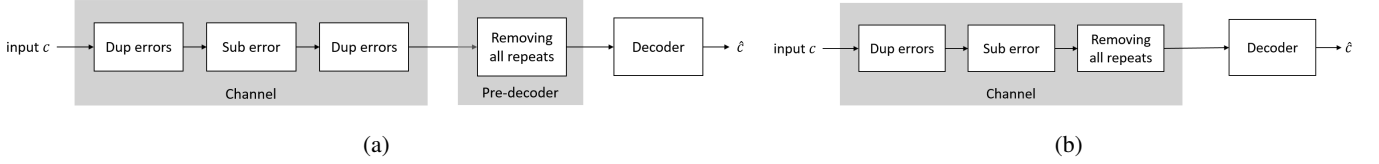


Figure 5: The duplication-substitution channel along with the decoder (a) and an equivalent representation of the end-to-end system (b).

Proof: Let the stored codeword be $\mathbf{x} = B_1\sigma \cdots \sigma B_N \in \mathcal{C}_{MDS}$, where $B_i = \zeta(c_i)$ for $i \in [N]$ and $\mathbf{c} \in C$, with C denoting an $MDS(N, N-4, 5)$ code. Based on the definitions of ζ and the set B_{σ}^m in Construction 10, the i th message block $B_i = \zeta(c_i)$ satisfies $\sigma B_i \sigma \in \text{Irr}(m+2l)$. Then, by Lemma 9, $\mathbf{x} \in \text{Irr}(N(m+l)-l)$ and so $\mathbf{x} \in \mathcal{C}_{\sigma}$. Therefore, $\mathcal{C}_{MDS} \subseteq \mathcal{C}_{\sigma}$. Suppose the retrieved word is \mathbf{y} . By Theorem 8, $\mathcal{D}_{\sigma}(\mathbf{y})$ suffers either at most two substitutions or at most four erasures of message blocks. Suppose the message block B_i is substituted by another string \mathbf{v} of length m . If $\zeta^{-1}(\mathbf{v})$ exists, this translates to a substitution of c_i . If not, we define $\zeta^{-1}(B_i)$ as an arbitrary element of \mathbb{F}_{2^t} , again leading to a possible substitution of c_i with another symbol. To decode, we can use the MDS decoder on $\zeta^{-1}(\mathcal{D}_{\sigma}(\mathbf{y}))$, which relative to \mathbf{c} suffers either ≤ 2 substitutions or ≤ 4 erasures. Given that the minimum Hamming distance of the MDS code is 5, the decoder can successfully recover \mathbf{c} . ■

B. Extension to edit errors

In this subsection, we extend Theorem 11 to include insertion and deletion errors in addition to substitution errors. We do so by showing an insertion can be viewed as a duplication plus a substitution and a deletion as a substitution plus a deduplication. The duplication-substitution channel discussed so far can be viewed as shown in Figure 5a, where in a pre-decoding step, the root of the retrieved string is found and then passed to the decoder. Recall that for ≤ 3 -TDs, applying a sequence of duplications to a string s and then removing all repeats is equivalent to removing all repeats from s . Hence, The process shown in Figure 5a is equivalent to the one shown in Figure 5b. The same equivalence holds if we replace the block representing a substitution error with a block representing an edit error. We can now prove the following corollary to Theorem 11.

Corollary 12. *If $m > \mathcal{L}$, then the error-correcting code \mathcal{C}_{MDS} in Construction 10 can correct any number of ≤ 3 -TDs and at most one edit error.*

Proof: Suppose that $\mathbf{c} \in \mathcal{C}_{MDS}$ suffers a sequence of errors consisting of duplications, an edit, and more duplications. Then all repeats are removed from the resulting string. This process is equivalent to applying the first set of duplications and the edit and then removing all the repeats. Denote this sequence of operations applied to \mathbf{c} by S and the final result by $S(\mathbf{c})$. We show that we can find a sequence S' consisting of duplications, at most one substitution, and removal of all repeats, such that $S(\mathbf{c}) = S'(\mathbf{c})$. If the edit in S is a substitution, then we let $S' = S$. If it is an insertion or a deletion, we again start by setting $S' = S$ and then modify S' as follows: i) If the edit is an insertion, we replace it in S' by a duplication and, if needed, a substitution. Namely, we duplicate the symbol before the insertion and then substitute the copy as needed. For example, if in S we have $abc \xrightarrow{ins} abxc$, we replace this step by $abc \xrightarrow{dup} abbc \xrightarrow{sub} abxc$, where *ins* stands for insertion, *dup* for duplication, and *sub* for substitution. The substitution is not necessary if $x = b$. ii) If the edit is a deletion, we replace it by a deduplication

that is preceded by a substitution if needed. Namely, we substitute the symbol that is deleted in S to be equal to the previous symbol and then deduplicate it. For example, if in S we have $abc \xrightarrow{\text{del}} ac$, we replace this step by $abc \xrightarrow{\text{sub}} aac \xrightarrow{\text{dedup}} ac$ in S' , where del stands for deletion and dedup for deduplication. If $b = a$, the substitution is not needed. Now, S' is a sequence consisting of duplications, at most one substitution, and then removal of all repeats, and furthermore $S(c) = S'(c)$. From Theorem 11 and the above discussion about the equivalence of Figures 5a and 5a, we find that the decoder can produce c given $S'(c)$. ■

C. Construction of message blocks

In this subsection, we study the set \mathcal{B}_σ^m of valid message blocks of length m with σ as the marker. Since in Construction 10, the markers σ do not contribute to the size of the code, to maximize the code rate, we set $l = |\sigma| = 5$, i.e., $\sigma \in \text{Irr}(5)$.

For a given σ , we need to find the set \mathcal{B}_σ^m . The first step in this direction is finding all irreducible sequences of length $m + 2l = m + 10$. We will then identify those that start and end with σ but contain no other σ s.

As shown in [3], the set of ≤ 3 -irreducible strings over an alphabet of size q is a regular language whose graph $G_q = (V_q, \xi_q)$ is a subgraph of the De Bruijn graph. The vertex set V_q consists of 5-tuples $a_1a_2a_3a_4a_5$ that do not have any repeats (of length at most 2). There is an edge from $a_1a_2a_3a_4a_5 \rightarrow a_2a_3a_4a_5a_6$ if $a_1a_2a_3a_4a_5a_6$ belongs to $\text{Irr}(6)$. The label for this edge is a_6 . The label for a path is the 5-tuple representing its starting vertex concatenated with the labels of the subsequent edges. In this way, the label of a path in this graph is an irreducible sequence and each irreducible sequence is the label of a unique path in the graph. The graph G_q , when $q = 3$, can be found in [3, Fig. 1].

The following theorem characterizes the set \mathcal{B}_σ^m and will be used in the next section to find the size of the code.

Theorem 13. *Over an alphabet of size q and for $\sigma \in \text{Irr}(5)$, there is a one-to-one correspondence between $B \in \mathcal{B}_\sigma^m$ and paths of length $m + 5$ in G_q that start and end in σ but do not visit σ in their interiors. Specifically, each sequence $B \in \mathcal{B}_\sigma^m$ corresponds to the path with the label $\sigma B \sigma$.*

Proof: Consider a path $p = v_1v_2 \cdots v_{k+1}$ where v_i are vertices of G_q and k is the length of the path. Denote the label of this path by $s = s_1s_2 \cdots s_{k+5}$. It can be shown by induction on k that $v_i = s_i s_{i+1} s_{i+2} s_{i+3} s_{i+4}$. Hence, the label of a path of length $m + 5$ that starts and ends in σ but does not visit σ otherwise is an irreducible sequence with exactly two occurrences of σ and is of the form $\sigma B \sigma$ where $B \in \mathcal{B}_\sigma^m$. Conversely, suppose $B \in \mathcal{B}_\sigma^m$. Then $\sigma B \sigma$ is an irreducible string of length $m + 10$ and thus the label of a unique path of length $m + 5$ in G_q . This path starts and ends in σ . But it does not visit σ in its interior since that would imply there are more than two occurrences of σ in $\sigma B \sigma$. ■

D. Code rate

We now turn to find the rate of the code introduced in this section. For a code C of length n and size $|C|$, the rate is defined as $R(C) = \frac{1}{n} \log |C|$. For the code of Construction 10,

$$R(\mathcal{C}_{\text{MDS}}) = \frac{N - 4}{Nm + (N - 1)l} \log(N + 1), \quad (9)$$

where N depends on the choice of $\sigma \in \text{Irr}(5)$. More specifically, $N \leq 2^{\lfloor \log M_{\sigma}^{(m)} \rfloor} - 1$. Choosing the largest permissible value for N implies that $N \geq (M_{\sigma}^{(m)} - 1)/2$ and

$$\begin{aligned} R(\mathcal{C}_{\text{MDS}}) &\geq \frac{1 - 4/N}{m + l} \log(N + 1) \\ &\geq \frac{1}{m + l} \left(1 - \frac{8}{M_{\sigma}^{(m)} - 1} \right) (\log M_{\sigma}^{(m)} - 1). \end{aligned} \quad (10)$$

If we let m and $M_{\sigma}^{(m)}$ grow large, the rate becomes

$$R(\mathcal{C}_{\text{MDS}}) = \frac{1}{m} \log M_{\sigma}^{(m)} (1 + o(1)). \quad (11)$$

For a given alphabet Σ_q , let A denote the adjacency matrix of G_q , where the rows and columns of A are indexed by $v \in V_q \subseteq \Sigma_q^5$. Furthermore, let $A_{(v)}$ be obtained by deleting the row and column corresponding to v from A and $c_{(v)}$ (resp. $r_{(v)}^T$) be the column (row) of A corresponding to v with the element corresponding to v removed. Recall that $M_{\sigma}^{(m)} = |\mathcal{B}_{\sigma}^m|$ and $l = 5$. From Theorem 13, $M_{\sigma}^{(m)}$ equals the number of paths of length $m + l$ in G_q that start and end with σ but do not visit σ in their interiors. The number of paths of length $m + l - 2$ from vertex u to vertex v in G_q is given by the (u, v) element of $(A_{(\sigma)})^{m+l-2}$. Noting that the paths start and end σ , we have

$$M_{\sigma}^{(m)} = r_{(\sigma)}^T (A_{(\sigma)})^{m+l-2} c_{(\sigma)}, \quad (12)$$

where $(\cdot)^T$ denotes matrix transpose. Here, multiplying by $r_{(\sigma)}^T$ from the left and $c_{(\sigma)}$ from the right allow us to sum over elements (u, v) of $(A_{(\sigma)})^{m+l-2}$ such that there is an edge from σ to u and an edge from v to σ . As $m \rightarrow \infty$, if $A_{(\sigma)}$ is primitive [23], we have

$$\frac{1}{m + l} \log M_{\sigma}^{(m)} \rightarrow \log(\lambda_{\sigma}), \quad (13)$$

where λ_{σ} is the largest eigenvalue of $A_{(\sigma)}$. Maximizing over $\sigma \in V_q$ yields the largest value for $M_{\sigma}^{(m)}$ in (12) and (13), and thus the highest code rate. This is possible to do computationally for small values of q and, in particular, for $q = 4$, which corresponds to data storage in DNA. In this case, $A_{(\sigma)}$ is primitive for all choices of $\sigma \in \text{Irr}(5)$ and the largest eigenvalue is obtained for $\sigma = 01201$ (and strings obtained from 01201 by relabeling the alphabet symbols). For this σ , we find $\lambda_{\sigma} = 2.6534$, leading to an asymptotic code rate of 1.4078 bits/symbol.

It was shown in [3] that the set of ≤ 3 -irreducible strings of length n is a code correcting any number of ≤ 3 -TDs. In [18], it was shown that the rate of this code, $\frac{1}{n} \log |\text{Irr}(n)|$, is asymptotically optimal. It is easy to see that $\frac{1}{n} \log |\text{Irr}(n)| \leq \log(q - 1)$ as no symbol can be repeated. For the case of $q = 4$, we have $\frac{1}{n} \log |\text{Irr}(n)| = \log 2.6590 = 1.4109$ bits/symbol. Therefore, the cost of protection against a single edit in our construction is only 0.003 bits/symbol. It should be noted, however, that here we have assumed m is large, thus ignoring the overhead from the MDS code and marker strings.

In addition to the computational rate obtained above for the important case of $q = 4$, we will provide analytical bounds on the code rate. An important quantity affecting the rate of the code is the number of outgoing edges from each vertex in G_q that do not lead to σ . The asymptotic rate of the code is bounded from below by the number of such edges. The next lemma, which establishes the number of outgoing edges for each vertex, will be useful in identifying an appropriate choice of σ , and

the following theorem provides a lower bound for $M_{\sigma}^{(m)}$ for such a choice.

Lemma 14. *For $q > 2$, a vertex $v = a_1a_2a_3a_4a_5$ in G_q has $q - 2$ outgoing edges if $a_3 = a_5$ or $a_1a_2 = a_4a_5$. Otherwise, it has $q - 1$ outgoing edges.*

Proof: Consider $v = a_1a_2a_3a_4a_5 \in \text{Irr}(5)$, and $w = a_2a_3a_4a_5a_6 \in \text{Irr}(5)$. There is an edge from v to w if $a_1a_2a_3a_4a_5a_6 \in \text{Irr}(6)$. The number of outgoing edges from v equals the number of possible values for a_6 such that this condition is satisfied. Clearly, $a_6 \neq a_5$. Furthermore, if $a_3 = a_5$, then $a_6 \neq a_4$ and if $a_1a_2 = a_4a_5$, then $a_6 \neq a_3$.

However, $a_3 = a_5$ and $a_1a_2 = a_4a_5$ cannot simultaneously hold, since that would imply $a_2 = a_3$, contradicting $v \in \text{Irr}(5)$. Hence, if either $a_3 = a_5$ or $a_1a_2 = a_4a_5$ holds, then there are $q - 2$ outgoing edges and if neither holds, there are $q - 1$ outgoing edges. ■

Since σ must also be excluded, it may seem that the number of outgoing edges may be as low as $q - 3$. But we show in the next theorem that with an appropriate choice of σ , we can have $q - 2$ as the lower bound.

Theorem 15. *Over an alphabet of size $q > 2$, there exists $\sigma \in \text{Irr}(5)$ such that $M_{\sigma}^{(m)} \geq (q - 2)^{m - c_q}$, where c_q is a constant independent from m .*

Proof: Recall that $M_{\sigma}^{(m)}$ is the number of paths of length $m + 5$ in G_q that start and end in σ but do not visit σ otherwise. Since the path must return to σ , we will show below that for an appropriate choice of σ , there is a path in G_q from any vertex to σ , and define c_q such that the length of this path is at most $c_q + 5$. Hence $M_{\sigma}^{(m)}$ is at least the number of paths of length $m - c_q$ from σ to another vertex that do not pass through σ .

As shown in Lemma 14, each vertex in G_q has at least $q - 2$ outgoing edges. We select σ such that this still holds even if edges leading to σ are excluded. We do so by ensuring that each vertex v with an outgoing edge to σ has $q - 1$ outgoing edges. Let $v = a_1a_2a_3a_4a_5$ and $\sigma = a_2a_3a_4a_5a_6$. Based on Lemma 14, if $a_2 \neq a_5$ and $a_3 \neq a_5$, then v has $q - 1$ outgoing edges. In particular, we can choose $\sigma = 01020$ since $q \geq 3$. With this choice, $M_{\sigma}^{(m)} \geq (q - 2)^{m - c_q}$.

To complete the proof, we need to show that there is a path in G_q from any vertex to $\sigma = 01020$. For $q = 3, 4, 5$, we have checked this claim computationally by explicitly forming G_q . Let us then suppose $q \geq 6$, where the alphabet Σ_q contains $\{3, 4, 5\}$. Let $v = a_1 \cdots a_5$ be some vertex in G_q . There is an edge from v to $a_2 \cdots a_6$ for some $a_6 \in \{3, 4, 5\}$ since, from Lemma 14, at most two elements of Σ_q are not permissible. Continuing in similar fashion, in 5 steps, we can go from v to some vertex $w = b_1 \cdots b_5$ whose elements b_i belong to $\{3, 4, 5\}$. We can then reach σ in 5 additional steps via the path $w \rightarrow b_2 \cdots b_4b_50 \rightarrow b_3b_4b_501 \rightarrow \cdots \rightarrow \sigma$, proving the claim. In particular, for $q \geq 6$, we have $c_q \leq 5$. ■

We can now find a lower bound on the asymptotic rate, based on (11) and the proceeding theorem:

Corollary 16. *For $q > 2$, as $m \rightarrow \infty$, $R(\mathcal{C}_{\text{MDS}}) \geq \log(q - 2)(1 + o(1))$.*

We note that this gives the lower bound of 1 bit/symbol for $q = 4$, which we can compare to the upper bound of $\log(q - 1) = 1.585$ for codes correcting only duplications and to the rate obtained computationally following (13), which was 1.4078 bits/symbol.

V. CONCLUSION

This paper considered constructing error-correcting codes for channels with many short duplications and one edit error. Because the channel allows an arbitrary number of duplications, a single edit may affect an unbounded segment of the output. For example, an inserted symbol may appear many times in different positions. However, with an appropriate construction of message blocks and processing of the output strings, the edit error leads to the erasure of at most 4 message blocks or substitution of at most 2. Therefore, a maximum distance separable (MDS) code with minimum Hamming distance 5 over message blocks can correct these errors. However, there is an additional requirement. Namely, the codewords must be irreducible. Separating the message blocks with a marker sequence σ of length at least 5 allows us to ensure that the codewords are repeat-free by guaranteeing that each message block is irreducible. The rate of the code is determined by the number of such blocks, which in turn depends on the marker sequence σ . We showed that permitted message blocks are paths in a modified De Bruijn graph and that choosing σ appropriately allows each vertex to have at least $q - 2$ outgoing edges, thus guaranteeing an asymptotic rate of at least $\log(q - 2)$. When $q = 4$, the case corresponding to DNA storage, a computational bound for the code rate shows that the asymptotic rate is only 0.003 bits/symbol smaller than that of the code that corrects short duplications but no edits.

The problem of correcting more edit errors is also of interest but left to future work. Another, possibly more challenging problem is correcting edits and duplications of length bounded by an arbitrary constant k . If k is larger than 3, the duplication root is no longer unique [3], which complicates the code design. Furthermore, a key feature of duplications of length at most 3 is that such duplications lead to regular languages. We used this fact to characterize the effect of the channel on the roots of sequences. However, if $k \geq 4$, then the language is not regular [24], leading to challenges in characterizing the channel. To guarantee a unique root under $\leq k$ duplications with $k \geq 4$, the work [25] constructed error-correcting codes by selecting a set of $(k + 1)$ -distinct strings with $q \geq (k + 1)$ and achieved an asymptotic code rate of $\frac{1}{k} \sum_{i=1}^k \log(q - i)$. This means that, despite the interesting approach, to correct duplications of length ≥ 4 , [25] requires $q \geq 5$, which is not suitable for the important application of data storage in DNA with $q = 4$.

Acknowledgement: The authors thank the anonymous reviewers of the paper for their insightful comments and helpful suggestions.

REFERENCES

- [1] 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.
- [2] 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.
- [3] S. Jain, F. Farnoud, M. Schwartz, and J. Bruck, "Duplication-correcting codes for data storage in the DNA of living organisms," *IEEE Transactions on Information Theory*, vol. 63, no. 8, pp. 4996–5010, 2017.
- [4] M. Kovačević and V. Y. Tan, "Asymptotically optimal codes correcting fixed-length duplication errors in DNA storage systems," *IEEE Communications Letters*, vol. 22, no. 11, pp. 2194–2197, 2018.
- [5] 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.

- [6] 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.
- [7] 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.
- [8] 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.
- [9] 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.
- [10] 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.
- [11] 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.
- [12] 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.
- [13] 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.
- [14] Y. Yehezkeally and M. Schwartz, "Uncertainty of reconstructing multiple messages from uniform-tandem-duplication noise," in *IEEE International Symposium on Information Theory (ISIT)*, 2020, pp. 126–131.
- [15] 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 International Symposium on Information Theory (ISIT)*. IEEE, 2020, pp. 694–699.
- [16] 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.
- [17] 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.
- [18] M. Kovačević, "Codes correcting all patterns of tandem-duplication errors of maximum length 3," *arXiv preprint arXiv:1911.06561*, 2019.
- [19] 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.
- [20] Y. Tang and F. Farnoud, "Error-correcting codes for noisy duplication channels," in *57th Annual Allerton Conference on Communication, Control, and Computing*, 2019, pp. 140–146.
- [21] D. Pumpernik, B. Oblak, and B. Borštnik, "Replication slippage versus point mutation rates in short tandem repeats of the human genome," *Molecular Genetics and Genomics*, vol. 279, no. 1, pp. 53–61, 2008.
- [22] F. Farnoud, M. Schwartz, and J. Bruck, "Estimation of duplication history under a stochastic model for tandem repeats," *BMC Bioinformatics*, vol. 20, no. 1, 2019. [Online]. Available: <https://doi.org/10.1186/s12859-019-2603-1>
- [23] D. Lind, B. Marcus, L. Douglas, and M. Brian, *An introduction to symbolic dynamics and coding*. Cambridge university press, 1995.
- [24] P. Leupold, V. Mitrana, and J. M. Sempere, "Formal languages arising from gene repeated duplication," in *Aspects of Molecular Computing*. Springer, 2003, pp. 297–308.
- [25] 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.