

Error-correcting Codes for Noisy Duplication Channels

Yuanyuan Tang and Farzad Farnoud (Hassanzadeh)

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

Abstract—Because of its high data density and longevity, DNA is emerging as a promising candidate for satisfying increasing data storage needs. Compared to conventional storage media, however, data stored in DNA is subject to wider range of errors resulting from various processes involved in the data storage pipeline. In this paper, we consider correcting duplication errors for both exact and noisy tandem duplications of a given length k . Specifically, we design codes that can correct any number of exact duplication and one noisy duplication errors, where in the noisy duplication case the copy is at Hamming distance 1 from the original. Our constructions rely upon recovering the duplication root of the stored codeword. We characterize the ways in which duplication errors manifest in the root of affected sequences and design efficient codes for correcting these error patterns. We show that the proposed construction is asymptotically optimal.

I. INTRODUCTION

The rapidly increasing amount of data and the need for long-term data storage have led to new challenges. In recent years, advances in DNA sequencing, synthesis, and editing technologies [13], [11] have made deoxyribonucleic acid (DNA) a promising alternative to conventional storage media. Compared to traditional media, DNA has several advantages, including high data density, longevity, and ease of copying information. For example, it may be possible to recover a DNA sequence after 10,000 years and a single human cell contains an amount of DNA that can ideally hold 6.4 Gb of information [13]. However, DNA storage technologies also encounter many challenges. One obvious challenge is that a diverse set of errors are possible, including substitution, duplication, insertion, and deletion. This paper focuses on error-correcting codes for noisy duplication channels. In such case, in addition to exact duplication, noisy duplication, where an approximate copy is inserted into the sequence, may occur.

In duplication channels, (tandem) duplication errors generate copies of substrings of the sequence and insert

each copy after the original substring [3]. This type of channel was first studied in the context of recovering from timing errors in communication systems that led to individual symbols being repeated [2]. The copying mechanism of DNA, however, allows multiple symbols being repeated, for example, via slipped-strand mispairings, where the slippage of the molecule copying DNA causes a substring to be repeated [3]. Properties of duplication in DNA have been studied from various vantage points, including the theory of formal languages and the entropy of DNA sequences (see, e.g., [7] and references therein). Codes for correcting duplication errors in the context of data storage in the DNA of living organisms, such as bacteria [9], were studied by [3], where optimal constructions for correcting exact duplications of constant length were presented. This and related problems were then further studied by a number of works including [4], [14], [5], [6], [1], [11]. Most related to this paper is [11], which studies error correction in duplication and substitution channels, when substitutions are independent from duplications and when they only occur in copies generated by duplications. The latter model, i.e., the *noisy duplication model*, which is motivated by the abundance of inexact copies in tandem repeat stretches in genomes [8], is the model studied in this work.

In the noisy duplication channel, two types of errors are possible: i) exact duplications, which insert an exact copy of a substring in tandem, such as $ACGTC \rightarrow ACGTCGTC$; and ii) noisy duplications, which insert approximate copies, e.g., $ACGTC \rightarrow ACGTC\underline{T}TC$. In both cases, the length of the duplication refers to the length of the duplicated substring (3 in our preceding examples). In this paper, we limit our attention to exact and noisy tandem duplications of length k , referred to as k -TDs and k -NDs, respectively. Furthermore, we only consider noisy duplications where the copy and the original substring differ in one position. In other words, each noisy duplication can be viewed as an exact duplication followed by a substitution in the inserted copy.

We will design codes that correct (infinitely) many

This work was supported in part by NSF grants under grant nos. 1816409 and 1755773.

k -TD and a single k -ND errors, as a step towards codes that can correct t_1 k -TDs and t_2 k -NDs, for given t_1 and t_2 . The proposed codes will rely on finding the duplication root of the stored codeword. The *duplication root* of a sequence \mathbf{x} is the sequence obtained from \mathbf{x} by removing all repeats of length k . While k -TDs do not alter the duplication root, k -NDs do. Thus, we will first analyze the effect of noisy duplications on the root of the sequence. We show that the root may change in a variety of ways, leading to several error patterns. We then design efficient error-correcting codes that correct these errors via a number of transforms that simplify the different error patterns.

It was shown in [3] that the rate of the optimal code capable of correcting many k -TDs is

$$1 - \frac{(q-1) \log_q e}{q^{k+2}} + o(1), \quad (1)$$

as the length n of the code grows, where q is the size of the alphabet. The question then arises as to whether it is possible to correct an additional noisy duplication without a rate penalty. It is worth noting that the best known code for correcting an additional unrestricted substitution, i.e., a substitution that can occur anywhere rather than in a copy generated by duplication, has rate that is bounded from below by [11]

$$1 - \frac{2}{k} \log_q \frac{q}{q-1} + o(1). \quad (2)$$

which indicates a rate penalty. In contrast, we show that the proposed codes have the same asymptotic rate as (1), and are thus asymptotically optimal.

This paper is organized as follows. The notation and preliminaries are given in Section II. In Section III, we analyze the error patterns that manifest as the result of passing through the noisy duplication channel. Finally, the code construction and the corresponding code size are presented in Section IV. Note that proofs of theorems are not presented because of the limited space.

II. NOTATION AND PRELIMINARIES

Throughout the paper, Σ_q represents a finite alphabet of size q , assumed without loss of generality to be $\{0, 1, \dots, q-1\}$. We use Σ_q^+ to denote the nonzero elements of Σ_q and Σ_q^* to denote all strings of finite length over Σ_q . In particular, Σ_q^* includes the empty string Λ . Furthermore, Σ_q^n represents the strings of length n over Σ_q . The set $\{1, \dots, n\}$ is represented by $[n]$.

We use bold symbols, such as \mathbf{x} and \mathbf{y}_j , to denote strings over Σ_q . The entries of strings are shown

with normal symbols, e.g., $\mathbf{x} = x_1 x_2 \dots x_n$ and $\mathbf{y}_j = y_{j1} y_{j2} \dots y_{jm}$, where $x_i, y_{ji} \in \Sigma_q$. The indices of elements of words over Σ_q^* start from 1, unless otherwise stated. For two words $\mathbf{x}, \mathbf{y} \in \Sigma_q^*$, their concatenation is denoted as \mathbf{xy} , and \mathbf{x}^m represents the concatenation of m copies of \mathbf{x} . Given a word $\mathbf{x} \in \Sigma_q^*$, the length of \mathbf{x} is represented as $|\mathbf{x}|$. In addition, for a word $\mathbf{x} \in \Sigma_q^*$, the Hamming weight $\text{wt}(\mathbf{x})$ denotes the number of non-zero symbols in \mathbf{x} . If a word $\mathbf{x} \in \Sigma_q^*$ can be expressed as $\mathbf{x} = \mathbf{uvw}$ with $\mathbf{u}, \mathbf{v}, \mathbf{w} \in \Sigma_q^*$, then \mathbf{v} is a substring of \mathbf{x} .

Given a word $\mathbf{x} \in \Sigma_q^*$, an (exact) *tandem duplication* of length k (k -TD) generates a copy of a substring \mathbf{v} of \mathbf{x} of length k and inserts the copy immediately after \mathbf{v} . More specifically, a k -TD can be expressed as [3]

$$T_{i,k}(\mathbf{x}) = \begin{cases} \mathbf{uvv} & \text{if } \mathbf{x} = \mathbf{uvw}, |\mathbf{u}| = i, |\mathbf{v}| = k \\ \mathbf{x} & \text{if } |\mathbf{x}| < i + k \end{cases} \quad (3)$$

For example, given the alphabet $\Sigma_3 = \{0, 1, 2\}$ and $k = 3$, a k -TD may result in

$$\mathbf{x} = 1201210 \rightarrow \mathbf{x}' = T_{1,3}(\mathbf{x}) = 1201\underline{201}210, \quad (4)$$

where the underlined substring 201 is the copy. We refer to \mathbf{x}' as a k -TD *descendant* of \mathbf{x} .

Given a word $\mathbf{x} \in \Sigma_q^n$, $n \geq k$, the *k -discrete-derivative* transform [3] is defined as $\phi(\mathbf{x}) = (\hat{\phi}(\mathbf{x}), \bar{\phi}(\mathbf{x}))$, where

$$\hat{\phi}(\mathbf{x}) = x_1 \dots x_k, \bar{\phi}(\mathbf{x}) = x_{k+1} \dots x_n - x_1 \dots x_{n-k}. \quad (5)$$

where the subtraction is performed entry-wise modulo q . Continuing the example given in (4),

$$\begin{aligned} \mathbf{x} = 1201210 &\rightarrow \mathbf{x}' = 1201\underline{201}210, \\ \phi(\mathbf{x}) = 120, 0012 &\rightarrow \phi(\mathbf{x}') = 120, \underline{0000}12. \end{aligned} \quad (6)$$

As seen in the example, after the k -TD in \mathbf{x} , $\bar{\phi}(\mathbf{x}')$ can be obtained by inserting 0^k into $\bar{\phi}(\mathbf{x})$, immediately after the i -th entry.

Copies generated by tandem duplications may not be always perfect. That is, the copy may not always be exact. Such a duplication is referred to as a *noisy duplication*. In this paper, we limit our attention to noisy duplications in which the copy is at Hamming distance 1 from the original. Continuing example (4), one symbol in the copy 201 may change,

$$\begin{aligned} \mathbf{x}' = 1201201210 &\rightarrow \mathbf{x}'' = 1201\underline{101}210, \\ \phi(\mathbf{x}') = 120, 000012 &\rightarrow \phi(\mathbf{x}'') = 120, \underline{0200}12. \end{aligned}$$

As seen in the example, a noisy duplication of length k (k -ND) can be regarded as an exact k -TD followed

by a substitution. Given a word $\mathbf{x} \in \Sigma_q^*$, the tandem duplication results in $\mathbf{x}' = T_{i,k}(\mathbf{x})$ and the following substitution results in $\mathbf{x}'' = T_{i,k}(\mathbf{x}) + a\mathbf{e}_j$, where $(i + k + 1) \leq j \leq (i + 2k)$, $a \in \Sigma_q^+$, and \mathbf{e}_j represents a unit vector with 1 in the j -th entry and 0 elsewhere. Note that the first k elements are not affected by exact or noisy duplications and $\hat{\phi}(\mathbf{x}) = \hat{\phi}(\mathbf{x}') = \hat{\phi}(\mathbf{x}'')$. Hence, we focus on changes in $\bar{\phi}(\cdot)$. The substitution changes at most two symbols of $\bar{\phi}(\mathbf{x}')$ and can be expressed as

$$\bar{\phi}(\mathbf{x}'') = \bar{\phi}(\mathbf{x}') + a\boldsymbol{\epsilon}_j, \quad (7)$$

where $\boldsymbol{\epsilon}_j = \mathbf{e}_{j-k} - \mathbf{e}_j$ if $(k + 1) \leq j \leq (|\mathbf{x}'| - k)$ and $\boldsymbol{\epsilon}_j = \mathbf{e}_{j-k}$ if $(|\mathbf{x}'| - k + 1) \leq j \leq |\mathbf{x}'|$. We refer to \mathbf{x}'' as a k -ND descendant of \mathbf{x} .

Since noisy duplications may occur at any position, the word \mathbf{x} can generate many descendants through noisy duplication errors. Let $D_k^{t(p)}(\mathbf{x})$ denote the *descendant cone* of \mathbf{x} obtained after t duplications, p of which are noisy, where $t \geq p$. Furthermore, the descendant cone with many exact k -TDs and at most P noisy duplications, i.e., at most P substitution errors, can be expressed as

$$D_k^{*(\leq P)}(\mathbf{x}) = \bigcup_{p=0}^{P} \bigcup_{t=p}^{\infty} D_k^{t(p)}(\mathbf{x}). \quad (8)$$

In this paper, we limit our attention to $P = 1$.

We define a mapping operation $\mu : \Sigma_q^* \rightarrow \Sigma_q^*$ by removing all runs of 0^k in $\mathbf{z} \in \Sigma_q^*$. More specifically, consider a string \mathbf{z} as

$$\mathbf{z} = 0^{m_0} w_1 0^{m_1} \dots w_t 0^{m_{t+1}},$$

where $t = \text{wt}(\mathbf{z})$, $w_1, \dots, w_t \in \Sigma_q^+$, and m_0, \dots, m_{t+1} are non-negative integers. The mapping $\mu(\mathbf{z})$ is defined as

$$\mu(\mathbf{z}) = 0^{m_0 \bmod k} w_1 0^{m_1 \bmod k} \dots w_t 0^{m_{t+1} \bmod k}.$$

Also, $\text{RLL}(m)$ denotes the set of strings of length m containing no 0^k . In other words, $\text{RLL}(m) = \{\mathbf{z} \in \Sigma_q^m \mid \mu(\mathbf{z}) = \mathbf{z}\}$.

According to [3], given a word $\mathbf{x} \in \Sigma_q^*$, after many (even infinite) k -TD errors, the string $(\hat{\phi}(\mathbf{x}), \mu(\bar{\phi}(\mathbf{x})))$ stays the same. To make use of this property, define the *duplication root* $\text{drt}(\mathbf{x})$ as the string obtained from \mathbf{x} after all copies of length k are removed. Note that we then have

$$\phi(\text{drt}(\mathbf{x})) = (\hat{\phi}(\mathbf{x}), \mu(\bar{\phi}(\mathbf{x}))). \quad (9)$$

If $\text{drt}(\mathbf{x}) = \mathbf{x}$, we call the word \mathbf{x} irreducible. The set of all irreducible words of length n can be written as $\text{Irr}(n) = \{\mathbf{x} \in \Sigma_q^n \mid \text{drt}(\mathbf{x}) = \mathbf{x}\}$. In other words, an irreducible word $\mathbf{x} \in \Sigma_q^n$ satisfies $\bar{\phi}(\mathbf{x}) \in \text{RLL}(n - k)$.

For a word $\mathbf{z} \in \Sigma_q^*$, we define its *indicator* $\Gamma(\mathbf{z}) : \Sigma_q^* \rightarrow \Sigma_2^*$ as $\Gamma(\mathbf{z}) = \Gamma_1(\mathbf{z}) \cdots \Gamma_{|\mathbf{z}|}(\mathbf{z})$, where

$$\Gamma_i(\mathbf{z}) = \begin{cases} 1, & \text{if } z_i \neq 0, \\ 0, & \text{otherwise.} \end{cases} \quad i = 1, \dots, |\mathbf{z}|. \quad (10)$$

Based on (7), the substitution in a noisy duplication alters two symbols in $\bar{\phi}(\mathbf{x}')$ at distance k . For the purpose of error correction, it will be helpful to rearrange the symbols into k strings such that the two symbols affected by the substitution appear next to each other in one of the strings. More precisely, for $j \in [k]$, we define a *splitting* operation that extracts entries whose position is equal to j modulo k . That is, for $\mathbf{u} \in \Sigma_q^n$ and $j \in [k]$, define $\mathbf{u}_j = (\mu_{ji})_i = \text{Sp}_k(\mathbf{u}, j)$ such that

$$\mu_{ji} = \mu_{j+(i-1)k}, \quad 1 \leq i \leq \left\lceil \frac{n-j}{k} \right\rceil + 1.$$

For $\mathbf{u} \in \Sigma_q^n$, we then define the *interleaving* operation $\text{IL} : \Sigma_q^n \rightarrow \Sigma_q^n$ as the concatenation of $\text{Sp}_k(\mathbf{u}, j)$, $j \in [k]$,

$$\text{IL}(\mathbf{u}) = \text{Sp}_k(\mathbf{u}, 1) \cdots \text{Sp}_k(\mathbf{u}, k).$$

Example 1. Given an alphabet $\Sigma_3 = \{0, 1, 2\}$, $k = 3$, and $\mathbf{u}' = \bar{\phi}(\mathbf{x}') = 221200012$, after splitting \mathbf{u}' , we obtain

$$\begin{aligned} \mathbf{u}'_1 &= \text{Sp}_3(\mathbf{u}', 1) = 220, \\ \mathbf{u}'_2 &= \text{Sp}_3(\mathbf{u}', 2) = 201, \\ \mathbf{u}'_3 &= \text{Sp}_3(\mathbf{u}', 3) = 102, \\ \text{IL}(\mathbf{u}') &= \mathbf{u}'_1 \mathbf{u}'_2 \mathbf{u}'_3 = 220201102. \end{aligned}$$

Based on (7), after one substitution error, we may obtain $\mathbf{u}'' = \bar{\phi}(\mathbf{x}'') = 22120\underline{1}011$. We then find

$$\begin{aligned} \mathbf{u}''_1 &= \text{Sp}_3(\mathbf{u}'', 2) = 201, \\ \mathbf{u}''_2 &= \text{Sp}_3(\mathbf{u}'', 1) = 220, \\ \mathbf{u}''_3 &= \text{Sp}_3(\mathbf{u}'', 3) = 1\underline{1}1, \\ \text{IL}(\mathbf{u}'') &= \mathbf{u}''_1 \mathbf{u}''_2 \mathbf{u}''_3 = 2202011\underline{1}1. \end{aligned}$$

We observe that the error is restricted to \mathbf{u}''_3 and that the two symbols changed by the substitution error are adjacent in $\text{IL}(\mathbf{u}'')$, while they are not so in \mathbf{u}'' .

Given a word $\mathbf{z} \in \Sigma_q^n$, we define the *cumulative-sum* operation $\text{CS} : \Sigma_q^n \rightarrow \Sigma_q^n$, as $\mathbf{r} = \text{CS}(\mathbf{z})$, where

$$r_i = \sum_{t=1}^i z_t \bmod q, \quad i = 1, \dots, n. \quad (11)$$

We further define the *odd subsequence* $\text{Od}(\mathbf{z})$ and the *even subsequence* $\text{Ev}(\mathbf{z})$ of a word $\mathbf{z} \in \Sigma_q^n$ as two sequences containing symbols in the odd and

even positions, respectively. More precisely, $\text{Od}(z) = \text{Sp}_2(z, 1)$ and $\text{Ev}(z) = \text{Sp}_2(z, 2)$.

Our results will rely on codes that can correct a single insertion or deletion. We thus recall the Varshamov-Tenengolts codes [10], [12], which are binary codes capable of correcting a single insertion or deletion (indel).

Construction 1. Given integers $m \geq 1$ and $0 \leq \alpha \leq (m-1)$, the binary Varshamov-Tenengolts (VT) code [10] $C_{VT}(\alpha, m)$ is given as

$$C_{VT}(\alpha, m) = \{z \in \Sigma_2^* \mid \sum_{i=1}^{|z|} iz_i = \alpha \pmod{m}\}. \quad (12)$$

Compared to the binary indel-correcting code, correcting indels in non-binary sequences is more challenging. We will use Tenengolts' q -ary single-indel-correcting code [12], which relies on the mapping $\zeta : \Sigma_q^* \rightarrow \Sigma_2^*$, where the i -th position of $\zeta(z)$ is

$$\zeta_i(z) = \begin{cases} 1, & \text{if } z_i \geq z_{i-1}, \\ 0, & \text{if } z_i < z_{i-1}. \end{cases} \quad i = 2, 3, \dots, |z|. \quad (13)$$

with $\zeta_1(z) = 1$.

Construction 2. Based on Tenengolts' q -ary code [12], given integers $m \geq 1$, $0 \leq \alpha \leq (q-1)$ and $0 \leq \beta \leq (m-1)$, we construct the code $C_{Tq}(\alpha, \beta, m)$ over Σ_q^* as

$$C_{Tq}(\alpha, \beta, m) = \left\{ z \in \Sigma_q^* \mid \sum_{j=1}^{|z|} z_j = \alpha \pmod{q}, \right. \\ \left. \sum_{i=1}^{|z|} (i-1)\zeta_i(z) = \beta \pmod{m} \right\}. \quad (14)$$

III. NOISY DUPLICATION CHANNELS

To enable designing error-correcting codes, in this section, we study the relation between the input and output sequences in *noisy duplication channels*. As before, we consider channels with many (possibly infinite) exact duplications and at most one noisy duplication in which one of the copied symbols is altered.

If a code $C \in \Sigma_q^n$ corrects many k -TD and one k -ND errors, then for any two distinct codewords $c_1, c_2 \in C$, we have

$$D_k^{*(\leq 1)}(c_1) \cap D_k^{*(\leq 1)}(c_2) = \emptyset. \quad (15)$$

This can be shown to be equivalent to

$$\begin{aligned} \text{drt}(c_2) &\neq \text{drt}(c_1), \\ \text{drt}(D_k^{*(\leq 1)}(c_1)) \cap \text{drt}(D_k^{*(\leq 1)}(c_2)) &= \emptyset. \end{aligned} \quad (16)$$

Since k -TDs do not alter the root of the sequence, $\text{drt}(c_2) \neq \text{drt}(c_1)$ ensures that k -TD errors can be corrected. Noisy tandem duplications however alter the roots. In fact, they may produce sequences with roots whose lengths are different from the roots of the stored sequences. Since the codewords have distinct roots, it suffices to recover the root of the retrieved word to correct any errors. We will restrict our constructions to codes whose codewords are irreducible, and thus are their own roots. While this is not necessary, it will simplify the code construction, as we will show, and does not incur a large penalty in terms of the size of the code.

For noisy duplication channels, given a codeword $x \in \Sigma_q^n$, the generation of descendants $x'' \in D_k^{*(\leq 1)}(x)$ includes three different cases: only k -TDs; k -TDs followed by one k -ND; and k -TDs, followed by a k -ND, followed by more k -TDs. Since the root is not affected by the k -TDs, to study $\text{drt}(D_k^{*(\leq 1)}(x))$, we only need to consider the second case, i.e., we focus on descendants x'' immediately after the noisy duplication.

Given an irreducible string $x \in \Sigma_q^n$ with $n > 2k$, our goal is to characterize $\text{drt}(D_k^{*(\leq 1)}(x))$. Based on (5), we have

$$\phi(x) = (\hat{\phi}(x), \bar{\phi}(x)) = (y, z), \quad (17)$$

where $y = \hat{\phi}(x) \in \Sigma_q^k$ and $z = \bar{\phi}(x) \in \Sigma_q^{n-k}$. Since x is an irreducible string, the string z contains no runs of 0^k , i.e. $z = \mu(z)$.

After many k -TDs and one k -ND, we have a descendant $x'' \in D_k^{*(\leq 1)}(x)$. Since the substitution only occurs in the copy, the first k symbols always stay the same. Thus x'' satisfies

$$\phi(x'') = (\hat{\phi}(x''), \bar{\phi}(x'')) = (\hat{\phi}(x), \bar{\phi}(x'')) = (y, z''). \quad (18)$$

Based on (9), it suffices to study the problem in the transform domain, i.e., we want to obtain all possible $(y, \mu(z''))$ derived from $(y, \mu(z))$. Our code constructions in the next section will also rely on certain sequences derived from $\mu(z)$. The next theorem characterizes how these sequences can be altered by k -TDs and one k -ND.

Theorem 1. Let $x \in \Sigma_q^n$ and let $x'' \in D_k^{*(\leq 1)}(x)$ be a descendent of x (produced by passing through the noisy duplication channel). Furthermore, let

$$\begin{aligned} z &= \bar{\phi}(x), & \mu &= \mu(z), \\ \mu_j &= \text{Sp}_k(\mu, j), & s_j &= \Gamma(\mu_j). \end{aligned}$$

We define $z'', \mu'', \mu_j'', s_j''$, similarly, based on x'' . The differences between sequences defined based on x and x'' are given in Table I and Table II.

Note that the length of μ can change by $-k$, 0 , k , or $2k$. This means that the noisy duplication may manifest as deletions, insertions, or substitutions in μ . Furthermore, the complex error patterns in μ are simplified when we consider $\mu_j, j \in [k]$. The errors marked by $(*)$ occur for at most one value of j . These correspond to positions affected by the substitution. (Rows marked by $(\$)$ relate to our error-correction strategy and are discussed in the next section.)

Now that we have determined all changes from (y, μ) to (y, μ'') resulting from passing through the noisy duplication channel, we consider the code design to correct many exact k -TDs and at most one noisy duplication in the next section.

IV. ERROR-CORRECTING CODES FOR NOISY DUPLICATION CHANNELS

Recall from Section III that we are interested in constructing a code $C \subseteq \text{Irr}(n) \cap \Sigma_q^n$ that can correct many exact k -TDs and at most one noisy duplication. Based on (16), for any code that corrects k -TDs, two distinct codewords must have distinct roots. Thus, for a stored codeword x and the retrieved word x'' , if we can recover the duplication root $\text{drt}(x)$ of x from x'' , we can recover the codeword x . But we have made a further simplifying assumption that $C \subseteq \text{Irr}(n)$ and thus $x = \text{drt}(x)$.

As shown in Theorem 1, duplication errors manifest in various ways in $\text{drt}(x'')$ and its counterpart in the μ -transform domain $\mu(\bar{\phi}(x''))$. Hence, for error correction, we utilize several sequences derived from x , including μ_j and $s_j, j \in [k]$, as defined in Theorem 1. Furthermore, we define $r = \text{CS}(\text{IL}(\mu))$ and $r'' = \text{CS}(\text{IL}(\mu''))$. We note that r (similarly r'') can be directly found by rearranging the elements $x_{k+1} \cdots x_n$.

The relationship between these mappings is illustrated in Figure 1. In the figure, solid edges represent invertible mappings. Since x is irreducible, the stored codeword can be recovered if any of $\mu, (\mu_j)_{j \in [k]}, \text{IL}(\mu)$ or r are recovered (note that $x_1 \cdots x_k$ are not affected by errors). We use these mappings to simplify and correct different error patterns described by Theorem 1 in an efficient manner.

The motivation behind defining $\mu_j, j \in [k]$, is to convert insertions and deletions of blocks of length k into simpler errors involving one or two symbols.

Some of the errors, marked by $(\$)$ in Tables I and II, involve 0s, which appear in the same positions in s_j and μ_j . Correcting these errors in s_j is more efficient since it will rely on binary codes rather than q -ary codes. We will first correct these errors in s_j and then correct the corresponding μ_j . Finally, the cumulative-sum mapping CS turns errors marked by $(*)$, e.g., $\Lambda \rightarrow a\bar{a}$ into a single q -ary insertion or substitution. Importantly, in each case there is only one such error. So if other errors are corrected, we can concatenate $\mu_j, j \in [k]$, and then correct the single occurrence of this error.

We will construct an error-correcting code that will allow us to recover μ from μ'' . As discussed, for certain errors occurring in μ_j , specifically those marked by $(\$)$ in Tables I and II, we may do so by correcting errors in s_j , via Construction 3 below.

The indicator vectors (s_1, \dots, s_k) are subject to several error patterns: insertion of 11; insertion of two 0s with distance at most 2; indel of 1 or 0; swaps of two adjacent elements; and substitution of one or two 0s with one or two 1s. The following code can correct a single occurrence of one of these errors, as shown in the next theorem. A slightly modified version of this code is used for the noisy duplication channel.

Construction 3. Given integers $0 \leq a \leq 2(n+1)$, $0 \leq b \leq 4$, and $0 \leq c \leq 2n$, we construct the code $C_{(a,b,c)}$ as

$$C_{(a,b,c)} = \{\mathbf{u} \in \Sigma_2^n \mid \mathbf{u} \in C_{VT}(a, 2n+3), \quad (19)$$

$$\sum_{i=1}^n u_i = b \pmod{5}, \quad (20)$$

$$\sum_{i=1}^n i \left(\sum_{j=1}^{j=i} u_j \right) = c \pmod{(2n+1)}, \quad (21)$$

where $n = |\mathbf{u}|$.

Theorem 2. The code $C_{(a,b,c)}$ can correct all error patterns shown in the s_j column of Tables I and II.

Since (s_1, \dots, s_k) are weight indicators of (μ_1, \dots, μ_k) , the 0s in (s_1, \dots, s_k) and (μ_1, \dots, μ_k) coincide. However, if a 1 is deleted from a run of 1s in s_j , we will not be able to identify which symbol is deleted from μ_j . This means that after recovering s_j from s_j'' we can recover μ_j only in certain cases, specifically, those marked by $(\$)$ in Table I and Table II. Interestingly, the errors not corrected by recovering $s_j, j \in [k]$ are marked by $(*)$, indicating that they occur only for a single value of j . Hence, to

Table I

THE CHANGES IN μ_j AND s_j , $j \in [k]$ AS A RESULT OF EXACT AND NOISY DUPLICATIONS, WHEN THE POSITION OF THE SUBSTITUTION IN \mathbf{x}'' SATISFIES $k < p \leq (|\mathbf{x}''| - k)$. HERE $a, b, c \in \Sigma_q$, $d \in \Sigma_2$, $\bar{a} = -a$, AND $a, b \neq 0$. FURTHERMORE, $\Lambda \rightarrow \mathbf{u}$ AND $\mathbf{u} \rightarrow \Lambda$ REPRESENT INSERTION AND DELETION OF THE STRING \mathbf{u} , RESPECTIVELY. ROWS MARKED BY (*) INDICATE THAT THIS TYPE OF ERROR OCCURS FOR AT MOST ONE VALUE OF $j \in [k]$. ROWS MARKED BY (\$), RELATED TO ERROR-CORRECTION CODE, ARE DISCUSSED IN THE NEXT SECTION

$ \mu'' - \mu $	$\mu \rightarrow \mu''$	$\mu_j \rightarrow \mu_j''$	$s_j \rightarrow s_j''$
$+2k$	insert $0^{j-1}a0^{k-j}$ and $0^{t-1}(0-a)0^{k-t}$	$\Lambda \rightarrow a\bar{a}$ (*) $\Lambda \rightarrow 00$ (\$) $c \rightarrow 0c0$ (\$)	$\Lambda \rightarrow 11$ $\Lambda \rightarrow 00$ $d \rightarrow 0d0$
$+k$	insert $0^{j-1}a0^{k-j}$ and substitute $b_i \rightarrow (b_i - a)$	$c \rightarrow a(c-a), c \neq a$ (*) $a \rightarrow a0$ ($\Lambda \rightarrow 0$) (\$) $\Lambda \rightarrow 0$ (\$)	$0 \rightarrow 11, 1 \rightarrow 11$ $1 \rightarrow 10$ ($\Lambda \rightarrow 0$) $\Lambda \rightarrow 0$
	substitute $0 \rightarrow a$ and insert $0^{t-1}(0-a)0^{k-t}$	$0 \rightarrow a\bar{a}$ (*) $\Lambda \rightarrow 0$ (\$)	$0 \rightarrow 11$ $\Lambda \rightarrow 0$
0	insert $0^{j-1}a0^{k-j}$ and delete $0^{t-1}a0^{k-t}$ with a at the same position	$b0 \rightarrow 0b$ (\$) stay same	$10 \rightarrow 01$ stay same
	substitute $0 \rightarrow a$ and $b_i \rightarrow (b_i - a)$ with distance k	$0c \rightarrow a(c-a)$ (*, \$) stay same	$00 \rightarrow 11, 01 \rightarrow 11, 01 \rightarrow 10$ stay same
$-k$	substitute $0 \rightarrow a$ and delete $0^{t-1}a0^{k-t}$	$0 \rightarrow \Lambda$ (\$)	$0 \rightarrow \Lambda$

Table II

THE CHANGES IN μ_j AND s_j , $j \in [k]$ AS A RESULT OF EXACT AND NOISY DUPLICATION, WHEN THE POSITION OF THE SUBSTITUTION IN \mathbf{x}'' SATISFIES $(|\mathbf{x}''| - k) < p \leq |\mathbf{x}''|$. HERE THE NOTATION IS THE SAME AS THAT OF TABLE I

$ \mu'' - \mu $	$\mu \rightarrow \mu''$	$\mu_j \rightarrow \mu_j''$	$s_j \rightarrow s_j''$
$+k$	insert $0^{j-1}a0^{k-j}$	$\Lambda \rightarrow a$ (*) $\Lambda \rightarrow 0$ (\$)	$\Lambda \rightarrow 1$ $\Lambda \rightarrow 0$
0	substitute $0 \rightarrow a$	$0 \rightarrow a$ (*, \$) stay same	$0 \rightarrow 1$ stay same

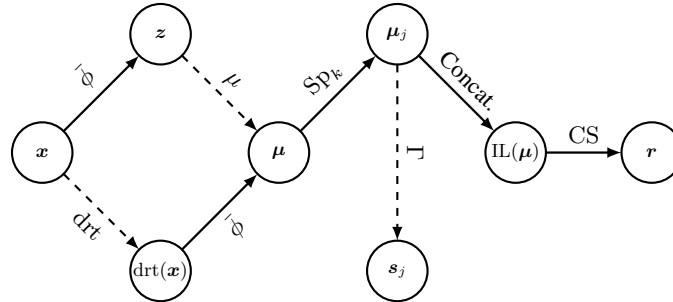


Figure 1. The various mapping used in the paper. “Concat.” stands for concatenation. Solid edges indicate invertible mappings, where we have assumed $x_1 \cdots x_k$ is known, since these symbols are not affected by the channel. The mapping μ is generally non-invertible, but in our constructions, since we assume \mathbf{x} is irreducible, if we recover $\mu = \mu(\mathbf{x})$, we can recover \mathbf{x} .

correct these errors, we apply the code constraints to the concatenation of $\boldsymbol{\mu}_j, j \in [k]$, rather than to each $\boldsymbol{\mu}_j$ separately.

Construction 4. Define $C_{nd} \subseteq \Sigma_q^n$ as

$$C_{nd} = \{\mathbf{x} \in \text{Irr}(n) \cap \Sigma_q^n \mid \boldsymbol{\mu} = \mu(\bar{\phi}(\mathbf{x})), \quad (22)$$

$$\boldsymbol{\mu}_j = \text{Sp}_k(\boldsymbol{\mu}, j), \mathbf{s}_j = \Gamma(\boldsymbol{\mu}_j), \quad (23)$$

$$\mathbf{s}_j \in C_{VT}(a_j, 2|\mathbf{s}_j| + 3), \quad (24)$$

$$\sum_{i=1}^{|\mathbf{s}_j|} i \left(\sum_{t=1}^{t=i} s_{jt} \right) = c_j \bmod (2|\mathbf{s}_j| + 1), \quad (25)$$

$$\sum_{j=1}^k \sum_{i=1}^{|\mathbf{s}_j|} s_{ji} = b \bmod 5, \quad (26)$$

$$\text{Od}(\text{IL}(\boldsymbol{\mu})) \in C_{Tq}(\bar{a}_1, \bar{b}_1, \lceil \frac{n-k}{2} \rceil), \quad (27)$$

$$\text{Ev}(\text{IL}(\boldsymbol{\mu})) \in C_{Tq}(\bar{a}_2, \bar{b}_2, \lceil \frac{n-k}{2} \rceil), \quad (28)$$

$$\text{CS}(\text{IL}(\boldsymbol{\mu})) \in C_{Tq}(\bar{a}_3, \bar{b}_3, n-k), \quad (29)$$

where $j, a_j, c_j, b, \bar{a}_i, \bar{b}_i$ are integers satisfying $j \in [k]$, $0 \leq a_j \leq 2(|\mathbf{s}_j| + 1)$, $0 \leq c_j \leq 2|\mathbf{s}_j|$, $0 \leq b \leq 4$, $0 \leq \bar{a}_1, \bar{a}_2, \bar{a}_3 < q$, $0 \leq \bar{b}_1, \bar{b}_2 \leq \lfloor \frac{n-k}{2} \rfloor$, and $0 \leq \bar{b}_3 < n-k$.

In Construction 4, the constraints (24), (25), and (26) play the same role as the code in Construction 3, and the constraints (27), (28), and (29) can correct the error patterns of $\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k\}$ not marked by (\$) in Table I and Table II. The constraint (24) corrects one insertion/deletion or two insertions of 0s or 1s in adjacent positions over Σ_2 . The constraint (25) corrects one transposition of $\{0, 1\}$ in two adjacent positions. The constraint (26) is a weight-indicating equation for $\{s_1, \dots, s_k\}$. The constraints (27), (28), and (29) can correct one insertion/deletion in $\text{Od}(\text{IL}(\boldsymbol{\mu}))$, $\text{Ev}(\text{IL}(\boldsymbol{\mu}))$ and $r = \text{CS}(\text{IL}(\boldsymbol{\mu}))$ over Σ_q , respectively.

Theorem 3. *The error-correcting code C_{nd} proposed in Construction 4 can correct infinitely many exact k -TD and up to one k -ND errors. There exists one such code with size*

$$\frac{|\text{Irr}(n)|}{5q^3 \lceil \frac{n-k}{2} \rceil 2(4 \lceil \frac{n}{k} \rceil^2 - 1)^k (n-k)} \leq |C_{nd}| \leq |\text{Irr}(n)|. \quad (30)$$

For a code $C \subseteq \Sigma_q^n$, define its rate $R_n(C)$ as $\frac{1}{n} \log_q |C|$. From (30),

$$\frac{1}{n} \log_q |\text{Irr}(n)| - \frac{(2k+3)}{n} \log_q n - \frac{2k}{n} \log_q 2 - \frac{3}{n} - \frac{1}{n} \log_q 5 \leq R_n(C_{nd}) \leq \frac{1}{n} \log_q |\text{Irr}(n)|.$$

It can then be shown that if $q+k \geq 4$, as $n \rightarrow \infty$,

$$\begin{aligned} R_n(C_{nd}) &= \frac{1}{n} \log_q |\text{Irr}(n)| + o(1) \\ &= 1 - \frac{(q-1) \log_q e}{q^{k+2}} + o(1). \end{aligned} \quad (31)$$

Since this is asymptotically the same as the rate of the code correcting only k -TDs [3], the code proposed here is asymptotically optimal. Furthermore, it outperforms the code proposed in [11] for correcting a single unrestricted substitution in addition to correcting many k -TDs.

REFERENCES

- [1] Y. M. Chee, J. Chrisnata, H. M. Kiah, and T. T. Nguyen, "Deciding the Confusability of Words under Tandem Repeats," *arXiv:1707.03956 [math]*, Jul. 2017.
- [2] L. Dolecek and V. Anantharam, "Repetition error correcting sets: Explicit constructions and prefixing methods," *SIAM Journal on Discrete Mathematics*, vol. 23, no. 4, pp. 2120–2146, 2010.
- [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] —, "Noise and uncertainty in string-duplication systems," in *2017 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2017, pp. 3120–3124.
- [5] M. Kovacevic 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.
- [6] A. Lenz, A. Wachter-Zeh, and E. Yaakobi, "Duplication-Correcting Codes," *arXiv:1712.09345 [cs, math]*, Dec. 2017.
- [7] H. Lou, M. Schwartz, and F. Farnoud, "Evolution of N-gram Frequencies under Duplication and Substitution Mutations," in *IEEE Int. Symp. Information Theory (ISIT)*, Jun. 2018.
- [8] 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.
- [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] N. J. Sloane, "On single-deletion-correcting codes," *Codes and designs*, vol. 10, pp. 273–291, 2000.
- [11] Y. Tang, Y. Yehezkeally, M. Schwartz, and F. Farnoud, "Single-error detection and correction for duplication and substitution channels," in *2019 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2019.
- [12] G. Tenengolts, "Nonbinary codes, correcting single deletion or insertion," *IEEE Transactions on Information Theory*, vol. 30, no. 5, pp. 766–769, 1984.
- [13] 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.
- [14] Y. Yehezkeally and M. Schwartz, "Reconstruction codes for DNA sequences with uniform tandem-duplication errors," in *2018 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2018, pp. 2535–2539.