

Asymptotically Optimal Codes Correcting One Substring Edit

Yuting Li*, Yuanyuan Tang[†], Hao Lou[†], Ryan Gabrys[‡], and Farzad Farnoud^{*†}

*Computer Science, University of Virginia, USA, mzy8rp@virginia.edu

[†]Electrical & Computer Engineering, University of Virginia, USA,
{yt5tz, hl2nu, farzad}@virginia.edu

[‡]Calit2, University of California-San Diego, USA, rgabrys@ucsd.edu

Abstract—The substring edit error is the operation of replacing a substring u of x with another string v , where the lengths of u and v are bounded by a given constant k . It encompasses localized insertions, deletions, and substitutions within a window. Codes correcting one substring edit have redundancy at least $\log n + k$. In this paper, we construct codes correcting one substring edit with redundancy $\log n + O(\log \log n)$, which is asymptotically optimal. The full version of this paper is available online.¹

I. INTRODUCTION

In data transmission and storage, especially in data storage in DNA [1] and other emerging media, the need for robust error correction for a diverse set of errors is an important and challenging problem. This paper presents a family of asymptotically optimal codes capable of correcting a burst of localized edits, encompassing combinations of insertions, deletions, and substitutions. In addition to their application in error-correction, the codes provide a promising approach to the challenge of synchronizing multiple copies of related files, where bursts of edits are common. The data synchronization problem, also known as document exchange, has been explored in various works demonstrating that error-correcting codes, especially those accommodating file edits, can substantially reduce the communication required to maintain file consistency [2].

Specifically, we focus on correcting a single burst of errors of length at most k , which we refer to as a *k -substring edit*. Such an error may be the result of insertions, deletions, or substitutions, occurring within a bounded interval of our strings. Formally, a k -substring edit in a string x is the operation of replacing a substring u of x with another string v , where $|u|, |v| \leq k$. An analysis of real and simulated data in [3] reveals that substring edits are common in file synchronization and DNA storage, where viewing errors as substring edits leads to smaller redundancies.

In this paper, we assume that k is a fixed constant. Notice that for the setting where we replace a string u with another string v of the same length (i.e., $|u| = |v|$), a substring edit is equivalent to a burst of substitutions. Furthermore, if v is the empty string, then a substring edit is equivalent to a burst of deletions. Analogously, if u is the empty string, then the substring edit is a burst of insertions. Despite the fact that

these specialized types of substring edit have each received significant attention in the past, the problem of constructing codes for substring edits has received relatively less attention, with the exception of [3] and [4]. A critical observation in this context is that the ability to correct an arbitrarily long burst of deletions does not enable correcting a k -substring edit [3, Lemma 1].

Using simple counting arguments, it can be shown that a code correcting one substring edit has redundancy at least $\log n + k$ and at most roughly $2 \log n$. In [3], the authors present an efficient construction for a code that corrects one substring edit with redundancy roughly $2 \log n$, which matches the existential upper bound of redundancy. However, the question of whether it is possible to construct a code with less than $2 \log n$ bits of redundancy remained open. In this paper, we provide an affirmative answer to this question and construct codes correcting one substring edit with redundancy roughly $\log n + O(\log \log n)$, i.e., at most $O(\log \log n)$ bits larger than the optimal redundancy. Hence, the codes we propose are asymptotically optimal in terms of redundancy and improve the state-of-the-art redundancy by a factor of 2.

The rest of the paper is organized as follows: In Section II, we introduce basic definitions and notation, review our techniques, and compare with related works. In Section III, we construct asymptotically optimal codes correcting one substring edit.

II. NOTATION, OVERVIEW, AND RELATED WORK

A. Basic Definitions and Notations

We will deal with two types of strings, strings over $\{0, 1\}$ and strings over the set of nonnegative integers. When necessary to identify the alphabet, we refer to a string as a *binary string* or an *integer string*, respectively. Strings are denoted by bold symbols, e.g., $\mathbf{u} = u_1 \cdots u_n$. Let $[i, j]$ represent the integers $i, i + 1, \dots, j$. Furthermore, the set $[1, j]$ is denoted as $[j]$. The substring of \mathbf{u} starting at position i and ending at position j is denoted by $\mathbf{u}_{[i, j]}$. We use $|\mathbf{u}|$ to denote the length of \mathbf{u} . For $\mathbf{z} \in \mathbb{N}^*$ and a positive integer A , we say that \mathbf{z} is A -bounded if each symbol of \mathbf{z} is at most A .

For binary strings $\mathbf{x}, \mathbf{y} \in \{0, 1\}^*$, we denote a k -substring edit by $\mathbf{x} \rightarrow \mathbf{y}$. If $\mathbf{x} \rightarrow \mathbf{y}$ and $\mathbf{x} \neq \mathbf{y}$, then we may also write $\mathbf{x} \not\rightarrow \mathbf{y}$. In our construction, substring edits over integer strings also arise. For clarity, we use a different notation for

¹<https://www.ece.virginia.edu/~ffh8x/papers/ecss.pdf>

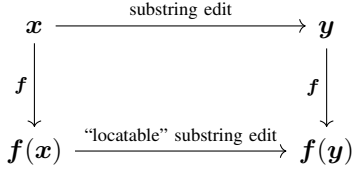


Figure 1: converting a substring edit to a “locatable” substring edit

these: For two integer strings $z, w \in \mathbb{N}^*$, we let $z \Rightarrow w$ denote a substring edit (not necessarily a k -substring edit) and let $z \xrightarrow{j:a,b} w$ denote the substring edit that deletes a substring of length a starting in position j , and inserts a substring of length b in its place. This operation is a $\max(a, b)$ -substring edit.

A partition \mathcal{P} of x is a rule under which one divides x into parts. We use $n_{\mathcal{P}}(x)$ to denote the number of parts and $x^{\mathcal{P}}$ to denote the vector $(x_1, x_2, \dots, x_{n_{\mathcal{P}}(x)})$, where x_i , $1 \leq i \leq n_{\mathcal{P}}(x)$, denotes a part. We call a string (\mathcal{P}, δ) -dense if all elements of $x^{\mathcal{P}}$ have length at most δ . Suppose f is a function on strings. Then we use $f^{\mathcal{P}}(x)$ to denote the string $f(x_1)f(x_2) \cdots f(x_{n_{\mathcal{P}}(x)})$. When the partition \mathcal{P} is clear from the context, we may simply use $\mathbf{f}(x)$ to denote $f^{\mathcal{P}}(x)$. (Note that $f(x)$, in contrast to $\mathbf{f}(x)$, is simply the result of applying f to x .) For a (binary or integer) string z , its VT sketch is defined as $\text{VT}(z) = \sum_{i=1}^{|z|} iz_i$.

B. Overview of the Techniques

The overview of our approach to correcting one substring edit error is given in Figure 1. Here, x is the input to the channel while y is the output. The first step in correcting the error is to approximately identify its position. To achieve this, we rely on a key insight that goes back to the pioneering work of Levenshtein [5] on correcting a burst of at most 2 deletions. Namely, while the VT sketch applied to x can locate only a single deletion, it can be more powerful when applied to carefully designed functions of x . Previous work has used this observation to correct deletions within a window, as well as to correct one deletion or one adjacent transposition [6][7][8]. Our main contribution in this direction is to introduce a novel mapping that, when used in conjunction with a VT sketch, can locate the position of any k -substring edit to within a small interval. In particular, this mapping allows us to locate a burst of substitutions, a challenging task since such an edit may leave the sequence composed of lengths between “markers” untouched.

To construct the aforementioned mapping, denoted in Figure 1 by \mathbf{f} , we identify a set of conditions on a substring edit over an integer string that, if satisfied, the location of the substring edit can be approximately identified using the VT sketch and some other small amount of information (Lemma 2). Such a substring edit is called “locatable”, a term that is defined precisely in Definition 1. Specifically, our approach will be to find a partition rule \mathcal{P} and a function f such that if $x \not\Rightarrow y$, then $\mathbf{f}(y) = \mathbf{f}^{\mathcal{P}}(y)$ is obtained from $\mathbf{f}(x) = \mathbf{f}^{\mathcal{P}}(x)$ by a locatable substring edit. Then we restrict

x to be (\mathcal{P}, δ) -dense. If $x \not\Rightarrow y$, we treat $\mathbf{f}^{\mathcal{P}}(x)$ and $\mathbf{f}^{\mathcal{P}}(y)$ as the z and w in Lemma 2, and determine the edit position of this locatable substring edit to within an interval of length $O(\delta) = O(\log n)$. Since we restrict x to be (\mathcal{P}, δ) -dense, which means $|x_i| \leq \delta$ for any i , we can determine the original edit position to within an interval of length $O(\delta^2)$ (a technique used in [7]), which is then corrected. The case in which $x = y$ will also be handled in Theorem 11. The VT sketch used in our construction will require a redundancy of about $\log n$ bits, while the other components will have negligible redundancy.

C. Related Works

Codes correcting localized errors have been a subject of extensive study over time. Fire codes [9] are designed to construct a burst of k substitutions with roughly $\log n$ redundancy. Several works, including [4], [10], [11], have construct burst-deletion correcting codes, where the length of the bursts is known in advance. More related to this work are codes that can correct a variable number of deletions. In a pioneering work by Levenshtein’s work [5], optimal codes are constructed to correct a burst of at most two consecutive deletions.

For longer deletions and non-consecutive but localized deletions, several works [3], [6]–[8] aim to first approximately locate the deletion. Their method uses the VT sketch of an integer string that is related to the original binary string to locate the approximate deletion position. Specifically, in [6], [7], and [3], each symbol of the integer string is the length of each part of the original binary string with respect to some partition. Using this method, Lenz et al. [6] propose codes correcting a variable-length burst of deletions with nearly optimal redundancy, but the deletions still need to be consecutive. Bitar et al. [7] later construct nearly optimal codes that remove the restriction of consecutive deletions, so that it could correct deletions within a window that do not necessarily occur consecutively. Also using VT sketch of the length integer string, Tang et al. [3] construct codes correcting a substring edit, but the VT sketch of the length integer string can locate the edit position only when the substring edit changes the length of the original binary string. So the construction from [3] needs another $\log n$ bits to handle the other cases, and the total redundancy is about $2 \log n$ bits, which is not optimal. By using a different integer string, Gabrys et al. [8] further construct codes correcting one deletion or one adjacent transposition with near-optimal redundancy, but their approach does not address the setup where additional edits may occur.

III. ASYMPTOTICALLY OPTIMAL CODES CORRECTING ONE SUBSTRING EDIT

In this section, we present the code construction, prove its correctness, and determine its redundancy. Based on our earlier discussion in Subsection II-B, the construction has three interrelated components: *a*) identifying a set of conditions that make a substring edit over an integer string “locatable” and showing that the position of a locatable substring edit can be found approximately; *b*) constructing a partition rule \mathcal{P} and a mapping f such that if $x \not\Rightarrow y$, then $\mathbf{f}^{\mathcal{P}}(y)$ is

obtained from $\mathbf{f}^P(\mathbf{x})$ through a “locatable” substring edit; and c) correcting the edit in \mathbf{y} and recovering \mathbf{x} . We discuss each in a subsection below. In the next subsection, we discuss each of these components. Then we prove the correctness of the overall approach.

A. Components of the Code Construction

a) *Locatable substring edits*: We will now define the notion of locatable substring edit in Definition 1, and prove Lemma 2, which roughly says that for two δ -bounded integer strings \mathbf{z} and \mathbf{w} , if \mathbf{w} is obtained from \mathbf{z} by a locatable substring edit at position j , then given \mathbf{w} , $\text{VT}(\mathbf{z})$, and some other little information, one can determine j to within an interval of length $O(\delta)$.

Definition 1. For two positive integer strings \mathbf{z} and \mathbf{w} , we say that \mathbf{w} is obtained from \mathbf{z} by a **locatable** K -substring edit at position j if $\mathbf{z} \xrightarrow{j:a,b} \mathbf{w}$, where $a, b \leq K$, and

- 1) The difference of the sum of \mathbf{z} and the sum of \mathbf{w} is bounded by all w_i , namely $\left| \sum_{j=1}^{|\mathbf{z}|} z_j - \sum_{j=1}^{|\mathbf{w}|} w_j \right| < w_i$ for all i 's.
- 2) $|\mathbf{z}| \neq |\mathbf{w}|$ or $\sum_{j=1}^{|\mathbf{z}|} z_j \neq \sum_{j=1}^{|\mathbf{w}|} w_j$.

Now we show that for A -bounded strings \mathbf{z} and \mathbf{w} , if \mathbf{w} is obtained from \mathbf{z} by a locatable K -substring edit at position j , then one can determine j to within an interval of length $O(A)$ given $\text{VT}(\mathbf{z}) \bmod O(m)$ (for some parameter m specified in the lemma statement) and some other little information of \mathbf{z} .

Lemma 2. For A -bounded strings \mathbf{z} and \mathbf{w} , if \mathbf{w} is obtained from \mathbf{z} by a locatable K -substring edit at position j , then one can determine j to within an interval of length $6K^2A$ given \mathbf{w} , $|\mathbf{z}| - |\mathbf{w}|$, $\sum_{i=1}^{|\mathbf{z}|} z_i - \sum_{i=1}^{|\mathbf{w}|} w_i$, and $\text{VT}(\mathbf{z}) \bmod m$, where $m > 8K^2 \left(\sum_{i=1}^{|\mathbf{z}|} z_i + A \right)$.

Proof. Suppose $\mathbf{z} \xrightarrow{j:a,b} \mathbf{w}$, where $a, b \leq K$ and $j \in [1, |\mathbf{z}| - a + 1]$. Let $\eta(v, \mathbf{z}, \mathbf{w}) := (|\mathbf{z}| - |\mathbf{w}|) \sum_{i=v}^{|\mathbf{w}|} w_i + (\sum_{i=1}^{|\mathbf{z}|} z_i - \sum_{i=1}^{|\mathbf{w}|} w_i) v$, where v is an integer. Note that since we know \mathbf{w} , $|\mathbf{z}| - |\mathbf{w}|$, and $\sum_{i=1}^{|\mathbf{z}|} z_i - \sum_{i=1}^{|\mathbf{w}|} w_i$, we know the expression of $\eta(\cdot, \mathbf{z}, \mathbf{w})$. Let $\Delta\text{VT} := \text{VT}(\mathbf{z}) - \text{VT}(\mathbf{w})$. We will prove the following.

- 1) $\eta(\cdot, \mathbf{z}, \mathbf{w})$ is strictly monotone on $[1, |\mathbf{w}| + 1]$.
- 2) $|\Delta\text{VT} - \eta(j, \mathbf{z}, \mathbf{w})| < 3K^2A$.
- 3) $|\Delta\text{VT}| \leq 4K^2 \left(\sum_{i=1}^{|\mathbf{z}|} z_i + A \right)$.

For 1), note that

$$\eta(v, \mathbf{z}, \mathbf{w}) - \eta(v-1, \mathbf{z}, \mathbf{w}) = (|\mathbf{z}| - |\mathbf{w}|)w_v + \sum_{i=1}^{|\mathbf{z}|} z_i - \sum_{i=1}^{|\mathbf{w}|} w_i.$$

If $|\mathbf{z}| \neq |\mathbf{w}|$, since \mathbf{w} is obtained from \mathbf{z} by a locatable K -substring edit, we have $\left| \sum_{i=1}^{|\mathbf{z}|} z_i - \sum_{i=1}^{|\mathbf{w}|} w_i \right| < w_v$. So $\eta(\cdot, \mathbf{z}, \mathbf{w})$ is strictly monotone on $[1, |\mathbf{w}| + 1]$. If $|\mathbf{z}| = |\mathbf{w}|$, since \mathbf{w} is obtained from \mathbf{z} by a locatable K -substring edit, we have $\left| \sum_{i=1}^{|\mathbf{z}|} z_i - \sum_{i=1}^{|\mathbf{w}|} w_i \right| \neq 0$. So, $\eta(\cdot, \mathbf{z}, \mathbf{w})$ is strictly monotone on $[1, |\mathbf{w}| + 1]$.

For 2), since $\mathbf{z} \xrightarrow{j:a,b} \mathbf{w}$, we have

$$\begin{aligned} \Delta\text{VT} &= f(j, \mathbf{z}, \mathbf{w}) + \sum_{i=1}^{a-1} iz_{j+i} - \sum_{i=1}^{b-1} iw_{j+i} \\ &\quad - (|\mathbf{z}| - |\mathbf{w}|) \sum_{i=j}^{j+b-1} w_i. \end{aligned}$$

Since \mathbf{z} and \mathbf{w} are A -bounded strings, we have $\left| \sum_{i=1}^{a-1} iz_{j+i} \right| \leq a^2A$, $\left| \sum_{i=1}^{b-1} iw_{j+i} \right| \leq b^2A$, and $\left| \sum_{i=j}^{j+b-1} w_i \right| \leq bA$. So we have

$$|\eta(j, \mathbf{z}, \mathbf{w}) - \Delta\text{VT}| \leq (a^2 + b^2)A + |a - b|bA < 3K^2A. \quad (1)$$

For 3), since $\eta(\cdot, \mathbf{z}, \mathbf{w})$ is monotone on $[1, |\mathbf{w}| + 1]$,

$$\begin{aligned} |\eta(j, \mathbf{z}, \mathbf{w})| &\leq |\eta(1, \mathbf{z}, \mathbf{w})| + |\eta(|\mathbf{w}| + 1, \mathbf{z}, \mathbf{w})| \\ &\leq K^2 \left(\sum_{i=1}^{|\mathbf{z}|} z_i + A \right). \end{aligned} \quad (2)$$

Thus, by (1) and (2), we have

$$\begin{aligned} |\Delta\text{VT}| &\leq |\Delta\text{VT} - \eta(j, \mathbf{z}, \mathbf{w})| + |\eta(j, \mathbf{z}, \mathbf{w})| \\ &\leq 4K^2 \left(\sum_{i=1}^{|\mathbf{z}|} z_i + A \right). \end{aligned}$$

Now since we know $\text{VT}(\mathbf{z}) \bmod m$, we can find ΔVT . Since $\eta(\cdot, \mathbf{z}, \mathbf{w})$ is strictly monotone on $[1, |\mathbf{w}| + 1]$ and $|\Delta\text{VT} - \eta(j, \mathbf{z}, \mathbf{w})| < 3K^2A$, we can determine j to within an interval of length $6K^2A$. \square

b) *Constructing a mapping to convert a general substring edit to a locatable edit*: Our goal here is to construct a partition \mathcal{P} and a function f over strings such that if $\mathbf{x} \xrightarrow{\cdot} \mathbf{y}$, then $\mathbf{f}(\mathbf{y}) = \mathbf{f}^P(\mathbf{y})$ is obtained from $\mathbf{f}(\mathbf{x}) = \mathbf{f}^P(\mathbf{x})$ by a locatable substring edit. We first give the partition rule \mathcal{P} . Let $\mathbf{p} = 0^k 1^k$, which we call a pattern.

We define the partition rule \mathcal{P} as follows. For $\mathbf{x} \in \{0, 1\}^*$, let $\mathbf{x}^P = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n_P(\mathbf{x})})$, where the starting position of \mathbf{x}_1 is the leftmost bit of \mathbf{x} and the starting position of \mathbf{x}_i ($i \geq 2$) is the leftmost bit of the $(i-1)$ th \mathbf{p} in \mathbf{x} .

Example 3. Let $k = 3$, $\mathbf{p} = 0^3 1^3$, and

$$\mathbf{x} = \overbrace{01001}^{\mathbf{x}_1} \underbrace{\overbrace{000111}^{\mathbf{x}_2}}_{\mathbf{p}} \underbrace{\overbrace{100}^{\mathbf{x}_3} \overbrace{000111}^{\mathbf{x}_4}}_{\mathbf{p}} \underbrace{\overbrace{101001}^{\mathbf{x}_5} \overbrace{000111}^{\mathbf{x}_6}}_{\mathbf{p}} 1.$$

To get \mathbf{x}^P , we first find all the patterns, and break it up at the left side of each pattern. So we get $\mathbf{x}^P = (01001, 000111100, 000111101001, 0001111)$.

We now explain the main idea behind the construction of the function f . For any lowercase function u over strings, we

use the uppercase function U to denote the sum of all the symbols of $\mathbf{u}(x)$, namely

$$U(\mathbf{x}) := \sum_{i=1}^{n_{\mathcal{P}}(\mathbf{x})} u(\mathbf{x}_i) = \sum_{i=1}^{|\mathbf{u}(\mathbf{x})|} u(\mathbf{x}_i),$$

and use ΔU to denote $U(\mathbf{x}) - U(\mathbf{y})$. For simplicity, we only consider the second condition of the locatable substring edit (see the full version for a complete treatment). Applied to the substring edit $f(\mathbf{x}) \Rightarrow f(\mathbf{y})$, the second condition of Definition 1 requires that if $\mathbf{x} \not\Rightarrow^{\neq} \mathbf{y}$, then $|\mathbf{f}(\mathbf{x})| \neq |\mathbf{f}(\mathbf{y})|$ or $\Delta F \neq 0$. Note that if the number of patterns changes ($n_{\mathcal{P}}(\mathbf{x}) \neq n_{\mathcal{P}}(\mathbf{y})$), then $|\mathbf{f}(\mathbf{x})| \neq |\mathbf{f}(\mathbf{y})|$. Hence, we assume that the number of patterns does not change. We want to ensure that $\Delta F = F(\mathbf{x}) - F(\mathbf{y}) \neq 0$. To better explain our approach, we first outline two different approaches that have increasing levels of complexity before arriving at our actual solution, which is described as the third attempt.

First attempt. For a string x , let $l(x)$ represent the length of x . Note that $L(x)$ is also the length of x . In our first attempt, we let $f(\cdot) = l(\cdot)$. Then $F(x) = L(x)$. In this case, if the lengths of the deleted and inserted strings differ, then $\Delta F \neq 0$. However, $\Delta F = 0$ occurs when the substring edit is a burst of substitutions. In fact, in [3], the authors use $f(\cdot) = l(\cdot)$, and use another $\log n$ bits of redundancy to handle a burst of substitutions separately. This leads to a redundancy of $2 \log n$ bits. In this work, we will introduce two other novel functions, d (second attempt), and n_{1^k} (third attempt), to resolve this issue.

Second attempt. Let $d : \{0, 1\}^n \rightarrow [2^k]$ be a function that can detect a burst of k substitutions, where the j th bit of $d(\mathbf{x})$ is the parity of every k th index of \mathbf{x} starting at position j , namely,

$$d(\mathbf{x}) := \left(\bigoplus_{i \equiv_k 1} x_i, \bigoplus_{i \equiv_k 2} x_i, \dots, \bigoplus_{i \equiv_k k} x_i \right),$$

where \oplus is addition in \mathbb{F}_2 . Note that $|\Delta D| \leq 2^k$. We let $f(\cdot) = d(\cdot) + Cl(\cdot)$, where C is a sufficiently large constant. Then, $\Delta F = \Delta D + C\Delta L$.

Now, if $\Delta L \neq 0$, then $\Delta F \neq 0$ provided that we choose $C > 2^k$. If $\Delta L = 0$, then the edit transforming \mathbf{x} into \mathbf{y} is a burst of substitutions. At first glance, it may seem that in this case $\Delta D \neq 0$, and thus $\Delta F \neq 0$. But while this is the case if this burst of substitutions does not involve any patterns, it may not be the case otherwise. The next example shows that if this burst of substitutions involves some patterns, then ΔD may be 0.

Example 4. Let $k = 3$, $p = 0^3 1^3$, and

$$x = \overbrace{000011}^{x_1} \underbrace{\overbrace{000111}^{x_2}}_p 001111,$$

$$y = \overbrace{000011000}^{y_1} \underbrace{\overbrace{01000111}^{y_2}}_p 1,$$

where the edited substring is in bold and red. Then $\mathbf{l}(\mathbf{x}) = \ell(\mathbf{x}_1)\ell(\mathbf{x}_2) = (6, 12)$, $\mathbf{l}(\mathbf{y}) = (11, 7)$. $\mathbf{d}(\mathbf{x}) = (011, 001) = (3, 1)$, $\mathbf{d}(\mathbf{y}) = (001, 011) = (1, 3)$. Then $\Delta L = 0$ and $\Delta D = 0$, so $\Delta F = 0$. Thus, $\mathbf{f}(\mathbf{x}) \Rightarrow \mathbf{f}(\mathbf{y})$ is not a locatable substring edit.

Third attempt. Recall that the method in the second attempt may not work only if the substring edit is a burst of substitutions that involves some patterns. By careful argument, one can find that the only case where $\Delta F = 0$ in the second attempt is when one pattern $0^k 1^k$ is shifted because of this burst of substitutions. So, we need a function, say g , such that $\Delta G \neq 0$ in this case, and thus our f will be of the form $f(\cdot) = d(\cdot) + Bg(\cdot) + Cl(\cdot)$, where B is a sufficiently large constant and C is much larger than B . This way, if $\Delta L \neq 0$, then $\Delta F \neq 0$. If $\Delta L = 0$, and $\Delta G \neq 0$, then $\Delta F \neq 0$. If $\Delta L = 0$ and $\Delta G = 0$, then $x \not\rightarrow y$ is a burst of substitutions that does not involve any pattern, then $\Delta D \neq 0$, thus $\Delta F \neq 0$.

Let $n_{1^k}(\mathbf{x})$ be the number of appearances of 1^k in \mathbf{x} . For example, if $\mathbf{x} = 111100111$, then $n_{1^2}(\mathbf{x}) = 5$. We will show that n_{1^k} is our desired g . The next example provides the intuition. Note that $N_{1^k}(\mathbf{x})$ is the sum of $n_{1^k}(\mathbf{x}_i)$ over $\mathbf{x}^{\mathcal{P}} = (\mathbf{x}_1, \dots, \mathbf{x}_{n_{\mathcal{P}}(\mathbf{x})})$.

Example 5. Let $N_{1^L}^L(\mathbf{x})$ and $N_{1^R}^R(\mathbf{x})$ denote the contribution of the left side and the right side with respect to the ‘|’ symbols below. Similar definitions apply to \mathbf{y} . Below, we show two typical cases in which $\mathbf{x} \not\rightarrow \mathbf{y}$ is a burst of substitutions that shifts a pattern. In both cases, we show that $N_{1^k}(\mathbf{x}) \neq N_{1^k}(\mathbf{y})$, implying that $\Delta N_{1^k} \neq 0$.

First, suppose $\mathbf{x} = \alpha 0^\epsilon |0^k 1^k \beta$ and $\mathbf{y} = \alpha |0^k 1^k 1^\epsilon \beta$, for some binary string α and β . Here $0 < \epsilon \leq k$. The burst of substitutions occurs at positions $[|\alpha| + k + 1, |\alpha| + k + \epsilon]$. Below is an example with $k = 10, \epsilon = 7$.

$$\begin{array}{l} x = \alpha \overbrace{0000000}^{0^\epsilon} \overbrace{000\textcolor{red}{0000000}}^{0^k} \overbrace{1111111111}^{1^k} \beta \\ y = \alpha \overbrace{00000000000}^{0^k} \overbrace{\textcolor{red}{1111111}111}^{1^k} \overbrace{11111111}^{1^\epsilon} \beta \end{array}$$

In this case, we can see $N_{1^k}^L(\mathbf{x}) = N_{1^k}^L(\mathbf{y})$, and because of the $1^{k+\epsilon}$ in \mathbf{y} compared to the 1^k in \mathbf{x} , we have $N_{1^k}^R(\mathbf{x}) = N_{1^k}^R(\mathbf{y}) - \epsilon$. So $N_{1^k}(\mathbf{x}) \neq N_{1^k}(\mathbf{y})$, and thus $\Delta N_{1^k} \neq 0$.

Second, suppose $\mathbf{x} = \alpha|0^k1^k0^{k-\epsilon}1^k\beta$ and $\mathbf{y} = \alpha0^{k1}k^{-\epsilon}|0^k1^k\beta$. Here $0 < \epsilon < k$ (when $\epsilon = k$, this edit is also described by the previous case). The burst of substitutions occurs at positions $[|\alpha| + 2k - \epsilon + 1, |\alpha| + 2k]$. An example with $k = 10, \epsilon = 7$ is shown below.

$$\begin{array}{l} x = \alpha \overbrace{0000000000}^{0^k} \overbrace{1111111111}^{1^k} \overbrace{000}^{0^{k-\epsilon}} \overbrace{1111111111}^{1^k} \beta \\ y = \alpha \overbrace{0000000000}^{0^k} \overbrace{111}^{1^{k-\epsilon}} \overbrace{0000000000}^{0^k} \overbrace{1111111111}^{1^k} \beta \end{array}$$

In this case, we can see $N_{1^k}^L(\mathbf{x}) = N_{1^k}^L(\mathbf{y})$, and because there are two 1^k 's between γ and β in \mathbf{x} compared to only one 1^k between γ and β in \mathbf{y} , we have $N_{1^k}^R(\mathbf{x}) = N_{1^k}^R(\mathbf{y}) + 1$. So $N_{1^k}(\mathbf{x}) \neq N_{1^k}(\mathbf{y})$, and thus $\Delta N_{1^k} \neq 0$.

It can be shown that the two cases in the above example are the only two key cases in which a burst of substitutions shifts a pattern (see full version). We can now formally construct f . For a string x , let both $l(x)$ and $L(x)$ represent the length of x , let $n_{1^k}(x)$ be the number of 1^k inside x , and let $d(x)$ be as in the second attempt. Let $B = 3(2^k)$, and $C = 40k(2^k)$, and define

$$f(x) := d(x) + Bn_{1^k}(x) + Cl(x). \quad (3)$$

If $x \not\rightarrow y$, then $f(x) \Rightarrow f(y)$ satisfies the second condition of locatability. It is not hard to prove that $f(x) \Rightarrow f(y)$ also satisfies the first condition of locatability (see full version). In summary, we have the following Lemma.

Lemma 6. *If $x \not\rightarrow y$, then $f(y)$ is obtained from $f(x)$ by a locatable 3-substring edit.*

Lemma 7. *If x is (\mathcal{P}, δ) -dense and $x \rightarrow y$, then $f(x)$ and $f(y)$ are $8\delta C$ -bounded strings.*

c) *Correcting an edit with known approximate location:*

We start by adapting a result on systematic codes correcting multiple edits [12] to a hash correcting a k -substring edit, since a k -substring edit can also be viewed as k edits. The adapted result states that a substring edit in a string of length n can be corrected using $4k \log n + o(\log n)$ bits of redundancy. Although our goal is to achieve about $\log n$ redundancy, this result will be useful in our construction.

Lemma 8 (c.f. [12]). *There exists a hash $\phi_n : \{0, 1\}^n \rightarrow \{0, 1\}^{4k \log n + o(\log n)}$ that can correct a single k -substring edit. Namely, if $r \in \{0, 1\}^n$ and $r \rightarrow s$, then r can be recovered from s and $\phi_n(r)$.*

Let $\psi_A(x) : \{0, 1\}^n \rightarrow [2^{4k \log A + o(\log A)}]$ be defined as

$$\psi_A(x) = \left(\left(\sum_{i \text{ odd}} \phi_A(x_i) \right) \bmod 2^{4k \log A + o(\log A)}, \left(\sum_{i \text{ even}} \phi_A(x_i) \right) \bmod 2^{4k \log A + o(\log A)} \right),$$

where the x_i result from partitioning x into parts of length A (note that in the definition of the map $\psi_A(x)$ we are not partitioning x according to $x^{\mathcal{P}}$).

Lemma 9 is an adaptation of a technique used in [7] to our substring edit scenario, which will allow us to correct a substring edit whose position is approximately known.

Lemma 9 (c.f. [7]). *If $x \rightarrow y$, then one can recover x given y and $\psi_A(x)$ and an interval of length A that contains the edit position.*

B. Substring edit-correcting code and its redundancy

We now put the components previously discussed together to construct the error-correcting code. Let \mathcal{P} be as in Sec-

tion III-A, and let f and C be the one in equation (3). Let

$$h(x) := \left(\text{VT}(f(x)) \bmod 2000Cn, \left(\sum_{i=1}^{n_{\mathcal{P}}(x)} f_i(x) \right) \bmod 10Ck, n_{\mathcal{P}}(x) \bmod 5 \right).$$

For $c_1, c_2 \in \mathbb{N}$, let

$$\mathcal{C} = \{x \text{ is } (\mathcal{P}, \delta)\text{-dense: } h(x) = c_1, \psi_{O(\delta^2)}(x) = c_2\}.$$

The next lemma will be used to show that the redundancy arising from restricting x to be (\mathcal{P}, δ) -dense is negligible.

Lemma 10. *Let $\delta = k2^{2k+3} \log n$. For a uniformly randomly chosen $x \in \{0, 1\}^n$,*

$$\Pr[x \text{ is not } (\mathcal{P}, \delta)\text{-dense}] < 1/n.$$

Theorem 11. *\mathcal{C} corrects one k -substring edit. Moreover, for $\delta = k2^{2k+3} \log n$, it has redundancy $\log n + O(\log \log n)$.*

Proof sketch: We first prove the redundancy. Note that $h(x)$ needs $\log n + O(1)$ redundancy and $\psi_{O(\delta^2)}(x)$ needs $O(\log \log n)$ redundancy. By Lemma 10, restricting x to be (\mathcal{P}, δ) -dense adds $O(1)$ redundancy. So \mathcal{C} has redundancy $\log n + O(\log \log n)$.

To prove the error-correction ability of the code, we first show that if $x \rightarrow y$, given y and $h(x)$, we can determine whether $x = y$. By Lemma 6 and the definition of locatable substring edit, $x = y$ if and only if $\sum_{i=1}^{n_{\mathcal{P}}(x)} f_i(x) \equiv \sum_{i=1}^{n_{\mathcal{P}}(y)} f_i(y) \bmod 10Ck$ and $n_{\mathcal{P}}(x) \equiv n_{\mathcal{P}}(y) \bmod 5$. These two equations can be checked given y and $h(x)$ to determine whether $x = y$. If $x = y$, then we just recover x by outputting y .

We will now assume $x \not\rightarrow y$ and show that we can determine the edit position to within an interval of length $O(\delta^2)$. By Lemmas 6 and 7, we can treat $f(x)$ and $f(y)$ as the z and w in Lemma 2, and determine the position of the locatable edit to within an interval of length $O(\delta)$. Since x is (\mathcal{P}, δ) -dense, we can further determine the position of the original edit to within an interval of length $O(\delta^2)$. In the final step, by Lemma 9, we can recover x by $\psi_{O(\delta^2)}(x)$. ■

REFERENCES

- [1] S. M. H. Tabatabaei Yazdi, Y. Yuan, J. Ma, H. Zhao, and O. Milenkovic, "A Rewritable, Random-Access DNA-Based Storage System," *Scientific Reports*, vol. 5, Sep. 2015. [Online]. Available: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4585656/>
- [2] A. Orlitsky, "Interactive communication: Balanced distributions, correlated files, and average-case complexity," in *[1991] Proceedings 32nd Annual Symposium of Foundations of Computer Science*, Oct. 1991, pp. 228–238.
- [3] Y. Tang, S. Motamen, H. Lou, K. Whitenour, S. Wang, R. Gabrys, and F. Farnoud, "Correcting a substring edit error of bounded length," in *2023 IEEE International Symposium on Information Theory (ISIT)*, 2023, pp. 2720–2725.
- [4] P. A. H. Bours, "Codes for correcting insertion and deletion errors," 1994.
- [5] V. Levenshtein, "Asymptotically optimum binary code with correction for losses of one or two adjacent bits," *Problemy Kibernetiki*, vol. 19, pp. 293–298, 1967.
- [6] A. Lenz and N. Polyanskii, "Optimal codes correcting a burst of deletions of variable length," in *2020 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2020, pp. 757–762.

- [7] R. Bitar, S. K. Hanna, N. Polyanskii, and I. Vorobyev, "Optimal codes correcting localized deletions," in *2021 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2021, pp. 1991–1996.
- [8] R. Gabrys, V. Guruswami, J. Ribeiro, and K. Wu, "Beyond single-deletion correcting codes: Substitutions and transpositions," *IEEE Transactions on Information Theory*, vol. 69, no. 1, pp. 169–186, 2023.
- [9] R. E. Blahut, *Algebraic codes for data transmission*. Cambridge university press, 2003.
- [10] L. Cheng, T. G. Swart, H. C. Ferreira, and K. A. S. Abdel-Ghaffar, "Codes for correcting three or more adjacent deletions or insertions," in *2014 IEEE International Symposium on Information Theory*, 2014, pp. 1246–1250.
- [11] C. Schoeny, A. Wachter-Zeh, R. Gabrys, and E. Yaakobi, "Codes correcting a burst of deletions or insertions," *IEEE Transactions on Information Theory*, vol. 63, no. 4, pp. 1971–1985, 2017.
- [12] J. Sima, R. Gabrys, and J. Bruck, "Optimal systematic t -deletion correcting codes," in *2020 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2020, pp. 769–774.