



Deciding the Confusability of Words under Tandem Repeats in Linear Time

YEOW MENG CHEE, National University of Singapore, Singapore
JOHAN CHRISNATA, HAN MAO KIAH, and TUAN THANH NGUYEN,
Nanyang Technological University, Singapore

Tandem duplication in DNA is the process of inserting a copy of a segment of DNA adjacent to the original position. Motivated by applications that store data in living organisms, Jain et al. (2016) proposed the study of codes that correct tandem duplications to improve the reliability of data storage. We investigate algorithms associated with the study of these codes.

Two words are said to be $\leq k$ -confusable if there exists a sequence of tandem duplications for each word, where each duplication is of length at most k , such that the resulting two words after duplications are equal. For $k = 3$, we demonstrate that the problem of deciding whether two words is ≤ 3 -confusable is linear-time solvable through a characterisation that can be checked efficiently. Combining with previous results, the decision problem is linear-time solvable for $k \leq 3$. We conjecture that this problem is undecidable for $k > 3$.

Using insights gained from the algorithm, we study the size of tandem-duplication codes. We improve the previous known upper bound and then construct codes with larger sizes as compared to the previous constructions. We determine the sizes of optimal tandem-duplication codes for lengths up to 20, develop recursive methods to construct tandem-duplication codes for all word lengths, and compute explicit lower bounds for the size of optimal tandem-duplication codes for lengths from 21 to 30.

CCS Concepts: • **Mathematics of computing** → **Combinatorics on words**; • **Applied computing** → **Bioinformatics**;

Additional Key Words and Phrases: Tandem duplications, DNA-based data storage

ACM Reference format:

Yeow Meng Chee, Johan Chrisnata, Han Mao Kiah, and Tuan Thanh Nguyen. 2019. Deciding the Confusability of Words under Tandem Repeats in Linear Time. *ACM Trans. Algorithms* 15, 3, Article 42 (July 2019), 22 pages. <https://doi.org/10.1145/3338514>

1 INTRODUCTION

At the beginning of the millenium, Lander et al. [9] published a draft sequence of the human genome and reported that more than 50% of the human genome consists of repeated substrings [9]. There are two types of common repeats: *interspersed* repeats and *tandem* repeats. Interspersed

This work is supported in part by the Singapore Ministry of Education under Grant No.: MOE2015-T2-2-086.

Authors' addresses: Y. M. Chee, National University of Singapore, Department of Industrial Systems Engineering and Management, 1 Engineering Drive 2, Singapore 117576; email: pvocym@nus.edu.sg; J. Chrisnata, H. M. Kiah, and T. T. Nguyen, Nanyang Technological University, School of Physical and Mathematical Sciences, 21 Nanyang Link, Singapore 637371; emails: {johanchr001, hmkiah, nguyentu001}@ntu.edu.sg.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2019 Association for Computing Machinery.

1549-6325/2019/07-ART42 \$15.00

<https://doi.org/10.1145/3338514>

repeats are caused by transposons when a segment of DNA is copied and pasted into new positions of the genome. In contrast, tandem repeats are caused by slipped-strand mispairings [13], and they occur when a pattern of one or more nucleotides is repeated and the repetitions are adjacent to each other. For example, consider the word AGTAGTCTGC. The substring AGTAGT is a tandem repeat, and we say that AGTAGTCTGC is generated from AGTCTGC by a *tandem duplication* of length three. Tandem repeats are believed to be the cause of several genetic disorders [3, 17, 18] and this motivated the study of tandem duplication mechanisms in a variety of contexts.

- **Formal languages:** Leopold et al. [10, 11] defined the *unbounded duplication language* and the *k-bounded duplication language* to be the set of words generated by seed word under tandem duplications of unbounded length and tandem duplications of length up to k , respectively. They investigated certain decidability problems involving unbounded duplication languages and showed that all k -bounded duplication languages are context free. Furthermore, they showed that k -bounded duplication language is always regular for any binary seed and $k \geq 1$. However, the k -bounded duplication language is not regular for any *square-free* seed word over an alphabet of at least three letters and $k \geq 4$. More recently, Jain et al. [6] completed this characterization and proved that k -bounded duplication languages are regular for $k \leq 3$.
- **Information theory:** Farnoud et al. [2] introduced the concept of *capacity* to determine average information content of a k -bounded duplication language. Later, Jain et al. [6] introduced the notion of *expressiveness* to measure a language's capability to generate words that contain certain desired substrings. A complete characterization of fully expressive k -bounded duplication languages was provided by Jain et al. for all alphabet sizes and all k .
- **Codes correcting tandem duplications:** Motivated by applications that store data in living organisms [1, 4, 12, 15], Jain et al. [7] proposed the study of codes that correct tandem duplications to improve the reliability of data storage. They investigated various types of tandem duplications and provided optimal code construction for uniform tandem duplication channel and in the case where tandem duplication length is at most two.

In this article, we study the last problem and investigate algorithms associated with these codes. In particular, given two words x and y , we look for efficient algorithms that answer the following question:

When are the words x and y *confusable* under tandem repeats? In other words, are there two sequences of tandem duplications such that the resulting words x' and y' are equal?

Interestingly, as we demonstrate, even for small duplication lengths, the solutions to this question are nontrivial. This is surprising as efficient algorithms are well known for analogous questions in other problems in string matching [14]. Before we give an account of these results, we introduce some necessary notations and provide a formal statement of our problem.

2 PRELIMINARIES

Let $\Sigma_q = \{0, 1, \dots, q-1\}$ be an alphabet of $q \geq 2$ symbols. For a positive integer n , let Σ_q^n denote the set of all strings or words of length n over Σ_q , and let Σ_q^* denote the set of all finite words over Σ_q , or the *Kleene closure* of Σ_q . Given two words $\mathbf{x}, \mathbf{y} \in \Sigma_q^*$, we denote their concatenation by \mathbf{xy} .

We state the *tandem duplication* rules. For nonnegative integers $k \leq n$ and $i \leq n-k$, we define $T_{i,k} : \Sigma_q^n \rightarrow \Sigma_q^{n+k}$ such that

$$T_{i,k}(\mathbf{x}) = \mathbf{u}\mathbf{v}\mathbf{w}, \text{ where } \mathbf{x} = \mathbf{u}\mathbf{v}\mathbf{w}, |\mathbf{u}| = i, |\mathbf{v}| = k.$$

If a finite sequence of tandem duplications of length k is performed to obtain \mathbf{y} from \mathbf{x} , then we say that \mathbf{y} is a k -descendant of \mathbf{x} , or \mathbf{x} is a k -ancestor of \mathbf{y} , and denote this relation by $\mathbf{x} \xrightarrow[k]{*} \mathbf{y}$.

Formally, $\mathbf{x} \xrightarrow[k]{*} \mathbf{y}$ means that for some positive t , there exist t non-negative integers i_1, i_2, \dots, i_t such that $\mathbf{y} = T_{i_t, k} \circ T_{i_{t-1}, k} \circ \dots \circ T_{i_1, k}(\mathbf{x})$.

Given a sequence \mathbf{x} and integer k , we consider the set of words that may be obtained from \mathbf{x} via a finite number of tandem duplications of length k . We define the k -descendant cone of \mathbf{x} to be the set of all k -descendants of \mathbf{x} and denote this cone by $D_k^*(\mathbf{x})$. In other words,

$$D_k^*(\mathbf{x}) \triangleq \left\{ \mathbf{y} \in \Sigma_q^* \mid \mathbf{x} \xrightarrow[k]{*} \mathbf{y} \right\}.$$

Our work studies tandem duplications whose lengths are upper bounded by an integer k , and we extend the previous definitions. Formally, if a finite sequence of tandem duplications of length up to k is performed to obtain \mathbf{y} from \mathbf{x} , then we say that \mathbf{y} is a $\leq k$ -descendant of \mathbf{x} and denote this relation by $\mathbf{x} \xrightarrow[\leq k]{*} \mathbf{y}$. We have analogous definitions of $\leq k$ -descendant, $\leq k$ -ancestor, and $\leq k$ -descendant cone $D_{\leq k}^*(\mathbf{x})$.

Example 2.1. Consider $\mathbf{x} = 01210$ over Σ_3 . Since $T_{1,3}(\mathbf{x}) = 01211210$ and $T_{0,2}(01211210) = 0101211210$, we write $01210 \xrightarrow[3]{*} 01211210$ and $01210 \xrightarrow[\leq 3]{*} 0101211210$. Alternatively, we have that $01211210 \in D_3^*(\mathbf{x})$ and $0101211210 \in D_{\leq 3}^*(\mathbf{x})$.

2.1 Problem Formulation

Motivated by applications that store data on living organisms, Jain et al. [7] looked at the $\leq k$ -descendant cones of a pair of words and asked whether the two cones have a nontrivial intersection. Specifically, we introduce the notion of confusability.

Definition 2.2 (Confusability). Two words \mathbf{x} and \mathbf{y} , are said to be k -confusable if $D_k^*(\mathbf{x}) \cap D_k^*(\mathbf{y}) \neq \emptyset$. Similarly, they are said to be $\leq k$ -confusable if $D_{\leq k}^*(\mathbf{x}) \cap D_{\leq k}^*(\mathbf{y}) \neq \emptyset$.

To design error-correcting codes that store information in the DNA of living organisms, Jain et al. then proposed the use of codewords that are not pairwise confusable.

Definition 2.3 ($\leq k$ -Tandem-Duplication Codes). A subset $C \subset \Sigma_q^n$ is a $\leq k$ -tandem-duplication code if for all $\mathbf{x}, \mathbf{y} \in C$ and $\mathbf{x} \neq \mathbf{y}$, we have that \mathbf{x} and \mathbf{y} are not $\leq k$ -confusable. We say that C is an $(n, \leq k; q)$ -tandem-duplication code or $(n, \leq k; q)$ -TD code.

Therefore, to determine if a set of words is a tandem-duplication code, we need to verify that all pairs of distinct words are not confusable. Hence, we state our problem of interest.

CONFUSABILITY PROBLEM

Instance: Two words \mathbf{x} and \mathbf{y} over Σ_q and an integer k

Question: Are \mathbf{x} and \mathbf{y} $\leq k$ -confusable?

While the confusability problem is a natural question, efficient algorithms are only known for the case where $k \in \{1, 2\}$. We review these results in the next subsection.

2.2 Previous Work

Confusability. We summarize known efficient algorithms that determine whether two words are k -confusable for all k , and whether they are $\leq k$ -confusable for $k \in \{1, 2\}$. We then highlight why these methods cannot be extended for the case $k \geq 3$. Crucial to the algorithms is the concept of irreducible words and roots.

Definition 2.4. A word \mathbf{x} is said to be k -irreducible if \mathbf{x} cannot be deduplicated into shorter words with deduplication of length k . In other words, if $\mathbf{y} \xrightarrow[k]{*} \mathbf{x}$, then $\mathbf{y} = \mathbf{x}$. The set of k -irreducible q -ary words is denoted by $\text{Irr}_k(q)$ and the set of k -irreducible words of length n is denoted by $\text{Irr}_k(n, q)$. The k -ancestors of $\mathbf{x} \in \Sigma_q^*$ that are k -irreducible words are called the k -roots of \mathbf{x} . The set of k -roots of \mathbf{x} is denoted by $R_k(\mathbf{x})$. In a similar fashion, we have the definitions of $\leq k$ -irreducible words and $\leq k$ -roots of \mathbf{x} , and the notation $\text{Irr}_{\leq k}(q)$, $\text{Irr}_{\leq k}(n, q)$, and $R_{\leq k}(\mathbf{x})$.

Example 2.5. For alphabet Σ_3 , we have $R_{\leq 3}(01012012) = \{012\}$, and $R_{\leq 2}(012012) = \{012012\}$. However, $R_{\leq 4}(012101212) = \{012, 0121012\}$.

As we see, it is possible for a word to have more than one root. Jain et al. [7] determined when a tandem duplication system has only one root, irrespective of the word, and provided efficient algorithms to compute this unique root.

PROPOSITION 2.6 (JAIN ET AL. [7]). *For any $\mathbf{x} \in \Sigma_q^*$, we have that $|R_k(\mathbf{x})| = 1$ for all k , and $|R_{\leq k}(\mathbf{x})| = 1$ for all $k \in \{1, 2, 3\}$. Furthermore, there exist linear-time algorithms to compute $R_k(\mathbf{x})$ for all k .*

One may easily derive linear-time algorithms to compute $R_{\leq k}(\mathbf{x})$ for all $k \in \{1, 2, 3\}$. For completeness, we formally describe these algorithms in Appendix A.

As it turns out, for certain cases, determining the confusability of two words is equivalent to computing the roots for the words.

PROPOSITION 2.7 (JAIN ET AL. [7]). *For all $\mathbf{x}, \mathbf{y} \in \Sigma_q^*$, we have that*

- (i) \mathbf{x} and \mathbf{y} are k -confusable if and only if $R_k(\mathbf{x}) = R_k(\mathbf{y})$ for all k ;
- (ii) \mathbf{x} and \mathbf{y} are $\leq k$ -confusable if and only if $R_{\leq k}(\mathbf{x}) = R_{\leq k}(\mathbf{y})$ for $k \in \{1, 2\}$.
- (iii) If $R_{\leq 3}(\mathbf{x}) \neq R_{\leq 3}(\mathbf{y})$, then \mathbf{x} and \mathbf{y} are not ≤ 3 -confusable.

In other words, when $k \in \{1, 2\}$, to determine \mathbf{x} and \mathbf{y} are $\leq k$ -confusable, it is both necessary and sufficient to compute the $\leq k$ -roots of \mathbf{x} and \mathbf{y} . Therefore, applying Proposition 2.6 and the algorithms in Appendix A, we are able to determine whether two words are $\leq k$ -confusable in linear time for $k \in \{1, 2\}$.

Unfortunately, when $k = 3$, it is no longer sufficient to compute the ≤ 3 -roots of \mathbf{x} and \mathbf{y} . Specifically, even though $R_{\leq 3}(\mathbf{x}) = R_{\leq 3}(\mathbf{y})$, it is possible that \mathbf{x} and \mathbf{y} are not ≤ 3 -confusable. We illustrate this in the next example.

Example 2.8. Consider $\mathbf{x} = 012012$ and $\mathbf{y} = 011112$ over Σ_3 . The words have the same root as $R_{\leq 3}(\mathbf{x}) = R_{\leq 3}(\mathbf{y}) = \{012\}$. However, \mathbf{x} and \mathbf{y} are not ≤ 3 -confusable, because any ≤ 3 -descendant of \mathbf{x} has a 2 to the left of a 0, whereas any ≤ 3 -descendant of \mathbf{y} does not.

Therefore, the next smallest open case is where $k = 3$. A polynomial-time algorithm is implied by the results of Leupold et al. [11] and Jain et al. [6].

PROPOSITION 2.9 (LEUPOLD ET AL. [11], JAIN ET AL. [6]). *Let $k \in \{2, 3\}$. Let $\mathbf{x} \in \Sigma_q^*$. Then $D_{\leq k}^*(\mathbf{x})$ is a regular language. Furthermore, if $|\mathbf{x}| = m$, then the deterministic finite automaton that generates $D_{\leq k}^*(\mathbf{x})$ has $O(m)$ vertices and $O(m)$ edges. The number of vertices and edges is independent of the alphabet size.*

Therefore, for two words \mathbf{x} and \mathbf{y} with $|\mathbf{x}| = m$ and $|\mathbf{y}| = n$, the language $D_{\leq 3}^*(\mathbf{x}) \cap D_{\leq 3}^*(\mathbf{y})$ is regular and the corresponding finite automaton has $O(mn)$ vertices and $O(mn)$ edges (see for example, Sipser [16, Footnote 3, Theorem 1.25]). Hence, to determine if \mathbf{x} and \mathbf{y} are confusable, it suffices to determine if the language $D_{\leq 3}^*(\mathbf{x}) \cap D_{\leq 3}^*(\mathbf{y})$ is nonempty. The latter can be done in $O(mn)$ time. We improve this running time by providing an algorithm that runs in $O(\max\{m, n\})$ time.

Code Construction. We are interested in determining the maximum possible size of an $(n, \leq k; q)$ -TD code. An $(n, \leq k; q)$ -TD code that achieves this maximum is said to be *optimal*.

Motivated by Proposition 2.7, Jain et al. used irreducible words to construct tandem-duplication codes.

CONSTRUCTION 1. For $k \in \{2, 3\}$ and $n \geq k$. An $(n, \leq k; q)$ -TD code $C_I(n, \leq k; q)$ is given by

$$C_I(n, \leq k; 3) = \bigcup_{i=1}^n \left\{ \xi_{n-i}(\mathbf{x}) \mid \mathbf{x} \in \text{Irr}_{\leq k}(i, q) \right\}.$$

Here, $\xi_i(\mathbf{x}) = \mathbf{x}z^i$, where z is the last symbol of \mathbf{x} . In other words, $\xi_i(\mathbf{x})$ is the sequence \mathbf{x} with its last symbol z repeated i more times. Furthermore, the size of $C_I(n, \leq k; q)$ is $\sum_{i=1}^n |\text{Irr}_{\leq k}(i, q)|$.

It then follows from Proposition 2.7 that $C_I(n, \leq 2; q)$ is an optimal $(n, \leq 2; q)$ -TD code. Unfortunately, the TD code $C_I(n, \leq 3; q)$ is not an optimal for $n \geq 6$, and we illustrate this via the following example.

Example 2.10. Consider $n = 6, q = 3$. The code $C_I(6, \leq 3; 3)$ from Construction 1 is given by

$$\begin{aligned} C_I(6, \leq 3; 3) = & \{aaaaaa \mid a \in \Sigma_3\} \cup \{abbbbbb \mid a, b \in \Sigma_3, a \neq b\} \\ & \cup \{abaaaa, abcccc \mid a, b, c \in \Sigma_3, a \neq b, b \neq c, a \neq c\} \\ & \cup \{abacce, abcaaa, abcbbb \mid a, b, c \in \Sigma_3, a \neq b, b \neq c, a \neq c\} \\ & \cup \{abacaa, abacbb, abcabb, abcacc, abcbaa \mid a, b, c \in \Sigma_3, a \neq b, b \neq c, a \neq c\} \\ & \cup \{abacab, abacba, abacbc, abcaba, abcacb, abcabab, abcabac \mid a, b, c \in \Sigma_3, a \neq b, b \neq c, a \neq c\}. \end{aligned}$$

Hence, $|C_I(6, \leq 3; 3)| = 111$. However, we may remove from $C_I(6, \leq 3; 3)$ the six codewords $\{abcccc \mid a, b, c \in \Sigma_3, a \neq b, b \neq c, a \neq c\}$ and augment the code with twelve more codewords $\{abcabc, abbbbc \mid a, b, c \in \Sigma_3, a \neq b, b \neq c, a \neq c\}$. Then, we can check that the new code has size 117, and we later verify using Proposition 2.11 that the new code is indeed optimal. For lengths at least six, we construct TD codes with strictly larger sizes.

Upper Bound. We also study upper bounds on the size of an optimal $(n, \leq 3; q)$ -TD code. By definition, an $(n, \leq 3; q)$ -TD code is also an $(n, \leq 2; q)$ -TD code. Since an optimal $(n, \leq 2; q)$ -TD code is provided by Construction 1, we have the following upper bound on the size of an optimal $(n, \leq 3; q)$ -TD code.

PROPOSITION 2.11. The size of an $(n, \leq 3; q)$ -TD code is at most $\sum_{i=1}^n |\text{Irr}_{\leq 2}(i, q)|$.

Proposition 2.11 implies that Construction 1 is tight for $k = 3$ and $n \leq 5$. Using a combinatorial characterization implied by our algorithm, we improve this upper bound for longer lengths in Section 4.

2.3 Our Contributions

- **Confusability.** In Section 3, we present sufficient and necessary conditions for two words to be ≤ 3 -confusable and propose a linear-time algorithm to solve the ≤ 3 -confusability problem.
- **Estimates on code sizes.** Using insights gained from Section 3, we study the size of tandem-duplication codes in Section 4. We first improve the upper bound given by Proposition 2.11 and then construct codes with larger sizes as compared to those given by Construction 1. We also provide certain explicit constructions and recursive constructions for tandem-duplication codes. Furthermore, we determine the sizes of optimal ≤ 3 -tandem-duplication codes for lengths up to 20.

3 DETERMINING ≤ 3 -CONFUSABILITY

We derive a linear-time algorithm to determine the ≤ 3 -confusability of two words \mathbf{x} and \mathbf{y} . For the sake of simplicity, we omit the use of " ≤ 3 " and assume confusable, descendant, and root to mean ≤ 3 -confusable, ≤ 3 -descendant, and ≤ 3 -root, respectively.

Prior to stating the algorithm, we make some technical observations on tandem duplications of length at most three. The first lemma states that given a sequence of tandem duplications, we may reorder the tandem duplications such that the lengths of the repeats are nonincreasing.

LEMMA 3.1. *Let $k_1 < k_2 \leq 3$ and $\mathbf{x} \in \Sigma_G^*$. Suppose that $T_{i_2, k_2} \circ T_{i_1, k_1}(\mathbf{x}) = \mathbf{x}'$ for some i_1, i_2 . Then there exist integers $3 \geq \ell_1 \geq \ell_2 \geq \dots \geq \ell_t$ and j_1, j_2, \dots, j_t such that*

$$T_{j_t, \ell_t} \circ T_{j_{t-1}, \ell_{t-1}} \circ \dots \circ T_{j_1, \ell_1}(\mathbf{x}) = \mathbf{x}'.$$

PROOF. Consider a sequence of two tandem duplications $T_{i_2, k_2} \circ T_{i_1, k_1}$ with $k_1 < k_2 \leq 3$. We replace the sequence of duplications according to the following rules.

- Suppose that $k_1 = 1$ and $k_2 = 3$.
 - (a) If $i_2 \leq i_1 - 2$, then substitute with $T_{i_1+3, 1} \circ T_{i_2, 3}$.
 - (b) If $i_2 = i_1 - 1$, then substitute with $T_{i_1+3, 1} \circ T_{i_1, 1} \circ T_{i_1-1, 2}$.
 - (c) If $i_2 = i_1$, then substitute with $T_{i_1+3, 1} \circ T_{i_1, 1} \circ T_{i_1, 2}$.
 - (d) If $i_2 \geq i_1 + 1$, then substitute with $T_{i_1, 1} \circ T_{i_2-1, 3}$.
- Suppose that $k_1 = 2$ and $k_2 = 3$.
 - (a) If $i_2 \leq i_1 - 1$, then substitute with $T_{i_1+3, 2} \circ T_{i_2, 3}$.
 - (b) If $i_2 = i_1$, then substitute with $T_{i_2, 1} \circ T_{i_1, 2} \circ T_{i_1, 2}$.
 - (c) If $i_2 = i_1 + 1$, then substitute with $T_{i_1+3, 1} \circ T_{i_1, 2} \circ T_{i_1, 2}$.
 - (d) If $i_2 \geq i_1 + 2$, then substitute with $T_{i_1, 2} \circ T_{i_2-2, 3}$.
- Suppose that $k_1 = 1$ and $k_2 = 2$.
 - (a) If $i_2 \leq i_1 - 1$, then substitute with $T_{i_1+2, 1} \circ T_{i_2, 2}$.
 - (b) If $i_2 = i_1$, then substitute with $T_{i_1, 1} \circ T_{i_1, 1} \circ T_{i_1, 1}$.
 - (c) If $i_2 \geq i_1 + 1$, then substitute with $T_{i_1, 1} \circ T_{i_2-1, 2}$.

Thus, given a sequence of tandem duplications, we may reorder the tandem duplications such that the lengths of the repeats are non-increasing. ■

The next lemma states that we may assume the duplications of length two and three are performed on segments whose symbols are distinct. This is because if a duplication is performed on a segment whose symbols are not distinct, then we may find an equivalent sequence of duplications of strictly shorter lengths.

LEMMA 3.2. *Let $k \in \{2, 3\}$. Suppose $T_{i, k}(\mathbf{x}) = \mathbf{u}\mathbf{v}\mathbf{v}\mathbf{w}$, where $\mathbf{x} = \mathbf{u}\mathbf{v}\mathbf{w}$, $|\mathbf{u}| = i$ and $|\mathbf{v}| = k$. If the symbols in \mathbf{v} are not pairwise distinct, then there exist integers $k > \ell_1 \geq \ell_2$ and j_1, j_2 such that*

$$T_{j_2, \ell_2} \circ T_{j_1, \ell_1}(\mathbf{x}) = T_{i, k}(\mathbf{x}).$$

PROOF. When $k = 3$, we consider the following three cases.

- Duplication of aba to $abaaba$ is equivalent to $T_{2, 1} \circ T_{1, 2}(aba) = T_{2, 1}(ababa) = abaaba$.
- Duplication of aab to $aabaab$ is equivalent to $T_{3, 1} \circ T_{1, 2}(aab) = T_{3, 1}(aabab) = aabaab$.
- Duplication of abb to $abbabb$ is equivalent to $T_{1, 1} \circ T_{0, 2}(abb) = T_{1, 1}(ababb) = abbaab$.
- Duplication of aaa to $aaaaaa$ is equivalent to $T_{0, 1} \circ T_{0, 2}(aaa) = T_{0, 1}(aaaaa) = aaaaaa$.

When $k = 2$, we consider the duplication of length two on aa to obtain $aaaa$. This is equivalent to $T_{0,1} \circ T_{0,1}(aa) = T_{0,1}(aaa) = aaaa$. Therefore, if a duplication is performed on a segment whose symbols are not distinct, then we may find an equivalent sequence of duplications of strictly shorter lengths. ■

Henceforth, we use the notation $\mathbf{x} \xrightarrow[k_d]{*} \mathbf{y}$ to denote that there is a sequence of tandem duplications of length k of *distinct* symbols from \mathbf{x} to \mathbf{y} . It follows from definition that $\mathbf{x} \xrightarrow[1]{*} \mathbf{y}$ is equivalent to $\mathbf{x} \xrightarrow[1_d]{*} \mathbf{y}$. The next lemma is immediate from Lemmas 3.1 and 3.2.

LEMMA 3.3. *Suppose that $\mathbf{x} \xrightarrow[\leq 3]{*} \mathbf{x}'$. Then there exist \mathbf{x}_1 and \mathbf{x}_2 such that $\mathbf{x} \xrightarrow[3_d]{*} \mathbf{x}_2 \xrightarrow[2_d]{*} \mathbf{x}_1 \xrightarrow[1_d]{*} \mathbf{x}'$. Furthermore, \mathbf{x}_1 and \mathbf{x}_2 are uniquely determined by $\mathbf{x}_1 = R_1(\mathbf{x}')$ and $\mathbf{x}_2 = R_{\leq 2}(\mathbf{x}')$.*

Equipped with this lemma, we state the following theorem that provides a necessary and sufficiency condition for two words to be confusable.

THEOREM 3.4. *Two words \mathbf{x} and \mathbf{y} are ≤ 3 -confusable if and only if there exist \mathbf{x}' and \mathbf{y}' such that*

$$\mathbf{x} \xrightarrow[3_d]{*} \mathbf{x}', \mathbf{y} \xrightarrow[3_d]{*} \mathbf{y}' \text{ and } R_{\leq 2}(\mathbf{x}') = R_{\leq 2}(\mathbf{y}').$$

PROOF. Suppose that \mathbf{x} and \mathbf{y} are confusable. Then there exists \mathbf{z} such that $\mathbf{x} \xrightarrow[\leq 3]{*} \mathbf{z}$ and $\mathbf{y} \xrightarrow[\leq 3]{*} \mathbf{z}$. By Lemma 3.3, there exist $\mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2$ such that $\mathbf{x} \xrightarrow[3_d]{*} \mathbf{x}_2 \xrightarrow[2_d]{*} \mathbf{x}_1 \xrightarrow[1_d]{*} \mathbf{z}$ and $\mathbf{y} \xrightarrow[3_d]{*} \mathbf{y}_2 \xrightarrow[2_d]{*} \mathbf{y}_1 \xrightarrow[1_d]{*} \mathbf{z}$. Hence, we set $\mathbf{x}' = \mathbf{x}_2$ and $\mathbf{y}' = \mathbf{y}_2$. Note that by Proposition 2.7, $R_{\leq 2}(\mathbf{x}') = R_{\leq 2}(\mathbf{y}')$, since \mathbf{x}' and \mathbf{y}' are ≤ 2 -confusable.

Conversely, suppose that there exist \mathbf{x}' and \mathbf{y}' such that $\mathbf{x} \xrightarrow[3_d]{*} \mathbf{x}'$, $\mathbf{y} \xrightarrow[3_d]{*} \mathbf{y}'$, and $R_{\leq 2}(\mathbf{x}') = R_{\leq 2}(\mathbf{y}')$. Since \mathbf{x}' and \mathbf{y}' have a common ≤ 2 -root, it follows from Proposition 2.7 that \mathbf{x}' and \mathbf{y}' are ≤ 2 -confusable. Thus, \mathbf{x} and \mathbf{y} are confusable. ■

In conclusion, for two words \mathbf{x} and \mathbf{y} to be confusable, it is equivalent to checking if there are tandem duplications of length three that make their descendants \mathbf{x}' and \mathbf{y}' have the same ≤ 2 -root. We make use of this important fact in the next section.

3.1 Strategy Behind Algorithm

We outline our strategy to determine the confusability of two words \mathbf{x} and \mathbf{y} . Recall from Proposition 2.7(iii), if \mathbf{x} and \mathbf{y} have different roots, then they are not confusable. Hence, our first step is to determine $R_{\leq 3}(\mathbf{x})$ and $R_{\leq 3}(\mathbf{y})$ and see if the roots are equal. We can do this in linear time using Algorithm 2. If the roots are not equal, then we immediately conclude that \mathbf{x} and \mathbf{y} are not confusable. Therefore, the nontrivial task is to determine confusability when $R_{\leq 3}(\mathbf{x}) = R_{\leq 3}(\mathbf{y})$. Henceforth, we assume that the *common root of \mathbf{x} and \mathbf{y} is \mathbf{r}* .

Suppose that \mathbf{r} contains at most two distinct symbols. Then from Lemma 3.3, we have that $\mathbf{r} \xrightarrow[\leq 2]{*} \mathbf{x}$ and $\mathbf{r} \xrightarrow[\leq 2]{*} \mathbf{y}$. In other words, \mathbf{r} is also the ≤ 2 -root of both \mathbf{x} and \mathbf{y} . Applying Proposition 2.7, we have that \mathbf{x} and \mathbf{y} are ≤ 2 -confusable and so, ≤ 3 -confusable. Therefore, it remains to consider the case where \mathbf{r} contains *at least three distinct symbols*.

Given \mathbf{r} , \mathbf{x} and \mathbf{y} , our next step is to compute certain proper prefixes $*\text{Pref}(\mathbf{r}, \mathbf{x})$ and $*\text{Pref}(\mathbf{r}, \mathbf{y})$ of \mathbf{x} and \mathbf{y} , respectively, such that *under certain conditions* the following holds:

$$(\mathbf{x} \text{ and } \mathbf{y} \text{ are confusable}) \text{ if and only if } (\mathbf{x} \setminus *\text{Pref}(\mathbf{r}, \mathbf{x}) \text{ and } \mathbf{y} \setminus *\text{Pref}(\mathbf{r}, \mathbf{y})) \text{ are confusable}.$$

Here, $\mathbf{x} \setminus z$ denotes the word we obtain by removing the prefix z from \mathbf{x} . Since the words on the righthand side are strictly shorter than the words on the lefthand side, we may repeat this process for the shorter words until the common root of the words has less than three distinct symbols.

Before we formally define the prefixes $\text{*Pref}(\mathbf{r}, \mathbf{x})$ and $\text{*Pref}(\mathbf{r}, \mathbf{y})$, we provide certain intuition behind our technical definitions via an example.

Example 3.5. Let $\mathbf{x} = 01102021020120111$ and $\mathbf{y} = 00101002110200120201$. Since both \mathbf{x} and \mathbf{y} have the same ≤ 3 -root $\mathbf{r} \triangleq 010201$, it is possible that \mathbf{x} and \mathbf{y} are ≤ 3 -confusable.

To do so, we make some observations on the root \mathbf{r} and look at how we obtain \mathbf{x} and \mathbf{y} from \mathbf{r} . Now, \mathbf{r} has two 3-substrings that contain three distinct symbols: namely, 102 and 201, and so, we divide \mathbf{r} two *overlapping* substrings $\mathbf{r}_1 = 0102$ and $\mathbf{r}_2 = 0201$.

We can also similarly divide each of \mathbf{x} and \mathbf{y} into two overlapping substrings. Namely,

$$\begin{aligned} \mathbf{x}_1 &= 0110202102, & \mathbf{x}_2 &= 020120111, \\ \mathbf{y}_1 &= 001010021102, & \mathbf{y}_2 &= 0200120201. \end{aligned}$$

We can easily check that $\mathbf{r}_1 \xrightarrow[\leq 3]{*} \mathbf{x}_1$, $\mathbf{r}_1 \xrightarrow[\leq 3]{*} \mathbf{y}_1$, $\mathbf{r}_2 \xrightarrow[\leq 3]{*} \mathbf{x}_2$, and $\mathbf{r}_2 \xrightarrow[\leq 3]{*} \mathbf{y}_2$. In what follows, we prove that \mathbf{x} and \mathbf{y} are ≤ 3 -confusable if and only if \mathbf{x}_i and \mathbf{y}_i are ≤ 3 -confusable for $i \in \{1, 2\}$. Furthermore, for $i \in \{1, 2\}$, we define a combinatorial property of \mathbf{x}_i and \mathbf{y}_i that allows one to easily determine that \mathbf{x}_i and \mathbf{y}_i are ≤ 3 -confusable.

Before this, we describe rules on how to break \mathbf{r} into the *regions*: \mathbf{r}_1 and \mathbf{r}_2 .

We now provide the formal definitions. Assume that \mathbf{r} contains at least three distinct symbols. Since $\mathbf{r} \in \text{Irr}_{\leq 3}(q)$, we easily check that either $r_1r_2r_3$ or $r_2r_3r_4$ is a substring with three distinct symbols. Define $\text{main}(\mathbf{r})$ to be the first substring in \mathbf{r} with three distinct symbols and we define the *first region of \mathbf{r}* , denoted by $\text{Reg}(\mathbf{r})$, to be a certain prefix of \mathbf{r} according to the following rule.

If $r_1 = r_3$, then we set

$$\text{main}(\mathbf{r}) \triangleq r_2r_3r_4, \text{ and } \text{Reg}(\mathbf{r}) \triangleq \begin{cases} r_1r_2r_3r_4 = r_1r_2r_1r_4, & \text{if } r_2 \neq r_5, \\ r_1r_2r_3r_4r_5 = r_1r_2r_1r_4r_2, & \text{if } r_2 = r_5, r_3 \neq r_6, \\ r_1r_2r_3r_4r_5r_6 = r_1r_2r_1r_4r_2r_1, & \text{otherwise.} \end{cases}$$

If $r_1 \neq r_3$, then we set

$$\text{main}(\mathbf{r}) \triangleq r_1r_2r_3, \text{ and } \text{Reg}(\mathbf{r}) \triangleq \begin{cases} r_1r_2r_3, & \text{if } r_1 \neq r_4, \\ r_1r_2r_3r_4 = r_1r_2r_3r_1, & \text{if } r_1 = r_4, r_2 \neq r_5, \\ r_1r_2r_3r_4r_5 = r_1r_2r_3r_1r_2, & \text{otherwise.} \end{cases}$$

Intuitively, the first region $\text{Reg}(\mathbf{r})$ of \mathbf{r} is defined as above so that the following properties of the region, or Lemma 3.6 hold.

LEMMA 3.6. *Suppose the root \mathbf{r} has three distinct symbols. Define $\text{Reg}(\mathbf{r})$ as above. Then $\text{Reg}(\mathbf{r})$ can be written as $\mathbf{w}(abc)^\ell ab$, where a, b, c are distinct symbols, $\ell \in \{0, 1\}$ and \mathbf{w} is prefix of length at most three over the alphabet $\{a, b, c\}$. Furthermore, if $\text{Reg}(\mathbf{r}) \xrightarrow[\leq 3]{*} z$, then z can be written as $\mathbf{w}(abc)^m ab$ with $m \geq \ell$. Also, the symbol in \mathbf{r} after the prefix $\text{Reg}(\mathbf{r})$ is not c .*

Example 3.7. Suppose that $\mathbf{r} = 010201$. Then $\text{main}(\mathbf{r}) = 102$ and $\text{Reg}(\mathbf{r}) = 0102$. Indeed, we can write $\text{Reg}(\mathbf{r})$ as $\mathbf{w}(abc)^\ell ab$ with $\mathbf{w} = 01$, $\ell = 0$, and $abc = 021$. Note that abc is not necessarily $\text{main}(\mathbf{r})$. We also check that the next symbol after $\text{Reg}(\mathbf{r})$ in \mathbf{r} is 0, which is not $c = 1$.

Remark 1. Given a root \mathbf{r} , we provide the rule to determine the *first region*. The subsequent regions of \mathbf{r} are then recursively defined by looking at the “first region” of a certain suffix of \mathbf{r} . Details are provided in Algorithm 1.

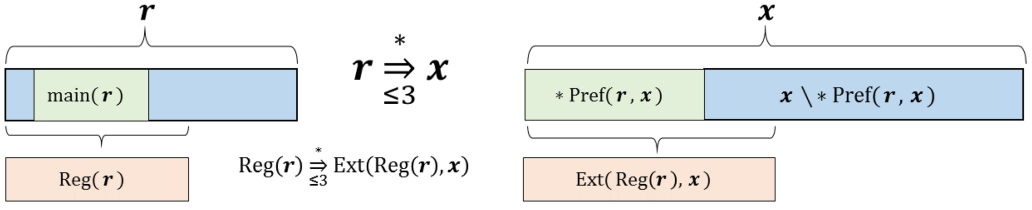


Fig. 1. Prefixes used to determine confusability.

Given $\text{Reg}(r) = w(abc)^\ell ab$, we consider the longest prefix of x that is generated from $\text{Reg}(r)$ through a sequence of tandem duplications of length at most three. We call this prefix the *first extended region of r in x* and denote it with $\text{Ext}(\text{Reg}(r), x)$. This prefix can be computed deterministically using the finite automaton that defines $D_{\leq 3}^*(\text{Reg}(r))$ (see Subsection 3.2 for details).

Finally, we define our special prefix $*\text{Pref}(r, x)$. Suppose that $\text{Ext}(\text{Reg}(r), x) = p_1 p_2 \cdots p_t = \mathbf{p}$ and p_i is the last appearance of a in \mathbf{p} , where a is the second last symbol of $\text{Reg}(r)$. The *special prefix* $*\text{Pref}(r, x)$ is given by $p_1 p_2 \cdots p_{i-1}$.

Example 3.8. Let $r = 010201$ and $x = 01102021020120111$. Recall that $\text{Reg}(r) = 0102$ and $a = 0$. Then,

$$\text{Ext}(0102, x) = 0110202102 \text{ and } *\text{Pref}(r, x) = 01102021.$$

Remark 2. Consider a word x whose root is r . Similar to the definition of the first region of r , we only provide the rule to determine the *first extended region of r in x* . The subsequent extended regions of r in x are then recursively defined by looking at the “first extended region” of a certain suffix of r in a certain suffix of x . As before, details are provided in Algorithm 1.

We summarize the relationships of r , x , $\text{Reg}(r)$, $\text{Ext}(\text{Reg}(r), x)$, and $*\text{Pref}(r, x)$ in Figure 1, and we are ready to state the main characterization theorem.

THEOREM 3.9. *Suppose that x and y are two words such that $r = R_{\leq 3}(x) = R_{\leq 3}(y)$ and r contains at least three distinct symbols. Set*

$$\begin{aligned} \mathbf{p} &= \text{Ext}(\text{Reg}(r), x), & \mathbf{q} &= \text{Ext}(\text{Reg}(r), y), \\ \mathbf{u} &= x \setminus *\text{Pref}(r, x), & \mathbf{v} &= y \setminus *\text{Pref}(r, y). \end{aligned}$$

Then x and y are confusable if and only if \mathbf{p} and \mathbf{q} are confusable, and \mathbf{u} and \mathbf{v} are confusable.

PROOF. From Lemma 3.6, we may set $\text{Reg}(r) = w(abc)^\ell ab$. Since $\text{Reg}(r) \xrightarrow[\leq 3]{*} \mathbf{p}$, Lemma 3.3 implies that $\text{Reg}(r) \xrightarrow[\leq 2]{*} R_{\leq 2}(\mathbf{p})$. So, $R_{\leq 2}(\mathbf{p}) = w(abc)^m ab$ for some $m \geq \ell$, and $R_{\leq 2}(\mathbf{q}) = w(abc)^j ab$ for some $j \geq \ell$.

Suppose that x and y are confusable. By Theorem 3.4, there exist 3_d -descendants of x and y namely x' and y' , such that $R_{\leq 2}(x') = R_{\leq 2}(y')$. Set

$$\begin{aligned} \mathbf{p}' &= \text{Ext}(\text{Reg}(r), x'), & \mathbf{q}' &= \text{Ext}(\text{Reg}(r), y'), \\ \mathbf{u}' &= x' \setminus *\text{Pref}(r, x'), & \mathbf{v}' &= y' \setminus *\text{Pref}(r, y'). \end{aligned}$$

Since the overlapping substring between \mathbf{p} and \mathbf{u} must be of the form $ab \cdots b$, each duplication of distinct triplets from x to x' must be either entirely in \mathbf{p} or in \mathbf{u} . In other words, $\mathbf{p} \xrightarrow[\leq 3]{*} \mathbf{p}'$ and $\mathbf{u} \xrightarrow[\leq 3]{*} \mathbf{u}'$, and likewise $\mathbf{q} \xrightarrow[\leq 3]{*} \mathbf{q}'$ and $\mathbf{v} \xrightarrow[\leq 3]{*} \mathbf{v}'$. Looking at the overlapping substring between \mathbf{p}'

and \mathbf{u}' , we can observe that $R_{\leq 2}(\mathbf{x}') = R_{\leq 2}(\mathbf{p}'\mathbf{u}')$. Then $\text{Reg}(\mathbf{r}) \xrightarrow[\leq 3]{*} \mathbf{p}'$ and Lemma 3.3 imply that $\text{Reg}(\mathbf{r}) \xrightarrow[\leq 3]{*} R_{\leq 2}(\mathbf{p})$ and so, $R_{\leq 2}(\mathbf{p}') = \mathbf{w}(abc)^i ab$ for some $i \geq m$.

Next, $\mathbf{w}(abc)^i R_{\leq 2}(\mathbf{u}')$ is a ≤ 2 -irreducible word, because $R_{\leq 2}(\mathbf{u}')$ starts with ab and by definition, $\mathbf{w}(abc)^i$ ends with c but not cac . Therefore, $\mathbf{w}(abc)^i R_{\leq 2}(\mathbf{u}')$ is the ≤ 2 -root of \mathbf{x}' , or, $R_{\leq 2}(\mathbf{x}') = \mathbf{w}(abc)^i R_{\leq 2}(\mathbf{u}')$ and likewise $R_{\leq 2}(\mathbf{y}') = \mathbf{w}(abc)^h R_{\leq 2}(\mathbf{v}')$ for some $h \geq j$.

Since the first symbol of $R_{\leq 2}(\mathbf{x}' \setminus \mathbf{p}')$ and $R_{\leq 2}(\mathbf{y}' \setminus \mathbf{q}')$ cannot be c , neither $R_{\leq 2}(\mathbf{u}')$ nor $R_{\leq 2}(\mathbf{v}')$ start with abc . Therefore, since $R_{\leq 2}(\mathbf{x}') = R_{\leq 2}(\mathbf{y}')$, we can deduce that $h = i$ and $R_{\leq 2}(\mathbf{u}') = R_{\leq 2}(\mathbf{v}')$. Since $h = i$, we have $R_{\leq 2}(\mathbf{p}') = R_{\leq 2}(\mathbf{q}')$. Together with $\mathbf{p} \xrightarrow[\leq 3]{*} \mathbf{p}'$, $\mathbf{q} \xrightarrow[\leq 3]{*} \mathbf{q}'$, by Theorem 3.4, we can have \mathbf{p} and \mathbf{q} are confusable. Since $\mathbf{u} \xrightarrow[\leq 3]{*} \mathbf{u}'$, $\mathbf{v} \xrightarrow[\leq 3]{*} \mathbf{v}'$, and $R_{\leq 2}(\mathbf{u}') = R_{\leq 2}(\mathbf{v}')$, we can also conclude by Theorem 3.4 that \mathbf{u} and \mathbf{v} are confusable.

Suppose that \mathbf{p} and \mathbf{q} are confusable, and \mathbf{u} and \mathbf{v} are confusable. Hence, \mathbf{pu} and \mathbf{qv} are confusable. Note that \mathbf{pu} is a descendant of \mathbf{x} and \mathbf{qv} is a descendant of \mathbf{y} , hence \mathbf{x} and \mathbf{y} are also confusable. \blacksquare

We can efficiently determine whether the prefixes \mathbf{p} and \mathbf{q} are confusable. When the substring $\text{main}(\mathbf{r})$ appears the same number of times in $R_{\leq 2}(\mathbf{p})$ as in $R_{\leq 2}(\mathbf{q})$, we may assume that $R_{\leq 2}(\mathbf{p}) = \mathbf{w}(abc)^m ab = R_{\leq 2}(\mathbf{q})$. Hence, \mathbf{p} and \mathbf{q} are confusable by Proposition 2.7. In particular, we have Corollary 3.10.

Consider $\mathbf{z} = abc$, where a, b, c are three distinct symbols. Let $\text{Count}(\mathbf{z}, \mathbf{x})$ denote the number of times \mathbf{z} appears as a substring in \mathbf{x} , and $\text{Count}(\text{rot}(\mathbf{z}), \mathbf{x})$ denote the number of times abc, bca , or cab appears as a substring in \mathbf{x} .

COROLLARY 3.10. *Suppose that \mathbf{x} and \mathbf{y} are two words such that $\mathbf{r} = R_{\leq 3}(\mathbf{x}) = R_{\leq 3}(\mathbf{y})$ and \mathbf{r} contains at least three distinct symbols. Set $\mathbf{p} = \text{Ext}(\text{Reg}(\mathbf{r}), \mathbf{x})$ and $\mathbf{q} = \text{Ext}(\text{Reg}(\mathbf{r}), \mathbf{y})$. Suppose further that $\text{Count}(\text{main}(\mathbf{r}), R_{\leq 2}(\mathbf{p})) = \text{Count}(\text{main}(\mathbf{r}), R_{\leq 2}(\mathbf{q}))$. Then \mathbf{x} and \mathbf{y} are confusable if and only if $\mathbf{x} \setminus * \text{Pref}(\mathbf{r}, \mathbf{x})$ and $\mathbf{y} \setminus * \text{Pref}(\mathbf{r}, \mathbf{y})$ are confusable.*

Example 3.11. Let $\mathbf{r} = 010201$, $\mathbf{x} = 01102021020120111$, and $\mathbf{y} = 00101002110200120201$. Note that $R_{\leq 3}(\mathbf{x}) = R_{\leq 3}(\mathbf{y}) = \mathbf{r}$. Recall that $\text{main}(\mathbf{r}) = 102$, $\text{Reg}(\mathbf{r}) = 0102$ and $a = 0$. Then,

$$\begin{aligned} \mathbf{p} &= \text{Ext}(0102, \mathbf{x}) = 0110202102, & * \text{Pref}(\mathbf{r}, \mathbf{x}) &= 01102021, & \mathbf{u} &= \mathbf{x} \setminus * \text{Pref}(\mathbf{r}, \mathbf{x}) = 020120111, \\ \mathbf{q} &= \text{Ext}(0102, \mathbf{y}) = 001010021102, & * \text{Pref}(\mathbf{r}, \mathbf{y}) &= 0010100211, & \mathbf{v} &= \mathbf{y} \setminus * \text{Pref}(\mathbf{r}, \mathbf{y}) = 0200120201. \end{aligned}$$

Note that $R_{\leq 2}(\mathbf{p}) = 0102102 = R_{\leq 2}(\mathbf{q})$, and $\text{Count}(\text{main}(\mathbf{r}), R_{\leq 2}(\mathbf{p})) = \text{Count}(\text{main}(\mathbf{r}), R_{\leq 2}(\mathbf{q})) = \text{Count}(102, 0102102) = 2$. Therefore, by Corollary 3.10, we have \mathbf{x} and \mathbf{y} are confusable if and only if \mathbf{u} and \mathbf{v} are confusable. Since $R_{\leq 2}(\mathbf{u}) = 0201201 = R_{\leq 2}(\mathbf{v})$, Proposition 2.7 implies that \mathbf{u} and \mathbf{v} are confusable. Therefore, \mathbf{x} and \mathbf{y} are confusable.

However, when the number of appearances is not equal, we may modify the proof to obtain the following corollary.

COROLLARY 3.12. *Suppose that \mathbf{x} and \mathbf{y} are two words such that $\mathbf{r} = R_{\leq 3}(\mathbf{x}) = R_{\leq 3}(\mathbf{y})$ and \mathbf{r} contains at least three distinct symbols. Set $\mathbf{p} = \text{Ext}(\text{Reg}(\mathbf{r}), \mathbf{x})$ and $\mathbf{q} = \text{Ext}(\text{Reg}(\mathbf{r}), \mathbf{y})$. Suppose further that $\text{Count}(\text{main}(\mathbf{r}), R_{\leq 2}(\mathbf{p})) < \text{Count}(\text{main}(\mathbf{r}), R_{\leq 2}(\mathbf{q}))$. Then \mathbf{x} and \mathbf{y} are confusable if and only if $\text{Count}(\text{rot}(\text{main}(\mathbf{r})), \mathbf{p}) > 0$ and $\mathbf{x} \setminus * \text{Pref}(\mathbf{r}, \mathbf{x})$ and $\mathbf{y} \setminus * \text{Pref}(\mathbf{r}, \mathbf{y})$ are confusable.*

PROOF. By Theorem 3.9, it is sufficient to show that $\text{Count}(\text{rot}(\text{main}(\mathbf{r})), \mathbf{p}) > 0$ iff \mathbf{p} and \mathbf{q} are confusable. Let $R_{\leq 2}(\mathbf{p}) = \mathbf{w}(abc)^k ab$ and $R_{\leq 2}(\mathbf{q}) = \mathbf{w}(abc)^j ab$ for some $k < j$.

ALGORITHM 1: CONFUSE(x, y)

Input: $x, y \in \Sigma_q^*$
Output: true if x and y are confusable, false otherwise.

```

1:  $r \leftarrow R_{\leq 3}(x)$ 
2:  $r' \leftarrow R_{\leq 3}(y)$ 
3: if  $r \neq r'$  then
4:   return false
5: else
6:   return CONFUSE-RECURSION( $x, y, r$ )

7: procedure CONFUSE-RECURSION( $x, y, r$ )
8:   if  $r$  has at most two distinct symbols then
9:     return true
10:  else
11:    Compute  $\text{main}(r)$  and  $\text{Reg}(r)$ 
12:     $p \leftarrow \text{Ext}(\text{Reg}(r), x), \quad q \leftarrow \text{Ext}(\text{Reg}(r), y)$ 
13:     $u \leftarrow x \setminus * \text{Pref}(r, x), \quad v \leftarrow y \setminus * \text{Pref}(r, y)$ 
14:     $C_p \leftarrow \text{Count}(\text{main}(r), R_{\leq 2}(\cdot, p)), \quad C_q \leftarrow \text{Count}(\text{main}(r), R_{\leq 2}(q))$ 
15:     $r^* \leftarrow r \setminus w(abc)^\ell$ , where  $\text{Reg}(r) = w(abc)^\ell ab$ 
16:    if  $C_p = C_q$  then
17:      return CONFUSE-RECURSION( $u, v, r^*$ )
18:    else if  $C_p < C_q$  and  $\text{Count}(\text{rot}(\text{main}(r)), p) > 0$  then
19:      return CONFUSE-RECURSION( $u, v, r^*$ )
20:    else if  $C_p < C_q$  and  $\text{Count}(\text{rot}(\text{main}(r)), q) > 0$  then
21:      return CONFUSE-RECURSION( $u, v, r^*$ )
22:    else
23:      return false

```

Suppose there exists at least one distinct triplet in p , we set p' to be the result of duplication of the distinct triplet in p as many as $j - k$ times and q' to be q . Then $R_{\leq 2}(p') = w(abc)^j ab = R_{\leq 2}(q')$, and hence by Theorem 3.4, p and q are confusable.

Suppose, however, p and q are confusable. By Theorem 3.4, there exist p' and q' such that $p \xrightarrow[3d]{*} p'$ and $q \xrightarrow[3d]{*} q'$, where $R_{\leq 2}(p') = R_{\leq 2}(q') = w(abc)^g ab$ for some $g \geq j > k$. Then there has to be at least one substring of length three with distinct symbols, or *distinct triplet*, in p . Since $p = \text{Ext}(\text{Reg}(r), x)$, then by definition, the only distinct triplet possible is from the set $\text{rot}(\text{main}(r))$. Therefore, $\text{Count}(\text{rot}(\text{main}(r)), p) > 0$. ■

Example 3.13. Let r and x be as defined in Example 3.11. We consider another word $y' = 0010100200120201$. As before, note that $R_{\leq 3}(x) = R_{\leq 3}(y') = r$. Recall that $\text{main}(r) = 102$, $\text{Reg}(r) = 0102$, and $a = 0$. Then p , $* \text{Pref}(r, x)$ and u are computed as before. In contrast, for y' , we have

$$q' = \text{Ext}(0102, y') = 00101002, \quad * \text{Pref}(r, y') = 001010, \quad v' = y' \setminus * \text{Pref}(r, y') = 0200120201.$$

Now, $R_{\leq 2}(q') = 0102$, and so, $\text{Count}(\text{main}(r), R_{\leq 2}(q')) = 1$. Hence, $\text{Count}(\text{main}(r), R_{\leq 2}(p)) \neq \text{Count}(\text{main}(r), R_{\leq 2}(q'))$. Furthermore, $\text{Count}(\text{rot}(\text{main}(r)), q) = \text{Count}(\text{rot}(102), 00101002) = 0$. Therefore, by Corollary 3.12, we have x and y are not confusable.

3.2 Confusability Algorithm

Based on the ideas in the previous subsection, we formally describe our algorithm CONFUSE in Algorithm 1 and demonstrate that the running time is linear. In the following analysis, we assume the definitions in Algorithm 1.

- Main(\mathbf{r}) and Reg(\mathbf{r}) can be computed in constant time. This is clear from definition.
- The prefixes \mathbf{p} and \mathbf{q} can be computed in time linear in $|\mathbf{p}|$ and $|\mathbf{q}|$, respectively. From Proposition 2.9, we have that $D_{\leq 3}^*(\text{Reg}(\mathbf{r}))$ is a regular language. Since the length of $\text{Reg}(\mathbf{r})$ is at most six, the finite automaton generating the language has a constant number of vertices. Hence, the longest prefixes of \mathbf{x} and \mathbf{y} that are accepted by the finite automaton are precisely \mathbf{p} and \mathbf{q} , respectively. In Appendix B, we show that the sum of the lengths of all prefixes that are computed in all recursive calls is at most $3(|\mathbf{x}| + |\mathbf{y}|)$. Therefore, the computation of all prefixes takes linear time.
- In the recursive calls, the root \mathbf{r}^* can be computed in constant time. As mentioned earlier, the first step is to compute the roots of \mathbf{x} and \mathbf{y} in linear time. When the roots are equal to \mathbf{r} and certain conditions are met, we make the recursive call CONFUSE-RECURSION. If the roots are recomputed for the shorter words \mathbf{u} and \mathbf{v} , then the overall running time for CONFUSE is quadratic. To avoid this, we observe that the roots for the shorter words are always equal and in fact, the common root \mathbf{r}^* of \mathbf{u} and \mathbf{v} may be easily inferred from \mathbf{r} . More precisely, if $\text{Reg}(\mathbf{r}) = \mathbf{w}(abc)^\ell ab$, then $\mathbf{r}^* = \mathbf{r} \setminus \mathbf{w}(abc)^\ell$.
- In the recursive calls, determining \mathbf{r} has at most two symbols can be done in constant time. If \mathbf{r} has length at least four, then \mathbf{r} necessarily has at least three distinct symbols. Otherwise, it contradicts the fact that \mathbf{r} is irreducible. Therefore, it suffices to check if the first four symbols \mathbf{r} contain three distinct symbols.

In summary, from the above observations, we conclude that the running time is $O(\max\{|\mathbf{x}|, |\mathbf{y}|\})$.

4 ≤ 3 -TANDEM-DUPLICATION CODES

We use insights gained from the previous section to construct $(n, \leq 3; q)$ -TD codes. Motivated by the concept of roots, we consider a ≤ 3 -irreducible word \mathbf{r} , and we say that a $(n, \leq 3; q)$ -TD code C is an $(n, \leq 3; \mathbf{r})$ -TD code if all words in C belong to $D_{\leq 3}^*(\mathbf{r})$.

For $\mathbf{r} \in \text{Irr}_{\leq 3}(i, q)$ with $i \leq n$, suppose that $C(n, \mathbf{r})$ is an $(n, \leq 3; \mathbf{r})$ -TD code. Then Proposition 2.7 implies that $\bigcup_{\mathbf{r} \in \text{Irr}_{\leq 3}(i, q), i \leq n} C(n, \mathbf{r})$ is an $(n, \leq 3; q)$ -TD code. Trivially, $\{\xi_{n-i}(\mathbf{r})\}$ is an $(n, \leq 3; \mathbf{r})$ -TD code for all $\mathbf{x} \in \text{Irr}_{\leq 3}(i, q)$. Taking the union of these codes, we recover Construction 1.

Therefore, in the rest of this section, our objective is to provide estimates on the size of an optimal $(n, \leq 3; \mathbf{r})$ -TD code for fixed ≤ 3 -irreducible word \mathbf{r} . To simplify our discussion, we focus on the case where $q = 3$ and let $T(n)$ and $T(n, \mathbf{r})$ to denote the sizes of an optimal $(n, \leq 3; 3)$ -TD code and an optimal $(n, \leq 3; \mathbf{r})$ -TD code, respectively.

Pick $\mathbf{x} \in D_{\leq 3}^*(\mathbf{r})$. Using concepts from Section 3, we define the following combinatorial characterisation of \mathbf{x} . Set $\mathbf{x}_1 = \mathbf{x}$ and $\mathbf{r}_1 = \mathbf{r}$. For $i \geq 2$, set $\mathbf{x}_i = \mathbf{x}_{i-1} \setminus * \text{Pref}(\mathbf{r}_{i-1}, \mathbf{x}_{i-1})$ and $\mathbf{r}_i = R_{\leq 3}(\mathbf{x}_i)$. We terminate this recursion when \mathbf{r}_{m+1} has less than three distinct symbols. If \mathbf{r}_m is the last root to have three distinct symbols, then we say that \mathbf{r} has m regions and for $1 \leq i \leq m$, define

$$\mathbf{p}_i = \text{Ext}(\text{Reg}(\mathbf{r}_i), \mathbf{x}_i), \mathbf{t}_i = \text{main}(\mathbf{r}_i), c_i = \text{Count}(\mathbf{t}_i, R_{\leq 2}(\mathbf{p}_i)), \text{ and } \delta_i = \begin{cases} +, & \text{if } \text{Count}(\text{rot}(\mathbf{t}_i), \mathbf{p}_i) > 0, \\ -, & \text{otherwise.} \end{cases}$$

For the word \mathbf{x} , we define its *label* by $\text{Label}(\mathbf{x}) = (\mathbf{r}, (c_1, \delta_1), (c_2, \delta_2), \dots, (c_m, \delta_m))$.

Example 4.1. Consider $\mathbf{r} = 01210$. Then \mathbf{r} has two regions with $\mathbf{t}_1 = 012$ and $\mathbf{t}_2 = 210$. Consider also the words 01210210, 01201210, and 01112110 that belongs to $D_{\leq 3}^*(01210)$. Then their labels

are as follows:

$$\text{Label}(01210210) = (01210, (1, +), (2, +)),$$

$$\text{Label}(01201210) = (01210, (2, +), (1, +)),$$

$$\text{Label}(01112110) = (01210, (1, -), (1, -)).$$

Then it follows from Theorem 3.9 and Corollaries 3.10 and 3.12 that 01210210 and 01201210 are confusable, while 01210210 and 01112110 are not confusable.

More generally, the next proposition is immediate from Theorem 3.9 and Corollaries 3.10 and 3.12.

PROPOSITION 4.2. *Let \mathbf{x} and \mathbf{y} belong to $D_{\leq 3}^*(\mathbf{r})$. Suppose that*

$$\text{Label}(\mathbf{x}) \triangleq (\mathbf{r}, (c_1^{(\mathbf{x})}, \delta_1^{(\mathbf{x})}), (c_2^{(\mathbf{x})}, \delta_2^{(\mathbf{x})}), \dots, (c_m^{(\mathbf{x})}, \delta_m^{(\mathbf{x})})),$$

$$\text{Label}(\mathbf{y}) \triangleq (\mathbf{r}, (c_1^{(\mathbf{y})}, \delta_1^{(\mathbf{y})}), (c_2^{(\mathbf{y})}, \delta_2^{(\mathbf{y})}), \dots, (c_m^{(\mathbf{y})}, \delta_m^{(\mathbf{y})})).$$

Then \mathbf{x} and \mathbf{y} are not ≤ 3 -confusable if and only if for some $1 \leq i \leq m$,

$$(c_i^{(\mathbf{x})} < c_i^{(\mathbf{y})} \text{ and } \delta_i^{(\mathbf{x})} = -) \text{ or } (c_i^{(\mathbf{x})} > c_i^{(\mathbf{y})} \text{ and } \delta_i^{(\mathbf{y})} = -).$$

Immediate from Proposition 4.2 are certain sufficient conditions for two words to be confusable.

COROLLARY 4.3. *Let $\mathbf{x}, \mathbf{y} \in D_{\leq 3}^*(\mathbf{r})$ and their labels $\text{Label}(\mathbf{x})$ and $\text{Label}(\mathbf{y})$ be as defined in Proposition 4.2.*

(i) *If $c_i^{(\mathbf{x})} = c_i^{(\mathbf{y})}$ for all i , then \mathbf{x} and \mathbf{y} are ≤ 3 -confusable.*

(ii) *If $\delta_i^{(\mathbf{x})} = \delta_i^{(\mathbf{y})} = +$ for all i , then \mathbf{x} and \mathbf{y} are ≤ 3 -confusable.*

From Corollary 4.3, we use the number of integer solutions to certain equations as an upper bound for $T(n, \mathbf{r})$. As this combinatorial argument is fairly technical, we state the upper bound and defer the proof to Appendix C.

PROPOSITION 4.4. *Let $i \leq n$. Suppose that $\mathbf{r} \in \text{Irr}_{\leq 3}(i, 3)$ has m regions. Then,*

$$T(n, \mathbf{r}) \leq U(n, i, m) \triangleq \begin{cases} \binom{(n-i)/3+m}{m} - \binom{(n-i)/3+m-1}{m-1} + 1, & \text{if } 3 \text{ divides } n-i, \\ \binom{\lfloor (n-i)/3 \rfloor + m}{m}, & \text{otherwise.} \end{cases}$$

4.1 Estimates of $T(|\mathbf{r}| + t, \mathbf{r})$ for $0 \leq t \leq 5$

Immediate from Proposition 4.4 is that $T(|\mathbf{r}|, \mathbf{r}) = T(|\mathbf{r}| + 1, \mathbf{r}) = T(|\mathbf{r}| + 2, \mathbf{r}) = 1$. Next, we provide lower bounds for $T(n, \mathbf{r})$ by constructing $(n, \leq 3; \mathbf{r})$ -TD codes. We first show that there exists an $(n, \leq 3; \mathbf{r})$ -TD code of size two.

CONSTRUCTION 2. *Let $\mathbf{r} \in \text{Irr}_{\leq 3}(i, 3)$ for some $i \geq 4$. Set*

$$C_2 \triangleq \begin{cases} \{r_1 r_2 r_3 r_1 r_2 r_3 r_4 \cdots r_i, & r_1 r_2 r_2 r_3 r_3 r_4 r_4 \cdots r_i\}, & \text{if } r_1 \neq r_3, \\ \{r_1 r_2 r_1 r_4 r_2 r_1 r_4 r_5 \cdots r_i, & r_1 r_2 r_1 r_1 r_4 r_4 r_5 r_5 \cdots r_i\}, & \text{if } r_1 = r_3 \text{ and } i \geq 5, \\ \{r_1 r_2 r_1 r_4 r_2 r_1 r_4, & r_1 r_2 r_1 r_1 r_4 r_4 r_4\}, & \text{if } r_1 = r_3 \text{ and } i = 4. \end{cases}$$

Then C_2 is an $(i + 3, \leq 3; \mathbf{r})$ -TD code of size two. By Proposition 4.4, C_2 is optimal.

PROOF. Let $C_2 = \{\mathbf{x}, \mathbf{y}\}$ and the labels of the words be as given in Proposition 4.2. Observe that in all cases, $(c_1^{(\mathbf{x})}, \delta_1^{(\mathbf{x})}) = (2, +)$ and $(c_1^{(\mathbf{y})}, \delta_1^{(\mathbf{y})}) = (1, -)$. Proposition 4.2 then implies that \mathbf{x} and \mathbf{y} are not confusable. ■

Following the above construction, we have an improvement to Construction 1.

COROLLARY 4.5. For $n \geq 6$, we have that

$$T(n) \geq \sum_{\substack{1 \leq i \leq 3, \text{ or} \\ n-2 \leq i \leq n}} \text{Irr}_{\leq 3}(i, 3) + 2 \sum_{i=4}^{n-3} \text{Irr}_{\leq 3}(i, 3).$$

Construction 2 also implies that $T(|\mathbf{r}| + 3, \mathbf{r}) = 2$ for all irreducible words \mathbf{r} . When $n \in \{|\mathbf{r}| + 4, |\mathbf{r}| + 5\}$, Proposition 4.4 states that $T(n, \mathbf{r}) \leq m + 1$, where m is the number of regions for \mathbf{r} . However, the upper bound can be reduced to a constant independent of m and we do this by analysing the possible labels of words in a tandem-duplication code.

PROPOSITION 4.6. $T(|\mathbf{r}| + 4, \mathbf{r}) \leq 4$.

PROOF. Let $T(|\mathbf{r}| + 4, \mathbf{r}) = b$, and suppose that $C(|\mathbf{r}| + 4, \mathbf{r})$ is the $(|\mathbf{r}| + 4, \leq 3; \mathbf{r})$ -TD code with size b . Let $C(|\mathbf{r}| + 4, \mathbf{r}) = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_b\}$. Let

$$\text{Label}(\mathbf{x}_i) \triangleq \left(\mathbf{r}, \left(c_1^{(\mathbf{x}_i)}, \delta_1^{(\mathbf{x}_i)} \right), \left(c_2^{(\mathbf{x}_i)}, \delta_2^{(\mathbf{x}_i)} \right), \dots, \left(c_m^{(\mathbf{x}_i)}, \delta_m^{(\mathbf{x}_i)} \right) \right).$$

Note that there can only be at most one word, $\mathbf{x}_i \in C(|\mathbf{r}| + 4, \mathbf{r})$, such that $c_t^{(\mathbf{x}_i)} = 1$ for all $1 \leq t \leq m$. We also know that $c_t^{(\mathbf{x}_i)} \leq 2$ for all $1 \leq i \leq b, 1 \leq t \leq m$. Otherwise, the length of \mathbf{x}_i is at least $|\mathbf{r}| + 6$. Furthermore, for each $\mathbf{x}_i \in C(|\mathbf{r}| + 4, \mathbf{r})$, there can only be at most one $1 \leq t \leq m$, such that $c_t^{(\mathbf{x}_i)} = 2$. So, we can define $C'(|\mathbf{r}| + 4, \mathbf{r}) = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{b'}\}$, where $b' \leq b$, to be the subset of $C(|\mathbf{r}| + 4, \mathbf{r})$ such that for all $1 \leq i \leq b'$, $c_t^{(\mathbf{x}_i)} = 2$ for exactly one $t \leq m$, and conclude that $b \leq 1 + b'$.

We say that \mathbf{x}_i loses to \mathbf{x}_j if there exists $1 \leq t \leq m$, such that $c_t^{(\mathbf{x}_i)} < c_t^{(\mathbf{x}_j)}$ and $\delta_t^{(\mathbf{x}_i)} = -$. Note that in $C'(|\mathbf{r}| + 4, \mathbf{r})$, for any distinct pair $1 \leq i < j \leq b'$, either \mathbf{x}_i loses to \mathbf{x}_j or \mathbf{x}_j loses to \mathbf{x}_i . Hence, in total, there is at least $b'(b' - 1)/2$ losses in $C'(|\mathbf{r}| + 4, \mathbf{r})$. Suppose that for some $1 \leq t \leq m$, we have $c_t^{(\mathbf{x}_i)} = 2$, then there can only be at most one $1 \leq t' \leq m, t' \neq t$, such that $\delta_{t'}^{(\mathbf{x}_i)} = -$. Otherwise, the length of \mathbf{x}_i is at least $|\mathbf{r}| + 5$. And since for distinct $i, j \leq b'$, we can never have $c_t^{(\mathbf{x}_i)} = c_t^{(\mathbf{x}_j)} = 2$. This implies that for each word $\mathbf{x}_i \in C'(|\mathbf{r}| + 4, \mathbf{r})$, the word \mathbf{x}_i can only lose at most once to another word in $C'(|\mathbf{r}| + 4, \mathbf{r})$. Therefore, there are at most b' losses in $C'(|\mathbf{r}| + 4, \mathbf{r})$. Hence, we have the inequality $b'(b' - 1)/2 \leq b'$, which implies that $b' \leq 3$, and therefore $b \leq 4$.

We conclude the proof with the labels of a possible set of four words in $C(|\mathbf{r}| + 4, \mathbf{r})$:

$$\begin{aligned} \text{Label}(\mathbf{x}_1) &= (\mathbf{r}, (2, +), (1, -), (1, +)), \\ \text{Label}(\mathbf{x}_2) &= (\mathbf{r}, (1, +), (2, +), (1, -)), \\ \text{Label}(\mathbf{x}_3) &= (\mathbf{r}, (1, -), (1, +), (2, +)), \\ \text{Label}(\mathbf{x}_4) &= (\mathbf{r}, (1, -), (1, -), (1, -)). \end{aligned} \quad \blacksquare$$

COROLLARY 4.7. $T(|\mathbf{r}| + 5, \mathbf{r}) \leq 6$.

PROOF. Following the proof of Proposition 4.6, there are at least $b'(b' - 1)/2$ losses in $C'(|\mathbf{r}| + 5, \mathbf{r})$. And for each word in $C'(|\mathbf{r}| + 5, \mathbf{r})$, it can have at most two losses. Hence, we have $b'(b' - 1)/2 \leq 2b'$, and so $T(|\mathbf{r}| + 5, \mathbf{r}) = b \leq 1 + b' \leq 1 + 5 = 6$.

As before, we provide the labels of a possible set of six words in $C(|\mathbf{r}| + 5, \mathbf{r})$:

$$\begin{aligned} \text{Label}(\mathbf{x}_1) &= (\mathbf{r}, (2, +), (1, -), (1, -), (1, +), (1, +)), \\ \text{Label}(\mathbf{x}_2) &= (\mathbf{r}, (1, +), (2, +), (1, -), (1, -), (1, +)), \\ \text{Label}(\mathbf{x}_3) &= (\mathbf{r}, (1, +), (1, +), (2, +), (1, -), (1, -)), \\ \text{Label}(\mathbf{x}_4) &= (\mathbf{r}, (1, -), (1, +), (1, +), (2, +), (1, -)), \\ \text{Label}(\mathbf{x}_5) &= (\mathbf{r}, (1, -), (1, -), (1, +), (1, +), (2, +)), \\ \text{Label}(\mathbf{x}_6) &= (\mathbf{r}, (1, -), (1, -), (1, -), (1, -), (1, -)). \end{aligned}$$

4.2 Exact Value of $T(n, \mathbf{r})$ When the Number of Regions is Exactly One

Proposition 4.2 suggests that we examine code constructions according to the number of regions of \mathbf{r} . When $|\mathbf{r}| \leq 2$, the number of regions is zero and hence $T(n, \mathbf{r}) = 1$. We consider the smallest nontrivial case where \mathbf{r} has exactly one region. Without loss of generality, assume that $\mathbf{r} \in R \triangleq \{012, 0120, 01201, 1012, 10120, 101201, 0121, 01202, 012010, 10121, 101202, 1012010\}$, and so, $\mathbf{t}_1 \triangleq 012$.

CONSTRUCTION 3. For $\mathbf{r} \in R$ and $\ell \geq 1$, define the words $\mathbf{x}(\mathbf{r}, \ell)$ and $\mathbf{z}(\mathbf{r}, \ell)$ with the following rule.

\mathbf{r}	$\mathbf{x}(\mathbf{r}, \ell)$	$\mathbf{z}(\mathbf{r}, \ell)$	\mathbf{r}	$\mathbf{x}(\mathbf{r}, \ell)$	$\mathbf{z}(\mathbf{r}, \ell)$
012	$0(112200)^{\ell-1}112$	$(012)^\ell$	1012	$10(112200)^{\ell-1}112$	$1(012)^\ell$
0120	$0(112200)^{\ell-1}11220$	$(012)^\ell 0$	10120	$10(112200)^{\ell-1}11220$	$1(012)^\ell 0$
01201	$0(112200)^{\ell-1}1122001$	$(012)^\ell 01$	101201	$10(112200)^{\ell-1}1122001$	$1(012)^\ell 01$
0121	$0(112200)^{\ell-1}1121$	$(012)^\ell 1$	10121	$10(112200)^{\ell-1}1121$	$1(012)^\ell 1$
01202	$0(112200)^{\ell-1}112202$	$(012)^\ell 02$	101202	$10(112200)^{\ell-1}112202$	$1(012)^\ell 02$
012010	$0(112200)^{\ell-1}11220010$	$(012)^\ell 010$	1012010	$10(112200)^{\ell-1}11220010$	$1(012)^\ell 010$

For $\mathbf{r} \in R$ and $n \geq |\mathbf{x}(\mathbf{r}, 2)|$, set $\ell_z = \lfloor (n - |\mathbf{r}|)/3 \rfloor + 1$ and

$$C_r(n) \triangleq \left\{ \xi_{n-|\mathbf{x}(\mathbf{r}, \ell)|}(\mathbf{x}(\mathbf{r}, \ell)) : |\mathbf{x}(\mathbf{r}, \ell)| \leq n \right\} \cup \left\{ \xi_{n-|\mathbf{z}(\mathbf{r}, \ell_z)|}(\mathbf{z}(\mathbf{r}, \ell_z)) \right\}.$$

Then $C_r(n)$ is an $(n, \leq 3; \mathbf{r})$ -TD code. Furthermore, $C_r(n)$ is optimal. Therefore, if we set $n_2 = |\mathbf{x}(\mathbf{r}, 2)|$, then we have

$$T(n, \mathbf{r}) = \begin{cases} \left\lfloor \frac{n-n_2}{6} \right\rfloor + 3, & \text{if } n \geq n_2, \\ 2, & \text{if } n_2 > n \geq |\mathbf{r}| + 3, \\ 1, & \text{if } |\mathbf{r}| + 3 > n \geq |\mathbf{r}|, \\ 0, & \text{otherwise.} \end{cases}$$

PROOF. Let $\mathcal{X} = \left\{ \xi_{n-|\mathbf{x}(\mathbf{r}, \ell)|}(\mathbf{x}(\mathbf{r}, \ell)) : |\mathbf{x}(\mathbf{r}, \ell)| \leq n \right\}$ and $M \geq 2$ is the biggest integer that satisfies $|\mathbf{x}(\mathbf{r}, M)| \leq n$. Observe that for all $\ell \geq 1$, we have $\text{Label}(\mathbf{x}(\mathbf{r}, \ell)) = (\mathbf{r}, (\ell, -))$ and $\text{Label}(\mathbf{z}(\mathbf{r}, \ell)) = (\mathbf{r}, (\ell, +))$. Furthermore, the length of $\mathbf{x}(\mathbf{r}, \ell)$ is at least $6(\ell - 1) + |\mathbf{r}| + 1$. Therefore, we have $6(M - 1) + |\mathbf{r}| + 1 \leq |\mathbf{x}(\mathbf{r}, M)| \leq n$. So, $\ell_z = \lfloor (n - |\mathbf{r}|)/3 \rfloor + 1 \geq \lfloor (6M - 5)/3 \rfloor + 1 \geq M + 1$ for $M \geq 2$.

Hence, for any word \mathbf{x} in \mathcal{X} , we have that $\text{Label}(\mathbf{x}) = (\mathbf{r}, (\ell', -))$ with $1 \leq \ell' \leq M < \ell_z$. Therefore, Proposition 4.2 implies that \mathbf{x} and $\mathbf{z}(\mathbf{r}, \ell_z)$ are not confusable. Also, Corollary 4.3(ii) implies that \mathbf{x} and \mathbf{x}' are not confusable for $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$. Therefore, $C_r(n)$ is an $(n, \leq 3; \mathbf{r})$ -TD code.

To demonstrate optimality, we observe that for $\mathbf{y} \in D_{\leq 3}^*(\mathbf{r})$, if $\text{Label}(\mathbf{y}) = (\ell, -)$, then $|\mathbf{y}| \geq |\mathbf{x}(\mathbf{r}, \ell)|$. Similarly, if $\text{Label}(\mathbf{y}) = (\ell, +)$, then $|\mathbf{y}| \geq |\mathbf{z}(\mathbf{r}, \ell)|$. Suppose that there is an $(n, \leq 3; \mathbf{r})$ -TD code C' with size $M + 2$. We know from Corollary 4.3(ii) that there is at most one codeword

in C' whose label δ_1 is +. So, there are at least $M + 1$ codewords in C' whose label δ_1 are -. Furthermore, Corollary 4.3(i), implies that no two codewords whose label δ_1 are -, have the same label c_1 . Therefore there exists a $\mathbf{y} \in C'$ whose label is $(c_1, -)$ with $c_1 \geq M + 1$. This implies that $n = |\mathbf{y}| \geq |\mathbf{x}(\mathbf{r}, M + 1)|$, which contradicts M being the biggest integer that satisfies $|\mathbf{x}(\mathbf{r}, M)| \leq n$.

Finally, the values of $T(n, \mathbf{r})$ follow from straightforward computations. \blacksquare

Hence, the value of $T(n, \mathbf{r})$ is completely determined whenever \mathbf{r} has at most one region. Combining this result with Proposition 4.4, we have the following theorem.

THEOREM 4.8. *Let $I(i, m)$ denote the number of irreducible words in $\text{Irr}_{\leq 3}(i, 3)$ with exactly m regions¹, and $U(n, i, m)$ be as defined in Proposition 4.4. Then,*

$$T(n) \leq \sum_{\mathbf{r} \in R} T(n, \mathbf{r}) + \sum_{i=5}^n \sum_{m=2}^i I(i, m)U(n, i, m). \quad (1)$$

4.3 Estimates of $T(n)$ for $n \leq 30$

In addition to the above constructions, we construct tandem-duplication codes for small lengths by searching for them exhaustively. Specifically, fix $\mathbf{r} \in \text{Irr}_{\leq 3}(i, 3)$ and $n \geq i$, we construct the graph $\mathcal{G}(n, \mathbf{r})$, whose vertices correspond to the set of all labels of descendants of \mathbf{r} of length n , or $\{L : L = \text{Label}(\mathbf{x}) \text{ and } \mathbf{x} \in \Sigma_3^n \cap D_{\leq 3}^*(\mathbf{r})\}$. Two vertices or labels L_1 and L_2 are connected in $\mathcal{G}(n, \mathbf{r})$ if and only if the words whose labels are L_1 and L_2 are not ≤ 3 -confusable. Hence, a clique of size M in the graph $\mathcal{G}(n, \mathbf{r})$ correspond to an $(n, \leq 3; \mathbf{r})$ -TD code of size M and $T(n, \mathbf{r})$ is the maximum size of a clique in $\mathcal{G}(n, \mathbf{r})$.

We use the exact algorithm MaxCliqueDyn [8] to determine the maximum size of the clique in these graph $\mathcal{G}(n, \mathbf{r})$ for $n \leq 20$. Since $T(n) = \sum_{\mathbf{r} \in \text{Irr}_{\leq 3}(i, 3), i \leq n} T(n, \mathbf{r})$, we tabulate the results $T(n)$ in Table 1.

Remark 3.

- Such a method to compute $T(n, \mathbf{r})$ is only possible because we have developed the algorithm to determine confusability in Section 3. Prior to this work, the necessary conditions for ≤ 3 -confusability was not known and hence, methods to compute $T(n, \mathbf{r})$ were not available.
- For a fixed length and root, even though the set of descendants is huge, the set of all labels are significantly smaller. Hence, the task of computing maximum cliques remains feasible. More concretely, when $n = 20$, the order of the largest graph $\mathcal{G}(20, \mathbf{r})$ is 366, despite the fact that the average size² of $\Sigma_3^{20} \cap D_{\leq 3}^*(\mathbf{r})$ is $3^{20}/27,687 \approx 125,935$.

Finally, we develop recursive constructions in the following proposition.

PROPOSITION 4.9. *Let $\mathbf{r} = r_1 r_2 \cdots r_i \in \text{Irr}_{\leq 3}(i, 3)$. Then the following holds:*

$$T(n, \mathbf{r}) \geq \begin{cases} T(n-1, \mathbf{r} \setminus r_1), & \text{if } r_1 = r_3, \\ \max\{2T(n-4, \mathbf{r} \setminus r_1), 3T(n-8, \mathbf{r} \setminus r_1)\}, & \text{if } r_1 \neq r_3, r_1 \neq r_4, \\ \max\{2T(n-5, \mathbf{r} \setminus r_1 r_2), 3T(n-10, \mathbf{r} \setminus r_1 r_2)\}, & \text{if } r_1 \neq r_3, r_1 = r_4, r_2 \neq r_5, \\ \max\{2T(n-6, \mathbf{r} \setminus r_1 r_2 r_3), 3T(n-12, \mathbf{r} \setminus r_1 r_2 r_3)\}, & \text{if } r_1 \neq r_3, r_1 = r_4, r_2 = r_5. \end{cases}$$

Furthermore, $T(n, \mathbf{r}) \geq T(n-1, \mathbf{r})$ and $T(n, \mathbf{r}) = T(n, \mathbf{r}^R)$, where \mathbf{z}^R denotes the reverse of word \mathbf{z} .

PROOF. First consider $r_1 = r_3$. Suppose that \mathcal{D} is an $(n-1, \leq 3; \mathbf{r} \setminus r_1)$ -TD code. To construct a code of length n , we simply prepend the prefix r_1 to all words in \mathcal{D} . For convenience, given a set of

¹We provide a formula in Appendix D.

²The number of ≤ 3 -irreducible ternary words of length at most 20 is 27,687.

Table 1. Lower, Upper Bounds and Exact Values for $T(n)$

Length	Constr. 1	This article	Equation (1)	Prop. 2.11	Length	Constr. 1	This article	Equation (1)	Prop. 2.11
1	3	3	3	3	16	5,985	10,641	12,267	15,495
2	9	9	9	9	17	8,781	16,287	19,479	25,077
3	21	21	21	21	18	12,879	25,005	30,957	40,581
4	39	39	39	39	19	18,885	38,223	49,245	65,667
5	69	69	69	69	20	27,687	57,957	78,417	106,257
6	111	117	117	117	21	40,587	83,619	125,001	171,933
7	171	195	195	195	22	59,493	116,145	199,467	278,199
8	261	315	315	321	23	87,201	166,761	318,621	450,141
9	393	495	495	525	24	127,809	249,159	509,457	728,349
10	585	777	777	855	25	187,323	375,129	815,361	1,178,499
11	867	1,221	1,227	1,389	26	274,545	558,573	1,306,107	1,906,857
12	1,281	1,887	1,941	2,253	27	402,375	813,771	2,093,967	3,085,365
13	1,887	2,913	3,075	3,651	28	589,719	1,164,309	3,359,685	4,992,231
14	2,775	4,527	4,875	5,913	29	864,285	1,675,935	5,394,369	8,077,605
15	4,077	6,969	7,731	9,573	30	1,266,681	2,464,419	8,667,075	13,069,845

The columns labelled Construction 1 and Proposition 2.11 correspond to the previous known lower and upper bounds for $T(n)$, respectively. The column labelled "This article" summarizes the best lower bounds that arise from this work, while the column labelled Equation (1) give the upper bound that result from Theorem 4.8. Optimal values of $T(n)$ are highlighted in **bold**.

words \mathcal{X} and a word \mathbf{p} , we use $\mathbf{p}\mathcal{X}$ to denote the set $\{\mathbf{p}\mathbf{x} : \mathbf{x} \in \mathcal{X}\}$. When $\mathbf{p} = p$ is of length one, we simply write $p\mathcal{X}$. Using this notation, we set $C = r_1\mathcal{D}$. We then apply Proposition 4.2 and verify that C is an $(n, \leq 3; \mathbf{r})$ -TD code whose size is given by $|\mathcal{D}|$.

Next, consider $r_1 \neq r_3$ and $r_1 \neq r_4$. Set $\mathbf{r}' = \mathbf{r} \setminus r_1$ and suppose that \mathcal{D} is an $(n-4, \leq 3; \mathbf{r}')$ -TD code. Let $\mathcal{D}_1 = r_1r_2r_2r_2\mathcal{D}$ and $\mathcal{D}_2 = r_1r_2r_3r_1\mathcal{D}$, and so, $\mathcal{D}_1 \cup \mathcal{D}_2 \subseteq D_{\leq 3}^*(\mathbf{r}) \cap \Sigma_3^n$.

For $\mathbf{x} \in \mathcal{D}$, let $\text{Label}(\mathbf{x}) = (\mathbf{r}', (c_1, \delta_1), \dots, (c_m, \delta_m))$. Then, we have that

$$\text{Label}(r_1r_2r_2r_2\mathbf{x}) = (\mathbf{r}, (1, -), (c_1, \delta_1), \dots, (c_m, \delta_m)),$$

$$\text{Label}(r_1r_2r_3r_1\mathbf{x}) = (\mathbf{r}, (2, +), (c_1, \delta_1), \dots, (c_m, \delta_m)).$$

Proposition 4.2 implies that $r_1r_2r_2r_2\mathbf{x}$ and $r_1r_2r_3r_1\mathbf{x}$ are not confusable. We can similarly check that any pair of distinct words in $\mathcal{D}_1 \cup \mathcal{D}_2$ are not confusable.

For the remaining cases, we choose the short code \mathcal{D} and prepend \mathcal{D} according to the rules below.

Conditions for \mathbf{r}	Short Code \mathcal{D}	$(n, \leq 3; \mathbf{r})$ -TD code
$r_1 = r_3$	$(n-1, \leq 3; \mathbf{r} \setminus r_1)$ -TD code	r_1
$r_1 \neq r_3, r_1 \neq r_4$	$(n-4, \leq 3; \mathbf{r} \setminus r_1)$ -TD code	$r_1r_2r_2r_2\mathcal{D}$ $\cup r_1r_2r_3r_1\mathcal{D}$
$r_1 \neq r_3, r_1 \neq r_4$	$(n-8, \leq 3; \mathbf{r} \setminus r_1)$ -TD code	$r_1r_2r_2r_2r_2r_2r_2\mathcal{D}$ $\cup r_1r_2r_2r_3r_3r_1r_1r_2\mathcal{D}$ $\cup r_1r_2r_3r_1r_2r_3r_1r_2\mathcal{D}$
$r_1 \neq r_3, r_1 = r_4, r_2 \neq r_5$	$(n-5, \leq 3; \mathbf{r} \setminus r_1r_2)$ -TD code	$r_1r_2r_2r_2r_3\mathcal{D}$ $\cup r_1r_2r_3r_1r_2\mathcal{D}$

Conditions for \mathbf{r}	Short Code \mathcal{D}	$(n, \leq 3; \mathbf{r})$ -TD code
$r_1 \neq r_3, r_1 = r_4, r_2 \neq r_5$	$(n - 10, \leq 3; \mathbf{r} \setminus r_1 r_2)$ -TD code	$r_1 r_2 r_2 r_3 r_3 r_3 r_3 r_3 r_3 \mathcal{D}$ $\cup r_1 r_2 r_2 r_3 r_3 r_1 r_1 r_2 r_2 r_3 \mathcal{D}$ $\cup r_1 r_2 r_3 r_1 r_2 r_3 r_1 r_2 r_3 r_3 \mathcal{D}$
$r_1 \neq r_3, r_1 = r_4, r_2 = r_5$	$(n - 6, \leq 3; \mathbf{r} \setminus r_1 r_2 r_3)$ -TD code	$r_1 r_2 r_2 r_3 r_3 r_1 \mathcal{D}$ $\cup r_1 r_2 r_3 r_1 r_2 r_3 \mathcal{D}$
$r_1 \neq r_3, r_1 = r_4, r_2 = r_5$	$(n - 12, \leq 3; \mathbf{r} \setminus r_1 r_2 r_3)$ -TD code	$r_1 r_2 r_2 r_3 r_3 r_1 r_1 r_1 r_1 r_1 r_1 \mathcal{D}$ $\cup r_1 r_2 r_2 r_3 r_3 r_1 r_1 r_2 r_2 r_3 r_1 \mathcal{D}$ $\cup r_1 r_2 r_3 r_1 r_2 r_3 r_1 r_2 r_3 r_1 r_1 \mathcal{D}$

To show that $T(n, \mathbf{r}) \geq T(n - 1, \mathbf{r})$, let \mathcal{D} be an $(n - 1, \leq 3; \mathbf{r})$ -TD code. Then $C = \{\xi_1(\mathbf{x}) : \mathbf{x} \in \mathcal{D}\}$ is an $(n, \leq 3; \mathbf{r})$ -TD code.

To show that $T(n, \mathbf{r}) = T(n, \mathbf{r}^R)$, let \mathcal{D} be an $(n, \leq 3; \mathbf{r})$ -TD code. Then $C = \{\mathbf{x}^R : \mathbf{x} \in \mathcal{D}\}$ is an $(n, \leq 3; \mathbf{r}^R)$ -TD code. ■

Using Proposition 4.9 with Constructions 2 and 3 and the values computed by MaxCliqueDyn, we derive lower bounds for $T(n)$ for $21 \leq n \leq 30$. The results are summarized in Table 1. In addition to the lower bounds for the code size $T(n)$, we also compare the upper bounds in Proposition 2.11 and (1). Observe that (1) is tight up to lengths at most ten and the constructions in this article improve the rates³ for Construction 1 by as much as 6.74%.

5 DISCUSSION

We studied the problem of determining the $\leq k$ -confusability of two words. Combining the results of this article, we have linear-time algorithms to solve the confusability problem for $k \in \{1, 2, 3\}$.

It remains open whether there exist efficient algorithms to determine ≤ 4 -confusability. One key obstacle is the fact that Proposition 2.7(iii) does not hold for $k = 4$. In particular, there exists \mathbf{x} and \mathbf{y} such that $R_{\leq 4}(\mathbf{x}) \neq R_{\leq 4}(\mathbf{y})$, but \mathbf{x} and \mathbf{y} are ≤ 4 -confusable. An example is provided by Jain et al. [7], where $\mathbf{x} = 012$ and $\mathbf{y} = 0121012$ belongs to $\text{Irr}_{\leq 4}(3)$, but 012101212 is a common descendant of \mathbf{x} and \mathbf{y} .

Another approach is to consider the intersection of the descendant cones $D_{\leq 4}^*(\mathbf{x}) \cap D_{\leq 4}(\mathbf{y})$ as formal languages. Unfortunately, Leupold et al. [11] demonstrated that while the language $D_{\leq 4}^*(\mathbf{x})$ is context free for all \mathbf{x} , the language is not regular in general. However, given two context free languages L_1 and L_2 , determining whether $L_1 \cap L_2$ is empty is an undecidable problem (see, for example, Sipser [16, Exercise 5.32]). While this does not imply that the confusability problem is undecidable, we nevertheless conjecture that it is undecidable for $k \geq 4$.

APPENDICES

A COMPUTING THE ROOT IN LINEAR TIME

We formally describe linear-time algorithms to compute all $\leq k$ -roots for a word for $k \in \{2, 3\}$. Given any $\mathbf{x} \in \Sigma_q^*$, since the root set of \mathbf{x} has size exactly one by Proposition 2.6, it suffices to find one $\mathbf{r} \in \text{Irr}_{\leq k}(q)$ so that $\mathbf{r} \xrightarrow[\leq k]{*} \mathbf{x}$. Lemma 3.3 then implies that we may reorder the duplications and assume that the tandem duplications are performed in decreasing lengths.

³The rate of a code C of length n is given by $(\log_3 |C|)/n$.

Therefore, to compute r , we simply remove duplicates in increasing length. Formally, given a word of the form $uvvw$ where $|v| = k'$, we say that v is a k' -duplicate. Removing the k' -duplicate v yields the word uvw . Therefore, to compute ≤ 2 -root, we simply remove all 1-duplicates from x and then all 2-duplicates from the result. If we want to compute ≤ 3 -root, then we remove all 3-duplicates from the ≤ 2 -root.

A formal description of the algorithm in Algorithm 2. Since removing all k' -duplicates may be performed in linear time, the algorithms run in linear-time.

ALGORITHM 2: Finding the $\leq k$ -root

```

1: procedure FIND- $\leq 2$ -ROOT( $x$ )
Input:  $x \in \Sigma_q^*$ 
Output:  $r$ , where  $R_{\leq 2}(x) = \{r\}$ .
2:    $r \leftarrow x$ 
3:   Remove all 1-duplicates in  $r$  using REMOVE_DUPLICATES( $r$ , 1)
4:   Remove all 2-duplicates in  $r$  using REMOVE_DUPLICATES( $r$ , 2)

5: procedure FIND- $\leq 3$ -ROOT( $x$ )
Input:  $x \in \Sigma_q^*$ 
Output:  $r$ , where  $R_{\leq 2}(x) = \{r\}$ .
6:    $r \leftarrow x$ 
7:   Remove all 1-duplicates in  $r$  using REMOVE_DUPLICATES( $r$ , 1)
8:   Remove all 2-duplicates in  $r$  using REMOVE_DUPLICATES( $r$ , 2)
9:   Remove all 3-duplicates in  $r$  using REMOVE_DUPLICATES( $r$ , 3)

10: procedure REMOVE_DUPLICATES( $r, k$ )
11:   while  $i + 2k - 1 \leq |r|$  do
12:     if the substring of length  $k$  starting at index  $i$  is equal to the substring of length  $k$  starting at
       index  $i + k$  then
13:       Remove the substring of length  $k$  starting at index  $i$ 
14:     else
15:        $i \leftarrow i + 1$ 

```

B ON THE SUM OF LENGTHS OF ALL PREFIXES

Recall in Section 3.2, the running time of Algorithm 1 is linear in the sum of the lengths of all prefixes that are computed in all recursive calls. Therefore, if this sum is linear in the lengths of the original words, then Algorithm 1 runs in linear time. Specifically, we establish the following proposition.

PROPOSITION B.1. *Let $x \in \Sigma_q^m$. Let p_1, p_2, \dots, p_s be all extended regions of r in x computed in Line 12 of Algorithm 1. Then $\sum_{i=1}^s |p_i| \leq 3m$.*

PROOF. Let $x = x_1x_2 \cdots x_m$. To establish the proposition, it suffices to show that each symbol x_j with $1 \leq j \leq m$ appears⁴ in at most three extended regions p_i, p_{i+1}, p_{i+2} for some $1 \leq i \leq s - 2$.

Let i be the smallest index such that x_j appears in the extended region p_i and let u_i be the suffix of x computed in Line 13 of Algorithm 1 in the same recursive call. Note that p_{i+1} is a prefix of u_i and u_{i+1} is a proper suffix of u_i . If x_j is still in p_{i+1} , which is a prefix of u_i , then x_j must belong

⁴In this proof, when we refer to a symbol x_j , we refer specifically to the symbol at index j . More formally, we may rewrite the word x as $(x_1, 1)(x_2, 2) \cdots (x_m, m)$.

to the overlapping substring of \mathbf{p}_i and \mathbf{u}_i (see Figure 1). Suppose that \mathbf{r}_i is the current root and $\text{Reg}(\mathbf{r}_i) = \mathbf{w}(abc)^\ell ab$. Then the overlapping substring is of the form ab^t for some $t \geq 1$, and we have the following cases.

- $x_j = a$. Then x_j is necessarily the first symbol of \mathbf{u}_i , and hence the first symbol of \mathbf{p}_{i+1} . Therefore, in the next recursive call, x_j is no longer in \mathbf{u}_{i+1} , which is a proper suffix of \mathbf{u}_i and hence not in \mathbf{p}_{i+2} . In other words, x_j appears only in \mathbf{p}_i and \mathbf{p}_{i+1} .
- $x_j = b$ for some b in the overlapping substring. To get \mathbf{u}_{i+1} , at least ab^{t-1} is removed from the prefix of \mathbf{u}_i . If exactly ab^{t-1} is removed from the prefix of \mathbf{u}_i to get \mathbf{u}_{i+1} , and x_j is the last b in the overlapping substring of \mathbf{p}_i and \mathbf{u}_i , then x_j still appears in \mathbf{u}_{i+1} and hence still appears in \mathbf{p}_{i+2} , but not in \mathbf{u}_{i+2} . In this case, x_j appears in \mathbf{p}_i , \mathbf{p}_{i+1} and \mathbf{p}_{i+2} , but not \mathbf{p}_{i+3} . Otherwise, if x_j is not the last b in the overlapping substring of \mathbf{p}_i and \mathbf{u}_i , x_j is no longer in \mathbf{u}_{i+1} , and hence x_j does not appear in \mathbf{p}_{i+2} . Hence, in this case, x_j appears in \mathbf{p}_i , \mathbf{p}_{i+1} , but not \mathbf{p}_{i+2} .

Therefore, from the two cases above, x_j appears in at most three consecutive extended regions. \blacksquare

C PROOF OF PROPOSITION 4.4

Let \mathbf{x} be a word of length n whose root is \mathbf{r} . Suppose that \mathbf{r} has length i and m regions and we set $\text{Label}(\mathbf{x}) = (\mathbf{r}, (c_1, \delta_1), (c_2, \delta_2), \dots, (c_m, \delta_m))$. To prove Proposition 4.4, we first show certain properties of the label of \mathbf{x} .

LEMMA C.1. *We have that*

$$i + 3(c_1 + c_2 + \dots + c_m) - m \leq n. \quad (2)$$

PROOF. Since $\mathbf{r} \xrightarrow[\leq 3]{*} \mathbf{x}$, then by Theorem 3.4, there exists \mathbf{x}' such that $\mathbf{r} \xrightarrow[3_d]{*} \mathbf{x}'$, and $\mathbf{x}' = R_{\leq 2}(\mathbf{x})$.

Let $1 \leq j \leq m$. Since there are c_j distinct triplets in the j th region of $R_{\leq 2}(\mathbf{x})$, we have to duplicate the distinct triplet in the j th region of \mathbf{r} by $c_j - 1$ times. Therefore, we have $|\mathbf{x}'| = |\mathbf{r}| + \sum_{j=1}^m 3(c_j - 1) = i + 3(c_1 + c_2 + \dots + c_m) - m$. Note that since $\mathbf{x}' = R_{\leq 2}(\mathbf{x})$, we have $|\mathbf{x}| \geq |\mathbf{x}'|$, and hence it yields Equation (2). \blacksquare

LEMMA C.2. *We have that $\delta_1 = \delta_2 = \dots = \delta_m = +$, whenever*

$$i + 3(c_1 + c_2 + \dots + c_m) - m = n. \quad (3)$$

PROOF. Since we know that $i + 3(c_1 + c_2 + \dots + c_m) - m$ is the length of $\mathbf{x}' = R_{\leq 2}(\mathbf{x})$, and the equality holds, that means there is no tandem duplication of length at most two in \mathbf{x} . Hence, a distinct triplet must remain in each of the extended regions of \mathbf{r} in \mathbf{x} and so, $\delta_j = +$ for all $1 \leq j \leq m$. \blacksquare

Finally, we complete the proof of Proposition 4.4.

Proposition 15. Let $i \leq n$. Suppose that $\mathbf{r} \in \text{Irr}_{\leq 3}(i, 3)$ has m regions. Then,

$$T(n, \mathbf{r}) \leq U(n, i, m) \triangleq \begin{cases} \binom{(n-i)/3+m}{m} - \binom{(n-i)/3+m-1}{m-1} + 1, & \text{if } 3 \text{ divides } n - i, \\ \binom{\lfloor (n-i)/3 \rfloor + m}{m}, & \text{otherwise.} \end{cases}$$

PROOF. Let C be an $(n \leq 3; \mathbf{r})$ -TD code. For $\mathbf{x} \in C$, set $\text{Label}(\mathbf{x}) = (\mathbf{r}, (c_1, \delta_1), (c_2, \delta_2), \dots, (c_m, \delta_m))$ and so, (c_1, c_2, \dots, c_m) is an integer solution to Equation (2) whose entries are all positive. Corollary 4.3(i) implies that these integer solutions are distinct for different \mathbf{x} chosen from C . Hence, the number of integer solutions to Equation (2) is an upper bound to the size of C . This number is well known (see Heubach and Mansour [5]) and is given by $\binom{\lfloor (n-i)/3 \rfloor + m}{m}$.

When $n - i$ is divisible by three, we improve this upper bound. Observe that the integer solutions to Equation (3) is a proper subset of Equation (2). Corollary 4.3(ii) then implies that there is at most one \mathbf{x} in C whose label yields a positive integer solution to Equation (3). Since the number of positive integer solutions to Equation (3) is $\binom{(n-i)/3+m-1}{m-1}$, we obtain the desired upper bound. ■

D NUMBER OF IRREDUCIBLE WORDS WITH CERTAIN NUMBER OF REGIONS

Recall that $I(i, m)$ denote the number of irreducible words in $\text{Irr}_{\leq 3}(i, 3)$ with exactly m regions. In this Appendix, we derive a recursive formula for $I(i, m)$ that allows us to efficiently compute Equation (1).

Let $\text{Irr}(aba, i, m)$ and $\text{Irr}(abc, i, m)$ denote the set of irreducible words of length i with exactly m regions that have two and three distinct symbols, respectively, in their prefixes of length three. Let $I(aba, i, m)$ and $I(abc, i, m)$ denote the sizes of $\text{Irr}(aba, i, m)$ and $\text{Irr}(abc, i, m)$, respectively.

We consider the maps,

$$\begin{aligned}\Phi_1 &: \text{Irr}(aba, i, m) \rightarrow \text{Irr}(abc, i-1, m), \\ \Phi_2 &: \text{Irr}(abc, i-1, m) \rightarrow \text{Irr}(aba, i, m), \\ \Psi_1 &: \text{Irr}(abc, i, m) \rightarrow \text{Irr}(aba, i-1, m-1) \cup \text{Irr}(aba, i-2, m-1) \cup \text{Irr}(aba, i-3, m-1), \\ \Psi_2 &: \text{Irr}(aba, i-1, m-1) \cup \text{Irr}(aba, i-2, m-1) \cup \text{Irr}(aba, i-3, m-1) \rightarrow \text{Irr}(abc, i, m),\end{aligned}$$

defined via the following rules. In what follows, we set $\mathbf{r} = r_1 r_2 \cdots r_{|\mathbf{r}|}$ and for distinct elements r_1, r_2 , set r^* to be the unique symbol distinct from r_1 and r_2 :

$$\begin{aligned}\Phi_1(\mathbf{r}) &= \mathbf{r} \setminus r_1, & \Phi_2(\mathbf{r}) &= r_2 \mathbf{r}, \\ \Psi_1(\mathbf{r}) &= \begin{cases} \mathbf{r} \setminus r_1, & \text{if } r_1 \neq r_4, \\ \mathbf{r} \setminus r_1 r_2, & \text{if } r_1 = r_4, r_2 \neq r_5, \\ \mathbf{r} \setminus r_1 r_2 r_3, & \text{if } r_1 = r_4, r_2 = r_5, \end{cases} & \Psi_2(\mathbf{r}) &= \begin{cases} r^* \mathbf{r}, & \text{if } |\mathbf{r}| = i-1, \\ r_2 r^* \mathbf{r}, & \text{if } |\mathbf{r}| = i-2, \\ r_1 r_2 r^* \mathbf{r}, & \text{if } |\mathbf{r}| = i-3. \end{cases}\end{aligned}$$

Then we check that the maps Φ_1, Φ_2, Ψ_1 , and Ψ_2 are well-defined and $\Phi_1 \circ \Phi_2, \Phi_2 \circ \Phi_1, \Psi_1 \circ \Psi_2$, and $\Psi_2 \circ \Psi_1$ are identity maps on their respective domains. Therefore, all four maps are bijections, and we establish the following recursion. For $m \geq 0$ and $i \geq 3$,

$$I(aba, i, m) = \begin{cases} I(abc, i-1, m), & i > 3, \\ 0, & \text{if } i = 3, m > 0, \\ 6, & \text{if } i = 3, m = 0. \end{cases}$$

$$I(abc, i, m) = \begin{cases} I(aba, i-1, m-1) + I(aba, i-2, m-1) + I(aba, i-3, m-1), & \text{if } i \geq 6, m \geq 1, \\ 0, & \text{if } i = 5, m > 2, \\ 6, & \text{if } i = 5, m = 2, \\ 12, & \text{if } i = 5, m = 1, \\ 0, & \text{if } i = 4, m > 1, \\ 12, & \text{if } i = 4, m = 1, \\ 0, & \text{if } i = 3, m > 1, \\ 6, & \text{if } i = 3, m = 1, \\ 0, & \text{if } m = 0. \end{cases}$$

Finally, to compute $I(i, m)$, we have that $I(i, m) = I(aba, i, m) + I(abc, i, m)$.

ACKNOWLEDGMENT

We are grateful for an anonymous reviewer who proposed the definition of label given in Section 4. This definition has led to a more succinct presentation and a spectrum of new results, including Proposition 4.4 and the computation of the size of optimal codes of lengths up to 20.

REFERENCES

- [1] Masanori Arita and Yoshiaki Ohashi. 2004. Secret signatures inside genomic DNA. *Biotechnol. Progr.* 20, 5 (2004), 1605–1607.
- [2] Farzad Farnoud, Moshe Schwartz, and Jehoshua Bruck. 2016. The capacity of string-duplication systems. *IEEE Trans. Info. Theory* 62, 2 (2016), 811–824.
- [3] John W. Fondon and Harold R. Garner. 2004. Molecular origins of rapid and continuous morphological evolution. *Proc. Natl. Acad. Sci. U.S.A.* 101, 52 (2004), 18058–18063.
- [4] Dominik Heider and Angelika Barnekow. 2007. DNA-based watermarks using the DNA-Crypt algorithm. *BMC Bioinform.* 8, 1 (2007), 176.
- [5] Silvia Heubach and Toufik Mansour. 2009. *Combinatorics of Compositions and Words*. CRC Press.
- [6] Siddharth Jain, Farzad Farnoud Hassanzadeh, and Jehoshua Bruck. 2017. Capacity and expressiveness of genomic tandem duplication. *IEEE Trans. Info. Theory* 63, 10 (2017), 6129–6138.
- [7] Siddharth Jain, Farzad Farnoud, Moshe Schwartz, and Jehoshua Bruck. 2017. Duplication-correcting codes for data storage in the DNA of living organisms. *IEEE Trans. Info. Theory* 63, 8 (2017), 4996–5010.
- [8] Janez Konc and Dušanka Janežic. 2007. An improved branch and bound algorithm for the maximum clique problem. *Proteins* 4, 5 (2007). Retrieved from <http://insilab.org/maxclique/>.
- [9] Eric S. Lander, Lauren M. Linton, Bruce Birren, Chad Nusbaum, Michael C. Zody, Jennifer Baldwin, Keri Devon, Ken Dewar, Michael Doyle, William FitzHugh et al. 2001. Initial sequencing and analysis of the human genome. *Nature* 409, 6822 (2001), 860–921.
- [10] Peter Leupold, Carlos Martin-Vide, and Victor Mitrana. 2005. Uniformly bounded duplication languages. *Discrete Appl. Math.* 146, 3 (2005), 301–310.
- [11] Peter Leupold, Victor Mitrana, and José M. Sempere. 2003. Formal languages arising from gene repeated duplication. In *Aspects of Molecular Computing*. Springer, 297–308.
- [12] Michael Liss, Daniela Daubert, Kathrin Brunner, Kristina Kliche, Ulrich Hammes, Andreas Leiberer, and Ralf Wagner. 2012. Embedding permanent watermarks in synthetic genes. *PLoS One* 7, 8 (2012), e42465.
- [13] Nicholas I. Mundy and Andreas J. Helbig. 2004. Origin and evolution of tandem repeats in the mitochondrial DNA control region of shrikes (*Lanius* spp.). *J. Molec. Evol.* 59, 2 (2004), 250–257.
- [14] Gonzalo Navarro. 2001. A guided tour to approximate string matching. *ACM Comput. Surveys* 33, 1 (2001), 31–88.
- [15] Seth L. Shipman, Jeff Nivala, Jeffrey D. Macklis, and George M. Church. 2017. CRISPR–Cas encoding of a digital movie into the genomes of a population of living bacteria. *Nature* 547, 7663 (2017), 345.
- [16] Michael Sipser. 2006. *Introduction to the Theory of Computation*, vol. 2. Thomson Course Technology, Boston.
- [17] Grant R. Sutherland and Robert I. Richards. 1995. Simple tandem DNA repeats and human genetic disease. *Proc. Natl. Acad. Sci. U.S.A.* 92, 9 (1995), 3636–3641.
- [18] Karen Usdin. 2008. The biological effects of simple tandem repeats: Lessons from the repeat expansion diseases. *Genome Res.* 18, 7 (2008), 1011–1019.

Received April 2018; revised May 2019; accepted May 2019