

Capacity and Expressiveness of Genomic Tandem Duplication

Siddharth Jain
Electrical Engineering
California Institute of Technology
Pasadena, CA 91125, U.S.A.
sidjain@caltech.edu

Farzad Farnoud (Hassanzadeh)
Electrical Engineering
California Institute of Technology
Pasadena, CA 91125, U.S.A.
farnoud@caltech.edu

Jehoshua Bruck
Electrical Engineering
California Institute of Technology
Pasadena, CA 91125, U.S.A.
bruck@caltech.edu

Abstract—To be considered for an 2015 IEEE Jack Keil Wolf ISIT Student Paper Award. The majority of the human genome consists of repeated sequences. An important type of repeats common in the human genome are tandem repeats, where identical copies appear next to each other. For example, in the sequence *AGTCTGTGC*, *TGTG* is a tandem repeat, namely, it was generated from *AGTCTGC* by tandem duplication of length 2. In this work, we investigate the possibility of generating a large number of sequences from a small initial string (called the seed) by tandem duplication of length bounded by a constant. Our results include *exact capacity* values for certain tandem duplication string systems with alphabet sizes 2, 3, and 4. In addition, motivated by the role of DNA sequences in expressing proteins via RNA and the genetic code, we define the notion of the *expressiveness* of a tandem duplication system, as the feasibility of expressing arbitrary substrings. We then *completely* characterize the expressiveness of tandem duplication systems for general alphabet sizes and duplication lengths. Noticing that a system with capacity = 1 is expressive, we prove that for an alphabet size ≥ 4 , the capacity is strictly smaller than 1, independent of the seed and the duplication lengths. The proof of this limit on the capacity (note that the genomic alphabet size is 4), is related to an interesting result by Axel Thue from 1906 which states that there exist arbitrary length sequences with no tandem repeats (square-free) for alphabet size ≥ 3 . Finally, our results illustrate that duplication lengths play a more significant role than the seed in generating a large number of sequences for these systems.

Index Terms—Expressiveness, tandem repeats, finite automata, irreducible strings.

I. INTRODUCTION

More than 50% of the human genome consists of repeated sequences [6]. These repeats are of two types i) interspersed repeats and ii) tandem repeats. Interspersed repeats are caused by transposons. A transposon (jumping gene) is a segment of DNA that can copy or cut and paste itself into new positions of the genome. Tandem repeats are thought to be caused by slipped-strand mispairings [10]. Slipped-strand mispairings are thought to occur when one DNA duplex becomes misaligned with the other.

Tandem Repeats are common in both prokaryote and eukaryote genomes. They are not only present in intergenic regions but also in both coding and non-coding regions. They are thought to be the cause of several genetic disorders. The effects of tandem repeats on several biological processes is understood by these disorders. They can result in generation of toxic or malfunctioning proteins, chromosome fragility, expansion

diseases, silencing of genes, modulation of transcription and translation [12] and rapid morphological changes [4].

The significance of sequences with tandem repeats and the fact that much of our unique DNA was likely originally a repeated sequence motivates us to study the *capacity* and *expressiveness* of string systems with tandem duplication. The model of a string duplication system consists of a starting string (seed) of finite length, a set of duplication rules and the set of all the sequences that can be obtained by applying the duplication rules on the seed a finite number of times. The notion of capacity was defined in [3]. It represents the average number of m -ary bits per symbol that are required asymptotically to encode a sequence in the string system (m is the alphabet size). The notion of expressiveness defined formally later answers the question whether each of the finite length sequence for a given alphabet can be obtained as a substring of some sequence in the string system. Expressiveness and capacity are closely related. More precisely, it is not difficult to show *if we have a system that is not expressive then capacity* < 1 [9].

A process that leads to tandem repeats is *tandem duplication* which allows substrings of certain lengths to be duplicated next to their original position. For example, from the sequence *AGTCGTCGCT*, a tandem duplication of length 2 can give *AGTCGTCGCGCT*, which if followed by a duplication of length 3 can give *AGTCGTCGTCGCGCT*. Tandem duplications have already been studied in [1], [2], [7], [8]. However the main concern of these works is to determine the place of tandem duplication rules in the Chomsky hierarchy of formal languages. A study related to our work can be found in [3]. In [3], the authors show that for a fixed duplication length the capacity is 0 in a tandem duplication string system. Further, they find a lower bound on the capacity of these systems, when duplications of all length are allowed. In this paper, we study tandem duplication string systems where we consider all the strings that can be obtained from a given starting string (seed) via a finite number of tandem duplications. More precisely, we consider tandem duplication string systems, where we restrict the maximum size of the block being tandemly duplicated to a certain *finite* length. In the rest of the paper, the term tandem duplication string system refers to these kind of string

duplication systems.

Example 1: To illustrate the notion of expressiveness and capacity for tandem duplication string systems, consider a string system on binary alphabet where the seed = 01 and the maximum allowable block size for duplication is 2. It is easy to check that the set of strings that can be generated by this system have to start with a 0 and end with a 1. In fact, it can be proved that all binary strings of length n which start with 0 and end with 1 can be generated by this system. The proof is based on the fact that every n length string which starts with 0 and ends with 1 can be rewritten as $0^{r_1}1^{r_2}\dots\dots\dots 0^{r_{m-1}}1^{r_m}$, where each $r_i \geq 1$ and m is even. Hence a natural way to generate such a string from seed = 01 is to repeat 01 $\frac{m}{2}$ times and then repeat each 0 or 1 at position i , r_i times.

Expressiveness: 11010 cannot be generated by this system. However, it can be generated as a substring of 0110101 in the following way:

$$01 \rightarrow 0101 \rightarrow 010101 \rightarrow 0110101.$$

If every binary string can be generated as a substring of some larger string in the duplication system, then we say that the system is expressive. In this case, since every binary string starting with 0 and ending with 1 can be generated, we can generate every binary string as a substring. Hence, the system is expressive.

Capacity: The number of n -length strings in this string system is 2^{n-2} and therefore the capacity is 1 bit/symbol.

Observing this fact for alphabet of size 2, one can ask related questions on expressiveness and capacity for higher alphabet sizes and duplication lengths. However, counting the number of n -length sequences for capacity calculation and characterizing expressive systems for higher alphabets is non-trivial for higher alphabets. In this paper, we study these questions and develop tools to answer them. It is interesting to observe that the string system over binary alphabet in the above example can be represented by a finite automata given in Figure 1. The regular expression for the language defined by the finite automata is given below which exactly represents all binary strings that start with 0 and end with 1.

$$R_{01} = (0^+1^+)^+ \quad (1)$$

One can use Perron-Frobenius theory [5], [9] to count the number of sequences which can be generated by a finite automata. In this paper, we use finite automata as a tool to calculate capacity for some string duplication systems with tandem repeats over higher alphabet. In our results, we find exact capacity of 0.876036 for a tandem duplication string system over ternary alphabet with seed = 012 and duplication size atmost 3. Furthermore, we show that for any tandem duplication string system over ternary alphabet with maximum duplication length 3, an expressive system does not exist. However, if the maximum duplication length is 4 and the seed is 012, then we get an expressive system. This shows that for such string duplication systems, the maximum duplication length plays a more significant role in generating larger number of strings than the seed. Further to emphasize this

fact we prove that over all tandem duplication string systems with a given alphabet size and maximum duplication length, an expressive tandem duplication string system has maximum capacity. We also find that for alphabet size > 3 , an expressive tandem duplication string system does not exist which shows that full capacity (i.e. capacity = 1) cannot be achieved by a tandem duplication string system for alphabet size > 3 .

It is easy to check that for binary alphabet any sequence of length ≥ 4 has a tandem repeat. The dependence of expressiveness and capacity on alphabet size is intuitively connected to a result by Thue [11] which states that for an alphabet size > 2 , there exists a square free sequence (sequence with no tandem repeat) for every length. In our proofs of results on expressiveness, we elaborate on this connection with Thue's result. The rest of the paper is organized as follows. In section II, we give some preliminary definitions and notation. In section III, we provide our results on capacity and expressiveness. We conclude the paper in section IV.

II. PRELIMINARIES

Let Σ be some finite alphabet. An n -string $x = x_1x_2\dots x_n \in \Sigma^n$ is a finite sequence where $x_i \in \Sigma$ and $|x| = n$. The set of all finite strings over alphabet Σ is denoted by Σ^* . For two strings $x \in \Sigma^n$ and $y \in \Sigma^m$, their concatenation is denoted by $xy \in \Sigma^{n+m}$. For a positive integer m and a string s , s^m denotes the concatenation of m copies of s . A string $v \in \Sigma^*$ is a substring of x if $x = uvw$, where $u, w \in \Sigma^*$.

A string system $S \subseteq \Sigma^*$ is represented as a tuple $S = (\Sigma, s, \mathcal{T})$, where $s \in \Sigma^*$ is a finite length string (seed) which is used to start the duplication process and \mathcal{T} is the duplication rule [3]. We denote by $N_S(n)$ the number of strings in S of length n . The *capacity* of the string system S is defined as:

$$cap(S) = \limsup_{n \rightarrow \infty} \frac{\log_{|\Sigma|} N_S(n)}{n}. \quad (2)$$

Next, we define the notion of *expressiveness*. A string system S is *expressive* if for each $y \in \Sigma^*$, there exists a $z \in S$, such that y is a substring of z .

Tandem Duplication of length k : $\mathcal{T}_k^{tan} : \Sigma^* \rightarrow \Sigma^*$, is defined as

$$\mathcal{T}_k^{tan}(x) = uvvw, \text{ where } x = uvw, |v| = k. \quad (3)$$

Furthermore, let $\mathcal{T}_{\leq k}^{tan}$ denote the set of tandem duplications of length at most k , i.e., $\mathcal{T}_{\leq k}^{tan} = \{\mathcal{T}_{k'}^{tan} | k' \leq k\}$. With this notation, for system considered in Example 1, $S = (\{0, 1\}, 01, \mathcal{T}_{\leq 2}^{tan})$.

III. RESULTS AND PROOFS

Our first result is on the capacity of a tandem duplication string system over ternary alphabet.

Theorem 1: For a tandem duplication string system $S = (\{0, 1, 2\}, 012, \mathcal{T}_{\leq 3}^{tan})$, $cap(S) = 0.876036$.

Proof: We prove this theorem by designing a finite automata for this system. Consider the finite automata given in Figure 2. The finite automata is designed in such a way so that it covers tandem duplications of length 1, 2 and 3. The self loops on each state cover duplication of length 1, connections

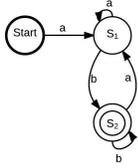


Fig. 1: Finite Automata for $S = (\{a, b\}, ab, \mathcal{T}_{\leq k}^{tan})$, where $k \geq 2$ (In binary alphabet, we do not gain anything by increasing k above 2), a and b are distinct symbols.

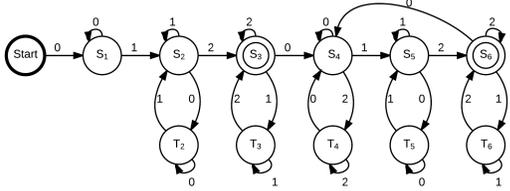


Fig. 2: Finite Automata for $S = (\{0, 1, 2\}, 012, \mathcal{T}_{\leq 3}^{tan})$. in between state pairs S_i and T_i cover duplications of length 2, and the connections in between states S_4, S_5 and S_6 covers duplications of length 3. In the rest of the proof, we show that the finite automata we construct in this way indeed represents S under consideration.

The regular expression for the language defined by this finite automata is given by

$$R_{012} = (0^+1^+)^+2^+(1^+2^+)^*[0^+(2^+0^+)^*1^+(0^+1^+)^*2^+(1^+2^+)^*]^* \quad (4)$$

Let $L_{R_{012}}$ denote the language defined by the regular expression R_{012} or finite automata in Figure 2. We claim that *Claim 1: $L_{R_{012}} \subseteq S$*

Before moving to the proof of Claim 1, we define the *de-duplication* process. Consider a de-duplication map $\mathcal{D}_{\leq k} : \Sigma^* \rightarrow P_{\leq k}^{\Sigma^*}$. Here $P_{\leq k}^{\Sigma^*}$ is the power set of strings in Σ^* which do not have a tandem repeat $\alpha\alpha$, where $|\alpha| \leq k$. For $x \in \Sigma^*$, $\mathcal{D}_{\leq k}(x)$ is the set of strings in $P_{\leq k}^{\Sigma^*}$ from which x can be obtained by tandem duplications of size at most k . For example, $\mathcal{D}_{\leq 3}(010100001) = \{01\}$, $\mathcal{D}_{\leq 3}(2122221212) = \{212\}$, $\mathcal{D}_{\leq 4}(012101212) = \{012, 0121012\}$.

Now, we show that for every string $x \in L_{R_{012}}$, $012 \in \mathcal{D}_{\leq 3}(x)$ or in other words every string $x \in L_{R_{012}}$ can be de-duplicated to 012 using $\mathcal{D}_{\leq 3}$.

The regular expression R_{012} in Eq. (4) can be represented as $R_{012} = B_1B_2^*$, where

$$B_1 = (0^+1^+)^+2^+(1^+2^+)^* \quad (5)$$

$$B_2 = 0^+(2^+0^+)^*1^+(0^+1^+)^*2^+(1^+2^+)^* \quad (6)$$

It is easy to check that de-duplication $\mathcal{D}_{\leq 3}$ converts $a^+ \rightarrow a$, $a^* \rightarrow \epsilon$ or a , $(ab)^+ \rightarrow ab$, $(ab)^* \rightarrow \epsilon$ or ab , $(abc)^+ \rightarrow abc$ and $(abc)^* \rightarrow \epsilon$ or abc , where a, b and c are distinct.

Next, we show that applying de-duplication $\mathcal{D}_{\leq 3}$ on B_1 gives 012 and on B_2 gives either 02012 or 012.

i) De-duplication $\mathcal{D}_{\leq 3}$ on B_1 : de-duplication

$$(0^+1^+)^+2^+(1^+2^+)^* \rightarrow 012(12)^* \rightarrow 012.$$

ii) De-duplication $\mathcal{D}_{\leq 3}$ on B_2 :

$$0^+(2^+0^+)^*1^+(0^+1^+)^*2^+(1^+2^+)^* \rightarrow 0(20)^*1(01)^*2(12)^*$$

$$\rightarrow 0(20)^*1(01)^*2 \rightarrow 0(20)^*12 \rightarrow 02012 \text{ or } 012.$$

Therefore, $B_1B_2^*$ can be de-duplicated to 012 by applying $\mathcal{D}_{\leq 3}$ since

$$B_1B_2^* \rightarrow 012(02012)^* \rightarrow 012$$

$$B_1B_2^* \xrightarrow{\text{or}} 012(012)^* \rightarrow 012.$$

Hence, any $x \in L_{R_{012}}$ can be de-duplicated to 012 by $\mathcal{D}_{\leq 3}$ or in other words, each $x \in L_{R_{012}}$ can be obtained by tandem duplications of length atmost 3 if the seed $s = 012$. Therefore, $L_{R_{012}} \subseteq S$.

Now, we claim that every $x \in S$ also belongs to $L_{R_{012}}$, i.e.

Claim 2: $S \subseteq L_{R_{012}}$

To prove this we need to show two things for the finite automata in Figure 2:

i) It can generate 012.

ii) If the automaton can generate pqr , with $p, q, r \in \Sigma^*$ and $|q| \leq 3$, it can also generate pq^2r .

(i) holds trivially. (See the path $Start - S_1 - S_2 - S_3$ in Figure 2)

Before moving further, we define:

- *Path Label (l_a)*: Given a path a in a finite automata, the path label $l_a \in \Sigma^*$ is defined as the concatenation of labels on the edges in the path.
- *Path Length* is the number of edges on the path.

For proving ii) we look at the adjacency matrix of the finite automata and show that for each state C all the 1, 2 and 3 length paths that end in C have a corresponding path with the same label which starts in C and ends in some state which is equivalent to C . The proof details are omitted here (See Appendix).

After proving $S = L_{R_{012}}$, we use Perron-Frobenius Theory [5], [9] to count the number of sequences which can be generated from this deterministic finite automata. We calculate the maximum absolute eigen value e^* of the adjacency matrix B of the strongly connected component of the finite automata in Figure 2 (i.e., $S_4, S_5, S_6, T_4, T_5, T_6$). B is given by

$$B = \begin{bmatrix} 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 \end{bmatrix}$$

$e^* = 2.618034$ for B above. By Perron-Frobenius Theory, $cap(S) = \log_3 e^* = 0.876036$ (upto 6 places of decimal). ■

Intuition as to why the capacity is less than 1: If we observe the regular expression for the finite automata, we see that we cannot generate a string which has 210, 021 or 102 as a substring (which also means that the system is not expressive). This further puts constraints on substrings of *size* > 3 that can be generated using this finite automata. Hence, we cannot achieve full capacity.

We know from Example 1 that the capacity for $S = (\{0, 1\}, 01, \mathcal{T}_{\leq 2}^{tan})$ is 1. Using this result, we can calculate the capacity for another tandem duplication string system given by $S_1 = (\{0, 1, 2\}, 012, \mathcal{T}_{\leq 2}^{tan})$. The main idea behind the proof here is the fact that S_1 can be decoupled into

two string systems equivalent to S . Similarly, we can calculate the capacity for $S_2 = (\{0, 1, 2, 3\}, 0123, \mathcal{T}_{\leq 3}^{tan})$ since it can be decoupled into two string systems equivalent to $S = (\{0, 1, 2\}, 012, \mathcal{T}_{\leq 3}^{tan})$. We omit the proof details here (See Appendix). Our capacity results are listed in Table I.

Σ	s	k	Capacity
$\{0, 1\}$	01	1	= 0
$\{0, 1\}$	01	≥ 2	= 1
$\{0, 1, 2\}$	012	2	= 0.630930
$\{0, 1, 2\}$	012	3	= 0.876036
$\{0, 1, 2, 3\}$	0123	3	= 0.694242

TABLE I: Capacity values for different string systems with starting string s that allow tandem duplications upto size k .

The next few theorems are on the expressiveness of tandem duplication string system,

Theorem 2: Consider $S = (\{0, 1, 2\}, s, \mathcal{T}_{\leq 3}^{tan})$, where s is any arbitrary starting string $\in \{0, 1, 2\}^*$. Then, S is not expressive.

Proof: Before, we move to the proof, let us define the notion of an *irreducible* string. A string $x \in \{0, 1, 2\}^*$ is irreducible if it does not have a tandem repeat $\alpha\alpha$, such that $|\alpha| \leq 3$. For example, 01201, 01210, 02101, 01210121 are irreducible strings. 01212, 021021, 01112 are not irreducible. To prove Theorem 2, we construct an irreducible string which cannot be generated by S as a substring of some $y \in S$.

At any stage of duplication in S , we can either do a tandem duplication of length 1 or 2 or 3. The string z on which the duplication is to be performed can be represented in the following way $z = uvw$, where $|v| \leq 3$ and v is the string that is to be tandemly duplicated. From tandem duplication of v in z , we get $z^* = uvvw$. We consider the following 3 cases and observe the irreducible substrings in z^* which do not possibly appear in z :

Case 1: $|v| = 1, v = a_1$.

Here $z = ua_1v$ and $z^* = ua_1a_1v$, the new substrings that we see in z^* are not irreducible. Since, they have a repeat a_1a_1 .

Case 2: $|v| = 2, v = a_1a_2$.

Here $z = ua_1a_2v$, and $z^* = ua_1a_2a_1a_2v$, the new possible irreducible substrings that we see in z^* of length ≥ 3 have either $a_1a_2a_1$ as suffix or $a_2a_1a_2$ as prefix, which means that if any new irreducible substring is generated in this step, either i) the letter on its first and third position is same or ii) the letter on its last and third last position are same.

Case 3: $|v| = 3, v = a_1a_2a_3$.

Here $z = ua_1a_2a_3v$, and $z^* = ua_1a_2a_3a_1a_2a_3v$, the new possible irreducible substrings that we see in z^* of length ≥ 4 have either $a_1a_2a_3a_1$ or $a_1a_2a_3a_1a_2$ as suffix or $a_3a_1a_2a_3$ or $a_2a_3a_1a_2a_3$ as prefix, which means that if any new irreducible substring is generated in this step, either i) the letter on its first and fourth position is same or ii) the letter on its last and fourth last position are same.

Consider, an arbitrary irreducible string $\in \{0, 1, 2\}^*$ of length ≥ 4 . Let $b_1b_2b_3b_4$ be its prefix and $c_4c_3c_2c_1$ be its suffix. From the 3 cases considered above, we have the following conditions, one of which has to be satisfied

by the irreducible substrings that can be generated by S , $b_1 = b_3$ or $b_1 = b_4$ or $c_1 = c_3$ or $c_1 = c_4$.

Now, we need to show that there are irreducible strings that do not satisfy any of the above 4 conditions. Consider irreducible strings of the form $t = (u_2u_1)^m u_2 a$ or $(u_1u_2)^m a$, where $u_1 = ab, u_2 = cb, m \geq 1$ and a, b and c are distinct symbols $\in \{0, 1, 2\}$. The 4-length suffix for strings of this form is $bcba$ and the 4-length prefix is either $abcb$ or $cbab$. None of these suffix or prefix satisfies any of the four conditions listed above. Hence, if not present in the seed s , irreducible substrings of this type cannot be generated by S . Since the seed s is of finite length, we have for some m , an irreducible string t with length $> |s|$ which cannot be generated as a substring of some string in S . Hence, S is not expressive. ■

Theorem 3: Consider $S = (\Sigma, s, \mathcal{T}_{\leq k}^{tan})$, where $|\Sigma| \geq 4, s$ is any arbitrary seed $\in \Sigma^*$ and k is some finite natural number, then S is not expressive, which also implies $cap(S) < 1$.

Proof: In this proof, we slightly change the definition of an *irreducible* string. We say that a string is irreducible if it does not have a tandem repeat. Now, we can extend and imitate the idea used in the proof of Theorem 2 (this time we have k cases). Consider an arbitrary irreducible string $\in \Sigma^*$ of length $\geq k+1$. Let $b_1b_2\dots b_k b_{k+1}$ be its prefix and $c_{k+1}c_k\dots c_1$ be its suffix. After considering k cases, we will get the following $2k - 2$ conditions, one of which has to be satisfied by the irreducible substrings that can be generated by S ,

$b_1 = b_{1+i}$ for some $i \in \{2, 3, 4, \dots, k\}$ or $c_1 = c_{1+j}$ for some $j \in \{2, 3, 4, \dots, k\}$.

Now, we show a construction of an irreducible string $\in \Sigma^*$ which does not satisfy any of the above listed conditions. Let $\Sigma = \{e_1, e_2, e_3, \dots, e_{|\Sigma|}\}$. Let $G = \{x : x \in \{e_2, \dots, e_{|\Sigma|}\}^*, |x| \geq k - 1, x \text{ is irreducible}\}$. Then, for any $y \in G$, it is easy to check that $t = e_1 y e_1$ does not satisfy any of the $2k - 2$ conditions listed above. By Thue [11], for alphabet size ≥ 3 , for any length there exists an irreducible string. Therefore, for each length $m \geq k - 1$, there exists a $y \in G$ with $|y| = m$. Since the seed s is of finite length, for some m we have an irreducible string t with length $> |s|$ which cannot be generated as a substring of some string in S . ■

Theorem 4 : Consider $S = (\{0, 1, 2\}, 012, \mathcal{T}_{\leq 4}^{tan})$, then S is an expressive string system.

Proof: Theorem 4 can be proved by an induction argument on the length of the substring that we want to generate. The system S considered in Theorem 4 clearly generates all the strings which can be generated using the system considered in Theorem 1. Looking at the finite automata in Figure 2 or R_{012} in Eq. (4) for string system considered in Theorem 1, it is easy to check that all possible 1 and 2 length strings over ternary alphabet can be obtained as substrings. Hence, all substrings of length 1 and 2 can be obtained using S considered in Theorem 4. Now to prove that all substrings of length 3 can also be obtained using S , it will be sufficient to prove that 021, 102 and 210 can be obtained as substrings using S , since other 3 length substrings can be obtained by system considered in Theorem 1 (again by observing the finite

automata in Figure 2 or R_{012} in Eq. (4)).

To generate 210,021 and 102 as substrings, here is the method:

$$\begin{aligned} &012 \rightarrow 01212 \rightarrow 012101212 \\ 012 &\rightarrow 012012 \rightarrow 01202012 \rightarrow 012021202012 \\ 012 &\rightarrow 012012 \rightarrow 01202012 \rightarrow 012020102012 \end{aligned}$$

For substring of length 4, we have the following 3 cases:

In case 1 and case 2 below, w is assumed to have at least one occurrence of each letter in the alphabet.

Case 1: The first 3 letters in the substring w are all distinct, i.e., if $w = w_1w_2w_3w_4$, then $w_1 \neq w_2 \neq w_3$. For generating such w as a substring, we first generate $w_1w_2w_3$ and then do a tandem duplication of w_3 if $w_4 = w_3$, of w_2w_3 if $w_4 = w_2$ and of $w_1w_2w_3$ if $w_4 = w_1$.

Case 2: Exactly 2 letters in the first 3 letters of w are same, i.e., if $w = w_1w_2w_3w_4$, then either $w_1 = w_2 \neq w_3$, or $w_1 = w_3 \neq w_2$, or $w_1 \neq w_2 = w_3$, if $w_1 = w_2$, then we first generate $w_1w_3w_4$ and then do a tandem duplication of w_1 to get $w = w_1w_1w_3w_4$. If $w_1 \neq w_2$, then we first generate $w' = w_4w_1w_2w_3$ as a substring, and then do a tandem duplication of w' to get w . (Note: w' is of type considered in *Case 1* since w_4 is different from both w_1 and w_2).

Case 3: w has ≤ 2 distinct letters, such a w has a tandem repeat. Therefore if $w = xyyz$, where either $|y| = 2$ and $|x| = |z| = 0$, or $|y| = 1$ and $|x| \leq 1$, $|z| = 2 - |x|$, then we first generate xyz and do a tandem duplication of y to get w . Till now, we have shown that all substrings w of length ≤ 4 can be generated.

For generating a substring w with $|w| > 4$, we use inductive argument. Let all substrings of length $\leq m$ can be generated (here $m \geq 4$), we need to prove that we can generate all substrings of length $m + 1$. Consider an arbitrary $w = a_1a_2\dots a_m a_{m+1}$, by induction assumption $w' = a_1a_2\dots a_m$ can be generated. Here, we have two cases: i) If all the three letters in the alphabet occur atleast once in $a_{m-3}a_{m-2}a_{m-1}a_m$, then w can be generated as a substring by a tandem duplication of some suffix of size ≤ 4 of w' . ii) If atleast one letter in the alphabet does not occur in $a_{m-3}a_{m-2}a_{m-1}a_m$, then $a_{m-3}a_{m-2}a_{m-1}a_m$ is a sequence over binary alphabet and hence is of the form $xyyz$ ($|y| = 1$ or 2), therefore w can be generated as a substring by tandem duplication of y on $a_1\dots a_{m-4}xyza_{m+1}$. (Note $|a_1\dots a_{m-4}xyza_{m+1}| \leq m$). Hence, we have proved Theorem 4. ■

Remark 1: It is important to note that in case (ii) above the binary sequence $a_{m-3}a_{m-2}a_{m-1}a_m$ has a tandem repeat. If the original alphabet size $|\Sigma|$ was > 3 , then this is not guaranteed over $|\Sigma| - 1$ -ary sequence of any length because of Thue's result [11].

Remark 2: For $S = (\{0, 1\}, s, \mathcal{T}_{\leq 1}^{tan})$, $(01)^m$ cannot be generated as a substring of any string $\in S$ for some m .

Table II gives a complete characterization of the expressiveness of tandem duplication string systems.

Theorem 5: For an expressive tandem duplication system $S = (\Sigma, s, \mathcal{T}_{\leq k}^{tan})$, $cap(S) \geq cap(S')$, where $S' = (\Sigma, s', \mathcal{T}_{\leq k}^{tan})$.

i.e. the capacity cannot be improved by only changing the seed if a tandem duplication string system is expressive for some seed s .

Proof: Since S is expressive, therefore s' can be generated as a substring of some string $z \in S$. Now, consider $S_z = (\Sigma, z, \mathcal{T}_{\leq k}^{tan})$, since $z \in S \Rightarrow S_z \subseteq S$, therefore $cap(S) \geq cap(S_z)$. Also $z = \alpha s' \beta$, where $\alpha, \beta \in \Sigma^*$. For any $x \in S'$, $\alpha x \beta \in S_z$ which implies that $N_{S'}(|x|) \leq N_{S_z}(|x| + |\alpha| + |\beta|)$, since $|\alpha|$ and $|\beta|$ are finite, we have $cap(S_z) \geq cap(S')$. Hence, $cap(S) \geq cap(S')$. ■

Remark 3: By Theorem 3, $|\Sigma| \leq 3$ in Theorem 5.

Σ	s	k	Expressive	Reason
$\{0\}$	0	≥ 1	Yes	Trivial
$\{0, 1\}$	arbitrary	1	No	Remark 2
$\{0, 1\}$	01	≥ 2	Yes	Example 1
$\{0, 1, 2\}$	arbitrary	≤ 3	No	Theorem 2
$\{0, 1, 2\}$	012	≥ 4	Yes	Theorem 4
Size ≥ 4	arbitrary	arbitrary	No	Theorem 3

TABLE II: Expressiveness of tandem duplication string systems where the maximum duplication length is k .

IV. CONCLUSION

It is still not clear to us what conditions on alphabet size, seed and maximum duplication length are needed for a tandem duplication string system to be representable as a regular language. As a future work, we would like to answer this question or develop other tools to calculate exact capacity or tight upper and lower bounds for such string duplication systems.

REFERENCES

- [1] J. Dassow, V. Mitrana, and G. Paun, "On the regularity of duplication closure," *Bulletin of the EATCS*, vol. 69, pp. 133-136, 1999.
- [2] J. Dassow, V. Mitrana, and A. Salomaa, "Operations and language generating devices suggested by the genome evolution," *Theoretical Computer Science*, vol. 270, no.1, pp. 701-738, 2002.
- [3] F. Farnoud, M. Schwartz, and J. Bruck, "The Capacity of String Duplication Systems," in *Proceedings of IEEE International Symposium on Information Theory*, pp. 1301-1305, 2014.
- [4] J. W. Fondon and H. R. Garner, "Molecular origins of rapid and continuous morphological evolution," *Proceedings of the National Academy of Sciences*, vol. 101, no. 52, pp. 18 058 – 18 063, 2004.
- [5] K. A. S. Imminck, *Codes for Mass Data Storage Systems*. Shannon Foundation Publishers, 2004.
- [6] E. S. Lander, L. M. Linton, B. Birren, C. Nusbaum, M. C. Zody, J. Baldwin, K. Devon, K. Dewar, M. Doyle, W. FitzHugh et al., "Initial sequencing and analysis of the human genome," *Nature*, vol. 409, no. 6822, pp. 860-921, 2001.
- [7] P. Leupold, C. Martin-Vide, and V. Mitrana, "Uniformly bounded duplication languages," *Discrete Applied Mathematics*, vol. 146, no. 3, pp. 301-310, 2005.
- [8] P. Leupold, V. Mitrana, and J. M. Sempere, "Formal languages arising from gene repeated duplication," in *Aspects of Molecular Computing*, Springer, 2004, pp. 297-308.
- [9] D. Lind and B. H. Marcus, *An Introduction to Symbolic Dynamics and Coding*. Cambridge University Press, 1985.
- [10] N. Mundy and A. J. Helbig, "Origin and evolution of tandem repeats in the mitochondrial DNA control region of shrikes (*Ianius* spp.)," *Journal of Molecular Evolution*, vol. 59, no. 2, pp. 250-257, 2004.
- [11] A. Thue, "über unendliche Zeichenreihen," *Kra. Vidensk. Selsk. Skrifter. I. Mat.-Nat. Kl., Cristiana* 7, 1906.
- [12] K. Usdin, "The biological effects of simple tandem repeats: lessons from the repeat expansion diseases," *Genome research*, vol. 18, no. 7, pp. 1011-1019, 2008.

APPENDIX

- *Claim 2 in Proof of Theorem 1:* We claim that every $x \in S = (\{0, 1, 2\}, 012, \mathcal{T}_{\leq 3}^{tan})$ also belongs to $L_{R_{012}}$, i.e.

$$S \subseteq L_{R_{012}}.$$

Proof: To prove this we need to show two things for the finite automata in Figure 2:

- It can generate 012.
 - If the automaton can generate pqr , with $p, q, r \in \Sigma^*$ and $|q| \leq 3$, it can also generate pq^2r .
- (i) holds trivially. (See the path $Start - S_1 - S_2 - S_3$ in Figure 2)

Before proving (ii), we define:

- *Path Label (l_a):* Given a path a in a finite automata, the path label $l_a \in \Sigma^*$ is defined as the concatenation of labels on the edges in the path.
- *Path Length* is the number of edges on the path.

For proving ii) we show that for each state C all the 1, 2 and 3 length paths that end in C have a corresponding path with the same label which starts in C and ends in some state which is equivalent to C . In order to prove (ii), we need to show the following condition holds for all states in Figure 2

Condition 1: For each state C in Figure 2, if q ($|q| \leq 3$) ends at C , then there is a path with the label q which starts in C and ends in a state equivalent to C .

If condition 1 is true for a given state C_o , it means that if pqr ($|q| \leq 3$) can be generated by finite automata in Figure 2 with q ending in C_o , then pq^2r can also be generated.

Below is a property which is sufficient for Condition 1 to hold for a given state.

Property 1 for a path of length j : Given a state u and a $j = 1, 2$ or 3 , let P_j^u be the set of all j -length paths ending in u and Q_j^u be the set of all j -length paths starting and ending in u . We say, property 1 holds if $\bigcup_{a \in P_j^u} l_a = \bigcup_{a \in Q_j^u} l_a$.

In rest of the proof we show that Condition 1 holds for all states in Figure 2.

Part 1 : We prove that Property 1 holds for all states $\{S_4, S_5, S_6, T_4, T_5, T_6\}$ in Figure 2 for paths of length 1, 2 and 3. This is shown by computing $\mathcal{A}_1, \mathcal{A}_1^2$ and \mathcal{A}_1^3 , where \mathcal{A}_1 is the adjacency matrix of the strongly connected component of the finite automata given in Figure 2. It is sufficient to consider only states in the strongly connected component here, since for every path starting from any state $\in \{S_1, S_2, S_3, T_2, T_3\}$ and ending in any state $\in \{S_4, S_5, S_6, T_4, T_5, T_6\}$, there is a corresponding path with the same labels starting from some state $\in \{S_4, S_5, S_6, T_4, T_5, T_6\}$. In fact, in finding the corresponding path, one can substitute S_1 by S_4 , S_2 by S_5 , S_3 by S_6 , T_2 by T_5 and T_3 by T_6 . For example, the path $S_2 - S_3 - S_4 - T_4$ has the same path label as $S_5 - S_6 - S_4 - T_4$.

The adjacency matrix \mathcal{A}_1 where x, y and z represent edges labelled by 0, 1 and 2 respectively is given by

$$\mathcal{A}_1 = \begin{bmatrix} x & y & 0 & z & 0 & 0 \\ 0 & y & z & 0 & x & 0 \\ x & 0 & z & 0 & 0 & y \\ x & 0 & 0 & z & 0 & 0 \\ 0 & y & 0 & 0 & x & 0 \\ 0 & 0 & z & 0 & 0 & y \end{bmatrix}$$

$$\mathcal{A}_1^2 = \begin{bmatrix} x^2+zx & y^2+xy & yz & z^2+zx & yx & 0 \\ zx & y^2+xy & z^2+yz & 0 & x^2+yx & zy \\ x^2+zx & xy & z^2+yz & xz & 0 & y^2+zy \\ x^2+zx & xy & 0 & z^2+zx & 0 & 0 \\ 0 & y^2+xy & yz & 0 & x^2+yx & 0 \\ zx & 0 & z^2+yz & 0 & 0 & y^2+zy \end{bmatrix}$$

By observing \mathcal{A}_1 and \mathcal{A}_1^2 , we can easily see that the diagonal entries in both the matrices is the union of the corresponding column. This means that Property 1 holds for 1 and 2 length paths. By computing \mathcal{A}_1^3 using computer it can be checked that Property 1 holds for all 3 length paths as well.

Part 2 : Here, we prove that Condition 1 holds for all states $\in \{S_1, S_2, S_3, T_2, T_3\}$. We first prove that property 1 holds for all states $\in \{S_1, S_2, T_2, T_3\}$ for paths of length 1, 2 and 3 and holds for S_3 for paths of length 1 and 2. Next, we show that Property 1 does not hold for paths of length 3 for S_3 , however Condition 1 still holds. Observe that there is no path of any length from any state $\in \{S_4, S_5, S_6, T_4, T_5, T_6\}$ to any state $\in \{S_1, S_2, S_3, T_2, T_3\}$, hence we only need the 5×5 adjacency matrix of $\{S_1, S_2, S_3, T_2, T_3\}$ represented by \mathcal{A}_2 which is given by

$$\mathcal{A}_2 = \begin{bmatrix} x & y & 0 & 0 & 0 \\ 0 & y & z & x & 0 \\ 0 & 0 & z & 0 & y \\ 0 & y & 0 & x & 0 \\ 0 & 0 & z & 0 & y \end{bmatrix}$$

$$\mathcal{A}_2^2 = \begin{bmatrix} x^2 & xy & yz & yx & 0 \\ 0 & y^2+xy & z^2+yz & x^2+yx & zy \\ 0 & 0 & z^2+yz & 0 & y^2+zy \\ 0 & y^2+xy & yz & x^2+yx & 0 \\ 0 & 0 & z^2+yz & 0 & y^2+zy \end{bmatrix}$$

By observing \mathcal{A}_2 and \mathcal{A}_2^2 , we see that the diagonal entries in both the matrices is the union of the corresponding column. This means that Property 1 holds for all states $\in \{S_1, S_2, S_3, T_2, T_3\}$ for paths of length 1 and 2. By computing \mathcal{A}_2^3 using computer, it can be checked that that Property 1 holds for all states $\in \{S_1, S_2, T_2, T_3\}$ for paths of length 3 as well.

For S_3 , there is a 3-length path $S_1 - S_1 - S_2 - S_3$ with label 012, for which there does not exist a corresponding path with the same label which starts and ends in S_3 due to which property 1 does not hold for S_3 for paths of length 3. But, for this 3-length path, we can traverse $S_3 - S_4 - S_5 - S_6$ which also has label 012, to get $p(012)^2$. Now, since S_3 and S_6 are equivalent, $p(012)^2r$ can be generated.

Hence, by Part 1 and 2, Condition 1 holds for all states and hence (ii) holds. Thus, we have proved Claim 2. ■

- *Calculation of Capacity for $S_1 = (\{0, 1, 2\}, 012, \mathcal{T}_{\leq 2}^{tan})$ and $S_2 = (\{0, 1, 2, 3\}, 0123, \mathcal{T}_{\leq 3}^{tan})$.*

Before, we calculate the capacities for S_1 and S_2 , we state and prove the following lemma:

Lemma 1: For a string duplication system $S = (\{a_1, a_2, \dots, a_{|\Sigma|}\}, a_1 a_2 \dots a_{|\Sigma|}, \mathcal{T}_{\leq k}^{tan})$, for each $x \in S$ if $i < j$, then the last occurrence of a_i in x is before the last occurrence of a_j .

Proof: We prove this using induction on the number of steps required to generate a $x \in S$.

Consider a string $x \in S$ that can be generated in one step, to generate such a x a substring of $s = a_1 a_2 \dots a_{|\Sigma|}$ of size $\leq k$ is chosen and is tandemly duplicated. The substring is of the form $a_p \dots a_q$, where $1 \leq p \leq q \leq \min\{p+k-1, |\Sigma|\}$. Then $x = a_1 a_2 \dots a_p \dots a_q a_p \dots a_q a_{q+1} \dots a_{|\Sigma|}$, hence Lemma 1 holds $\forall x \in S$ which can be generated in one step.

Now, suppose Lemma 1 holds $\forall x \in S$ that can be generated in $\leq m$ steps. We need to prove that Lemma 1 holds $\forall x \in S$ that can be generated in $m+1$ steps. Let $x' = uyvw$ be a string $\in S$ that can be generated in m steps. Hence, Lemma 1 holds for x' . Here w starts with a_1 and has only one occurrence of a_1 , i.e. the last occurrence of a_1 in x' is the starting position of w . Now to generate x from x' , we consider two cases:

Case 1: If a substring y of uyv is chosen and tandemly duplicated, then $x = uyyvw$ and Lemma 1 holds for x .

Case 2: If a substring z of w is chosen and tandemly duplicated. Let $w = w_1 z w_2$, and a_p, \dots, a_q have their last occurrences ($1 \leq p \leq q \leq |\Sigma|$) in z (in the same order), then a_1, \dots, a_{p-1} have their last occurrences in w_1 (in the same order) and $a_{q+1}, \dots, a_{|\Sigma|}$ have their last occurrences in w_2 (in the same order). After repeating z , we have $x = uyvw_1 z z w_2$. For x , a_1, \dots, a_{p-1} have their last occurrences in w_1 (in the same order), a_p, \dots, a_q have their last occurrences in zz (in the same order) and $a_{q+1}, \dots, a_{|\Sigma|}$ (in the same order) have their last occurrences in w_2 . Hence, Lemma 1 holds for x .

Since, there was no restriction on the choice of x' except for that it can be generated in m steps, Lemma 1 holds for $\forall x \in S$ that can be generated in $m+1$ steps. Thus, we have proved Lemma 1. ■

Now, using Lemma 1 and our results on capacity of $S = (\{0, 1\}, 01, \mathcal{T}_{\leq 2}^{tan})$ and $S = (\{0, 1, 2\}, 012, \mathcal{T}_{\leq 3}^{tan})$, we calculate the capacity of S_1 and S_2 respectively.

Example 2: Consider a string system $S = (\{0, 1, 2\}, 012, \mathcal{T}_{\leq 2}^{tan})$. For this system $\forall x \in S$, last occurrence of 0 is always before the first occurrence of 2, since by Lemma 1, 0 and 2 are always separated by atleast one occurrence of 1 and the maximum duplication size is 2, hence the regular expression is given by:

$$R = R_{01}^+ R_{12}^+ + 0^+ R_{12}^+ + R_{01}^+ 2^+. \quad (7)$$

Here R_{01} is given by Eq. (1) and R_{12} is same as R_{01} with 0 substituted by 1 and 1 substituted by 2. The finite automata for this system is given in Figure 3. The number

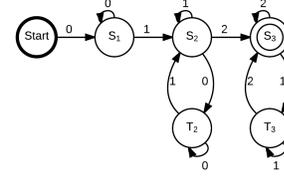


Fig. 3: Finite Automata for $S = (\{0, 1, 2\}, 012, \mathcal{T}_{\leq 2}^{tan})$. The regular expression $R = R_{01}^+ R_{12}^+ + 0^+ R_{12}^+ + R_{01}^+ 2^+$.

of strings that can be generated by this new system is of the order $n * 2^n$, which is asymptotically the same (to the first order in the exponent) to that of the string duplication system $S = (\{0, 1\}, 01, \mathcal{T}_{\leq 2}^{tan})$ and hence the capacity for this new system is given by $\log_3 2 = 0.630930$.

Example 3: Consider a string system $S = (\{0, 1, 2, 3\}, 0123, \mathcal{T}_{\leq 3}^{tan})$. For this system $\forall x \in S$, last occurrence of 0 is always before the first occurrence of 3, since by Lemma 1, 0 and 3 are always separated by atleast one occurrence of 1 and one occurrence of 2 and the maximum duplication size is 3, hence the regular expression is given by:

$$R = R_{012}^+ R_{123}^+ + 0^+ R_{123}^+ + R_{012}^+ 3^+. \quad (8)$$

Here R_{012} is defined in Eq. (4) and R_{123} is same as R_{012} with 0 substituted by 1, 1 substituted by 2 and 2 substituted by 3. The number of strings that can be generated by this new system is of the order $n * 3^{0.876036n}$ which is asymptotically the same (to the first order in the exponent) to that of the string duplication system $S = (\{0, 1, 2\}, 012, \mathcal{T}_{\leq 3}^{tan})$ and hence the capacity for this new system is given by $\log_4 3^{0.876036} = 0.694242$.

- Not Expressive \Rightarrow capacity < 1 .

Proof: Since S is not expressive therefore there exists a $z \in \Sigma^*$ such that z does not appear as a substring of any $y \in S$. Let $|z| = m$, then for $n = m\lambda + \mu$, where $\lambda \in \{0, 1, 2, \dots\}$ and $\mu \in \{0, 1, \dots, m-1\}$,

$$N_S(n) \leq (|\Sigma|^m - 1) \lfloor \frac{n}{m} \rfloor |\Sigma|^\mu.$$

Since m is finite, therefore $cap(S) < 1$. ■