

# A Stochastic Model for Genomic Interspersed Duplication

Farzad Farnoud  
Electrical Engineering  
California Institute of Technology  
Pasadena, CA 91125, USA  
Email: farnoud@caltech.edu

Moshe Schwartz  
Electrical and Computer Engineering  
Ben-Gurion University of the Negev  
Beer Sheva 8410501, Israel  
Email: schwartz@ee.bgu.ac.il

Jehoshua Bruck  
Electrical Engineering  
California Institute of Technology  
Pasadena, CA 91125, USA  
Email: bruck@paradise.caltech.edu

**Abstract**—Mutation processes such as point mutation, insertion, deletion, and duplication (including tandem and interspersed duplication) have an important role in evolution, as they lead to genomic diversity, and thus to phenotypic variation. In this work, we study the expressive power of interspersed duplication, i.e., its ability to generate diversity, via a simple but fundamental stochastic model, where the length and the location of the subsequence that is duplicated and the point of insertion of the copy are chosen randomly. In contrast to combinatorial models, where the goal is to determine the set of possible outcomes regardless of their likelihood, in stochastic systems, we investigate the properties of the set of high-probability sequences. In particular we provide results regarding the asymptotic behavior of frequencies of symbols and short words in a sequence evolving through interspersed duplication. The study of such a systems is an important step towards the design and analysis of more realistic and sophisticated models of genomic mutation processes.

## I. INTRODUCTION

It is estimated that there are about 8.7 million species on earth [5]. Of course, individuals within each species are also different from each other. Thus there is a vast amount of biological diversity on Earth. This diversity is, for the most part, the result of genomic mutation. The types of mutation include point mutation, insertion/deletion, and duplication. Duplication mutations, where a segment of DNA is copied and inserted elsewhere in the genome, may in turn be of the tandem or interspersed type. In tandem duplication, the copy is inserted immediately after the original, while in interspersed repeats, the copy may be inserted far from the original DNA segment. Interspersed duplications are caused by transposons, which are segments of DNA that can copy and insert themselves into new positions of the genome.

The general goal of this work is to move towards better understanding the effects of genomic interspersed duplication on generating novel sequences and creating biological diversity. Among the aforementioned mutation processes, interspersed duplication is of particular interest as it leads to interspersed repeated sequences, which form 45% of the human genome [4]. Here, we take a *probabilistic* approach to model interspersed duplications and investigate their ability to generate novelty and diversity. This complements our previous work [2], [3], in which we considered the same problem from a *combinatorial* point of view, with the goal of studying the

set of *possible* sequences arising from duplication systems. In contrast, the probabilistic view is concerned with identifying the *probable* outcomes of stochastic duplication systems. Particularly, we seek to find certain properties which the outcome of an interspersed-duplication system will possess with high probability.

In our interspersed-duplication model, a string evolves through random interspersed-duplication events, i.e., in each step, a random segment of the string is duplicated and then inserted in a random position in the string, independent of the position of the original segment. To avoid complications arising from boundary cases, we consider *circular strings*. It is worth noting that in fact many bacteria have circular chromosomes. While in practice, different mutation processes work together to create novel sequences, the scope of this work is limited to analyzing interspersed duplications in isolation. This helps us to obtain a better understanding of the properties of this type of mutation. We leave the study of more complex systems that evolve through more than one mutation type to future work.

Our analysis starts by considering how the frequencies (multiplicities divided by the length of the evolving string) of the alphabet symbols change as duplications occur. We show that under general conditions, the frequencies are martingales and thus converge almost surely. The same argument does not apply to the frequencies of strings of length larger than one. To analyze such frequencies, we use the stochastic approximation method which enables modeling of a discrete dynamic system by a corresponding continuous model described by ordinary differential equations. We show then that for interspersed-duplication systems, the frequencies of strings of length larger than one are, in the limit, consistent with those of iid sequences; implying that in a certain sense, a sequence evolving through interspersed duplication is unrecognizable from an iid sequence. Note that an iid sequence has the maximum entropy among sequences with a given symbol distribution.

The rest of the paper is organized as follows. Notation and preliminaries are given in the next section. Section III contains the analysis of the evolution of symbol frequencies in the evolving string. In Section IV, we present the necessary background and preliminaries for the use of stochastic approximation in this work. Section V is devoted to the analysis of



#### IV. STOCHASTIC APPROXIMATION FOR DUPLICATION SYSTEMS

In this section, we present a brief overview of the stochastic approximation method adapted to duplication systems. For an ordered set  $U$ , let  $\boldsymbol{\mu}_n = (\mu_n^u)_{u \in U}$  be a vector representing the number of appearances of objects  $u \in U$  in the string  $s$  at time  $n$  and let  $\mathbf{x}_n = \frac{\boldsymbol{\mu}_n}{L_n}$  be the normalized version of  $\boldsymbol{\mu}_n$ . For example,  $U$  can be the set of all strings over  $\mathcal{A}$  with length at most three. We also let  $\{\mathcal{F}_n\}$  be the filtration generated by the random variables  $\{\mathbf{x}_n, L_n\}$ . Our goal is to find out how  $\mathbf{x}_n$  changes with  $n$  by finding a differential equation whose solution approximates  $\mathbf{x}_n$ .

We state a set of conditions that must be satisfied for our analysis. Let  $\mathbb{E}_\ell[\cdot]$  denote the expected value conditioned on the fact that the length of the duplicated substring is  $\ell$  and let  $\boldsymbol{\delta}_\ell = \mathbb{E}_\ell[\boldsymbol{\mu}_{n+1}|\mathcal{F}_n] - \boldsymbol{\mu}_n$ . We consider the following conditions. Among them, we assume (A1) and, for now, accept the others without a proof.

- (A1) There exists  $K \in \mathbb{N}$  such that  $q_i = 0$  for  $i = 0$  or  $i \geq K$ .
- (A2)  $\boldsymbol{\mu}_{n+1} - \boldsymbol{\mu}_n$ , and thus  $\boldsymbol{\delta}_\ell$ , are bounded.
- (A3)  $\mathbf{x}_n$  is bounded.
- (A4) For each  $\ell$ ,  $\boldsymbol{\delta}_\ell$  is a function of  $\mathbf{x}_n$  only, so we can write  $\boldsymbol{\delta}_\ell = \boldsymbol{\delta}_\ell(\mathbf{x}_n)$ .
- (A5) The function  $\boldsymbol{\delta}_\ell(\mathbf{x}_n)$  is Lipschitz.

To understand how  $\mathbf{x}_n$  varies, our starting point is its difference sequence  $\mathbf{x}_{n+1} - \mathbf{x}_n$ . We note that

$$\mathbf{x}_{n+1} - \mathbf{x}_n = \mathbb{E}[\mathbf{x}_{n+1} - \mathbf{x}_n|\mathcal{F}_n] + (\mathbf{x}_{n+1} - \mathbb{E}[\mathbf{x}_{n+1}|\mathcal{F}_n]). \quad (1)$$

For the first term of the right side of (1), we have<sup>1</sup>

$$\begin{aligned} \mathbb{E}[\mathbf{x}_{n+1} - \mathbf{x}_n|\mathcal{F}_n] &= \sum_{\ell} q_\ell (\mathbb{E}_\ell[\mathbf{x}_{n+1}|\mathcal{F}_n] - \mathbf{x}_n) \\ &= \sum_{\ell} q_\ell \left( \frac{\boldsymbol{\mu}_n + \boldsymbol{\delta}_\ell(\mathbf{x}_n)}{L_n + \ell} - \frac{\boldsymbol{\mu}_n}{L_n} \right) \\ &= \frac{1}{L_n} \sum_{\ell} q_\ell \mathbf{h}_\ell(\mathbf{x}_n) (1 + O(L_n^{-1})) \\ &= \frac{1}{L_n} \mathbf{h}(\mathbf{x}_n) (1 + O(L_n^{-1})), \end{aligned} \quad (2)$$

where  $\mathbf{h}_\ell(\mathbf{x}) = \boldsymbol{\delta}_\ell(\mathbf{x}) - \ell \mathbf{x}$ ,  $\mathbf{h}(\mathbf{x}) = \sum_{\ell} q_\ell \mathbf{h}_\ell(\mathbf{x})$ , and where we have used  $1/(L_n + \ell) = (1 + O(L_n^{-1}))/L_n$  which follows from the boundedness of  $\ell$  (see (A1)).

Furthermore, for the second term of the right side of (1), we have

$$\begin{aligned} \mathbf{x}_{n+1} - \mathbb{E}[\mathbf{x}_{n+1}|\mathcal{F}_n] &= \frac{\boldsymbol{\mu}_{n+1}}{L_{n+1}} - \mathbb{E} \left[ \frac{\boldsymbol{\mu}_{n+1}}{L_{n+1}} \middle| \mathcal{F}_n \right] \\ &= \frac{1 + O(L_n^{-1})}{L_n} (\boldsymbol{\mu}_{n+1} - \mathbb{E}[\boldsymbol{\mu}_{n+1}|\mathcal{F}_n]) \\ &= \frac{1}{L_n} (1 + O(L_n^{-1})) M_{n+1} \end{aligned} \quad (3)$$

<sup>1</sup>Here, we abuse the notation by using  $\ell$ , which is a random variable, as an index of summations when the goal is finding expected values by conditioning on  $\ell$ .

where  $M_{n+1} = \boldsymbol{\mu}_{n+1} - \mathbb{E}[\boldsymbol{\mu}_{n+1}|\mathcal{F}_n]$ . Note that  $M_n$  is a bounded martingale difference sequence.

From (1), (2), and (3), we find  $\mathbf{x}_{n+1} = \mathbf{x}_n + \frac{1}{L_n} (\mathbf{h}(\mathbf{x}_n) + M_{n+1} + O(L_n^{-1}))$ , where we have used the fact that  $\mathbf{h}(\mathbf{x}_n)(1 + O(L_n^{-1})) = \mathbf{h}(\mathbf{x}_n) + O(L_n^{-1})$ . This follows from the boundedness of  $\mathbf{h}(\mathbf{x}_n)$ , which in turn follows from the boundedness of  $\boldsymbol{\delta}(\mathbf{x}_n)$ .

The following theorem relates the discrete system describing  $\mathbf{x}_n$  to a continuous system.

**Theorem 3.** [1, Theorem 2] *The sequence  $\{\mathbf{x}_n\}$  converges almost surely to a compact connected internally chain transitive invariant set of the ode  $d\mathbf{x}_t/dt = \mathbf{h}(\mathbf{x}_t)$ .*

Note the dual use of the symbol  $\mathbf{x}$  in the theorem; the meaning is however clear from the subscript. Recall that a set  $A$  is an invariant set of an ode  $d\mathbf{z}_t/dt = \mathbf{f}(\mathbf{z}_t)$  if it is closed and  $\mathbf{z}_{t'} \in A$  for some  $t' \in \mathbb{R}$  implies that  $\mathbf{z}_t \in A$  for all  $t \in \mathbb{R}$ . The invariant set  $A$  is internally chain transitive with respect to the ode  $d\mathbf{z}_t/dt = \mathbf{f}(\mathbf{z}_t)$ , provided that for every  $\mathbf{y}, \mathbf{y}' \in A$  and positive reals  $T$  and  $\epsilon$ , there exist  $N \geq 1$  and a sequence  $\mathbf{y}_0, \dots, \mathbf{y}_N$  with  $\mathbf{y}_i \in A$ ,  $\mathbf{y}_0 = \mathbf{y}$ , and  $\mathbf{y}_N = \mathbf{y}'$  such that for  $0 \leq i < n$ , if  $\mathbf{z}_0 = \mathbf{y}_i$ , then for some  $t \geq T$ ,  $\mathbf{z}_t$  is in the  $\epsilon$ -neighborhood of  $\mathbf{y}_{i+1}$ .

#### V. INTERSPERSED DUPLICATION

Now, we use the technique presented in Section IV to extend the results of Section III by analyzing the frequencies of strings  $u \in \mathcal{A}^*$  in an interspersed-duplication system.

Let  $\ell'_{max} \in \mathbb{N}$  and let  $U$  be an ordered set consisting of all strings of length at most  $\ell'_{max}$ . The vectors  $\mathbf{x}_n$  and  $\boldsymbol{\mu}_n$  are defined as before using  $U$ . Consider  $u \in U$ . To illustrate, we assume  $\mathcal{A} = \{A, C, G, T\}$  and will use  $u = ACT$  and  $\ell = 1$ . In an interspersed-duplication system, for  $\ell < |u|$ , we have

$$\begin{aligned} \delta_\ell^u &= -(|u| - 1)x_n^u + \sum_{i=1}^{\ell} x_n^{u_1, i} x_n^{u_{i+1}, |u|-i} \\ &\quad + \sum_{i=1}^{\ell} x_n^{u_1, |u|-i} x_n^{u_{|u|-i+1}, i} + \sum_{i=1}^{|u|-\ell-1} x_n^{u_1, i} x_n^{u_{i+\ell+1}, |u|-\ell-i} x_n^{u_{i+1}, \ell}. \end{aligned}$$

Here, the term  $-(|u| - 1)x_n^u$  accounts for the expected number of lost occurrences of  $u$  in  $s$  as a result of inserting the duplicate substring. For example, an occurrence of  $u = ACT$  will be lost if the symbol G is duplicated and inserted after A in this occurrence of  $u$ , since it becomes AGCT. The probability that a certain occurrence is lost equals  $\frac{|u|-1}{L_n}$ . Since there are  $\mu_n^u$  such occurrences, the expected number of lost occurrences of  $u$  equals  $\mu_n^u \frac{|u|-1}{L_n} = x_n^u (|u| - 1)$ . Note that if the symbol T is duplicated and inserted after C in an occurrence of ACT, we still count the original occurrence as lost, but count a new occurrence in the resulting ACTT, as seen in what follows. We now explain the first summation above. This summation represents the newly created occurrences of  $u$  where the first  $i$  symbols come from the duplicate and the next  $|u| - i$  are from the substring that starts after the point of insertion of the duplicate. There are  $\mu_n^{u_1, i}$  occurrences of  $u_{1, i}$ . The duplicate

starts with one of these with probability  $\frac{\mu_{n,1,i}}{L_n} = x_n^{u_{1,i}}$ . Furthermore, the duplicate is inserted before an occurrences of  $u_{i+1,|u|-i}$  with probability  $x_n^{u_{i+1,|u|-i}}$ . Hence, the probability of a new occurrence created in this way is  $x_n^{u_{1,i}} x_n^{u_{i+1,|u|-i}}$ , and so is the expected number of such new occurrences. The role of the second summation is similar, except that the duplicate provides the second part of  $u$ . The last summation accounts for new occurrences of  $u$  in which the duplicate substring forms a middle part of  $u$  of length  $\ell$  and previously existing substrings contribute a prefix of length  $i$  and a suffix of length  $|u| - \ell - i$ . In terms of our running example with  $u = \text{ACT}$  and  $\ell = 1$ , one such new occurrence is created if **C** is duplicated and inserted after **A** in an occurrence of **AT**. The probability of such an event is  $x_n^{u_{1,i} u_{i+\ell+1,|u|-\ell-i}} x_n^{u_{i+1,\ell}} = x_n^{\text{AT}} x_n^{\text{C}}$ , where  $i = 1$ .

For  $\ell \geq |u|$ , we have

$$\begin{aligned} \delta_\ell^u &= -(|u| - 1)x_n^u + \sum_{i=1}^{|u|-1} x_n^{u_{1,|u|-i}} x_n^{u_{|u|-i+1,i}} \\ &\quad + \sum_{i=1}^{|u|-1} x_n^{u_{1,i}} x_n^{u_{i+1,|u|-i}} + (\ell - |u| + 1)x_n^u \end{aligned}$$

where the first two summations are similar to the first two summations for the case of  $\ell < |u|$ , but a term corresponding to the third summation is not present. The term  $(\ell - |u| + 1)x_n^u$  corresponds to the cases in which a new occurrence of  $u$  is created as a substring of the duplicate substring. Note that  $\delta_\ell^u$  is bounded, depends only on  $x_n$ , and is Lipschitz since  $x_n \in [0, 1]^{|U|}$ .

Since  $h_\ell^u(x) = \delta_\ell^u(x) - \ell x^u$ , we have

$$\begin{aligned} h_\ell^u(x) &= -(\ell + |u| - 1)x^u + \sum_{i=1}^{\ell} x_n^{u_{1,i}} x_n^{u_{i+1,|u|-i}} \\ &\quad + \sum_{i=1}^{\ell} x_n^{u_{1,|u|-i}} x_n^{u_{|u|-i+1,i}} + \sum_{i=1}^{|u|-\ell-1} x_n^{u_{1,i} u_{i+\ell+1,|u|-\ell-i}} x_n^{u_{i+1,\ell}} \end{aligned} \quad (4)$$

for  $\ell < |u|$ , and

$$h_\ell^u(x) = -2(|u| - 1)x^u + 2 \sum_{i=1}^{|u|-1} x_n^{u_{1,i}} x_n^{u_{i+1,|u|-i}} \quad (5)$$

for  $\ell \geq |u|$ . Recall that  $\mathbf{h}_\ell(x) = (h_\ell^v(x))_{v \in U}$ . So from (4) and (5), we can find the ode  $dx_t/dt = \mathbf{h}(x_t) = \sum_{\ell} q_\ell \mathbf{h}_\ell(x_t)$ . As an example, if  $\ell_{max}^u = 2$  and  $\mathcal{A} = \{A, C\}$ , then  $U = (A, C, AA, AC, CA, CC, AAA, \dots, CCC)$  and some of the equations of the ode system are

$$\begin{aligned} \frac{d}{dt} x_t^A &= \frac{d}{dt} x_t^C = 0, \\ \frac{d}{dt} x_t^{AA} &= -2x_t^{AA} + 2(x_t^A)^2, \quad \frac{d}{dt} x_t^{AC} = -2x_t^{AC} + 2x_t^A x_t^C, \\ \frac{d}{dt} x_t^{AAC} &= -(4 - q_1)x_t^{AAC} + 2x_t^A x_t^{AC} + (2 - q_1)x_t^C x_t^{AA}. \end{aligned} \quad (6)$$

For a vector  $x$  that contains the elements  $(x^a)_{a \in \mathcal{A}}$  and for  $v \in \mathcal{A}^*$ , define  $p(v, x) = \prod_{a \in \mathcal{A}} (x^a)^{n_v(a)}$  and note that  $p(vw, x) = p(v, x)p(w, x)$ . We now turn to find the solutions to the ode  $dx_t/dt = \mathbf{h}(x_t)$ .

**Lemma 4.** Consider the ode  $dx_t/dt = \mathbf{h}(x_t)$  where  $\mathbf{h}(x) = \sum_{\ell} q_\ell \mathbf{h}_\ell(x)$  and the elements of  $\mathbf{h}_\ell(x)$  are given by (4) and (5). The solution to this ode is

$$x_t^v = p(v, x_0) + \sum_i b_i^v e^{-d_i^v t}, \quad v \in U, \quad (7)$$

where  $x_0 = x_t|_{t=0}$ ; the range of  $i$  in the summation is finite; and  $b_i^v$  and  $d_i^v$  are constants with  $d_i^v > 0$ .

*Proof:* We prove the lemma by induction. The claim (7) holds for  $v \in \mathcal{A}$ , since the equations for  $x_t^a$ ,  $a \in \mathcal{A}$ , are of the form  $dx_t^a/dt = 0$  and so  $x_t^a = x_0^a$ . Fix  $u \in U$  such that  $|u| > 1$ , and assume that (7) holds for all  $v \in U$  such that  $|v| < |u|$ . We show that it also holds for  $u$ , i.e.,  $x_t^u = p(u, x_0) + \sum_i b_i^u e^{-d_i^u t}$ . Using the assumption, we rewrite (4) and (5) as

$$h_\ell^u(x_t) = -(\ell + |u| - 1)(x_t^u - p(u, x_0)) + \sum_i b_i' e^{-d_i' t}$$

for  $\ell < |u|$ , and

$$h_\ell^u(x_t) = -2(|u| - 1)(x_t^u - p(u, x_0)) + \sum_i b_i'' e^{-d_i'' t}$$

for  $\ell \geq |u|$ , where  $b_i', d_i', b_i'', d_i''$  are constants with  $d_i', d_i'' > 0$ . Hence,  $h^u(x_t)$  can be written as

$$h^u(x_t) = -c^u(x_t^u - p(u, x_0)) + \sum_i b_i''' e^{-d_i''' t},$$

where  $c^u = 2|u| - 2 - \sum_{\ell=1}^{|u|-1} q_\ell(|u| - 1 - \ell)$ , and  $b_i''', d_i'''$  are constants with  $d_i''' > 0$ . Thus the solution to the ode  $dx_t^u/dt = h^u(x_t)$  is

$$\begin{aligned} x_t^u &= e^{-c^u t} \int e^{c^u t'} \left( c^u p(u, x_0) + \sum_i b_i''' e^{-d_i''' t'} \right) dt' + \bar{b} e^{-c^u t} \\ &= p(u, x_0) + \sum_i b_i^u e^{-d_i^u t}, \end{aligned}$$

where  $\bar{b}, b_i^u, d_i^u$  are some constants, with  $d_i^u > 0$  (note that  $c^u > 0$  since  $|u| > 1$ ). This completes the proof. ■

For example, the solutions to (6) with  $q_1 = 0$  are

$$\begin{aligned} x_t^A &= x_0^A, \quad x_t^C = x_0^C, \\ x_t^{AA} &= (x_0^A)^2 + b_1^{AA} e^{-2t}, \quad x_t^{AC} = x_0^A x_0^C + b_1^{AC} e^{-2t}, \\ x_t^{AAC} &= (x_0^A)^2 x_0^C + b_1^{AAC} e^{-2t} + b_2^{AAC} e^{-4t}, \end{aligned}$$

where  $b_1^{AAC} = x_0^A b_1^{AC} + x_0^C b_1^{AA}$ .

In the next theorem, we use Lemma 4 to characterize the limits of the frequencies of substrings in interspersed-duplication systems.

**Theorem 5.** Let  $U$  be an ordered set consisting of all strings over the alphabet  $\mathcal{A}$  of a certain maximum length and let

$\mathbf{x}_n = (x_n^u)_{u \in U}$  be the vector of frequencies of these strings at time  $n$  in an interspersed-duplication system. The vector  $\mathbf{x}_n$  converges almost surely. Furthermore, its limit  $\mathbf{x}_\infty$  satisfies

$$x_\infty^u = \prod_{a \in \mathcal{A}} (x_\infty^a)^{n_u(a)}, \quad \text{for all } u \in U.$$

*Proof:* From Theorem 3, we know that the limit set of  $\mathbf{x}_n$  is an internally chain transitive invariant set of the ode described by (4) and (5). Let this set, which consists of points of the form  $\mathbf{y} = (y^v)_{v \in U}$ , be denoted by  $A$ . Since for each  $u \in U$ ,  $x_n^u \in [0, 1]$ , we have that  $A \subseteq [0, 1]^{|U|}$ . We now use these facts to show that for each  $\mathbf{y} \in A$  and  $u \in U$ , we have  $y^u = p(u, \mathbf{y})$ .

Suppose to the contrary that there exist  $\mathbf{y} \in A$  and  $u \in U$  such that  $y^u \neq p(u, \mathbf{y})$ . Among all possible choices of such  $\mathbf{y}$  and  $u$ , choose the ones where the length  $|u|$  of  $u$  is minimum. Hence,  $y^u \neq p(u, \mathbf{y})$  but  $z^v = p(v, \mathbf{z})$  for all  $v \in \mathcal{A}^*$  with  $|v| < |u|$ , and all  $\mathbf{z} \in A$ . Then, similar to the proof of Lemma 4, one can show that if  $\mathbf{x}_0 = \mathbf{z} \in A$ , then  $x_t^u = p(u, \mathbf{z}) + be^{-c^u t}$ , where  $b = z^u - p(u, \mathbf{z})$  and  $c^u \geq |u|$ .

By the definition of internal chain transitivity, for any  $\epsilon > 0$  and  $T > 0$ , there exist  $N \geq 1$  and a sequence  $\mathbf{y}_0, \dots, \mathbf{y}_N$  with  $\mathbf{y}_i \in A$ ,  $\mathbf{y}_0 = \mathbf{y}_N = \mathbf{y}$  such that for  $0 \leq i < n$ , if  $\mathbf{x}_0 = \mathbf{y}_i$ , then there exists  $t \geq T$  such that  $\mathbf{x}_t$  is in the  $\epsilon$ -neighborhood of  $\mathbf{y}_{i+1}$ . Suppose  $\mathbf{x}_0 = \mathbf{y}_i$  and suppose for  $t' \geq T$ ,  $\mathbf{x}_{t'}$  is in the  $\epsilon$ -neighborhood of  $\mathbf{y}_{i+1}$ . We have

$$\begin{aligned} y_{i+1}^u &\leq x_{t'}^u + \epsilon \leq p(u, \mathbf{y}_i) + (y_i^u - p(u, \mathbf{y}_i))e^{-c^u t'} + \epsilon \\ &\leq p(u, \mathbf{y}_i) + e^{-c^u T} + \epsilon, \end{aligned}$$

where we have used the fact that  $y_i^u \leq 1$ . Furthermore, since  $y_{i+1}^a \leq y_i^a + \epsilon$  for  $a \in \mathcal{A}$ , we find

$$p(u, \mathbf{y}_{i+1}) - p(u, \mathbf{y}_i) \geq (1 + \epsilon)^{|u|} - 1,$$

where we have again used the fact that  $y_i^a \leq 1$  for  $a \in \mathcal{A}$ . It thus follows that

$$y_{i+1}^u - p(u, \mathbf{y}_{i+1}) \leq e^{-c^u T} + \epsilon + (1 + \epsilon)^{|u|} - 1.$$

In particular, this holds for  $i = n - 1$ , i.e.,

$$y^u - p(u, \mathbf{y}) \leq e^{-c^u T} + \epsilon + (1 + \epsilon)^{|u|} - 1.$$

But we can make the right side of the above inequality arbitrary small by choosing  $T$  large enough and  $\epsilon$  small enough. Thus  $y^u = p(u, \mathbf{y})$ , which is a contradiction. Hence, for each  $\mathbf{y} \in A$  and  $u \in U$ , we have  $y^u = p(u, \mathbf{y})$ , and the theorem follows.  $\blacksquare$

In words, the theorem shows that for  $u \in \mathcal{A}^*$ , the frequency of  $u$  converges to the frequency of same in an iid sequence where the probability of  $a \in \mathcal{A}$  equals  $x_\infty^a$ . Figure 2 illustrates an example, obtained via simulation, where the system starts with  $s^{(0)} = \text{AGCGTATGCG}$  and duplications of lengths 4 and 6 occur with equal probability. As the number  $n$  of duplications increases, the frequency vector  $\mathbf{x}_n$  becomes more compatible with that of an iid sequence. For example for  $n = 15000$ , we have  $x_n^{\text{AC}} = 0.0251 \simeq x_n^{\text{A}}x_n^{\text{C}} = 0.0266$ ,  $x_n^{\text{GT}} = 0.0872 \simeq x_n^{\text{G}}x_n^{\text{T}} = 0.0880$ , and  $x_n^{\text{GGG}} = 0.0992 \simeq (x_n^{\text{G}})^3 = 0.1084$

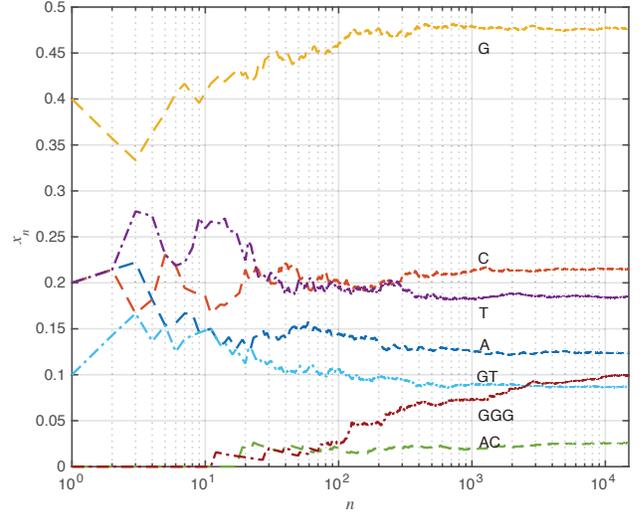


Figure 2. Symbol frequencies vs the number of duplications in an interspersed-duplication system, with  $s^{(0)} = \text{AGCGTATGCG}$ , and  $q_4 = q_6 = 1/2$ .

## VI. CONCLUSION

We studied the limiting behavior of stochastic interspersed-duplication systems in order to evaluate their ability in creating biological diversity. We showed that the composition of long sequences does not vary greatly in random duplication systems, given that all substrings of the same length are duplicated with equal probability. We also established that frequencies of sequences in interspersed-duplication systems tend to the corresponding probabilities in sequences generated by iid sources, which have the highest possible entropy for given symbol probabilities. It thus seems plausible that diversity may arise from random interspersed-duplication events. Since this work was limited to the asymptotic analysis of these systems, further research is required to quantify their finite-time behavior. Furthermore, in this work we did not consider the more realistic scenarios in which different duplications have different probabilities.

## REFERENCES

- [1] V. S. Borkar, "Stochastic approximation," *Cambridge Books*, 2008.
- [2] F. Farnoud, M. Schwartz, and J. Bruck, "The capacity of string-duplication systems," in *Proc. IEEE Int. Symp. Information Theory*, Honolulu, HI, USA, Jun. 2014, pp. 1301–1305.
- [3] —, "The capacity of string-replication systems," *Submitted to IEEE Trans. Inf. Theory*, *arXiv preprint: http://arxiv.org/abs/1401.4634*, 2014.
- [4] 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.
- [5] C. Mora, D. P. Tittensor, S. Adl, A. G. B. Simpson, and B. Worm, "How many species are there on earth and in the ocean?" *PLoS Biol*, vol. 9, no. 8, Aug. 2011. [Online]. Available: <http://dx.doi.org/10.1371/journal.pbio.1001127>