

Research Article

Optimal Control of Gene Mutation in DNA Replication

Juanyi Yu, Jr-Shin Li, and Tzyh-Jong Tarn

Department of Electrical and Systems Engineering, Washington University in St. Louis, Green Hall, Campus Box 1042, One Brookings Drive, St. Louis, MO 63130, USA

Correspondence should be addressed to Juanyi Yu, juanyi.yu@wustl.edu

Received 30 June 2011; Accepted 3 September 2011

Academic Editor: T. Akutsu

Copyright © 2012 Juanyi Yu et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

We propose a molecular-level control system view of the gene mutations in DNA replication from the finite field concept. By treating DNA sequences as state variables, chemical mutagens and radiation as control inputs, one cell cycle as a step increment, and the measurements of the resulting DNA sequence as outputs, we derive system equations for both deterministic and stochastic discrete-time, finite-state systems of different scales. Defining the cost function as a summation of the costs of applying mutagens and the off-trajectory penalty, we solve the deterministic and stochastic optimal control problems by dynamic programming algorithm. In addition, given that the system is completely controllable, we find that the global optimum of both base-to-base and codon-to-codon deterministic mutations can always be achieved within a finite number of steps.

1. Introduction

Systems biology is an emerging academic field aiming at system-level understanding of biological systems. The early development of systems biology started in the late 1940s [1]. Recent progress in molecular biology has enabled us to gain information on the interactions among the underlying molecules from comprehensive experimental data sets. In general, a system-level understanding of a biological system can be derived from insight into four key properties: (1) the system's structure, (2) the system dynamics, (3) the control method, and (4) the design method [2]. Equivalently, identifying related components and their interactions, gathering qualitative and quantitative information about the system's evolution under different circumstances, achieving the desired outputs by controlling the input with appropriate definitions of inputs and outputs of the system, and reconstructing analogous systems by eliminating the undesired properties are four essential steps in systems biology done by collaboration among engineers, biologists, and doctors. Figure 1 shows a typical method of system construction and verification commonly applied currently. Control engineers construct models, run simulations, and predict the system behaviors. Biologists design and carry out the experiments

and measure the output data. Control engineers revise and verify the models by comparing the predictions and experimental results.

Systems biology is a cross-cutting research area connecting control engineering, biology, and medical science, as shown in Figure 2. It provides a systematic view of the biological system and related medical interventions. It aims at understanding the bare function and integration function of the cell to reconstruct the biological systems with desired features. Control and automation play critical roles in this novel field not only by providing new technology and equipment for biologists to design and perform meticulous experiments, to take high-throughput measurements, and to analyze experimental data efficiently, but also by offering doctors new medical applications and improving the precision of medical manipulations. The wide range of aspects which control and automation have been applied to include, but are not limited to, gene regulation [3, 4], drug delivery [2, 5], and neuron networks [6, 7]. The equipment provided by control engineers includes, but is not limited to, nanodevices, biochips, cuvettes for electroporation, and gene guns. Biologists perform various biological experiments, such as protein synthesis and virus DNA modifications, to gather measurements for model revisions and verifications,

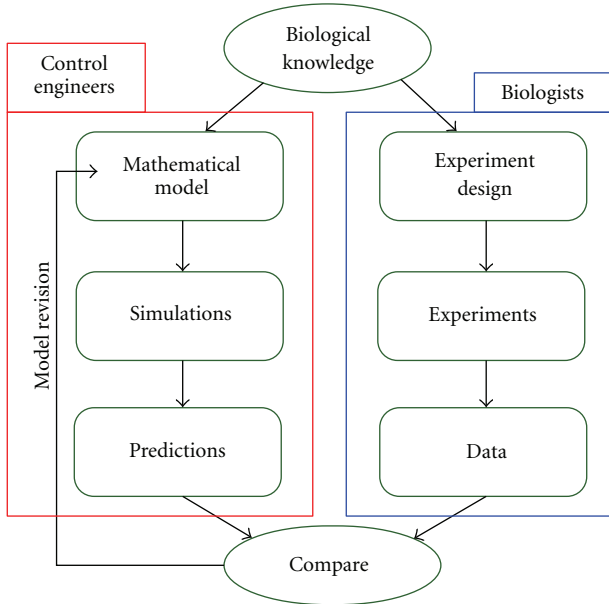


FIGURE 1: Typical analysis of biological systems.

to conclude theoretical and practical results from evidence, and to help medical practice. Doctors use both theoretical and practical results from biologists to perform tissue engineering, such as organ transplants and artificial tissue construction.

According to their scales, biological systems can be divided into three levels: the molecular level (nm), cellular level (μm), and tissue level (cm), analogous to the part, individual, and group, respectively. Molecular-level research focuses on how, when, where, and to what extent [10] a gene is expressed. The essential goal is to sketch a complete blueprint of genes by identifying the control sequences of coding DNA segments and their interactions. Cellular level research, in general, treats one cell as a plant in classical control theory and investigates the reactions of the cell to the changing environment, for instance, concentration changes of related chemicals. State-of-the-art medical therapies are primarily based on experimental results at the cellular level. Tissue-level research mainly concerns tissue reconstruction, artificial tissue substitutes, or tissue function recovery. The cell differentiation process is an important topic at the tissue level. Typical biological systems are collaboratively controlled at all three levels. Most current research work focuses on either cellular level or tissue level systems. Not much work has been done at the molecular-level. In contrast, understanding biological systems at the molecular level is crucial, since species have the same basic inheritance, DNA macromolecules, and follow a common rule in gene expression, the central dogma in molecular biology. Molecular level understanding of biological systems provides instrumental information about radical causes of many diseases and the genetic evidence of evolution. It also helps biologists to gain a better understanding of molecular level interactions, draw a complete blueprint of gene networks, improve existing

means, create novel means to cure genetic diseases, and to elaborate on the theory of evolution. Recent technology in gene sequencing makes it possible to conquer the difficulty in measurement at the molecular level and to identify the nucleotide sequences of a particular DNA segment. Targeted sequencing is the most promising step toward maximizing the efficiency of the next-generation sequencing technology using polymerase chain reaction. The availability of DNA microarray makes it possible to accomplish tens of thousands of genetic tests for picomoles (10^{-12}) of a specific DNA sequence.

Researchers have applied various methods to model, simulate, and control the gene regulation processes. Early attempts to model and simulate gene regulatory systems are summarized in [10], including direct graphs, Bayesian networks, Boolean networks, ordinary and partial differential equations, qualitative differential equations, quantitative differential equations, stochastic equations, and role-based formalisms [10]. Other approaches include Petri nets [11], transformational grammars [12, 13], and process algebra [14]. Three important modeling methods in recent work are gene regulatory units viewed under compound control [4, 15–18], logic network models [19, 20], and base-to-base molecular-level formulation [21, 22]. The first modeling method quantitatively describes chemical concentration variations corresponding to external environmental changes at the cellular level. The second qualitatively illustrates the interactions among operons in gene regulatory units. The last modeling converts DNA segments to discrete vectors. In our paper, we adapt the base-to-base molecular level formulation to express state variables.

Current obstacles in systems biology are obvious. The structure and dynamics of biological systems are sometimes unclear. Most existing models are constructed by data-driven or hypothesis-driven methods, with only partial information available. Due to the complexity of the systems and incomplete information, the mathematical models are usually formulated by modifying empirical equations or proposing heuristic equations. The parameters of proposed models are obtained by estimation methods. Although those models can disclose significant details of the system's structure and dynamics, the inconsistency between theoretical and experimental results creates difficulties for control engineers to verify the models, develop optimal controls, and reconstruct systems with desired properties.

In this paper, we use a novel approach to build up abstract mathematical models at the molecular level in Section 2.2, based directly on biological theory. With reasonable assumptions, we can avoid the conventional obstacles mentioned above. Different from existing methods, focusing on a gene or changes in chemical concentration, we emphasize the base change in the nucleotide bases. The cost function, a summation of costs for applying mutagens and the off-trajectory penalty, together with the system equations, formulates the optimal control problem. The optimal control is then solved by dynamic programming algorithm in Section 2.3. Section 3 shows simulation results of the optimal control problem at different scales and is followed by several important propositions.

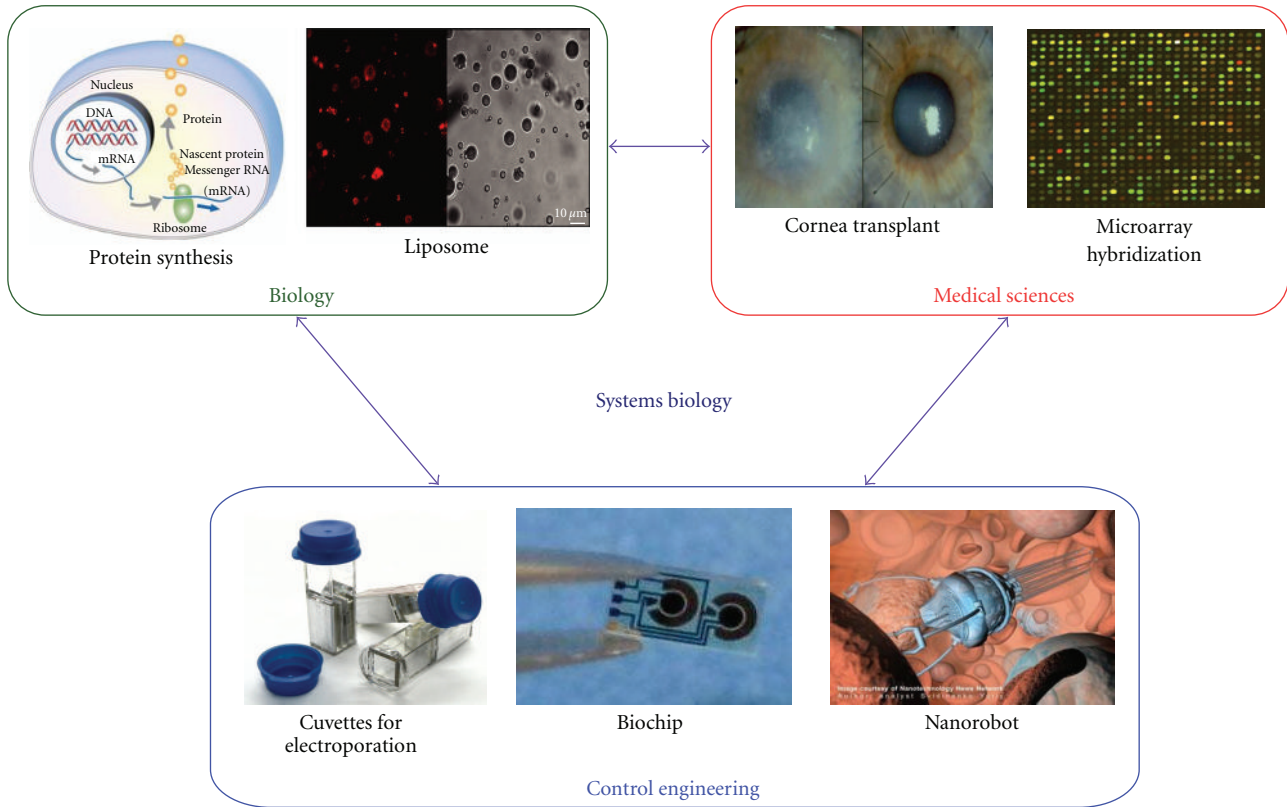


FIGURE 2: Systems biology is a cross-cutting research area connecting control engineering, biology, and medical science. Sources: protein synthesis <http://www.anticancer.de/>, liposome [8], corneal transplant <http://www.avclinic.com/>, microarray hybridization [9], cuvettes for electroporation <http://www.en.wikipedia.org/>, biochip <http://www.clemson.edu/>, nano robot <http://www.molecularlab.it/>.

2. System Equations and Generalized Optimal Control Problem Formulation

The central dogma of molecular biology, first elaborated in [23] and restated in [24], illustrates the detailed residue-by-residue transfer of genetic sequential information. Nowadays, it is widely recognized as the backbone of molecular biology. It describes the genetic information flow among three kinds of biopolymers: DNA, RNA, and protein. In most living organisms, genetic information transfers from DNA to RNA, and then to protein. This process is usually irreversible, thus protein always acts as the sink of information flow. A codon consists of three consecutive nucleotide bases, corresponding to one amino acid according to the genetic codes. Since there are only 20 kinds of amino acids and 64 combinations of codons, there exists redundancy in genetic codes.

We are particularly interested in mutations that happen during the process of DNA replication, as DNA serves as long-term genetic information storage and is the basis of genetic inheritance, the accuracy of which is particularly important to ensure the correct expression of genes. DNA molecules consist of four kinds of nucleotide acids, *adenine* (A), *thymine* (T), *guanine* (G), and *cytosine* (C), and a backbone made of sugars and phosphate. In 1953, James D.

Watson and Francis Crick found the double helix structure of DNA and the rule of basepairing, known as Watson-Crick basepairing [25, 26]. A always pairs with T, G always pairs with C, and vice versa. In nature, replication errors occur at a very low rate, one error for every 10^7 nucleotides added [27]. The redundancy of information caused by the double-helix structure ensures the accuracy of DNA replication. Some DNA self-repair mechanisms, listed in [28], such as proofreading, also help to eliminate errors during the replication process.

Gene mutations are changes in the nucleotide sequence of DNA or RNA. Usually, we focus only on mutations occurring in coding DNA sequences and RNA. Mutations are caused by various reasons. Induced mutations are caused by either chemical mutagens or radiation. In general, radiation induces higher randomness than chemical mutagens. Point mutation is the simplest form of mutation, involving only one base. Point mutations can be further divided into transitions ($A \leftrightarrow G$ or $C \leftrightarrow T$) and transversions ($A/G \leftrightarrow C/T$). Transversions are theoretically expected to be twice as frequent as transitions, but transitions may be favored over transversions in coding DNA because they usually result in a more conserved polypeptide sequence [29].

In this section, we first give the problem statement in Section 2.1, and then we construct system equations for both

deterministic and stochastic mutations in Section 2.2. At last, in Section 2.3, we formulate the optimal control problem and apply dynamic programming algorithm to solve it.

2.1. Problem Statement. Figure 3 shows the system diagram of restoring an abnormal DNA segment back to a normal sequence by applying mutagens during the process of DNA replication. After we obtain a patient's genome, we compare the coding DNA segments with normal DNA segments in our database to figure out the possible range of mutated segments. Due to the redundancy in genetic codes, as long as any two DNA segments can be transcribed and then translated to the same amino acid sequence, the distance reference between them is considered to be zero. Therefore, instead of having a predetermined final state or a neighborhood of a final state, our final state lies in a set where the distance reference between any sequence in this set and the desired sequence is zero. We name this set the final desired set. The prescription is then determined by comparing the current measurement and every sequence in the desired set. Both internal noises and external disturbance can be eliminated by the measurement. We treat DNA sequences as state variables, the ON/OFF controls of all available mutagens at every spot on the given DNA segment as inputs, the measurements as the outputs, and one cell cycle as the step increment in our system equations.

The objective function is defined as a summation of the costs (including risks) of applying mutagens and the off-trajectory penalty. The optimal control sequences are computed beforehand to let doctors make treatment plans according to the patient's condition. In general, the optimal control sequence and the corresponding optimal trajectory are not unique because the bases mutate independently in most cases and the order of mutating different bases does not matter if the number of medical treatment sessions is not under a tight restriction. Additional measurements are taken before and after each treatment, if necessary. In deterministic cases, the purpose of taking additional measurements is to check the current sequence and to eliminate both internal and external disturbances. The treatment plan is adjusted if the measurement is not the same as expected. In stochastic settings, the measurements are taken to conquer the randomness caused by both mutagens and other noises. The treatment is then updated accordingly.

To sum up, our system is a discrete-time dynamic system with finite state space and output space, and a set of ON/OFF switches as controls. Our goal is to optimally drive this system from a given initial state to a desired final set at the lowest cost.

2.2. System Equations Formulation. We mainly focus on applying chemical mutagens and radiation to restore the original amino acid sequence during the process of DNA replication. Other factors that may affect the gene mutation, including temperature and electroporation, are not within our consideration. In addition, we assume that chemical mutagens or radiation target one and only one nucleotide base at any preset site, despite the technical limitation, and

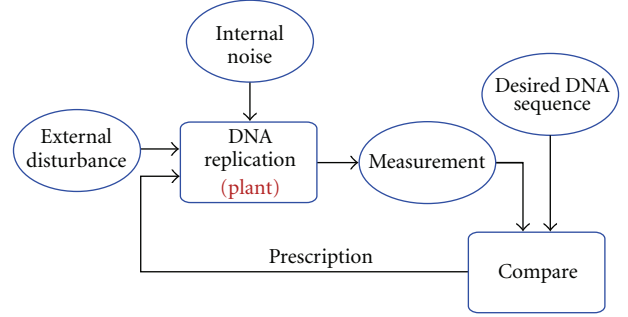


FIGURE 3: System diagram of restoring an abnormal DNA segment back to a normal sequence by applying mutagens during the process of DNA replication.

the results of applying chemical mutagens and radiation are independent. For simplicity, we normalize the dose of mutagens to transfer one nucleotide base to another in one step to 1. In most cases, nucleotide bases mutate independently, therefore there is no chain effect caused by mutagens. To avoid reactions among different mutagens, we require that at most one chemical mutagen and one radiation be applied in each cell cycle. While constructing a generalized model, since the order of applying chemical mutagens and radiation does not affect the results, without loss of generality, we require they be applied in the order shown in Figure 4. That is, chemical mutagens are always applied before the duplication process starts, radiation is always applied in the middle of the cell cycle, and the measurements are taken before every replication starts. Lastly, we assume that the measurements are always correct, and DNA replication error, background mutation rate, and other random noise can be eliminated from measurements by considering them as spontaneous mutation.

2.2.1. Base-to-Base Deterministic Mutations. Denote the targeted DNA segment with n nucleotide bases at k th step by a column vector x_k , as shown in Figure 5. x_k^i is the i th element of x_k . Let P be the transfer matrix from x_k to x_{k+1} , for all k , $k \in \mathbb{Z}^+ \cup \{0\}$, without mutation. Then, the perfect DNA replication process can be expressed as

$$x_{k+1} = Px_k. \quad (1)$$

Proposition 1. $P = -I$.

Proof. As no mutation occurs, x_{k+1} is completely complementary to x_k by Watson-Crick base pairing rule, and x_{k+2} is completely complementary to x_{k+1} . Therefore, x_{k+2} is exactly the same as x_k , thus,

$$x_{k+2} = Px_{k+1} = P^2x_k \implies P^2 = I. \quad (2)$$

Since every base mutates independently, every element of x_{k+1}^i only depends on the corresponding element of x_k^i , thus P is diagonal. In addition, $x_{k+1} \neq x_k$, we conclude $P = -I$. \square

Based on Proposition 1, we assign values to nucleotide bases set $\{A, G, C, T, O\}$, where O is an artificial non-sense base. Define an equivalence relationship between

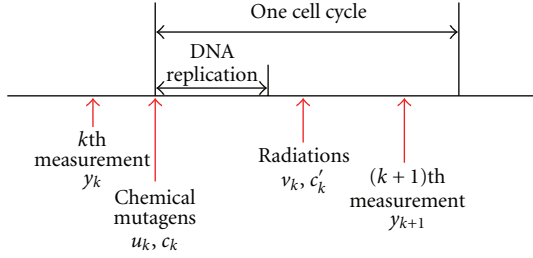


FIGURE 4: The order of taking measurements and applying chemical mutagens and radiation in a cell cycle.

$\{A, G, C, T, O\}$ and $\{1, 2, -2, -1, 0\}$, that is, $\{A, G, C, T, O\} \Leftrightarrow \{1, 2, -2, -1, 0\}$, with

$$x_k^i = \begin{cases} 1, & \text{if } A, \\ 2, & \text{if } G, \\ -2, & \text{if } C, \\ -1, & \text{if } T, \\ 0, & \text{if } O. \end{cases} \quad (3)$$

Proposition 2. $\{1, 2, -2, -1, 0\}$ is a field under proper definitions of addition and multiplication.

Proof. Defining the addition table and multiplication table as in Tables 1 and 2, we check if the set $\{1, 2, -2, -1, 0\}$ satisfies the definition of field. \square

Closed under Addition and Multiplication. Satisfied obviously from Tables 1 and 2.

Associativity of Addition and Multiplication. Implicitly satisfied by integer addition and multiplication.

Commutativity of Addition and Multiplication. Satisfied as Tables 1 and 2 are symmetric according to the diagonal.

Additive and Multiplicative Identity. Additive identity is 0, and multiplicative identity is 1.

Additive and Multiplicative Inverses. Additive inverses pair: $1 \leftrightarrow -1, 2 \leftrightarrow -2, 0 \leftrightarrow 0$.

Multiplicative inverses pair: $1 \leftrightarrow 1, 2 \leftrightarrow -2, -1 \leftrightarrow -1$.

Distributivity of Multiplication over Addition. Implicitly satisfied by integer addition and multiplication.

We conclude $\{0, 1, 2, -2, -1\}$ is a field under addition and multiplication defined by Tables 1 and 2.

From now on, we use \mathcal{F} to denote the field $\{0, 1, 2, -2, -1\}$. And $x_k \in \mathcal{F}^n$ is the state vector representing a DNA segment with n nucleotide bases, where \mathcal{F}^n is the set of \mathcal{F} -valued vectors of dimension n .

TABLE 1: Addition table for $\{1, 2, -2, -1, 0\}$.

+	1	2	-2	-1	0
1	2	-2	-1	0	1
2	-2	-1	0	1	2
-2	-1	0	1	2	-2
-1	0	1	2	-2	-1
0	1	2	-2	-1	0

TABLE 2: Multiplication table for $\{1, 2, -2, -1, 0\}$.

\times	1	2	-2	-1	0
1	1	2	-2	-1	0
2	2	-1	1	-2	0
-2	-2	1	-1	2	0
-1	-1	-2	2	1	0
0	0	0	0	0	0

TABLE 3: Possible values of Δs and Δw . The corresponding values of Δs and Δw are obtained by substituting the value of x_k and x_{k+1} into (4), with $\Delta w \neq 0$ only if $x_k = 0$.

kth	(k+1)th				
	A	G	C	T	O
A	2	-2	-1	0	1
G	-1	2	0	-2	1
C	-2	0	2	-1	1
T	0	-1	-2	2	1
O	1	2	-2	-1	0

$\left. \begin{array}{l} \Delta s \\ \Delta w \end{array} \right\}$

We start with the simplest form of mutations, point mutation. Suppose there is a point mutation, we write it mathematically as

$$x_{k+1} = (-I + \Delta s)x_k + \Delta w, \quad (4)$$

where $x_{k+1}, x_k \in \mathcal{F}$, and $-I$ reduces to -1 as only one base is involved. The corresponding values of Δs and Δw , obtained by reverse engineering with all possible pairs of x_k and x_{k+1} , are listed in Table 3.

Here, Δs represents the mutation from four normal nucleotide bases, and Δw corresponds to mutation from nonsense base, that is, $\Delta w \neq 0$ only if $x_k = 0$.

Rewriting (4) by collecting all values of Δs and Δw in Table 3, we get

$$x_{k+1} = \left(-I + \sum_{j=0}^4 u_k^j s_j \right) x_k + \sum_{j=0}^4 c_k^j w_j \quad (5a)$$

$$= (-I + u_k s) x_k + c_k w, \quad (5b)$$

where $\{s_0, s_1, s_2, s_3, s_4\} = \{w_0, w_1, w_2, w_3, w_4\} = \{0, 1, 2, -2, -1\}$, $u_k^j, c_k^j \in \{0, 1\}$, representing the on/off controls, $u_k = [u_k^0 \ u_k^1 \ u_k^2 \ u_k^3 \ u_k^4]$, $c_k = [c_k^0 \ c_k^1 \ c_k^2 \ c_k^3 \ c_k^4]$, and $s = w = [0 \ 1 \ 2 \ -2 \ -1]^T$. In (5a), s_j and w_j are constants for all k and j . u_k^j and c_k^j , the inputs of the system, are the on/off controls for chemical mutagens or radiation. Clearly,

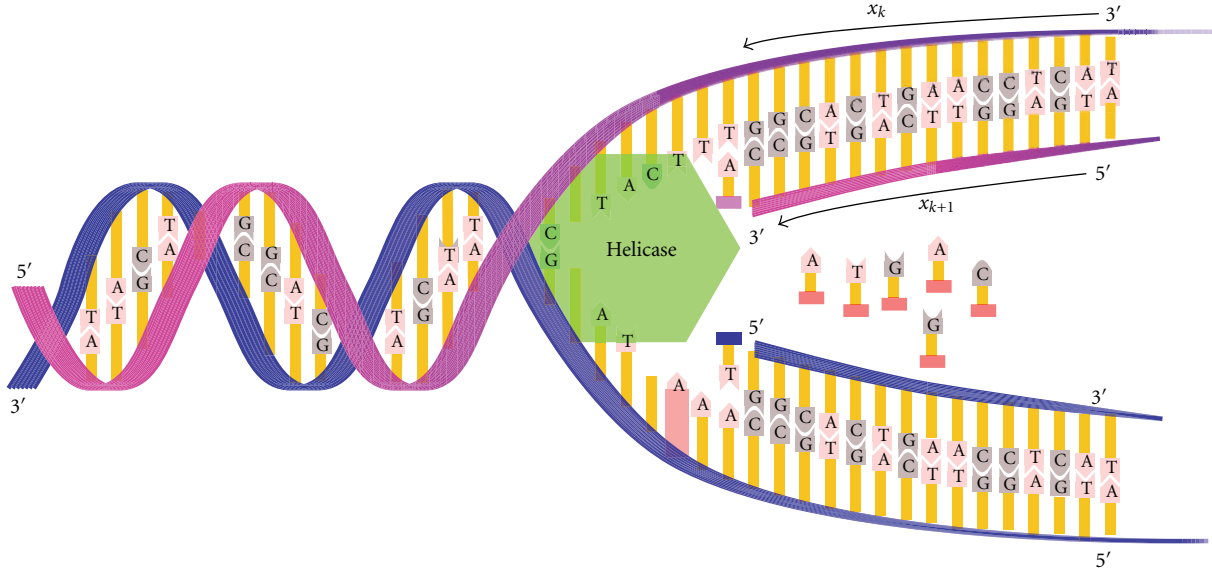


FIGURE 5: Biological information flow in central dogma of molecular biology.

$\sum_{j=0}^4 c_k^j = 1$ only if $x_k = 0$. Equation (5b) is a simplified version of (5a) as we put u_k^j, s_j, c_k^j, w_j into vector form u_k, s, c_k, w . s and w serve as vector basis for base-to-base deterministic model. u_k and c_k are now multi-input controls; each of them contains 5 on/off controls, corresponding to all possible transfer patterns. For a particular k , at most one of u_k^j 's and c_k^j 's can be 1, as stated in Proposition 3. This is consistent with the fact that every state can be transferred to only one of the five states in the state space \mathcal{F} with corresponding mutagens available.

Proposition 3. *It is always 1 – 1 transfer when mutation occurs, that is, one nucleotide base can only transfer to another one, therefore*

- (i) if $x_k = 0$ and $c_k = 0$, or $c_k = [1 \ 0 \ 0 \ 0 \ 0]$, then $x_{k+1} = 0$,
- (ii) if $x_k \neq 0$, then $c_k = 0$ and u_k is either 0 or a unit row vector,
- (iii) if $x_k = 0$, then $u_k = 0$ and c_k is either 0 or a unit row vector,
- (iv) $u_k + c_k$ is either 0 or a unit row vector, for all $k \in \mathbb{Z}^+ \cup \{0\}$.

Now, suppose for some reason we need to take an addition, measurement in the middle of the cell cycle, after the completion of the k th duplication and before the start of the $(k+1)$ th. We name this kind of measurement an intermediate state, and denote by it x'_k . Then, we have

$$x_{k+1} = (I + \Delta s') x'_k + \Delta w', \quad (6)$$

where the values of $\Delta s'$ and $\Delta w'$, listed in Table 4, are obtained in the same way as getting Δs and Δw in Table 3.

Comparing Tables 3 and 4, we find the collection of Δs and $\Delta s'$, Δw and $\Delta w'$, form the same set, respectively. Thus,

TABLE 4: Possible values of $\Delta s'$ and $\Delta w'$. The corresponding values of $\Delta s'$ and $\Delta w'$ are obtained by substituting the value of x'_k and x_{k+1} into (4), with $\Delta w' \neq 0$ only if $x'_k = 0$.

k th	$(k+1)$ th				
	A	G	C	T	O
A	0	1	2	-2	-1
G	2	0	-2	1	-1
C	1	-2	0	2	-1
T	-2	2	1	0	-1
O	1	2	-2	-1	0

we continue using s and w when rewriting (6) in the form of (5a) and (5b), that is,

$$x_{k+1} = (I + v_k s) x'_k + c'_k w, \quad (7)$$

where v_k, c'_k are the counterparts of u_k, c_k , respectively, and s, w are the same as in (5b).

Similar to Proposition 3, we get Proposition 4.

Proposition 4. v_k and c'_k in (7) need to satisfy the following conditions.

- (i) If $x_k = 0$ and $c'_k = 0$, or $c'_k = [1 \ 0 \ 0 \ 0 \ 0]$, then $x_{k+1} = 0$.
- (ii) If $x_k \neq 0$, then $c'_k = 0$ and v_k is either 0 or a unit row vector.
- (iii) If $x_k = 0$, then $v_k = 0$ and c_k is either 0 or a unit row vector.
- (iv) $v_k + c'_k$ is either 0 or a unit row vector, for all $k \in \mathbb{Z}^+ \cup \{0\}$.

Now take, both chemical mutagens and radiative rays under our consideration and apply them in the order as shown in Figure 4. Then, we can express our system equation as

$$x'_k = \left(-I + \underbrace{u_k s}_{\text{mutations caused by chemical mutagens from normal bases}} \right) x_k + \underbrace{c_k w}_{\text{mutations caused by chemical mutagens from } O}, \quad (8a)$$

$$x_{k+1} = \left(I + \underbrace{v_k s}_{\text{mutations caused by radiative rays from normal bases}} \right) x'_k + \underbrace{c'_k w}_{\text{mutations caused by radiative rays from } O}, \quad (8b)$$

$$y_k = x_k, \quad (8c)$$

where u_k and v_k are the inputs of the system and y_k is the measurement. Obviously, (8a) is modified from (5b) and (8b) from (7).

The two-step mutation and the intermediate state x'_k avoid the case x_k is changed to different bases by radiation and chemical mutagens simultaneously, which causes confusion. Substituting (8a) and (8b), we get

$$x_{k+1} = (I + v_k s)(-I + u_k s)x_k + (I + v_k s)c_k w + c'_k w, \quad (9a)$$

$$y_k = x_k. \quad (9b)$$

Obviously, Proposition 3 still holds for u_k and c_k , and Proposition 4 holds for v_k and c'_k for (9a).

For point mutations, we have 20 on/off controls in total for every step k , 10 for chemical mutagens as described before, and the rest for radiation.

2.2.2. Gene-to-Gene Deterministic Mutations. In general, mutations involve multiple bases. Therefore, large-scale deterministic model is necessary. Now, we show how to extend our model to large-scale systems.

Suppose we have a coding DNA segment with length n , then $x_k \in \mathcal{F}^n$. Since a coding DNA segment usually contains integer number of codons, which is made of three consecutive bases, n is a multiple of 3. Let x_k^i denote the i th component of x_k . This notation is consistent with the one in Section 2.2.1. Initiated by the base-to-base deterministic

model from Section 2.2.1, we write our system equation for large-scale system as

$$x'_k = \left(-I + \underbrace{\sum_{i=1}^n u_k^i S_k^i}_{\text{mutations caused by chemical mutagens from normal bases}} \right) x_k + \underbrace{\sum_{i \in \mathcal{O}_k} c_k^i W_k^i}_{\text{mutations caused by chemical mutagens from } O}, \quad (10)$$

$$x_{k+1} = \left(I + \underbrace{\sum_{i=1}^n v_k^i Q_k^i}_{\text{mutations caused by radiative rays from normal bases}} \right) x'_k + \underbrace{\sum_{i \in \mathcal{O}'_k} b_k^i R_k^i}_{\text{mutations caused by radiative rays from } O},$$

$$y_k = x_k,$$

where $u_k^i, v_k^i, c_k^i, b_k^i$ are on/off controls of the i th element, S_k^i, Q_k^i are $n \times n$ square matrices corresponding to the mutations between normal bases or from normal bases by chemicals and radiation, respectively, W_k^i, R_k^i are n -dimensional column vectors representing mutations from nonsense bases by chemicals and radiation, respectively, and $\mathcal{O}_k = \{i : x_k^i = 0, 1 \leq i \leq n\}$, $\mathcal{O}'_k = \{i : x_k^i = 0, 1 \leq i \leq n\}$.

S_k^i and Q_k^i are diagonal matrices since each base mutates independently. The values in the first four rows of Tables 3 and 4 correspond to the diagonal elements of S_k^i and Q_k^i , respectively. The last rows of Tables 3 and 4 are assigned to W_k^i and R_k^i , n -dimensional vectors, at nonsense base's spots for x_k .

Define $\mathcal{S} = \{s_j e_i e_i^T, \forall i, j, 0 \leq j \leq 4, 1 \leq i \leq n\}$, a collection of $n \times n$ matrices, where s_j is the same as in (5a) and (5b), e_i is the unit column vector of length n with i th component equal to 1 and all other components equal to 0, and $e_i e_i^T$ is the square matrix with only the i th element on the diagonal equals to 1, and 0 otherwise. Then, S_k^i, Q_k^i can be written as linear combinations of all elements from \mathcal{S} , with the coefficient of each element either 0 or 1 corresponding to the on/off control $u_k^{(i,j)}$ and $v_k^{(i,j)}$, respectively.

Similarly, define $\mathcal{W} = \{w_j e_i, \forall i, j, 0 \leq j \leq 4, 1 \leq i \leq n\}$, where w_j is the same as (5a) and (5b). W_k^i, R_k^i can be written as linear combinations of all components from \mathcal{W} , with coefficient of every component either 0 or 1 corresponding to the on/off control $c_k^{(i,j)}$ and $c_k'^{(i,j)}$, respectively.

Therefore, instead of using step-varying $S_k^i, S_k^j, W_k^i, R_k^i$, we find matrix basis for those four square matrices to make the

controls to be the only variables depending on k , as we did for single-base cases. Then, we can, write (10) as

$$x'_k = \left(-I + \sum_{i=1}^n \sum_{j=0}^4 u_k^{(i,j)} s_j e_i e_i^T \right) x_k + \sum_{i \in \mathcal{O}_k} \sum_{j=0}^4 c_k^{(i,j)} w_j e_i, \quad (11a)$$

$$x_{k+1} = \left(I + \sum_{i=1}^n \sum_{j=0}^4 v_k^{(i,j)} s_j e_i e_i^T \right) x'_k + \sum_{i \in \mathcal{O}'_k} \sum_{j=0}^4 c_k^{(i,j)} w_j e_i, \quad (11b)$$

$$y_k = x_k, \quad (11c)$$

where $u_k^{(i,j)}, v_k^{(i,j)}, c_k^{(i,j)}, c_k^{(i,j)'} \in \{0, 1\}$.

As shown in (11a), (11b), and (11c), multisites mutations contain $20n$ controls in total for every step k , where n is the number of nucleotide bases on the targeted gene. Similar to point mutations, every single site has 20 controls in each step, 10 for chemical mutagens and 10 for radiation.

We can view u_k, v_k, c_k, c_k' as binary matrices of dimension $n \times 5$, and $u_k^{(i,j)}, v_k^{(i,j)}, c_k^{(i,j)}, c_k^{(i,j)'}$ are the corresponding element of i th row and j th column. Use $u_k^i, v_k^i, c_k^i, b_k^i$, binary row vectors of dimension 5, to denote the i th row of u_k, v_k, c_k, c_k' , respectively. Again, $s = w = [0 \ 1 \ 2 \ -2 \ -1]^T$.

Combining (11a) and (11b), and writing control variables in vector forms, we get

$$x_{k+1} = \left(I + \sum_{i=1}^n v_k^i s e_i e_i^T \right) \left(-I + \sum_{i=1}^n u_k^i s e_i e_i^T \right) x_k + \left(I + \sum_{i=1}^n v_k^i s e_i e_i^T \right) \sum_{i \in \mathcal{O}_k} c_k^i w e_i + \sum_{i \in \mathcal{O}'_k} b_k^i w e_i, \quad (12)$$

$$y_k = x_k.$$

Proposition 5. For large-scale deterministic system, u_k, v_k, c_k, c_k' satisfy conditions below.

- (i) If $e_i^T x_k = 0$, then $i \in \mathcal{O}_k$.
- (ii) If $e_i^T x_k = 0, c_k^i = 0$ or $c_k^i = [1 \ 0 \ 0 \ 0 \ 0]$, then $i \in \mathcal{O}'_k$.
- (iii) For all $i \notin \mathcal{O}_k, u_k^i$ is either 0 or a row unit vector and $c_k^i = 0$.
- (iv) For all $i \in \mathcal{O}_k, c_k^i$ is either 0 or a row unit vector and $u_k^i = 0$.
- (v) For all $i \notin \mathcal{O}'_k, v_k^i$ is either 0 or a row unit vector and $b_k^i = 0$.
- (vi) For all $i \in \mathcal{O}'_k, b_k^i$ is either 0 or a row unit vector and $v_k^i = 0$.
- (vii) For all $i, k, 1 \leq i \leq n, k \in \mathbb{Z}^+ \cup \{0\}, u_k^i + c_k^i$ is either 0 or a unit row vector and $v_k^i + b_k^i$ is either 0 or a unit row vector.

The mathematical model (12) is quite flexible and can be easily extended to many cases, such as transcription process, multiple spot mutations within one-step or broken DNA strands.

Take broken DNA strands as an example. DNA strand breaks due to various reasons. Our system equation can represent this phenomenon by dividing the whole system into small subsystems. Significant brokage of DNA strands is simply eliminated by cell mechanism to ensure the accuracy to DNA replication. Equation (13) shows the case of one single DNA strand breaking into two segments by chemical mutagens

$$\begin{aligned} \begin{pmatrix} x'_k(1) \\ x'_k(2) \end{pmatrix} &= \begin{pmatrix} -I_m + \sum_{i=1}^m u_k^i s e_i e_i^T & 0 \\ 0 & -I_{n-m} + \sum_{i=m+1}^n u_k^i s e_i e_i^T \end{pmatrix} \\ &\times \begin{pmatrix} x_k(1) \\ x_k(2) \end{pmatrix} + \begin{pmatrix} \sum_{i \in \mathcal{O}_k, 1 \leq i \leq m} c_k^i w e_i \\ \sum_{i \in \mathcal{O}_k, (m+1) \leq i \leq n} c_k^i w e_i \end{pmatrix}, \\ x_{k+1}(1) &= \left(I_m + \sum_{i=1}^m v_k^i s e_i e_i^T \right) x'_k(1) + \sum_{i \in \mathcal{O}'_k, 1 \leq i \leq m} b_k^i w e_i, \\ x_{k+1}(2) &= \left(I_{n-m} + \sum_{i=m+1}^n v_k^i s e_i e_i^T \right) x'_k(2) \\ &+ \sum_{i \in \mathcal{O}'_k, (m+1) \leq i \leq n} b_k^i w e_i. \end{aligned} \quad (13)$$

2.2.3. Gene-to-Gene Stochastic Mutations. In reality, mutagens, no matter chemical or radiative, always cause randomness in mutation. Therefore, we need to derive the model for gene-to-gene stochastic mutations.

Introduce new random variables, $h_{k,l_1}^{(i,j)}, r_{k,l_2}^{(i,j)}, h_{k,l_3}^{(i,j)'} , r_{k,l_4}^{(i,j)'} \in \{0, 1\}$, associated with probability $p_{l_1,j}^{(h)}, p_{l_2,j}^{(r)}, p_{l_3,j}^{(h')}, p_{l_4,j}^{(r')}$, for all $i, k, 1 \leq i \leq n, k \in \mathbb{Z}^+ \cup \{0\}$, respectively, where k is the step index, l_1, l_2 are indices for chemical mutagens inducing mutation from normal bases and from O , respectively, l_3, l_4 are indices for radiation inducing mutation from normal bases and from O , respectively, i is the index of DNA segment, and the value of j corresponds to the transfer pattern, which can be found in Tables 3 and 4. Note different mutagens have different probability assignments, the probability assignments are only related to the type of mutagens, and the probability associated with every kind of mutagens sums up to 1, that is,

$$\begin{aligned} \sum_{j=0}^4 p_{l_1,j}^{(h)} &= 1, \quad \forall l_1, 1 \leq l_1 \leq l, \\ \sum_{j=0}^4 p_{l_2,j}^{(r)} &= 1, \quad \forall l_2, 1 \leq l_2 \leq m, \\ \sum_{j=0}^4 p_{l_3,j}^{(h')} &= 1, \quad \forall l_3, 1 \leq l_3 \leq l', \\ \sum_{j=0}^4 p_{l_4,j}^{(r')} &= 1, \quad \forall l_4, 1 \leq l_4 \leq m'. \end{aligned} \quad (14)$$

The controls are $u_{k,l_1}^i, c_{k,l_2}^i, v_{k,l_3}^i, b_{k,l_4}^i \in \{0, 1\}$, with the fact that 1 representing mutagen with corresponding index is applied at i th spot of DNA segment at k th generation, and 0 representing mutagen with corresponding index is not applied at spot i at k th step, similar to Sections 2.2.1 and 2.2.2. The mutagen indices l_1, l_2, l_3, l_4 can be omitted in deterministic mutations since given the current state and control, the next state is unique. However, they are necessary for stochastic mutations, because there exist multiple possible states for the next stage given the control. In other words, the next state is determined by random variables $h_{k,l_1}^{(i,j)}, r_{k,l_2}^{(i,j)}, h'_{k,l_3}{}^{(i,j)}, r'_{k,l_4}{}^{(i,j)}$, given the values of $u_{k,l_1}^i, c_{k,l_2}^i, v_{k,l_3}^i, b_{k,l_4}^i$, and x_k .

Suppose we have $(l + m)$ kinds of chemical mutagens available, with l kinds to induce mutations from normal bases and m kinds to induce mutations from O . And we have $(l' + m')$ kinds of radiation available, with l' kinds to induce mutations from normal bases and m' kinds to induce mutations from O . Therefore, we have total $(l + m + l' + m')$ controls for each spot i at each step k . We can write our system equation as

$$x'_k = \left(-I + \underbrace{\sum_{l_1=1}^l \sum_{i=1}^n u_{k,l_1}^i \sum_{j=0}^4 h_{k,l_1}^{(i,j)} s_j e_i e_i^T}_{\text{mutations caused by chemical mutagens from normal bases}} \right) x_k \quad (15a)$$

$$+ \underbrace{\sum_{l_2=1}^m \sum_{i \in \mathcal{O}_k} c_{k,l_2}^i \sum_{j=0}^4 r_{k,l_2}^{(i,j)} w_j e_i}_{\text{mutations caused by chemical mutagens from O}}$$

$$x_{k+1} = \left(I + \underbrace{\sum_{l_3=1}^{l'} \sum_{i=1}^n v_{k,l_3}^i \sum_{j=0}^4 h'_{k,l_3}{}^{(i,j)} s_j e_i e_i^T}_{\text{mutations caused by radiative rays from normal bases}} \right) x'_k \quad (15b)$$

$$+ \underbrace{\sum_{l_4=1}^{m'} \sum_{i \in \mathcal{O}'_k} b_{k,l_4}^i \sum_{j=0}^4 r'_{k,l_4}{}^{(i,j)} w_j e_i}_{\text{mutations caused by radiative rays from O}}$$

$$y_k = x_k. \quad (15c)$$

Again, we define $h_{k,l_1}^{(i,j)}, r_{k,l_2}^{(i,j)}, h'_{k,l_3}{}^{(i,j)}, r'_{k,l_4}{}^{(i,j)}$ the elements at i th row and j th column of $n \times 5$ binary matrices $h_{k,l_1}, r_{k,l_2}, h'_{k,l_3}, r'_{k,l_4}$, respectively. $h_{k,l_1}^i, r_{k,l_2}^i, h'_{k,l_3}{}^i, r'_{k,l_4}{}^i$ and binary row vectors of dimension 5 denote the i th row of $h_{k,l_1}, r_{k,l_2}, h'_{k,l_3}, r'_{k,l_4}$, respectively. Then, we can simplify (15a), (15b), and (15c) and combine (15a) and (15b) as

$$x_{k+1} = \left(I + \sum_{l_3=1}^{l'} \sum_{i=1}^n v_{k,l_3}^i h'_{k,l_3}{}^i s e_i e_i^T \right)$$

$$\begin{aligned} & \times \left(-I + \sum_{l_1=1}^l \sum_{i=1}^n u_{k,l_1}^i h_{k,l_1}^i s e_i e_i^T \right) x_k \\ & + \left(I + \sum_{l_2=1}^m \sum_{i=1}^n c_{k,l_2}^i r_{k,l_2}^i s e_i e_i^T \right) \sum_{l_2=1}^m \sum_{i \in \mathcal{O}_k} c_{k,l_2}^i r_{k,l_2}^i w e_i \\ & + \sum_{l_4=1}^{m'} \sum_{i \in \mathcal{O}'_k} b_{k,l_4}^i r'_{k,l_4}{}^i w e_i, \\ & y_k = x_k. \end{aligned} \quad (16)$$

Proposition 6. For large-scale stochastic system, $u_{k,l_1}^i, h_{k,l_1}^i, c_{k,l_2}^i, r_{k,l_2}^i, v_{k,l_3}^i, h'_{k,l_3}{}^i, b_{k,l_4}^i, r'_{k,l_4}{}^i$ follow the rules below.

- (i) If $e_i^T x_k = 0$, then $i \in \mathcal{O}_k$.
- (ii) If $e_i^T x_k = 0$ and $\sum_{l_2=1}^m c_{k,l_2}^i = 0$, then $i \in \mathcal{O}'_k$.
- (iii) If $e_i^T x_k = 0$, $\sum_{l_2=1}^m c_{k,l_2}^i = 1$ and $r_{k,l_2}^i = [1 \ 0 \ 0 \ 0 \ 0]$, then $i \in \mathcal{O}'_k$.
- (iv) For all $i, k, l_1, 1 \leq i \leq n$, $k \in \mathbb{Z}^+ \cup \{0\}, 1 \leq l_1 \leq l$, if $u_{k,l_1}^i = 1$, then h_{k,l_1}^i is a unit row vector.
- (v) For all $i, k, l_2, 1 \leq i \leq n$, $k \in \mathbb{Z}^+ \cup \{0\}, 1 \leq l_2 \leq m$, if $c_{k,l_2}^i = 1$, then r_{k,l_2}^i is a unit row vector.
- (vi) For all $i, k, l_3, 1 \leq i \leq n$, $k \in \mathbb{Z}^+ \cup \{0\}, 1 \leq l_3 \leq l'$, if $v_{k,l_3}^i = 1$, then $h'_{k,l_3}{}^i$ is a unit row vector.
- (vii) For all $i, k, l_4, 1 \leq i \leq n$, $k \in \mathbb{Z}^+ \cup \{0\}, 1 \leq l_4 \leq m'$, if $b_{k,l_4}^i = 1$, then $r'_{k,l_4}{}^i$ is a unit row vector.
- (viii) For all $i \notin \mathcal{O}_k$, $\sum_{l_1=1}^l u_{k,l_1}^i = 0$ or 1 and $c_{k,l_2}^i = 0$, for all $l_2, 1 \leq l_2 \leq m$.
- (ix) For all $i \in \mathcal{O}_k$, $\sum_{l_2=1}^m c_{k,l_2}^i = 0$ or 1 and $u_{k,l_1}^i = 0$, for all $l_1, 1 \leq l_1 \leq l$.
- (x) For all $i \notin \mathcal{O}'_k$, $\sum_{l_3=1}^{l'} v_{k,l_3}^i = 0$ or 1 and $b_{k,l_4}^i = 0$, for all $l_4, 1 \leq l_4 \leq m'$.
- (xi) For all $i \in \mathcal{O}'_k$, $\sum_{l_4=1}^{m'} b_{k,l_4}^i = 0$ or 1 and $v_{k,l_3}^i = 0$, for all $l_3, 1 \leq l_3 \leq l'$.
- (xii) For all $i, k, 1 \leq i \leq n$, $k \in \mathbb{Z}^+ \cup \{0\}, \sum_{l_1=1}^l u_{k,l_1}^i + \sum_{l_2=1}^m c_{k,l_2}^i = 0$ or 1 and $\sum_{l_3=1}^{l'} v_{k,l_3}^i + \sum_{l_4=1}^{m'} b_{k,l_4}^i = 0$ or 1.

We close this section with the definition of controllability to the system equations proposed above. DNA replication systems with system equations proposed as (9a), (9b), (12), and (16) are *completely controllable* if and only if for all $x_0, x_{2k_1}, x_{2k_2+1} \in \mathcal{F}$, $k_1, k_2 \in \mathbb{Z}^+ \cup \{0\}, \exists$ at least one path from x_0 to x_{2k_1} and at least one path from x_0 to x_{2k_2+1} by applying proper mutagens in the correct order, with k_1, k_2 finite.

2.3. Generalized Optimal Control Problem Formulation. We first define our objective function that can be adapted to all kinds of systems proposed in Section 2.2 with minor changes. Mathematically, in systems where the controllable parameters of interest are discrete, the objective function is

usually a weighted sum representing the number of times that a piece of equipment is turned “on” or “off” or the number of resources needed to execute certain tasks in the frequent cases [30]. In our case, this summation is the total number of times that different mutagens are applied weighted by the corresponding cost (including the risk). Another key factor of objective function is the off-trajectory penalty. Designing a trajectory beforehand is necessary to avoid other hidden risks. If the measurement indicates that the current state is off the predefined trajectory, we include a distance reference between current state and desired state as penalty and change the treatment plan accordingly.

Therefore, our objective function can be expressed as

$$\begin{aligned}
J_0(x_0) &= \min_{u,c,v,c',h,h',r,r'} \mathbb{E} \left[\underbrace{\sum_{k=0}^{N-1} \sum_{l_1=1}^l \sum_{i=1}^n \alpha_{l_1} u_{k,l_1}^i + \sum_{k=0}^{N-1} \sum_{l_2=1}^m \sum_{i=1}^n \beta_{l_2} c_{k,l_2}^i}_{\text{cost of applying chemical mutagens}} \right. \\
&\quad + \underbrace{\sum_{k=0}^{N-1} \sum_{l_3=1}^{l'} \sum_{i=1}^n \alpha'_{l_3} v_{k,l_3}^i + \sum_{k=0}^{N-1} \sum_{l_4=1}^{m'} \sum_{i=1}^n \beta'_{l_4} b_{k,l_4}^i}_{\text{cost of applying radiative rays}} \\
&\quad \left. + \underbrace{\sum_{k=0}^N d(x_k, \{x_k^d\})}_{\text{tracing cost}} \right], \tag{17}
\end{aligned}$$

with $x_0, x_k^d \in \mathcal{F}^n, 1 \leq k \leq N, n \equiv 0 \pmod{3}$ given. The physical meaning of $u_{k,l_1}^i, c_{k,l_2}^i, v_{k,l_3}^i, b_{k,l_4}^i, l_1, l_2, l_3, l_4$ is the same as in Section 2.2.3. $\alpha_{l_1}, \beta_{l_2}, \alpha'_{l_3}, \beta'_{l_4} \in \mathbb{R}$, for all $l_1, l_2, l_3, l_4, 1 \leq l_1 \leq l, 1 \leq l_2 \leq m, 1 \leq l_3 \leq l', 1 \leq l_4 \leq m'$, are the corresponding cost of applying chemical mutagens and radiative rays indexed l_1, l_2, l_3, l_4 , respectively. $\{x_k^d\} : \mathcal{F}^n \times \mathcal{F}^n \rightarrow \mathbb{R}^+ \cup \{0\}$ denotes the desired set at k th stage, generated by the DNA sequences representing the same amino acid sequence as x_k^d , the desired state at k th step. And $d(x_k, \{x_k^d\})$ is the distance reference of the current state x_k to the desired set $\{x_k^d\}$ at k th step. The final penalty, the distance reference from the final state to the desired set at $k = N$, is included in the last term.

In general, $\beta_{l_2}, \beta'_{l_4} \ll \alpha_{l_1}, \alpha'_{l_3}$, for all $l_1, l_2, l_3, l_4, 1 \leq l_1 \leq l, 1 \leq l_2 \leq m, 1 \leq l_3 \leq l', 1 \leq l_4 \leq m'$, because physically O is a set of nonsense bases and more details are necessary to convert an O back to normal bases, for instance, the cost to identify the exact element in the set O . Our goal is to drive our system optimally from initial state x_0 to the desired final set $\{x_N^d\}$ by applying a sequence of mutagens indexed with $\{l_1, l_2, l_3, l_4\}$, at problematic positions i , and in a correct order k .

In (17), the first four terms inside the expectation do not depend on random variables $h_{k,l_1}^i, r_{k,l_2}^i, h_{k,l_3}^i$, and r_{k,l_4}^i ,

for all i, k, l_1, l_2, l_3, l_4 as the treatment plan is computed based on the initial state x_0 . Given y_k , the updated treatment plan is computed accordingly but still not related to random variables. The last term inside expectation, $\sum_{k=0}^N d(x_k, \{x_k^d\})$, is the only term in summation that depends on the distribution of the random variables.

The constraint of the optimal control problem, in general, is the system equation. We choose multidimensional stochastic system equation as the generalized constraints as it can be degenerated to one-dimensional and multidimensional deterministic cases by proper modifications.

Therefore, we can rewrite our objective function and formulate our optimal control problem as

$$\begin{aligned}
J_0(x_0) &= \min_{\{u,c,v,c'\}_{0,1,\dots,N-1}} \left[\sum_{k=0}^{N-1} \sum_{l_1=1}^{l'} \sum_{i=1}^n \alpha_{l_1} u_{k,l_1}^i + \sum_{k=0}^{N-1} \sum_{l_2=1}^m \sum_{i=1}^n \beta_{l_2} c_{k,l_2}^i \right. \\
&\quad + \sum_{k=0}^{N-1} \sum_{l_3=1}^{l'} \sum_{i=1}^n \alpha'_{l_3} v_{k,l_3}^i + \sum_{k=0}^{N-1} \sum_{l_4=1}^{m'} \sum_{i=1}^n \beta'_{l_4} b_{k,l_4}^i \\
&\quad \left. + \sum_{k=0}^N \mathbb{E}_{\{h,r,h',r'\}_{0,1,\dots,N-1}} \left[d(x_k, \{x_k^d\}) \right] \right], \tag{18}
\end{aligned}$$

subject to

$$\begin{aligned}
x_{k+1} &= \left(I + \sum_{l_3=1}^{l'} \sum_{i=1}^n v_{k,l_3}^i h_{k,l_3}^i s e_i e_i^T \right) \\
&\quad \times \left(-I + \sum_{l_1=1}^l \sum_{i=1}^n u_{k,l_1}^i h_{k,l_1}^i s e_i e_i^T \right) x_k \\
&\quad + \left(I + \sum_{l_3=1}^{l'} \sum_{i=1}^n v_{k,l_3}^i h_{k,l_3}^i s e_i e_i^T \right) \sum_{l_2=1}^m \sum_{i \in \mathcal{O}_k} c_{k,l_2}^i r_{k,l_2}^i w e_i \\
&\quad + \sum_{l_4=1}^{m'} \sum_{i \in \mathcal{O}'_k} b_{k,l_4}^i r_{k,l_4}^i w e_i, \\
y_k &= x_k. \tag{19}
\end{aligned}$$

We need to choose a proper distance reference to quantitatively describe the relationship between DNA segments of same length. We first define the distance reference between codons, and the distance reference between DNA segments is a weighted sum of distance reference between every pair of codons.

The distance reference between codons, $d(\varphi_1, \varphi_2), \varphi_1, \varphi_2 \in \mathcal{F}^3$, needs to fulfill the biological requirements as below.

- (1) Nonnegativity: the distance reference between any two codons is either positive or zero. Mathematically, $d : \mathcal{F}^3 \times \mathcal{F}^3 \rightarrow \mathbb{R}^+ \cup \{0\}, d(\varphi_1, \varphi_2) \geq 0$.

- (2) The distance reference between two codons corresponding to the same amino acid is zero.
- (3) Symmetry: the distance reference from codon φ_1 to codon φ_2 equals the distance reference from codon φ_2 to codon φ_1 , that is, $d(\varphi_1, \varphi_2) = d(\varphi_2, \varphi_1)$.
- (4) The distance reference between two codons corresponding to different amino acids should reveal the chemical and physical differences between two amino acids.
- (5) The distance reference from stop codons to all other codons is much larger than those between other codons as early termination of amino acid sequences is more harmful than other forms of mutations.

All the existing metric defined on the finite field cannot achieve all the requirements above. The second requirement violates the identity of indiscernible, that is, $d(\varphi_1, \varphi_2) = 0$ if and only if $\varphi_1 = \varphi_2$. The redundancy in genetic codes implies $d(\varphi_1, \varphi_2) = 0$ if those two amino acids, φ_1 and φ_2 , are translated into the same amino acids. In addition, the triangular inequality is not necessarily true, according to the underlying physical meanings. We take the assumption that the stop codons are of the same distance reference from and to all other codons.

Important physical and chemical properties are listed in Table 5. We ignore codons containing O since their chemical and physical properties cannot be found in literature.

From Table 5, we can see all codons are divided into different sets with each set corresponding to one amino acid. The size and the elements in one codon set vary from one amino acid to another. This implies that the costs of driving one codon to the desired final set generated by the desired final state might be different from the costs of driving the complementary codon to the desired final set generated by the complementary of desired final state. More discussions about this issue are presented in Section 3.

The distance reference between any two codons can be defined by a weighted sum of the differences between physical and chemical properties or other reasonable functions. And the distance reference between two DNA sequences is defined as the sum of distance reference between the corresponding pair of codons. The biological statics plays a crucial rule to define this distance function in practical.

An example of the distance function can be expressed as

$$\begin{aligned}
 d(\xi_1, \xi_2) &= \zeta_{\text{polarity}} \text{polarity}(\xi_1, \xi_2) \\
 &\quad + \zeta_{\text{PH}} \text{PH}(\xi_1, \xi_2) + \zeta_{\text{size}} \text{size}(\xi_1, \xi_2), \\
 \text{polarity}(\xi_1, \xi_2) &= \begin{cases} 0 & \text{if } \xi_1, \xi_2 \text{ are both polar or non-polar,} \\ 1 & \text{if one of } \xi_1, \xi_2 \text{ is polar, and the other non-polar,} \end{cases}
 \end{aligned}$$

$$\begin{aligned}
 \text{PH}(\xi_1, \xi_2) &= |\text{PH value of } \xi_1 - \text{PH value of } \xi_2|, \\
 \text{size}(\xi_1, \xi_2) &= \begin{cases} 0, & \text{if } \xi_1, \xi_2 \text{ are both tiny, small, or normal,} \\ \sigma_1, & \text{if one of } \xi_1, \xi_2 \text{ is tiny, and the other small,} \\ \sigma_2, & \text{if one of } \xi_1, \xi_2 \text{ is tiny, and the other normal,} \\ \sigma_3, & \text{if one of } \xi_1, \xi_2 \text{ is small, and the other normal,} \end{cases} \quad (20)
 \end{aligned}$$

where ξ_1, ξ_2 are two amino acids.

$d(\xi_1, \xi_2)$ is then assigned to $d(\varphi_1, \varphi_2)$ with φ_1, φ_2 corresponding to amino acids ξ_1, ξ_2 , respectively.

Since the generalized optimal control problem in (18) and (19) is a multistage problem that can be broken down into simpler steps at different time points. Therefore, we can solve it by dynamic programming.

For dynamic programming, the optimal control policy is constructed backward. And Bellman's principle of optimality states that the optimal policy for x_0 to $\{x_N^d\}$ is also the optimal policy for the tail problem, from x_q to $\{x_N^d\}$.

The tail problem is defined as

$$\begin{aligned}
 J_q(x_q) &= \min_{\{u, c, v, c'\}_{q, q+1, \dots, N-1}} \left\{ \sum_{k=q}^{N-1} \sum_{l_1=1}^l \sum_{i=1}^n \alpha_{l_1} u_{k, l_1}^i + \sum_{k=q}^{N-1} \sum_{l_2=1}^m \sum_{i=1}^n \beta_{l_2} c_{k, l_2}^i \right. \\
 &\quad + \sum_{k=q}^{N-1} \sum_{l_3=1}^{l'} \sum_{i=1}^n \alpha'_{l_3} v_{k, l_3}^i + \sum_{k=q}^{N-1} \sum_{l_4=1}^{m'} \sum_{i=1}^n \beta'_{l_4} b_{k, l_4}^i \\
 &\quad \left. + \sum_{k=q}^N \mathbb{E}_{\{h, r, h', r'\}_{q, q+1, \dots, N-1}} [d(x_k, \{x_k^d\})] \right\}. \quad (21)
 \end{aligned}$$

The iterative update equation to find optimal policy can be expressed by (22), according to the dynamic programming algorithm in [31].

$$J_N(x_N) = d(x_N, \{x_N^d\}),$$

$$\begin{aligned}
 J_q(x_q) &= \min_{u_q, c_q, v_q, c'_q, h_q, r_q, h'_q, r'_q} \mathbb{E} \left[\sum_{l_1=1}^l \sum_{i=1}^n \alpha_{l_1} u_{q, l_1}^i + \sum_{l_2=1}^m \sum_{i=1}^n \beta_{l_2} c_{q, l_2}^i + \sum_{l_3=1}^{l'} \sum_{i=1}^n \alpha'_{l_3} v_{q, l_3}^i \right. \\
 &\quad \left. + \sum_{l_4=1}^{m'} \sum_{i=1}^n \beta'_{l_4} c'_{q, l_4}^i + d(x_q, \{x_q^d\}) + J_{q+1}(x_{q+1}) \right]
 \end{aligned}$$

TABLE 5: Properties of amino acids.

Amino Acid	Abbrev.	Codon(s)	Polarity	PH	Size
Alanine	Ala	<i>GCT, GCC, GCA, GCG</i>	Nonpolar	6.01	Tiny
Arginine	Arg	<i>CGA, CGG, CGC, CGT, AGA, AGG</i>	Polar	10.76	Normal
Asparagine	Asn	<i>AAC, AAT</i>	Polar	5.41	Small
Aspartic acid	Asp	<i>GAT, GAC</i>	Polar	2.85	Small
Cysteine	Cys	<i>TGT, TGC</i>	Nonpolar	5.05	Small
Glutamine	Gln	<i>CAA, CAG</i>	Polar	5.65	Normal
Glutamic acid	Glu	<i>GAA, GAG</i>	Polar	3.15	Normal
Glycine	Gly	<i>GGA, GGG, GGC, GGT</i>	Nonpolar	6.06	Tiny
Histidine	His	<i>CAC, CAT</i>	Polar	7.60	Normal
Isoleucine	Ile	<i>ATA, ATC, ATT</i>	Nonpolar	6.05	Normal
Leucine	Leu	<i>TTA, TTG, CTA, CTG, CTC, CTT</i>	Nonpolar	6.01	Normal
Lysine	Lys	<i>AAA, AAG</i>	Polar	9.60	Normal
Methionine	Met	<i>ATG</i>	Nonpolar	5.74	Normal
Phenylalanine	Phe	<i>TTC, TTT</i>	Nonpolar	5.49	Normal
Proline	Pro	<i>CCA, CCG, CCC, CCT</i>	Nonpolar	6.30	Small
Serine	Ser	<i>TCA, TCG, TCC, TCT, AGT, AGC</i>	Polar	5.68	Tiny
Threonine	Thr	<i>ACT, ACC, ACA, ACG</i>	Polar	5.60	Small
Tryptophan	Trp	<i>TCC</i>	Nonpolar	5.89	Normal
Tyrosine	Tyr	<i>TAC, TAT</i>	Polar	5.64	Normal
Valine	Val	<i>GTA, GTG, GTC, GTT</i>	Nonpolar	6.00	Small
Stop codon	Term	<i>TAA, TAG, TGA</i>	—	—	—

$$\begin{aligned}
&= \min_{u_q, c_q, v_q, c'_q} \left\{ \sum_{l_1=1}^l \sum_{i=1}^n \alpha_{l_1} u_{q,l_1}^i + \sum_{l_2=1}^m \sum_{i=1}^n \beta_{l_2} c_{q,l_2}^i + \sum_{l_3=1}^{l'} \sum_{i=1}^n \alpha'_{l_3} v_{q,l_3}^i \right. \\
&\quad \left. + \sum_{l_4=1}^{m'} \sum_{i=1}^n \beta'_{l_4} c'_{q,l_4}{}^i \right. \\
&\quad \left. + \mathbb{E}_{h_q, r_q, h'_q, r'_q} \left[d(x_q, \{x_q^d\}) + J_{q+1}(x_{q+1}) \right] \right\}, \\
&\quad q = 0, 1, \dots, N-1.
\end{aligned} \tag{22}$$

3. Results and Discussion

In the following examples, we consider applying chemical mutagens only because the randomness of applying radiation is much larger and more difficult to control. We also omit the mutations between a normal base and O because of high-cost β_{l_2} and the unavailable chemical and physical properties for codons containing O .

The distance reference between codons used in Sections 3.2 and 3.3 is computed by (20) with $\zeta_{\text{polarity}} = 8$, $\zeta_{\text{PH}} = 3$, $\zeta_{\text{size}} = 1$, $\sigma_1 = 2$, $\sigma_2 = 5$ and $\sigma_3 = 3$. We only keep the final penalty but omit the off-trajectory penalty along the trajectory.

3.1. Base-to-Base, Deterministic Optimal Control Problem. We define the distance reference between bases as

$$d(\psi_1, \psi_2) = \begin{cases} 0, & \text{if } x_N = x_N^d, \\ \infty, & \text{if } x_N \neq x_N^d, \end{cases} \tag{23}$$

with $\psi_1, \psi_2 \in \mathcal{F}_{\setminus\{0\}}$, where $\mathcal{F}_{\setminus\{0\}}$ denotes the set \mathcal{F} excluding the element 0.

Our optimal control problem for point mutations is

$$J_0(x_0) = \min_{u_{k,l_1}, 0 \leq k \leq N, 1 \leq l_1 \leq l} \left\{ \sum_{k=0}^{N-1} \sum_{l_1=1}^l \alpha_{l_1} u_{k,l_1} \right\}, \tag{24}$$

subject to

$$x_{k+1} = \left(-I + \sum_{l_1=1}^l u_{k,l_1} s \right) x_k, \tag{25}$$

$$x_N = x_N^d,$$

with x_0 given, $x_k \in \mathcal{F}_{\setminus\{0\}}$.

Suppose that there are 12 kinds of mutagens ($l_1 = 12$), each corresponding to a specific transfer pattern as in Table 6, all available controls and the respective costs can be immediately listed as in Tables 7 and 8.

The elements along the antidiagonal of Table 7, u_{AT} , u_{GC} , u_{CG} , and u_{TA} , are artificially added, because the complementary transfers naturally happen and no mutagen is necessary. Thus, the costs along the antidiagonal of Table 8 are all zero, that is, $\alpha_{AT} = \alpha_{GC} = \alpha_{CG} = \alpha_{TC} = 0$. The equivalence relationship between subscription in two nucleotide bases and subscription in integer $l_1(\psi_1\psi_2) : \{A, T, G, C\} \times \{A, T, G, C\} \rightarrow \{\text{integers from 1 to 12}\}$ is defined by Table 6.

TABLE 6: An example of chemical mutagens and their corresponding transfer patterns in deterministic mutations.

Index (l_1)	1	2	3	4	5	6
Transfer pattern	$A \rightarrow A$	$A \rightarrow G$	$A \rightarrow C$	$G \rightarrow A$	$G \rightarrow G$	$G \rightarrow T$
Index (l_1)	7	8	9	10	11	12
Transfer pattern	$C \rightarrow A$	$C \rightarrow C$	$C \rightarrow T$	$T \rightarrow G$	$T \rightarrow C$	$T \rightarrow T$

TABLE 7: Controls corresponding to transfer between bases within one step. The leftmost column denotes the state k th step, and the upmost row denotes the $(k + 1)$ th state.

kth	(k + 1)th			
	A	G	C	T
A	u_{AA}	u_{AG}	u_{AC}	u_{AT}
G	u_{GA}	u_{GG}	u_{GC}	u_{GT}
C	u_{CA}	u_{CG}	u_{CC}	u_{CT}
T	u_{TA}	u_{TG}	u_{TC}	u_{TT}

TABLE 8: Corresponding step cost of controls as shown in Table 7.

kth	(k + 1)th			
	A	G	C	T
A	α_{AA}	α_{AG}	α_{AC}	α_{AT}
G	α_{GA}	α_{GG}	α_{GC}	α_{GT}
C	α_{CA}	α_{CG}	α_{CC}	α_{CT}
T	α_{TA}	α_{TG}	α_{TC}	α_{TT}

Under the above assumptions, we can write update equation for optimal policy explicitly as

$$J_q(x_q) = \min_{u_{q,j_1}} \left\{ \alpha_{x_q \psi} + J_{q+1}(\psi), \forall \psi \in \{A, T, G, C\} \iff \mathcal{F}_{\setminus \{0\}} \right\}. \quad (26)$$

Proposition 7. For the same x_N^d , $J_q(\psi) \leq J_{q+1}(\bar{\psi})$, for all q , $0 \leq q \leq N - 1$, for all $\psi \in \{A, T, G, C\}$, where $\bar{\psi}$ denotes the complementary base of ψ . If, in addition, the system is completely controllable, $\exists M$, s.t. $J_M(\psi)$ is the global minimum and for all $q \leq M$, $J_q(\psi) = J_M(\bar{\psi})$ if $M - q \equiv 1 \pmod{2}$, and $J_q(\psi) = J_M(\psi)$ if $M - q \equiv 0 \pmod{2}$. In our example, $M \geq N - 6$.

Proof. This first part is due to the zero cost for the transfers between complementary bases in the consecutive steps.

For any $0 \leq q \leq N - 1$, the relationship between minimal costs in consecutive steps is shown in (26). Since $\psi \in \{A, T, G, C\}$, $\alpha_{\psi \bar{\psi}} + J_{q+1}(\bar{\psi})$ is one of the four elements in the set from which the $J_q(\psi)$ is picked. Moreover, $\alpha_{\psi \bar{\psi}} = 0$. Therefore, $J_{q+1}(\bar{\psi})$ is one of the four elements in the set. Since $J_q(\psi)$ is the minimum picking for a set containing $J_{q+1}(\bar{\psi})$, we conclude that $J_q(\psi) \leq J_{q+1}(\bar{\psi})$.

The M value in our example is proved by brute force method, that is, $J_{N-6}(\psi)$ is a guaranteed global minimum. For completely controllable systems, this M always exists.

TABLE 9: Sample step costs.

x_k	x_{k+1}			
	A	G	C	T
A	5.21	6.60	2.33	0
G	6.15	8.95	0	3.82
C	4.61	0	9.17	7.24
T	0	0.64	5.09	10.28

The existence of M implies that for without limitation in the number of steps, we can reach the global optimal in $N - M$ steps, 6 steps in our example.

Suppose that $q = M$, $J_M(\psi)$ is the global minimum, thus $J_{M-1}(\bar{\psi}) \geq J_M(\psi)$. However, $J_{M-1}(\bar{\psi}) \leq J_M(\psi)$ according to the first part of the proposition. Therefore, $J_{M-1}(\bar{\psi}) = J_M(\psi)$ for the same x_N^d . Therefore, $J_{M-1}(\bar{\psi})$ is also a global minimum.

By backward induction, suppose for $q = q_1$, the statement is true, that is, $J_{q_1-1}(\bar{\psi}) = J_{q_1}(\psi)$ is the global optimal either from $x_{q_1-1} = \bar{\psi}$ or $x_{q_1} = \psi$ to x_N^d . Obviously, for $q = q_1 - 1$, the statement is still true. Therefore, $J_q(\psi) = J_{q-2}(\psi) = J_{q-1}(\bar{\psi})$, $\psi \in \{A, T, G, C\}$, for all q , $2 \leq q \leq M$. \square

In the proof of global minimum that can be reached in the finite step in Proposition 7, we also discover Proposition 8. Here, $J_q(x_q, x_N^d)$ denotes the optimal cost from x_q to x_N^d .

Proposition 8. Given two single base mutation optimal control problems, with the same fixed N , with and desired final states complementary to each other. If $J_M(\psi, x_N^d)$ is the global minimum, then $J_M(\bar{\psi}, \bar{x}_N^d)$ is also the global minimum, that is, the global minimum of both systems is reach at the same stage M . Moreover, for all q , $0 \leq q \leq M$,

$$J_q(\psi, x_N^d) = J_q(\bar{\psi}, \bar{x}_N^d), \quad \psi, x_N^d \in \{A, T, G, C\}. \quad (27)$$

Physically, Proposition 8 states that the optimal can be achieve at the same step from a pair of complementary bases to another pair of complementary bases at the same cost. However, this fact is true only for base-to-base deterministic mutations, because the distance reference is well defined by (23).

Now, we show an example with simulation results.

The costs of applying different mutagens are listed in Table 9. It is a numerical assignment to Table 8. Since we apply mutagens before the replication starts, u_{AA} actually transfer A to T and then to A by replication. For simplicity, we just use the k th and $(k + 1)$ th step states as subscripts to represent the corresponding control and cost. The costs of transitions are lower than the costs of transversions. Therefore, $\alpha_{AC}, \alpha_{CA}, \alpha_{GT}, \alpha_{TG}$ is smaller than other mutagens, except artificial ones.

If we use χ to denote the costs of mutagens as listed in Table 9, then

$$\chi = \begin{bmatrix} 5.21 & 6.60 & 2.33 & 0 \\ 6.15 & 8.95 & 0 & 3.82 \\ 4.61 & 0 & 9.17 & 7.24 \\ 0 & 0.64 & 5.09 & 10.28 \end{bmatrix} = \begin{bmatrix} \alpha_{AA} & \alpha_{AG} & \alpha_{AC} & \alpha_{AT} \\ \alpha_{GA} & \alpha_{GG} & \alpha_{GC} & \alpha_{GT} \\ \alpha_{CA} & \alpha_{CG} & \alpha_{CC} & \alpha_{CT} \\ \alpha_{TA} & \alpha_{TG} & \alpha_{TC} & \alpha_{TT} \end{bmatrix}$$

$$\Leftrightarrow \begin{bmatrix} \alpha_1 & \alpha_2 & \alpha_3 & 0 \\ \alpha_4 & \alpha_5 & 0 & \alpha_6 \\ \alpha_7 & 0 & \alpha_8 & \alpha_9 \\ 0 & \alpha_{10} & \alpha_{11} & \alpha_{12} \end{bmatrix} \quad (28)$$

Running the dynamic programming for every pair of $(x_q, x_N^d) \in \{A, T, G, C\} \times \{A, T, G, C\}$, $N = 9$. Here, we slightly modify our notation. We use $J_q(x_q, x_N^d)$ to denote the optimal cost from x_q to x_N^d . Then,

$$J_q = \begin{bmatrix} J_q(A, A) & J_q(A, G) & J_q(A, C) & J_q(A, T) \\ J_q(G, A) & J_q(G, G) & J_q(G, C) & J_q(G, T) \\ J_q(C, A) & J_q(C, G) & J_q(C, C) & J_q(C, T) \\ J_q(T, A) & J_q(T, G) & J_q(T, C) & J_q(T, T) \end{bmatrix} \quad (29)$$

The path to reach the optimal cost is denoted by

$$P_q = \begin{bmatrix} P_q(A, A) & P_q(A, G) & P_q(A, C) & P_q(A, T) \\ P_q(G, A) & P_q(G, G) & P_q(G, C) & P_q(G, T) \\ P_q(C, A) & P_q(C, G) & P_q(C, C) & P_q(C, T) \\ P_q(T, A) & P_q(T, G) & P_q(T, C) & P_q(T, T) \end{bmatrix}, \quad (30)$$

where $P_q(x_q, x_N^d)$ is the $(q + 1)$ th state from x_q to x_N^d , that is, $x_{q+1} = P_q(x_q, x_N^d)$.

The simulation results are shown as below, including optimal costs for all possible transfer pairs $(x_q, x_N^d) \in$

$\{A, T, G, C\} \times \{A, T, G, C\}$, J_q , $0 \leq q \leq 8$, graphical representation in Figure 6, and optimal path P_q , $0 \leq q \leq 7$.

$$J_0 = \begin{bmatrix} 5.21 & 5.09 & 0.64 & 0 \\ 6.15 & 6.79 & 0 & 3.82 \\ 3.82 & 0 & 6.79 & 6.15 \\ 0 & 0.64 & 5.09 & 5.21 \end{bmatrix},$$

$$J_1 = \begin{bmatrix} 0 & 0.64 & 5.09 & 5.21 \\ 3.82 & 0 & 6.79 & 6.15 \\ 6.15 & 6.79 & 0 & 3.82 \\ 5.21 & 5.09 & 0.64 & 0 \end{bmatrix},$$

$$J_2 = \begin{bmatrix} 5.21 & 5.09 & 0.64 & 0 \\ 6.15 & 6.79 & 0 & 3.82 \\ 3.82 & 0 & 6.79 & 6.15 \\ 0 & 0.64 & 5.09 & 5.21 \end{bmatrix},$$

$$J_3 = \begin{bmatrix} 0 & 0.64 & 5.09 & 5.21 \\ 3.82 & 0 & 6.79 & 6.15 \\ 6.15 & 6.79 & 0 & 3.82 \\ 5.21 & 5.09 & 0.64 & 0 \end{bmatrix},$$

$$J_4 = \begin{bmatrix} 5.21 & 5.09 & 0.64 & 0 \\ 6.15 & 6.79 & 0 & 3.82 \\ 3.82 & 0 & 6.79 & 6.15 \\ 0 & 0.64 & 5.09 & 5.21 \end{bmatrix}, \quad (31)$$

$$J_5 = \begin{bmatrix} 0 & 0.64 & 5.09 & 5.21 \\ 3.82 & 0 & 6.79 & 6.15 \\ 6.15 & 6.79 & 0 & 3.82 \\ 5.21 & 5.09 & 0.64 & 0 \end{bmatrix},$$

$$J_6 = \begin{bmatrix} 5.21 & 5.09 & 0.64 & 0 \\ 6.15 & 6.79 & 0 & 3.82 \\ 3.82 & 0 & 7.88 & 6.15 \\ 0 & 0.64 & 5.09 & 5.21 \end{bmatrix},$$

$$J_7 = \begin{bmatrix} 0 & 0.64 & 5.09 & 5.21 \\ 3.82 & 0 & 8.48 & 6.15 \\ 6.15 & 7.88 & 0 & 3.82 \\ 5.21 & 5.09 & 0.64 & 0 \end{bmatrix},$$

$$J_8 = \begin{bmatrix} 5.21 & 6.60 & 2.33 & 0 \\ 6.15 & 8.95 & 0 & 3.82 \\ 4.61 & 0 & 9.17 & 7.24 \\ 0 & 0.64 & 5.09 & 10.28 \end{bmatrix},$$

$$\begin{aligned}
P_0 &= \begin{bmatrix} A, T & T & T & T \\ A, C & A, C & C & C, T \\ G & G & G & G \\ A & A, G & A, C & A \end{bmatrix}, \\
P_1 &= \begin{bmatrix} T & T & T & A, T \\ C, T & C & A, C & A, C \\ G & G & G & G \\ A & A, C & A, G & A \end{bmatrix}, \\
P_2 &= \begin{bmatrix} A, T & T & T & T \\ A, C & A, C & C & C, T \\ G & G & G & G \\ A & A, G & A, C & A \end{bmatrix}, \\
P_3 &= \begin{bmatrix} T & T & T & A, T \\ C, T & C & A, C & A, C \\ G & G & G & G \\ A & A, C & A, G & A \end{bmatrix}, \\
P_4 &= \begin{bmatrix} A, T & T & T & T \\ A, C & A, C & C & C, T \\ G & G & G & G \\ A & A, G & A, C & A \end{bmatrix}, \\
P_5 &= \begin{bmatrix} T & T & T & A, T \\ C, T & C & A & A, C \\ G & G & G & G \\ A & A, C & A, G & A \end{bmatrix}, \\
P_6 &= \begin{bmatrix} A, T & T & T & T \\ A, C & A & C & C, T \\ G & G & T & G \\ A & A, G & A, C & A \end{bmatrix}, \\
P_7 &= \begin{bmatrix} T & T & T & A \\ T & C & A & A \\ G & T & G & G \\ A & C & G & A \end{bmatrix}.
\end{aligned} \tag{32}$$

For simplicity, we use 1 to represent A, 2 to G, 3 to C, and 4 to T in graphical interpretation. From Figure 6, we can see clearly that the optimal cost decreases as q decreases in the first few steps, and then optimal cost remains at the global minimum. This phenomenon obeys Proposition 7. In this example, global optimal is reached at $M = 5$ for all pairs of initial and final states as $J_7 \neq J_5 = J_3$ and $J_6 \neq J_4 = J_2$. So, the global minimum is achieved before we reach $N - 5 = 4$ in this particular case. This also implies that with N free we can reach desired final state in 4 steps from given initial state.

Observing closely to J_q , $0 \leq q \leq 5$, J_{q-1} equals to J_q by exchanging the first and the last columns, and the second and the third columns, which is consistent with Proposition 8. Or we can exchange the first and the last rows, and the second and the third rows of J_q to obtain J_{q-1} . J_{q_1} and J_{q_2} are the same for $q_1, q_2 \leq M = 5$ for $q_1 - q_2 = 0 \pmod{2}$. This obeys Proposition 7.

The optimal trajectories are generated from $P_q(x_q, x_N^d)$. For example, given $x_2 = T$, and the final state $x_9^d = G$, we want to generate the optimal trajectories.

$$x_3 = P_2(T, G) = A, G.$$

If $x_3 = A$, $x_4 = P_3(A, G) = T$; if $x_3 = G$, $x_4 = P_3(A, G) = C$.

If $x_4 = T$, $x_5 = P_4(T, G) = A, G$; if $x_4 = C$, $x_5 = P_4(C, G) = G$.

If $x_5 = A$, $x_6 = P_5(A, G) = T$; if $x_5 = G$, $x_6 = P_5(G, G) = C$.

If $x_6 = T$, $x_7 = P_6(T, G) = A, G$; if $x_6 = C$, $x_7 = P_6(C, G) = G$.

If $x_7 = A$, $x_8 = P_7(A, G) = T$; if $x_7 = G$, $x_8 = P_7(G, G) = C$.

So, the optimal routes are

$$T \rightarrow A \rightarrow T \rightarrow A \rightarrow T \rightarrow A \rightarrow T \xrightarrow[\alpha_{TG}]{u_{TG}} G;$$

$$T \rightarrow A \rightarrow T \rightarrow A \rightarrow T \xrightarrow[\alpha_{TG}]{u_{TG}} G \rightarrow C \rightarrow G;$$

$$T \rightarrow A \rightarrow T \xrightarrow[\alpha_{TG}]{u_{TG}} G \rightarrow C \rightarrow G \rightarrow C \rightarrow G;$$

$$T \xrightarrow[\alpha_{TG}]{u_{TG}} G \rightarrow C \rightarrow G \rightarrow C \rightarrow G \rightarrow C \rightarrow G.$$

Consequently, the optimal cost is $J_2(T, G) = 0.64 = \alpha_{TG}$.

Optimal trajectories for other pairs of initial and final states can be obtained in the same manner.

It takes less than 1 second to generate optimal trajectories for all pairs of initial and final states with $N = 9$ on a regular desktop. Since we have already proven by the brute force method that the global optimal can be achieved with $M \leq 6$, the computation time can be further reduced by taking $N = 6$ with all the results necessary for this example.

3.2. Codon-to-Codon, Deterministic Optimal Control Problem.

For codon-to-codon deterministic mutations, we formulate our optimal control problem as

$$J_0(x_0) = \min_{\substack{u_{k,l}^i, 0 \leq k \leq N, \\ 1 \leq l_1 \leq l, 1 \leq i \leq 3}} \left\{ \sum_{k=0}^{N-1} \sum_{l_1=1}^l \sum_{i=1}^3 \alpha_{l_1} u_{k,l_1}^i + d(x_N, \{x_N^d\}) \right\}, \tag{33}$$

subject to

$$x_{k+1} = \left(-I + \sum_{l_1=1}^l \sum_{i=1}^3 u_{k,l_1}^i s e_i e_i^T \right) x_k, \tag{34}$$

with $x_0, x_N^d \in \mathcal{F}_{\setminus\{0\}}^3$ given, $x_k \in \mathcal{F}_{\setminus\{0\}}^3$, for all k , $0 \leq k \leq N$, and $d(\varphi_1, \varphi_2)$, $\varphi_1, \varphi_2 \in \mathcal{F}_{\setminus\{0\}}^3$ as defined in Section 2.3.

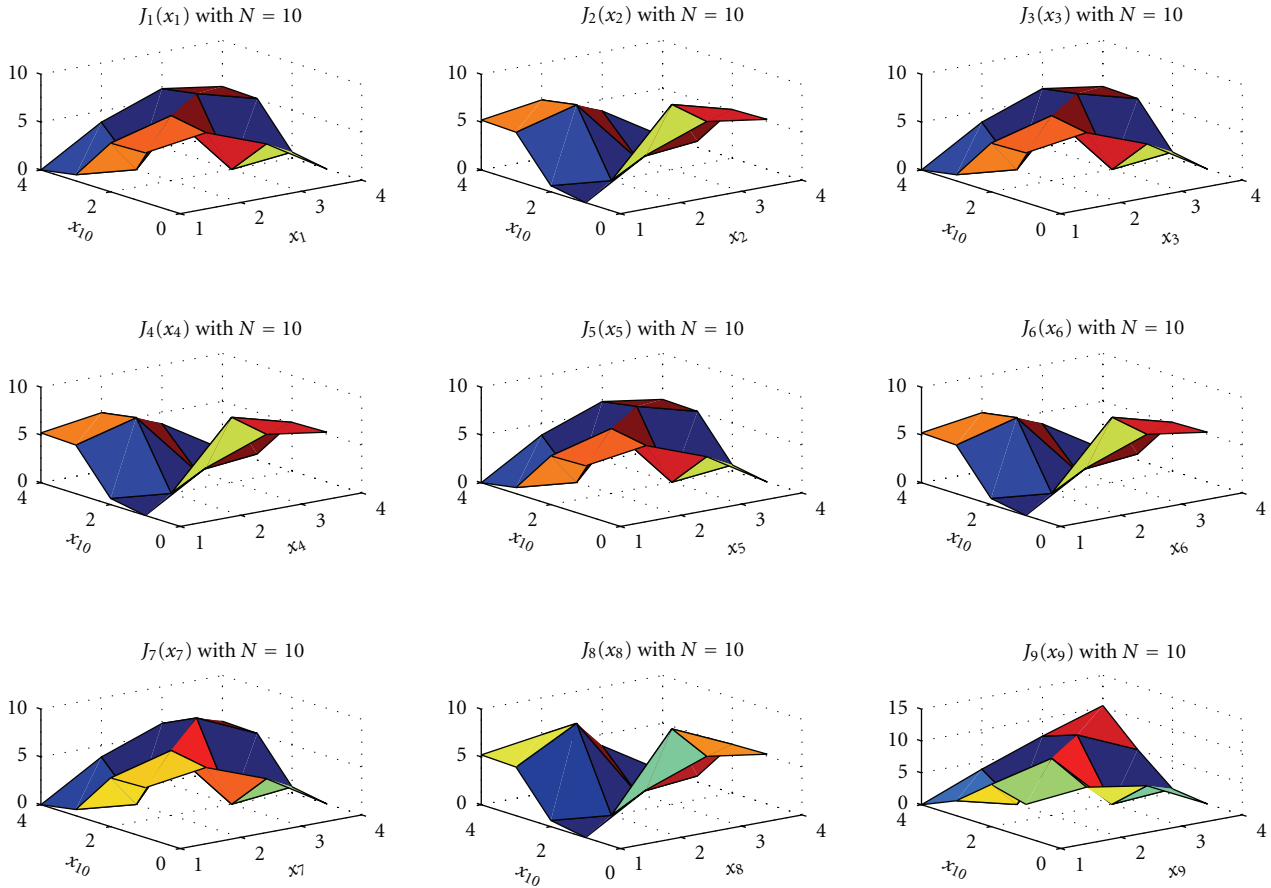


FIGURE 6: Graphical representation of J_q , $0 \leq q \leq 8$, $N = 9$, in single-base deterministic mutation example. The x -axis and y -axis represent x_q and x_N^d , respectively. $J_q(x_q, x_N^d)$ are represented by 16 isolated points. Those discrete points are connected together to show the surface.

If we take the same assumption on available mutagens as in Section 3.1, then we can write update equation for optimal control policy explicitly as

$$\begin{aligned}
 & J_q(x_q) \\
 = & \min_{u_{q,l}^i, 1 \leq l \leq 3, 1 \leq i \leq 3} \left\{ \alpha_{x_q^1 \psi_1} + J_{q+1} \left(\begin{bmatrix} \psi_1 \\ \overline{x_q^2} \\ \overline{x_q^3} \end{bmatrix} \right), \right. \\
 & \left. \alpha_{x_q^2 \psi_2} + J_{q+1} \left(\begin{bmatrix} \overline{x_q^1} \\ \psi_2 \\ \overline{x_q^3} \end{bmatrix} \right), \alpha_{x_q^3 \psi_3} + J_{q+1} \left(\begin{bmatrix} \overline{x_q^1} \\ \overline{x_q^2} \\ \psi_3 \end{bmatrix} \right), \right. \\
 & \left. \psi_1, \psi_2, \psi_3 \in \{A, T, G, C\} \Leftrightarrow \mathcal{F}_{\setminus\{0\}} \right\}, \quad (35)
 \end{aligned}$$

with $x_q^i \in \{A, T, G, C\} \Leftrightarrow \mathcal{F}_{\setminus\{0\}}$, $1 \leq i \leq 3$ denotes the i th element of $x_q \in \mathcal{F}_{\setminus\{0\}}^3$, and $\overline{x_q^i}$ denotes the complementary base of x_q^i .

The optimal control sequences depend on the numerical values of α_i s and $d(\varphi_1, \varphi_2)$, $\varphi_1, \varphi_2 \in \mathcal{F}_{\setminus\{0\}}^3$. Though we do not have real values of α_i s and $d(\varphi_1, \varphi_2)$, we can always obtain simulation results to compare the result differences by assigning different numerical values to those parameters.

Therefore, we use three different assignments of α_i s and the same $d(\varphi_1, \varphi_2)$ to generate our simulation results. Those three assignments of α_i s are χ , 5χ , and 0.5χ , respectively, with χ the same as assigned in Section 3.1. In every particular example, it takes approximately 2 seconds on a regular desktop to generate the optimal path table for all pairs of initial and final states with $N = 19$, and the dynamic programming algorithm ensures that the optimal control for tail problems is generated at the same time.

The graphical interpretation of three assignments are shown in Figures 7, 8, and 9, respectively. The x -axis and y -axis denote x_q and x_N^d , respectively. For a codon $[\psi_1 \ \psi_2 \ \psi_3]^T$, $\psi_1, \psi_2, \psi_3 \in \{A, T, G, C\} \Leftrightarrow \mathcal{F}_{\setminus\{0\}}$, its index is calculated by

$$4^2(\psi_1 - 1) + 4(\psi_2 - 1) + \psi_3, \quad (36)$$

where $\psi_i = 1$ if A, $\psi_i = 2$ if G, $\psi_i = 3$ if C, and $\psi_i = 4$ if T, $1 \leq i \leq 3$, for the simplicity of graphical interpretation. A codon has $4^3 = 64$ combinations. Thus, there are 64^2 pairs

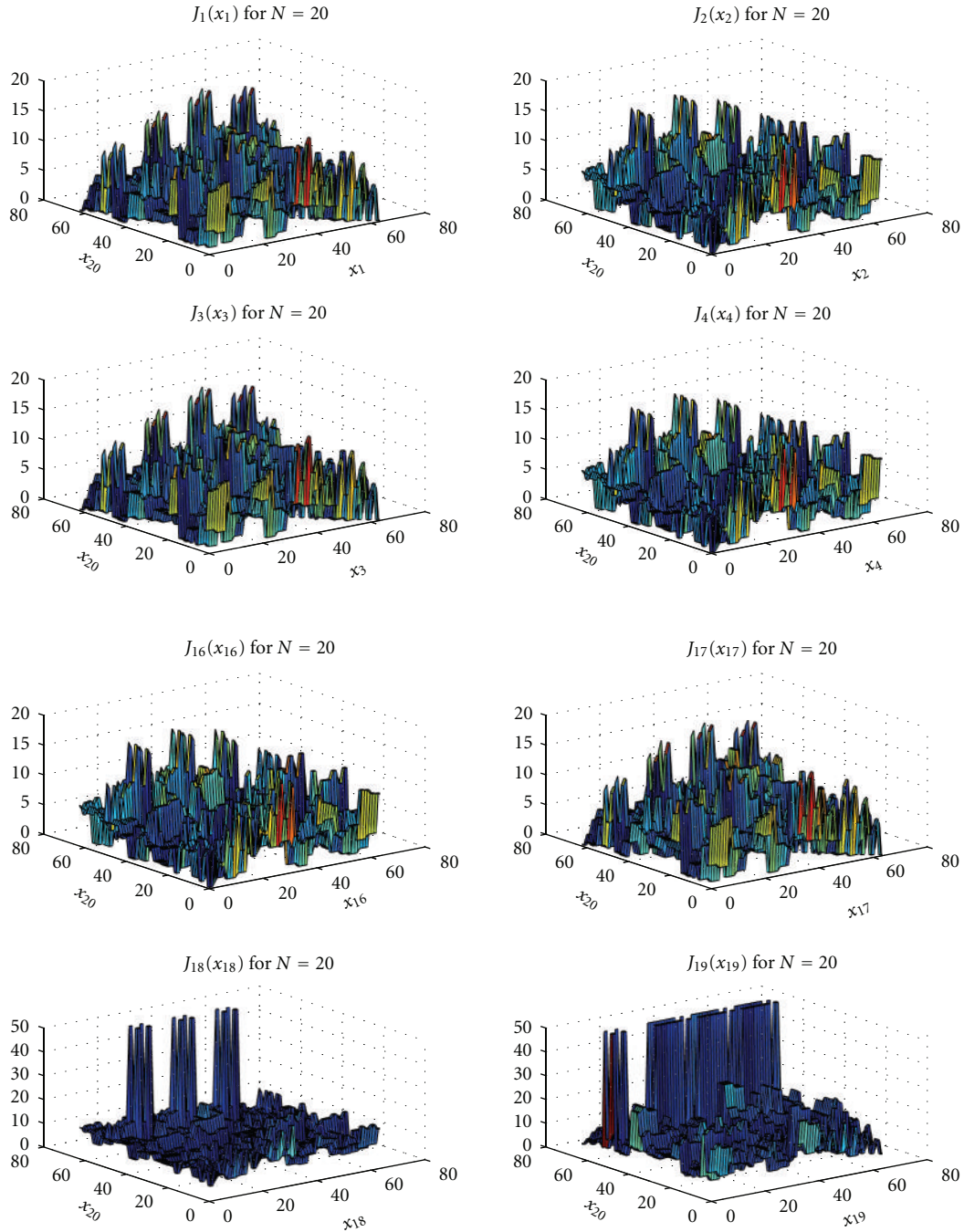


FIGURE 7: Graphical representation of J_q , $q = 0, 1, 2, 3, 15, 16, 17, 18$ for codon-to-codon deterministic mutations, with $\alpha_i = \chi$, $N = 19$.

of initial and final desired states, and there are 64×21 pairs of initial state and final desired set. The surface is generated by connecting 64×64 discrete points together. J_q is calculated following the same procedure as in base-to-base deterministic cases. The value of optimal cost can be read directly from graphical interpretation, and the optimal path can be generated from path matrix P_q , similar to base-to-base

deterministic case. Both J_q and P_q , for all q , $0 \leq q \leq N$, are of 64×64 dimension.

From the graphical interpretation and Table 10, we find that the value of q where the global minimum is reached at the first time decreases as α_i decreases. And J_0 is more similar to J_{18} with $\alpha = 5\chi$ than with $\alpha = \chi$ or $\alpha = 0.5\chi$. This implies that if $d(\varphi_1, \varphi_2)$ are the deterministic term in our objective

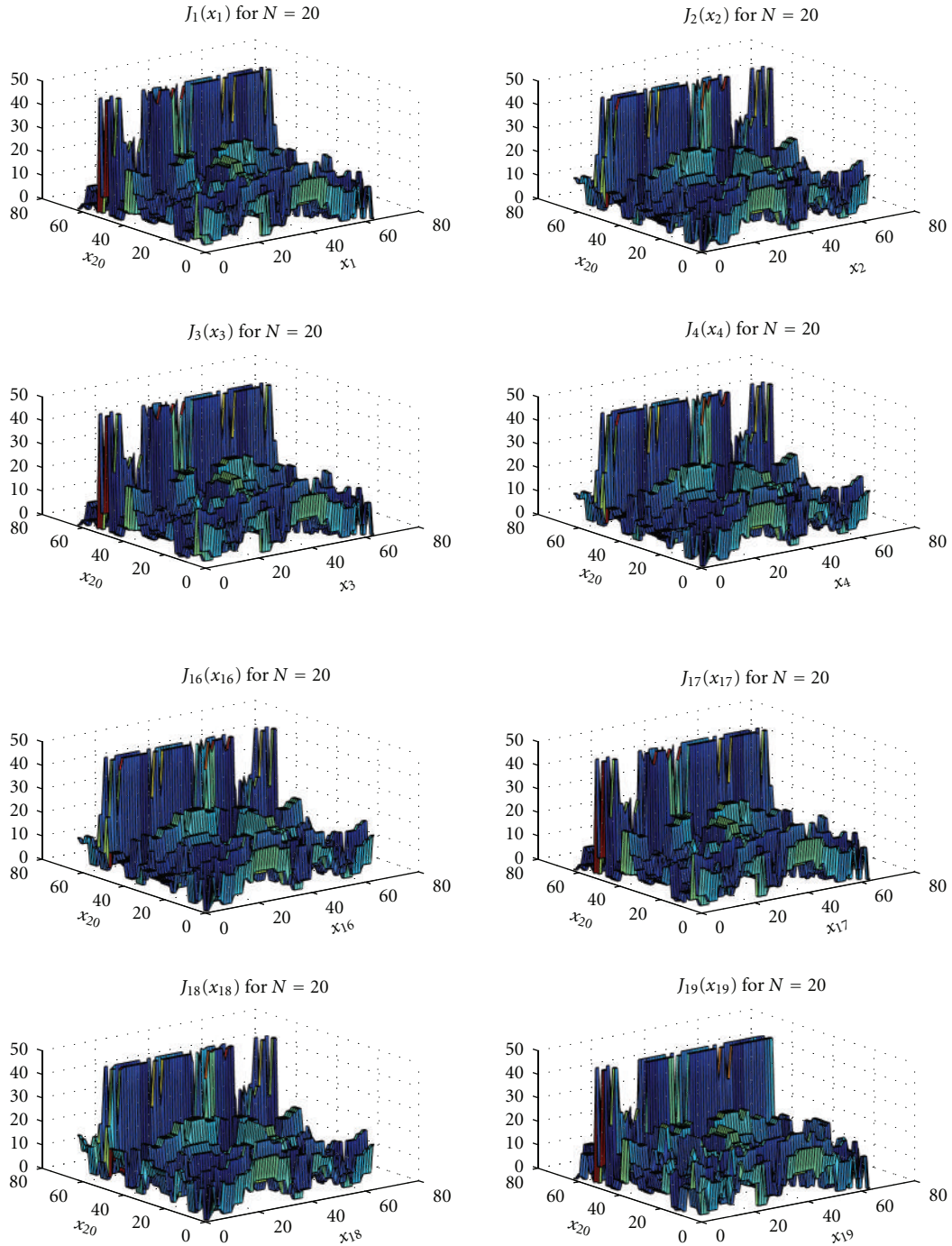


FIGURE 8: Graphical representation of J_q , $q = 0, 1, 2, 3, 15, 16, 17, 18$ for codon-to-codon deterministic mutations, with $\alpha_i = 5\chi$, $N = 19$.

TABLE 10: Simulation results with different α_i assignments and the first q where the global optimal is reached.

N	α_i	M (first q when the global minimum is reached)	Global minimum
19	0.5χ	$q = 12$	$J_{12}(x_{12})$
19	χ	$q = 13$	$J_{13}(x_{13})$
19	5χ	$q = 15$	$J_{15}(x_{15})$

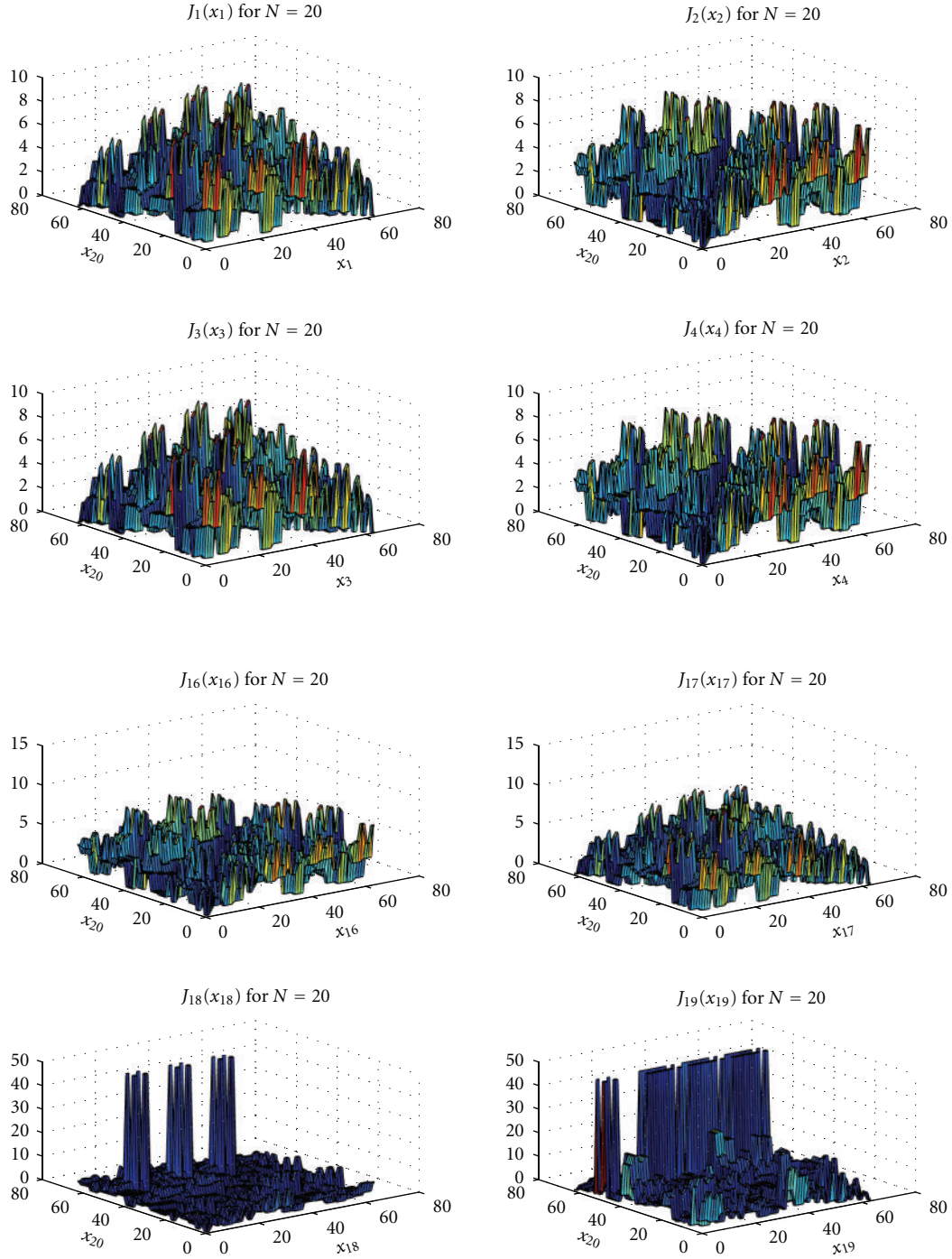


FIGURE 9: Graphical representation of J_q , $q = 0, 1, 2, 3, 15, 16, 17, 18$ for codon-to-codon deterministic mutations, with $\alpha_i = 0.5\chi$, $N = 19$.

function, then the treatment plan is made to drive the final state as close to the desired set as possible; if the costs of applying mutagens is the deterministic term in the objective function, then the treatment plan tends to stay in the original state and applies less mutagens; if they are of equal weight, then the treatment plan deals with this tradeoff.

Moreover, no matter how the numerical values of final penalty and the costs of applying mutagens changes in our objective function as shown in (33), there is always a M ,

$M \leq N - 18$, $J_M(x_M)$ is global minimum. Proposition 7 can be extended to codon-to-codon deterministic mutations as stated in Proposition 9.

Proposition 9. *Given an optimal control problem with objective function in the form of (33), constraints in the form of (34), and all available chemical mutagens and their corresponding transfer pairs and costs as listed in Tables 6, 7, and 8, $J_q(\varphi) \leq J_{q+1}(\varphi)$, $\varphi \in \mathcal{F}_{\setminus\{0\}}^3$. If, in addition, the system is*

TABLE 11: 12 kinds of mutagens, each corresponding to major transfer patterns, and the probabilities of different mutagens on different transfer patterns.

Index (l_1)	From	To				Major transfer pattern
		A	G	C	T	
1	A	$P_{1,AA}^{(h)}$	$P_{1,AG}^{(h)}$	$P_{1,AC}^{(h)}$	$P_{1,AT}^{(h)}$	A → A
2	A	$P_{2,AA}^{(h)}$	$P_{2,AG}^{(h)}$	$P_{2,AC}^{(h)}$	$P_{2,AT}^{(h)}$	A → G
3	A	$P_{3,AA}^{(h)}$	$P_{3,AG}^{(h)}$	$P_{3,AC}^{(h)}$	$P_{3,AT}^{(h)}$	A → C
4	G	$P_{4,GA}^{(h)}$	$P_{4,GG}^{(h)}$	$P_{4,GC}^{(h)}$	$P_{4,GT}^{(h)}$	G → A
5	G	$P_{5,GA}^{(h)}$	$P_{5,GG}^{(h)}$	$P_{5,GC}^{(h)}$	$P_{5,GT}^{(h)}$	G → G
6	G	$P_{6,GA}^{(h)}$	$P_{6,GG}^{(h)}$	$P_{6,GC}^{(h)}$	$P_{6,GT}^{(h)}$	G → T
7	C	$P_{7,CA}^{(h)}$	$P_{7,CG}^{(h)}$	$P_{7,CC}^{(h)}$	$P_{7,CT}^{(h)}$	C → A
8	C	$P_{8,CA}^{(h)}$	$P_{8,CG}^{(h)}$	$P_{8,CC}^{(h)}$	$P_{8,CT}^{(h)}$	C → C
9	C	$P_{9,CA}^{(h)}$	$P_{9,CG}^{(h)}$	$P_{9,CC}^{(h)}$	$P_{9,CT}^{(h)}$	C → T
10	T	$P_{10,TA}^{(h)}$	$P_{10,TG}^{(h)}$	$P_{10,TC}^{(h)}$	$P_{10,TT}^{(h)}$	T → G
11	T	$P_{11,TA}^{(h)}$	$P_{11,TG}^{(h)}$	$P_{11,TC}^{(h)}$	$P_{11,TT}^{(h)}$	T → C
12	T	$P_{12,TA}^{(h)}$	$P_{12,TG}^{(h)}$	$P_{12,TC}^{(h)}$	$P_{12,TT}^{(h)}$	T → T

TABLE 12: Sample probabilities with respect to different mutagens and different transfer patterns.

Index (l_1)	From	To			
		A	G	C	T
1	A	0.90	0.05	0.03	0.02
2	A	0.11	0.58	0.21	0.10
3	A	0.14	0.16	0.42	0.28
4	G	0.85	0.07	0.03	0.05
5	G	0.02	0.02	0.92	0.04
6	G	0.10	0.09	0.22	0.59
7	C	0.79	0.13	0.04	0.04
8	C	0.01	0.02	0.87	0.10
9	C	0.04	0.12	0.09	0.75
10	T	0.13	0.76	0.05	0.06
11	T	0.07	0.03	0.62	0.28
12	T	0.08	0.04	0.25	0.63

completely controllable, $\exists M$, s.t. $J_M(\varphi)$ is the global minimum and for all $q \leq M$, $J_q(\varphi) = J_M(\bar{\varphi})$ if $M - q \equiv 1 \pmod{2}$, and $J_q(\varphi) = J_M(\varphi)$ if $M - q \equiv 0 \pmod{2}$. In our example, $M \geq N - 18$.

Proof. We only prove that $M \geq N - 18$ here since the rest is similar to the proof of Proposition 7.

The objective function in (33) can be written as the summation of three separate single-base mutation systems and the distance reference between final states and the final desired set, that is,

$$J_q(x_q)$$

$$= \min_{\substack{n_1, n_2, n_3 \geq 0 \\ 2N \leq n_1 + n_2 + n_3 \leq 3N - 1}} \left\{ \underbrace{J_{n_1}(x_{n-1}) \left(x_q^1, \psi_1 \right)}_{\text{optimal costs of base-to-base deterministic optimal control problem formed by the 1st base}} \right.$$

$$+ \underbrace{J_{n_2}(x_{n-2}) \left(x_q^2, \psi_2 \right)}_{\text{optimal costs of base-to-base deterministic optimal control problem formed by the 2nd base}}$$

$$+ \underbrace{J_{n_3}(x_{n-3}) \left(x_q^3, \psi_3 \right)}_{\text{optimal costs of base-to-base deterministic optimal control problem formed by the 3rd base}}$$

$$+ d \left(\begin{bmatrix} \psi_1 \\ \psi_2 \\ \psi_3 \end{bmatrix}, \{x_N^d\} \right) \Bigg\},$$

(37)

where $N - q = (N - n_1) + (N - n_2) + (N - n_3)$.

According to Proposition 7, $J_{N-6}(x_{N-6})$ is guaranteed to be the global optimal for single-base mutations. Therefore, optimal costs corresponding to three single-base mutation

systems, $J_{n_1}(x_{n-1})(x_q^1, \psi_1)$, $J_{n_2}(x_{n-2})(x_q^2, \psi_2)$, $J_{n_3}(x_{n-3})(x_q^3, \psi_3)$ are guaranteed to reach their own global optimal at $n_1 = n_2 = n_3 = N - 6$ with all possible combinations of $\psi_1, \psi_2, \psi_3 \in \{A, T, G, C\}$. Therefore, $q = N - 18$ is a guaranteed global optimal. \square

Indeed, Proposition 9 is a quite loose condition. In real examples above, the M value where the first global optimal is reached at earlier stage as listed in Table 10.

Graphically, the indices of complementary codons $\varphi_1 = [\psi_1 \ \psi_2 \ \psi_3]^T$ and $\varphi_2 = [\bar{\psi}_1 \ \bar{\psi}_2 \ \bar{\psi}_3]^T$ sum up to 65, that is,

$$\begin{aligned} & (16(\psi_1 - 1) + 4(\psi_2 - 1) + \psi_3) \\ & + (16(\bar{\psi}_1 - 1) + 4(\bar{\psi}_2 - 1) + \bar{\psi}_3) \\ & = (16(\psi_1 - 1) + 4(\psi_2 - 1) + \psi_3) \\ & + (16((5 - \psi_1) - 1) + 4((5 - \psi_2) - 1) + (5 - \psi_3)) \\ & = 65. \end{aligned} \quad (38)$$

Therefore, J_q and J_{q-1} , $1 \leq q \leq M$ are symmetric about the plane $x = 32.5$, J_{q-2} and J_q , $2 \leq q \leq M$ are the same, as shown in Figures 7, 8, and 9.

However, Proposition 8 cannot be extended to codon-to-codon deterministic case due to the redundancy of genetic codes, that is, the set of codons translated to the same amino acid varies from one amino acid to another as shown in Table 5. The simulation results show that the costs, a pair of complementary codons, to two final desired set generated by a pair of complementary final desired codons are different, that is, $J_q(x_q, \{x_N^d\}) \neq J_q(\bar{x}_q, \{\bar{x}_N^d\})$, in general, for any q . Graphically, the optimal cost profile J_q is not symmetric about the plane $y = 32.5$ for J_{q-1} , $1 \leq q \leq M$. Therefore, the doctors need to pick the strand with lower cost to make the treatment plan. This also implies that in large-scale cases, for instance, a gene containing hundreds of nucleotide bases, the doctors should make the treatment plan based on the strand the total cost of which is lower than the other.

3.3. Codon-to-Codon, Stochastic Optimal Control Problem. The optimal control problem of codon-to-codon stochastic mutations can be written as

$$\begin{aligned} J_0(x_0) = \min_{\substack{u_{k,l}^i, 0 \leq k \leq N-1 \\ 1 \leq l_1 \leq l, 1 \leq i \leq 3}} & \left\{ \sum_{k=0}^{N-1} \sum_{l_1=1}^l \sum_{i=1}^3 \alpha_{l_1} u_{k,l_1}^i \right. \\ & \left. + \mathbb{E}_{\substack{h_{k,l_1}^i, 0 \leq k \leq N-1 \\ 1 \leq l_1 \leq l, 1 \leq i \leq 3}} [d(x_N, \{x_N^d\})] \right\}, \end{aligned} \quad (39)$$

subject to

$$x_{k+1} = -Ix_k + \sum_{l_1=1}^l \sum_{i=1}^3 u_{k,l_1}^i h_{k,l_1}^i s e_i e_i^T x_k, \quad (40)$$

with $x_0, x_N^d \in \mathcal{F}_{\setminus\{0\}}^3$ given, $x_k \in \mathcal{F}_{\setminus\{0\}}^3$.

The major difference between deterministic and stochastic systems is that we impose the random binary vector, h_{k,l_1}^i , in our system equation (40). We denote the probability associated with $h_{k,l_1}^{i,j}$ to be $p_{l_1, \psi_1 \psi_2}^{(h)}$ with $\psi_1, \psi_2 \in \{A, T, G, C\}$. The equivalence relationship between j and $\psi_1 \psi_2$ can be found in Table 3.

Again, we assume that we have $l_1 = 12$ kinds of mutagens, each corresponding to one major transfer pattern, associated with probability assignments, as listed in Table 11.

Then, we can write updated formula for optimal control policy explicitly as

$$\begin{aligned} J_q(x_q) & = \min_{u_{q,l_1}^i, 1 \leq l_1 \leq l, 1 \leq i \leq 3} \left\{ \alpha_{x_q^1 \psi_1} + h_{q,l_1}^1 \mathbb{E}_{(x_q^1 \psi_1)} \left[J_{q+1} \left(\begin{bmatrix} \cdot \\ \bar{x}_q^2 \\ \bar{x}_q^3 \end{bmatrix} \right) \right], \right. \\ & \alpha_{x_q^2 \psi_2} + h_{q,l_1}^2 \mathbb{E}_{(x_q^2 \psi_2)} \left[J_{q+1} \left(\begin{bmatrix} x_q^1 \\ \cdot \\ \bar{x}_q^3 \end{bmatrix} \right) \right], \\ & \alpha_{x_q^3 \psi_3} + h_{q,l_1}^3 \mathbb{E}_{(x_q^3 \psi_3)} \left[J_{q+1} \left(\begin{bmatrix} x_q^1 \\ \bar{x}_q^2 \\ \cdot \end{bmatrix} \right) \right], \\ & \left. \psi_1, \psi_2, \psi_3 \in \{A, T, G, C\} \iff \mathcal{F}_{\setminus\{0\}} \right\} \end{aligned} \quad (41)$$

where

$$\begin{aligned} h_{q,l_1}^1 \mathbb{E}_{(x_q^1 \psi_1)} \left[J_{q+1} \left(\begin{bmatrix} \cdot \\ \bar{x}_q^2 \\ \bar{x}_q^3 \end{bmatrix} \right) \right] & = p_{l_1(x_q^1 \psi_1), x_q^1 A}^{(h)} \left[J_{q+1} \left(\begin{bmatrix} A \\ \bar{x}_q^2 \\ \bar{x}_q^3 \end{bmatrix} \right) \right] \\ & + p_{l_1(x_q^1 \psi_1), x_q^1 G}^{(h)} \left[J_{q+1} \left(\begin{bmatrix} G \\ \bar{x}_q^2 \\ \bar{x}_q^3 \end{bmatrix} \right) \right] \\ & + p_{l_1(x_q^1 \psi_1), x_q^1 C}^{(h)} \left[J_{q+1} \left(\begin{bmatrix} C \\ \bar{x}_q^2 \\ \bar{x}_q^3 \end{bmatrix} \right) \right] \\ & + p_{l_1(x_q^1 \psi_1), x_q^1 T}^{(h)} \left[J_{q+1} \left(\begin{bmatrix} T \\ \bar{x}_q^2 \\ \bar{x}_q^3 \end{bmatrix} \right) \right], \end{aligned} \quad (42)$$

where $x_q^i \in \{A, T, G, C\} \Leftrightarrow \mathcal{F}_{\setminus\{0\}}$, $0 \leq q \leq N - 1$, $1 \leq i \leq 3$ denotes the i th element of $x_q \in \mathcal{F}_{\setminus\{0\}}^3$, \bar{x}_q^i

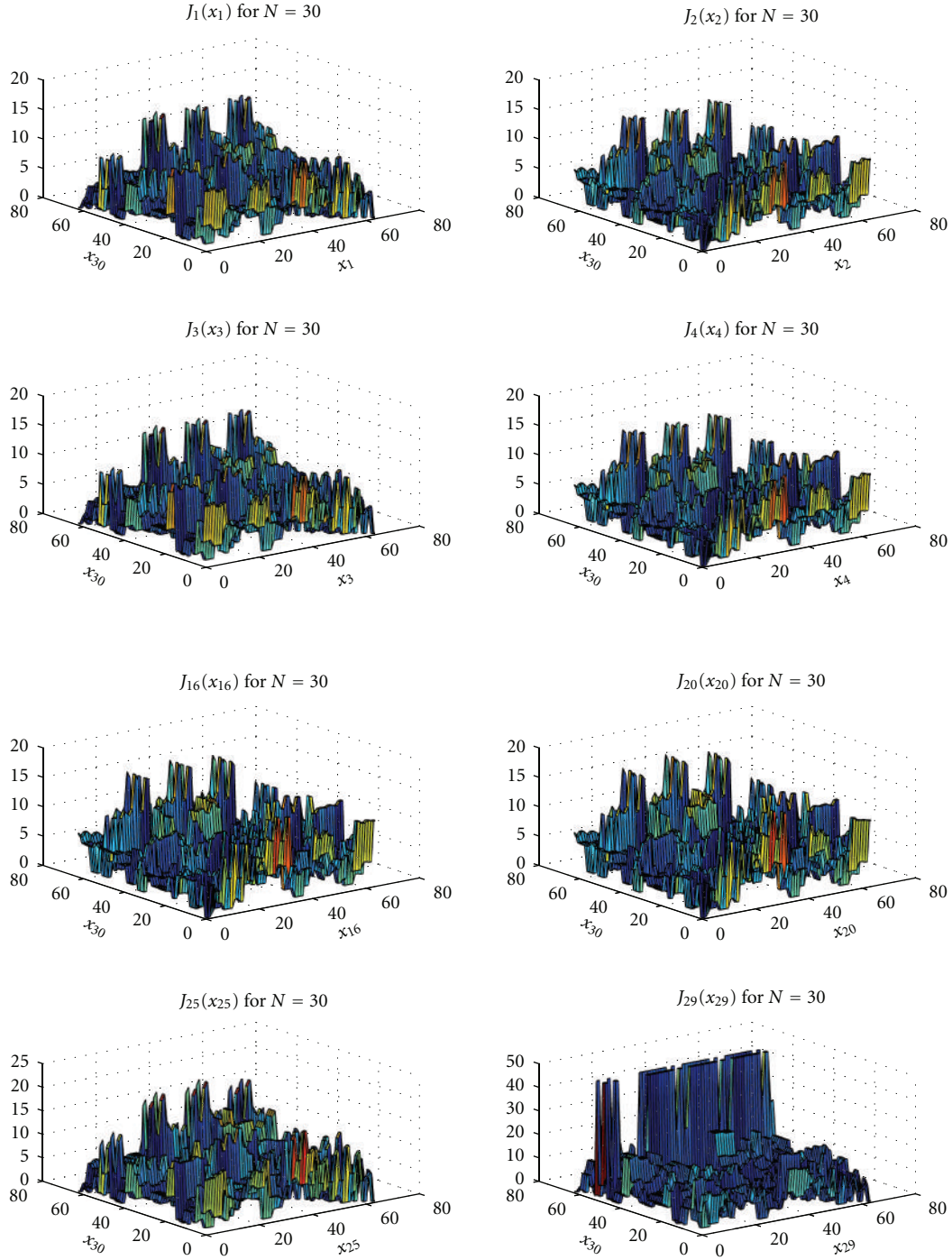


FIGURE 10: Graphical representation of $J_q(x_q)$, $q = 0, 1, 2, 3, 15, 19, 24, 28$ for codon-to-codon stochastic mutations, with $\alpha_{l_1} = \chi$, probability assignment as in Table 12, $N = 29$.

denotes the complementary base of x_q^i , and $l_1 : \psi_1 \psi_2 \in \{A, T, G, C\} \times \{A, T, G, C\} \rightarrow \{\text{integers from 1 to 12}\}$, the mapping from major transfer pattern $\psi_1 \rightarrow \psi_2$ to mutagen index, as shown in Table 11. The mathematical expression of $\mathbb{E}_{h^2} [J_{q+1}([\overline{x}_q^1 \cdot \overline{x}_q^3]^T)]$ and $\mathbb{E}_{h^3} [J_{q+1}([\overline{x}_q^1 \overline{x}_q^2 \cdot]^T)]$ is similar to $\mathbb{E}_{h^1} [J_{q+1}([\cdot \overline{x}_q^2 \overline{x}_q^3]^T)]$ as shown above.

In order to run the simulation, we assign numerical values to probabilities in Table 11, as illustrated in Table 12.

As in Section 3.2, we use three different assignments for α_{l_1} s, χ , 5χ , and 0.5χ , respectively. The optimal cost profile J_q with selected q values, for every pair of (x_q, x_N^d) , is graphically interpreted in Figures 10, 11, and 12, respectively, with $N = 29$, with computation time of approximately 7 seconds on a regular desktop.

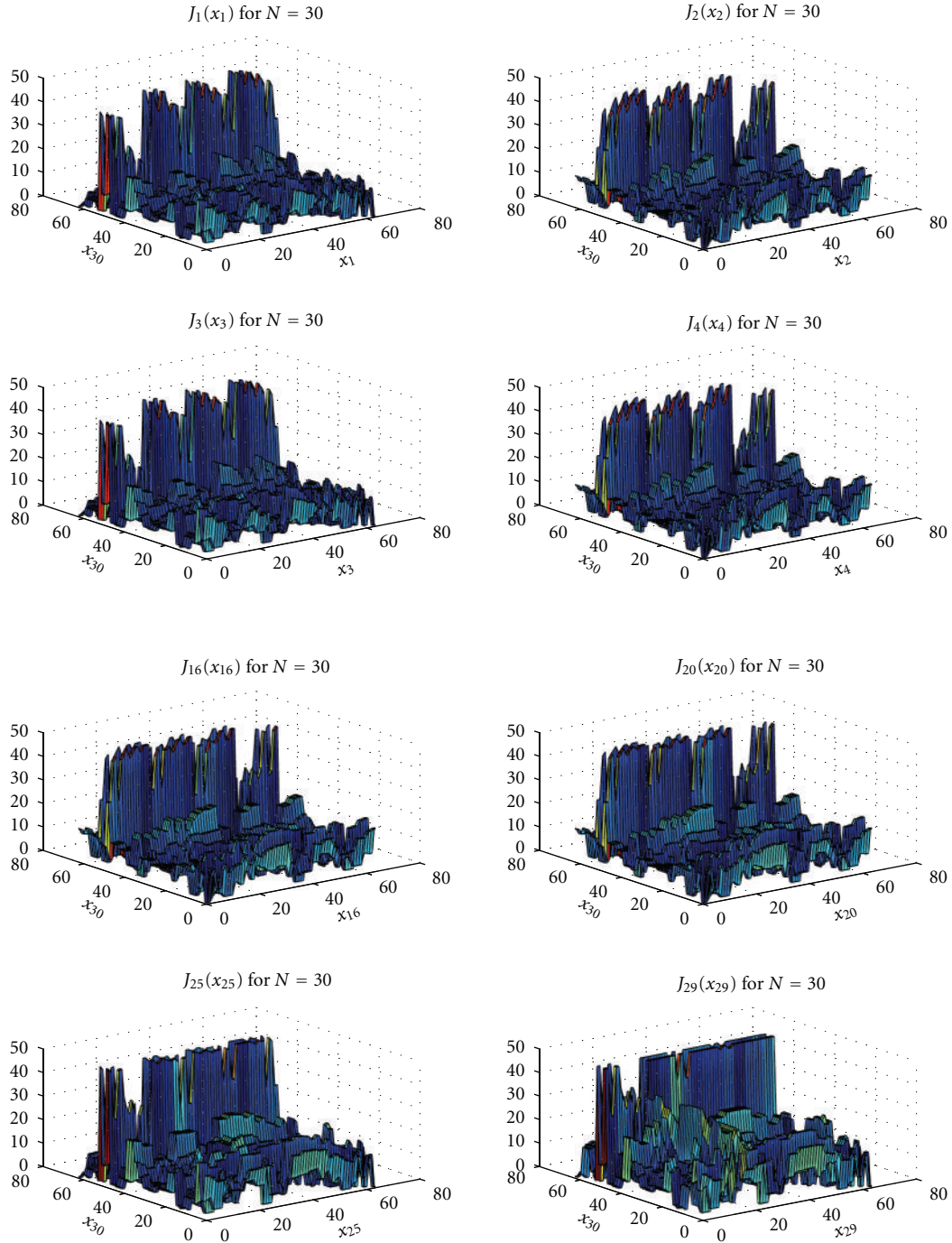


FIGURE 11: Graphical representation of $J_q(x_q)$, $q = 0, 1, 2, 3, 15, 19, 24, 28$ for codon-to-codon stochastic mutations, with $\alpha_{i_1} = 5\chi$, $N = 29$.

The simulation results are similar to those in Section 3.2. The profile of J_0 is more similar to J_{29} when α_{i_1} s are assigned 5χ than χ or 0.5χ . This implies in codon-to-codon stochastic mutations; the optimal control sequence behaves as codon-to-codon deterministic cases, that is, the system tends to getting as close as possible to the final desired set if α_{i_1} s are much smaller than $d(\varphi_1, \varphi_2)$, and the system tends to remain in the same state with minor mutations when α_{i_1} s are relatively larger than $d(\varphi_1, \varphi_2)$, $\varphi_1, \varphi_2 \in \mathcal{F}_{\setminus\{0\}}^3$.

And $J_q(\bar{\varphi}) \leq J_{q+1}(\varphi)$, $\varphi \in \mathcal{F}_{\setminus\{0\}}^3$ is still valid in codon-to-codon stochastic case.

However, for stochastic cases, we cannot reach a global minimum because of the randomness caused by mutagens. Since, in usual cases, there exists no stationary global minimum, we need to define error tolerance ϵ , that is, if $|J_q(\psi) - J_{q-1}(\bar{\psi})| \leq \epsilon$ with the same $\{x_N^d\}$, then we can stop at $J_q(x_q)$. Otherwise, we need to proceed to calculate $J_{q-2}(\psi)$. The value of ϵ is decided based on the doctors, experience. Obviously,

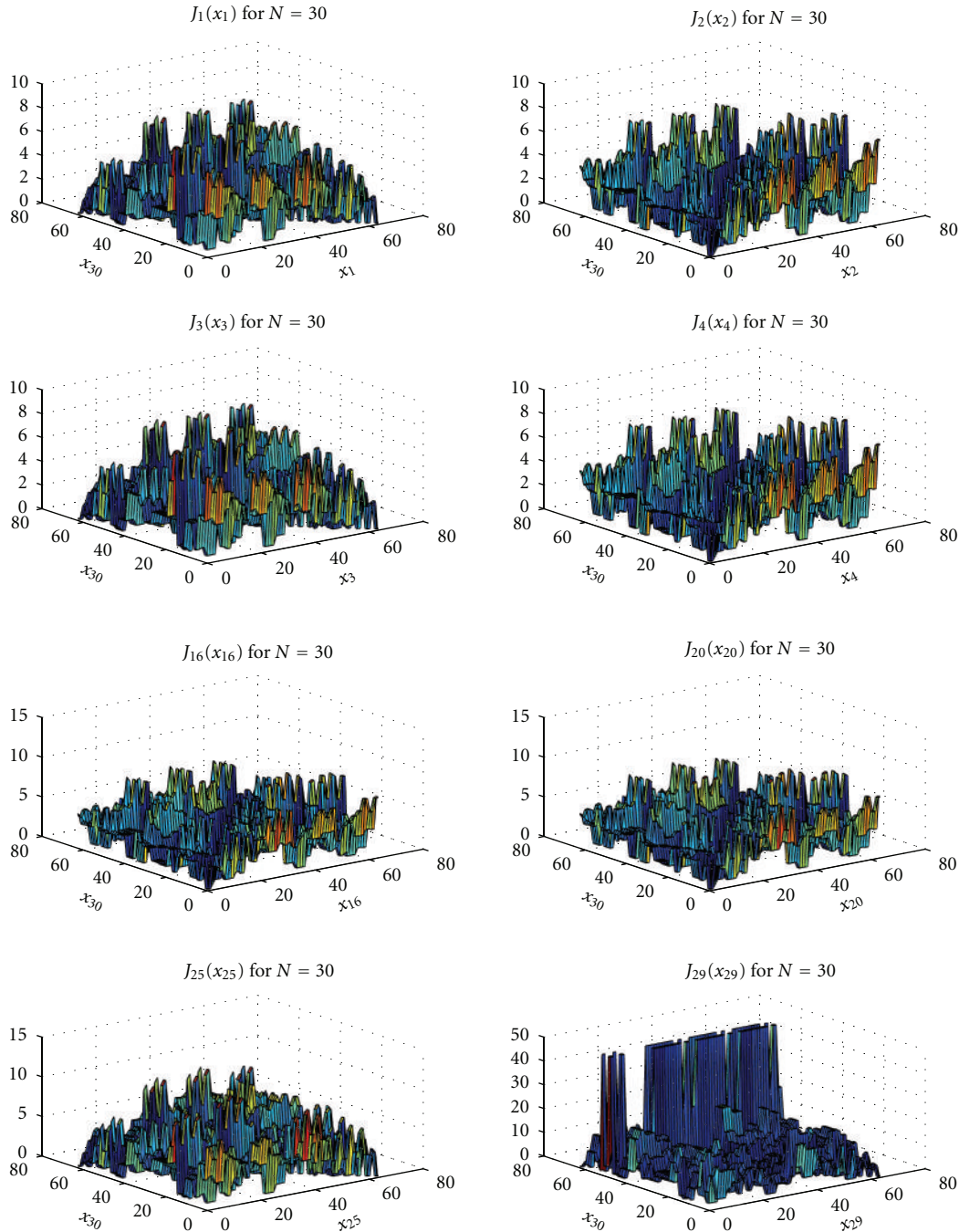


FIGURE 12: Graphical representation of $J_q(x_q)$, $q = 0, 1, 2, 3, 15, 19, 24, 28$ for codon-to-codon stochastic mutations, with $\alpha_{t_1} = 0.5\chi$, $N = 29$.

the smaller ϵ is, the better treatment plan is. However, we can still observe Figures 10, 11, and 12 to conclude that J_0 and J_2 are almost of the same shape in all three different parameter assignments. Higher dimensional optimal control problems, gene-to-gene stochastic mutations, can be solved as a series of cascade codon-to-codon stochastic problems.

4. Conclusion

In this paper, we present a mathematical model to deal with mutations in the process of DNA replication in the

view of control systems. Different from the existing models, our model is constructed directly from the basic biological theories, the central dogma in molecular biology, and the complementary base pair for DNA molecules with double helix structure. It precisely describes how the induced mutations affect the targeted DNA segments at molecular level. It provides instrumental information of molecule interactions in gene mutation for biologists and doctors to gain a better understanding of cellular and tissue level systems' behavior. Our model is adaptive to point and multisite, deterministic and stochastic mutations. Though we emphasize that we

target at induced mutations during the process of DNA replication in our work, this model can be extended to other biological processes at molecular level, such as transcription process and DNA brokage.

In our optimal control problem, the objective function includes two factors: the risk/cost of applying mutagens and the off-trajectory penalty. Under optimal control policy, the summation of those two factors are minimized, by dynamic programming, to propose a low-risk treatment plan. We define the distance reference following the chemical and physical properties of amino acids, representing the penalty. Our objective is to drive the system from given initial state to the final desired set generated by the final desired state at the lowest cost. We define the final desired set since redundancy in genetic codes gives us additional options of final desired state to further lower the cost. Dynamic programming algorithm ensures the optimality of the solution. We also discuss three different small-scale system, and show the simulation results of examples. The optimal control problems of base-to-base deterministic mutations and codon-to-codon deterministic mutations are of theoretical importance. As shown in the propositions, the global optimal can be reached within finite steps. If the step limit is larger than the number of steps that global optimal can achieve, then we have some flexibility in our treatment plan. In addition, there exist multiple optimal paths with the same total cost, given the initial state and the final desired set. The optimal control problem of codon-to-codon stochastic mutations is of practical importance, since codon is the basic component forming long DNA sequences. The step limit N is decided by doctors according to patients' conditions, and the treatment plan is made according to the initial state, the final desired set, and the step limit. Since the doctors constantly take measurement to see the result of treatments at current stage, the treatment plan is updated accordingly. Solving codon-to-codon stochastic optimal control problem is a key step to realize the optimal control to gene-to-gene stochastic mutations in the real world.

Our work contributes to several aspects in systems biology. The optimal control sequences generated by dynamic programming make it possible for biologists and doctors to mutate certain sections of genes on purpose in laboratory, at a relative low cost and low risk, which is an essential step to identify the functional units, to examine the interactions among different segments, and to find healthy, harmful, and lethal nucleotide sequences. All those results are beneficial in gene network construction. In addition, the fundamental details of gene mutations at molecular level help biologists to elaborate on the biological theories at the cellular and tissue levels, such as the theory of evolution. Moreover, by our method, biologists can distinguish the harmful and beneficial mutations and induce beneficial mutations during the evolution process in a proper way, which greatly helps to save rare species in danger. Furthermore, our solution to the optimal control problem proposed provides a new medical intervention to genetic diseases. Compared to the existing gene therapy, treatments by mutagens are safer by avoiding the side effect of virus infection. Lastly, our work also contributes to the construction of DNA computers.

Calculation errors, the mispairings in the process of two single-stranded DNA segments, can be corrected at lowest cost by applying a correct mutagen sequence.

Further work can be done by extending codon-to-codon stochastic optimal control problem to gene-to-gene stochastic mutations. The distance reference between DNA segment with equal length can be defined as a weighted sum of the distance references between codons. Since certain combinations of amino acids are lethal, those high-risk states should be avoided. This goal can be achieved by either defining a preset trajectory or eliminating high-risk sequences in the state space. Another possibility is to examine system's behavior under noisy measurements. Under this situation, the spontaneous mutations can be modeled as an additional random factor in our state updating equations, and another random noises should be added to the output equation.

References

- [1] N. Wiener, *Cybernetics*, J. Wiley, 1948.
- [2] H. Kitano, "Systems biology: a brief overview," *Science*, vol. 295, no. 5560, pp. 1662–1664, 2002.
- [3] R. J. Tanaka, H. Okano, and H. Kimura, "Mathematical description of gene regulatory units," *Biophysical Journal*, vol. 91, no. 4, pp. 1235–1247, 2006.
- [4] N. Yildirim and M. C. Mackey, "Feedback regulation in the lactose operon: a mathematical modeling study and comparison with experimental data," *Biophysical Journal*, vol. 84, no. 5, pp. 2841–2851, 2003.
- [5] R. Yang, T. J. Tarn, and M. Zhang, "Data-driven feedforward control for electroporation mediated gene delivery in gene therapy," *IEEE Transactions on Control Systems Technology*, vol. 18, no. 4, Article ID 5286233, pp. 935–943, 2010.
- [6] J. Feng and H. C. Tuckwell, "Optimal control of neuronal activity," *Physical Review Letters*, vol. 91, no. 1, pp. 018101/1–018101/4, 2003.
- [7] J. Moehlis, E. Shea-Brown, and H. Rabitz, "Optimal inputs for phase models of spiking neurons," *Journal of Computational and Nonlinear Dynamics*, vol. 1, no. 4, pp. 358–367, 2006.
- [8] I. Lentacker, B. Geers, J. Demeester, S. C. De Smedt, and N. N. Sanders, "Design and evaluation of doxorubicin-containing microbubbles for ultrasound-triggered doxorubicin delivery: cytotoxicity and mechanisms involved," *Molecular Therapy*, vol. 18, no. 1, pp. 101–108, 2010.
- [9] V. Reinke, *Germline Genomics*, WormBook, 2006.
- [10] H. De Jong, "Modeling and simulation of genetic regulatory systems: a literature review," *Journal of Computational Biology*, vol. 9, no. 1, pp. 67–103, 2002.
- [11] H. Matsuno, A. Doi, M. Nagasaki, and S. Miyano, "Hybrid Petri net representation of gene regulatory network," in *Pacific Symposium on Biocomputing*, vol. 5, p. 87, 2000.
- [12] J. Collado-Vides, "A transformational-grammar approach to the study of the regulation of gene expression," *Journal of Theoretical Biology*, vol. 136, no. 4, pp. 403–425, 1989.
- [13] J. Collado-Vides, R. M. Gutiérrez-Ríos, and G. Bel-Enguix, "Networks of transcriptional regulation encoded in a grammatical model," *BioSystems*, vol. 47, no. 1-2, pp. 103–118, 1998.
- [14] A. Regev, W. Silverman, and E. Shapiro, "Representation and simulation of biochemical processes using the pi-calculus

- process algebra,” in *Pacific Symposium on Biocomputing*, vol. 6, pp. 459–470, 2001.
- [15] E. M. Ozbudak, M. Thattai, H. H. Lim, B. I. Shraiman, and A. Van Oudenaarden, “Multistability in the lactose utilization network of *Escherichia coli*,” *Nature*, vol. 427, no. 6976, pp. 737–740, 2004.
- [16] M. Santillán and M. C. Mackey, “Influence of catabolite repression and inducer exclusion on the bistable behavior of the *iac* operon,” *Biophysical Journal*, vol. 86, no. 3, pp. 1282–1292, 2004.
- [17] Y. Setty, A. E. Mayo, M. G. Surette, and U. Alon, “Detailed map of a cis-regulatory input function,” *Proceedings of the National Academy of Sciences of the United States of America*, vol. 100, no. 13, pp. 7702–7707, 2003.
- [18] R. J. Tanaka and H. Kimura, “Mathematical classification of regulatory logics for compound environmental changes,” *Journal of Theoretical Biology*, vol. 251, no. 2, pp. 363–379, 2008.
- [19] L. Glass and S. A. Kauffman, “The logical analysis of continuous, non-linear biochemical control networks,” *Journal of Theoretical Biology*, vol. 39, no. 1, pp. 103–129, 1973.
- [20] R. Thomas, “Boolean formalization of genetic control circuits,” *Journal of Theoretical Biology*, vol. 42, no. 3, pp. 563–585, 1973.
- [21] L. M. Adleman, “Molecular computation of solutions to combinatorial problems,” *Science*, vol. 266, no. 5187, pp. 1021–1024, 1994.
- [22] M. Zhang, M. X. Cheng, and T. J. Tarn, “A mathematical formulation of DNA computation,” *IEEE Transactions on Nanobioscience*, vol. 5, no. 1, pp. 32–40, 2006.
- [23] F. H. Crick, “On protein synthesis,” *Symposia of the Society for Experimental Biology*, vol. 12, pp. 138–163, 1958.
- [24] F. Crick, “Central dogma of molecular biology,” *Nature*, vol. 227, no. 5258, pp. 561–563, 1970.
- [25] J. D. Watson and F. H. C. Crick, “Molecular structure of nucleic acids: a structure for deoxyribose nucleic acid,” *Nature*, vol. 171, no. 4356, pp. 737–738, 1953.
- [26] J. D. Watson and F. H. C. Crick, “A structure for deoxyribose nucleic acid,” in *A Century of Nature: Twenty-One Discoveries that Changed Science and the World*, p. 82, University Of Chicago Press, 2003.
- [27] S. D. McCulloch and T. A. Kunkel, “The fidelity of DNA synthesis by eukaryotic replicative and translesion synthesis polymerases,” *Cell Research*, vol. 18, no. 1, pp. 148–161, 2008.
- [28] E. C. Friedberg, G. C. Walker, and W. Siede, *DNA Repair and Mutagenesis*, ASM Press, 1995.
- [29] T. Strachan and A. P. Read, *Human Molecular Genetics*, vol. 3, Garland Science, 2004.
- [30] D. Hristu-Varsakelis and W. S. Levine, *Handbook of Networked and Embedded Control Systems*, Birkhauser, 2005.
- [31] D. P. Bertsekas, *Dynamic Programming and Optimal Control*, vol. 1,2, Athena Scientific, 1995.