



Missouri University of Science and Technology  
**Scholars' Mine**

---

Computer Science Faculty Research & Creative Works

Computer Science

---

26 Jun 2006

## Genetic Code Based Coding and Mathematical Formulation for DNA Computation

Mingjun Zhang

Maggie Xiaoyan Cheng

*Missouri University of Science and Technology*, [chengm@mst.edu](mailto:chengm@mst.edu)

Tzyh-Jong Tarn

Follow this and additional works at: [https://scholarsmine.mst.edu/comsci\\_facwork](https://scholarsmine.mst.edu/comsci_facwork)

 Part of the [Computer Sciences Commons](#)

---

### Recommended Citation

M. Zhang et al., "Genetic Code Based Coding and Mathematical Formulation for DNA Computation," *Proceedings of the 2006 IEEE International Conference on Robotics and Automation*, Institute of Electrical and Electronics Engineers (IEEE), Jun 2006.

The definitive version is available at <https://doi.org/10.1109/ROBOT.2006.1642256>

This Article - Conference proceedings is brought to you for free and open access by Scholars' Mine. It has been accepted for inclusion in Computer Science Faculty Research & Creative Works by an authorized administrator of Scholars' Mine. This work is protected by U. S. Copyright Law. Unauthorized use including reproduction for redistribution requires the permission of the copyright holder. For more information, please contact [scholarsmine@mst.edu](mailto:scholarsmine@mst.edu).

# Genetic Code Based Coding and Mathematical Formulation for DNA Computation

Mingjun Zhang  
Agilent Technologies  
Email: mingjunzhang@ieee.org

Maggie X. Cheng  
University of Missouri  
Email: chengm@umr.edu

T. J. Tarn  
Washington University  
Email: tarn@wuauto.wustl.edu

**Abstract**—DNA computation is to use DNA molecules for information storing and processing. Challenges currently faced by DNA computation are (1) lack of theoretical computational models for applications, and (2) high error rate for implementation. This paper attempts to address these problems from genetic coding and mathematical modeling aspects. The proposed genetic coding approach provides a promising alternative to reduce high error rate. The mathematical formulation lays down groundwork for studying theoretical aspects of DNA computation.

## I. INTRODUCTION

In the late 1950's, the Nobel laureate Richard Feynman first introduced the idea of computation at a molecular level. In 1994, the concept of DNA computation was demonstrated using experiments to solve a directed Hamiltonian Path Problem (HPP) by Adleman [1]. Since then, the possibility of DNA computation has attracted many researchers' attention.

DNA computation is to use DNA molecules for information storing and processing by encoding and interpreting DNA molecules in suspended solutions before and after DNA complementary binding reactions. The central idea of DNA computation is the Watson-Crick model of DNA structure, which specifies complementary binding properties of DNA molecules. DNA computation involves to use single-stranded DNA segments to code the problem, let the single-stranded DNA segments react in test tubes or substrate surfaces, and then to find DNA binding strands and interpret the results by applying bio-molecular techniques.

DNA computation is attractive mainly for three reasons. First, the computation realizes fast parallel information processing. Second, the process is remarkably energy efficient. Finally, DNA molecules have very high storing capacity.

Unfortunately, DNA computation currently is too error-prone to achieve its great potential. Many ideas of DNA computation assume a zero error rate. In reality, errors appear at every stage. In [1], [7], the problem of high error rates was identified as the most challenging problem for the success of DNA computation. High error resistant method is needed for DNA computation. One open question is whether the error rates in DNA manipulations can be adequately controlled [4], [8]. Some algorithms have been proposed to handle a few of the apparently crippling errors. Paper [6] proposed a surface-based DNA computation algorithm to solve the minimal set cover problem. The technique decreases errors caused by potential DNA strand lost by affixing the DNA onto a silicon

surface. In [3], a DNA computation model has been developed that uses dynamic programming and large size of memory available to DNA computers. The goal is to reduce error rates by increasing DNA strands. A more thorough study of decreasing error rates can be found in [4], where methods for making volume decreasing algorithms (the number of strands decreases as the algorithm executes) more resistant to certain types of errors are proposed. One effort in the paper is to convert the decreasing volume problem to a constant volume problem (the number of strands remains the same throughout the computation). The basic idea is to add DNA strand redundancy by increasing solution volume. The technique requires to increase steps of operations and cannot be applied to an algorithm that has constant volume to begin with. The other effort proposed in the paper is to reduce the false negative error rate in the bead separation procedure by double encoding DNA bases. The idea is to have each DNA-encoded base appear twice in separate locations in the strand to increase the possibility of being extracted. However, it is still not clear at present stage whether error rates can be reduced sufficiently to allow a general-purpose DNA computation.

This paper will propose a genetic code based approach to reduce error rate for DNA computation. The genetic code has great quality assurance, because of its redundancy. Some mis-binding results in no change in the coding words.

To better understand theoretical aspects of DNA computation, a mathematical model of DNA computation is also proposed. The model is useful to apply mathematical tools to solve DNA computation problems. Based on the formulation, character-based DNA computation is converted into a numerical computation problem. Propositions based on the formulation are also presented.

## II. GENETIC CODE BASED DNA COMPUTATION

Errors in DNA computation usually come from (1) DNA strand extraction, (2) Random errors, or (3) PCR errors. These errors can be reduced by either designing high error resistant coding approach or developing better molecular techniques for later DNA strand extraction. To avoid difficulty in reducing errors at later stages, a method to address the problem at the early phase of DNA computation is preferred. The proposed genetic code based approach targets earlier coding stage of DNA computation.

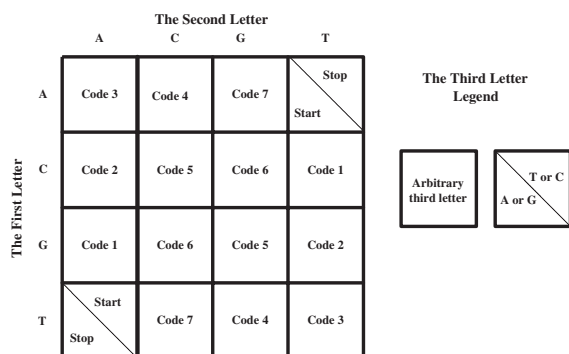


Fig. 1. Coding set for DNA computation

### A. Genetic code based DNA computation

If a DNA computation problem is coded using redundant genetic codes, the code will be highly error resistant and the error rates will be low. Different from biological systems that are highly diverging and require a large set of DNA codons to code various genetic information, the set of “codons” for DNA computation will be small. To take advantage of chemical properties of the molecular structures and obtain a reduced set of “codons” (called coding set for DNA computation), we will reduce the regular codons of biological systems to obtain smaller coding sets. One additional concern is that all coding sets should be closed, i.e., elements of a coding set and their complementaries (anti-codon) are within the same coding set. This is based on the concern that DNA computation does not use the same mechanism as biological systems for recognizing DNA strands. We expect to have a robust coding system. Second, the triple code mechanism will still be used. Since there may be biochemistry and stability reasons for the triple codes, though it is not completely clear at the present time. The following coding sets are proposed as shown in Fig. 1.

- Start coding set: ATG, ATA, TAC, or TAT.
- Set 1: GAA, GAG, GAT, GAC, CTT, CTC, CTA, CTG.
- Set 2: CAA, CAG, CAT, CAC, GTT, GTC, GTA, GTG.
- Set 3: AAA, AAG, AAC, AAT, TTT, TTC, TTG, TTA.
- Set 4: ACC, ACA, ACG, ACT, TGG, TGT, TGC, TGA.
- Set 5: CCT, CCC, CCA, CCG, GGA, GGG, GGT, GGC.
- Set 6: GCT, GCC, GCA, GCG, CGA, CGG, CGT, CGC.
- Set 7: TCT, TCC, TCA, TCG, AGT, AGC, AGA, AGG.
- Stop coding set: TAA, TAG, ATT, or ATC.

All the above coding sets are closed with respect to the DNA complementary operation. In addition, the above coding scheme allows significant amount of overlapping, which leads to highly error resistant. The coding sets are enough to code significantly large problems by varying the length of wording. By coding this way, many DNA base mutations may not cause changes in word meaning for DNA computation.

### B. Genetic Code Based DNA Computation to Solve the Hamiltonian Path Problem (HPP)

The HPP problem is to find (if there is) a Hamiltonian path for a given graph. A Hamiltonian path is a sequence of

compatible one-way edges of a directed graph that begins and ends at a specified vertex and enters every other vertex exactly once. Known algorithms for this problem have exponential worst-case complexity. The problem has been proved to be NP-complete.

Assume that the graph has  $n > 0$  vertices (cities) and  $i$  is the index of a vertex [1]. The following steps are usually followed.

- 1) Associate the start vertex  $i = 1$  with one code (three bases) from the Start coding set. Associate the end vertex  $i = n$  with a code from the Stop coding set.
- 2) Associate each of the other vertices  $i$  ( $1 < i < n$ ) with a  $3m$ -mer sequence generated by  $m$  codes (usually an even number to keep the left and right side edges of a city with equal length), and denote it by  $O_i$ . For each edge  $i \rightarrow j$ , an oligonucleotide  $O_{i \rightarrow j}$  is created, which is the 3'  $3m/2$ -mer of  $O_i$  followed by the 5'  $3m/2$ -mer of  $O_j$ .
- 3) Keep only those paths that begin with codes from the Start coding set, and end with codes from the Stop coding set. This can be done by PCR amplifying products of the Step 1) using primers starting with codes from the start coding set or the Stop coding set.
- 4) Keep only those paths that enter exactly  $n$  vertices. The product of Step 2) is run on an agarose gel and the  $3m$  base pair band (corresponding to dsDNA encoding paths entering exactly  $n$  vertices) is excised and soaked to extract DNA.
- 5) Keep only those paths that enter all vertices of the graph exactly once. This can be done by first generating single-stranded DNA sequences from the dsDNA product of Step 4) and then incubating the ssDNA with  $\bar{O}_2$  conjugated to magnetic beads. Only those ssDNA molecules containing  $\bar{O}_2$  (and hence encoded paths which enter vertex 2 at least once) anneal to the bound  $\bar{O}_2$  and were retained. The process repeat successively with  $\bar{O}_3, \bar{O}_4, \dots, \bar{O}_{n-1}$  and  $\bar{O}_n$ .
- 6) The remaining DNA sequences (paths) in the test tube represent solutions.

To further illustrate the idea, consider an  $n = 7$  vertex HPP graph as given in [1], we use 12 DNA bases (4 codes from the coding sets) to uniquely code each of the 7 cities. Fig. 2 shows the final coding graph for each edge and vertex. The next step is to apply DNA molecular techniques to obtain biological solutions. Similar to the work in [1], same conclusion will be obtained. Compared with the brute-force approach used for solving HPP problem in [1], the proposed approach can save 40% of the coding characters, which will eventually speed up the problem solving time. More importantly, the proposed approach has a much lower error rate.

To demonstrate that the error rates can be reduced, assume that the following stochastic transfer matrix of DNA sequences

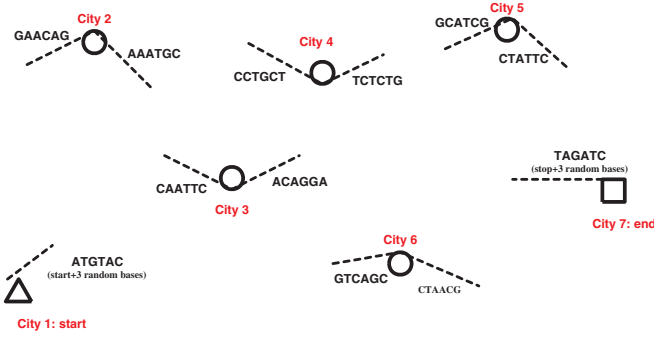


Fig. 2. Coding for the seven cities of the HPP problem

are held for DNA bases  $\{A, T, G, C\}$ .

$$\Gamma = \begin{bmatrix} 0.9990 & 0.0003 & 0.0004 & 0.0003 \\ 0.0003 & 0.9990 & 0.0003 & 0.0004 \\ 0.0004 & 0.0003 & 0.9990 & 0.0003 \\ 0.0003 & 0.0004 & 0.0003 & 0.9990 \end{bmatrix} \quad (1)$$

After one transformation, a DNA sequence  $X = GATCAG$  coded by codons from the coding sets 1 and 2 can be expressed in numerical values as

$$\Gamma_1 = \begin{bmatrix} 0.0004 & 0.0003 & 0.9990 & 0.0003 \\ 0.9990 & 0.0003 & 0.0004 & 0.0003 \\ 0.0003 & 0.9990 & 0.0003 & 0.0004 \\ 0.0003 & 0.0004 & 0.0003 & 0.9990 \\ 0.9990 & 0.0003 & 0.0004 & 0.0003 \\ 0.0004 & 0.0003 & 0.9990 & 0.0003 \end{bmatrix} \quad (2)$$

where each element shows the probability that the corresponding base is obtained after molecular manipulation.

It can be concluded that the original sequence  $X = GATCAG$  is still preserved very well. After about 1000 transformations, the sequence turns to be ambiguous as follows.

$$\Gamma_2 = \begin{bmatrix} 0.2021 & 0.1748 & 0.4484 & 0.1748 \\ 0.4484 & 0.1748 & 0.2021 & 0.1748 \\ 0.1748 & 0.4484 & 0.1748 & 0.2021 \\ 0.1748 & 0.2021 & 0.1748 & 0.4484 \\ 0.4484 & 0.1748 & 0.2021 & 0.1748 \\ 0.2021 & 0.1748 & 0.4484 & 0.1748 \end{bmatrix} \quad (3)$$

It seems that the original sequence may be turned into a different format. However, by following the proposed coding scheme, the original word coded in  $X$  is still well preserved. This is because the redundancy of the coding scheme. There are  $8 \times 8 = 64$  different combinations of codings for the six-base DNA sequence  $X$ , and they are all coded for the same word. Even mutation occurs after multiple transformations, the original information remains intact. This is the beauty of the genetic coding based method. However, if a brute force fixed length coding approach is used, i.e., each combination represents one scheme only, the original coding meaning cannot still be kept intact. If any other redundancy coding approach that has less than 64 combinations to represent one coding scheme is used, the proposed genetic coding scheme still performs the best.

### III. A MATHEMATICAL FORMULATION OF DNA COMPUTATION

Define the following notations:

- Let  $X = x_i x_{i+1} \dots x_j$  and  $Y = y_i y_{i+1} \dots y_j$  represent single-stranded DNA segments, where  $i, j \in N$  and  $i \leq j$ .  $N$  is the natural number.  $x_i, y_i \in \{A, T, G, C\}$ .
- The complementary sequence of  $X$  is defined as  $\bar{X}$ .
- Let  $T_i$  represents the  $i$ -th test tube, where  $i \in N$ .  $T_i(+X)$  means that the test tube containing DNA segment  $X$ .  $T_i(-X)$  means that the test tube  $T_i$  does not contain DNA segment  $X$ .

DNA computation involves many bio-molecular operations including hybridization, separation, cutting and pasting DNA strands at desired locations. These operations can be generalized at DNA strand level as follows, where “ $\rightarrow$ ” represents “the result of” reactions from the left hand side operations.

- Ligation: plus “+” operation. Ligation concatenates segments of DNA. Biochemically, it is often invoked after an annealing operation. For two single-stranded DNA segments  $X$  and  $Y$ , a ligation operation can be expressed as “ $X + Y \rightarrow [XY]$ ”, where  $[XY]$  represents a newly created single-stranded DNA segment.
- Cut: minus “-” operation. Restriction enzymes can cut a strand of DNA at a specific address. Some restriction enzymes only cleave single-stranded DNA, while others only cleave double-stranded DNA segments. If a single-stranded DNA  $X$  is cut at position  $n$  from the 3’ end of a DNA segment, the process can be expressed as “ $-X(n) \rightarrow Y + Z$ ”, where  $Y$  and  $Z$  are newly generated single-stranded DNA segments and  $Y$  has a length of  $n - 1$ .
- Hybridization: multiplication “ $\bullet$ ” operation. It is a process when single-stranded complementary DNA segments spontaneously form a double-stranded DNA. For single-stranded DNA  $X$  and  $\bar{X}$ , the binding process can be described as “ $X \bullet \bar{X} \rightarrow (X\bar{X})$ ”, where  $(X\bar{X})$  is a newly created double-stranded DNA. DNA strands enclosed by brackets “(” and “)” are double-stranded DNA and cannot be bonded with other strands unless further melting operation is applied.
- Melting: division “ $\backslash$ ” operation. Heating can be selectively used to melt apart short double-stranded DNA segments while leaving longer double-stranded segment intact. For example, “ $\backslash(X\bar{X}) \rightarrow X + \bar{X}$ ” means melting the double-stranded DNA  $(X\bar{X})$  as two complementary single-stranded DNA segments  $X$  and  $\bar{X}$ .

All the above operations are single step DNA molecule reactions, which are extremely fast compared with conventional silicon computation. In addition to strand level operations, the DNA computation may use the following operations at test tube level, which are performed using sets of DNA segments.

- Merge: union operation “ $\cup$ ”. The operation means two test tubes can be combined, usually by pouring one test tube into the other. For example, “ $T_1 \cup T_2 \rightarrow T$ ” means

melting two test tubes  $T_1$  and  $T_2$  together to produce a new test tube  $T$ .

- Separation or extraction: difference operation “-”. The expression “ $-(T_i, X) \rightarrow T_j(+X) \cup T_k(-X)$ ” represents a separation operation applied to the test tube  $T_i$  on DNA segment  $X$ . The operation produces two test tubes, where the tube  $T_j$  contains a string  $X$  and the tube  $T_k$  does not contain the DNA segment  $X$ , where  $i, j, k \in N$ . Either  $T_j$  and  $T_k$  could be empty set  $\phi$ . This step is done using gel electrophoresis. It requires the DNA strands to be extracted from the gel once the DNA segments of different length have been identified.
- Amplification: product operation “ $\times$ ”. Given a test tube containing DNA strands, the operation is to make multiple copies of a subset of the strands presented. Copies are usually made with PCR. For example, “ $\times T \rightarrow T_1 \cup T_2$ ” means two test tubes  $T_1$  and  $T_2$  containing copies of a subset of DNA strands are produced from the test tube  $T$ .
- Detection: question operation “?”. This operation means that gel electrophoresis is applied to see if anything of the appropriate length is left within a test tube after PCR amplification. For example, “ $T?X \rightarrow \text{True}$ ” means that the test tube  $T$  contains at least one string  $X$ . Otherwise, “ $T?X \rightarrow \text{False}$ ”.
- Destroy: intersection operation “ $\cap$ ” with an empty set. Subsets of strands can be systematically destroyed or “digested” by enzymes that preferentially break apart nucleotides in either single- or double-stranded DNA segments. The process can be expressed as “ $T \cap \phi \rightarrow \phi$ ”.

To further investigate DNA computation as a computational problem, the following concepts are developed.

#### A. Conversion of character-based DNA sequences to numerical sequences

Three original methods are proposed to convert character-based DNA sequences into numerical sequences. One method is to use complex numbers. The second method is to use integer numbers. The third method is to convert DNA sequence into vectors.

1) *Complex number representation*: Define a function  $f(x) : \{A, T, G, C\} \rightarrow \{1, -1, i, -i\}$  as

$$f(x) = \begin{cases} 1, & x = A; \\ -1, & x = T; \\ i, & x = G; \\ -i, & x = C. \end{cases} \quad (4)$$

where  $x$  is one of the four nucleotides.

The complementary base of each DNA base  $x$  can then be calculated by the following inverse function

$$\bar{x} = f^{-1}(-f(x)) = \begin{cases} T, & x = A; \\ G, & x = C; \\ C, & x = G; \\ A, & x = T. \end{cases} \quad (5)$$

By definitions (4) and (5), complementary DNA sequences (either numerical or character-based) can be easily obtained.

This means only single-stranded DNA segments need to be specified. The complementary strands can be easily generated using the above functions in either character-based or numerical format.

2) *Integer number representation*: DNA bases may be mapped as integer numbers as well. Define a function  $f(x) : \{A, T, G, C\} \rightarrow \{0, 1, 2, 3\}$  as

$$f(x) = \begin{cases} 0, & x = A; \\ 1, & x = C; \\ 2, & x = G; \\ 3, & x = T. \end{cases} \quad (6)$$

Similarly, the complementary base of  $x$  can be determined by the following inverse function

$$\bar{x} = f^{-1}(\{3\} - f(x)) = \begin{cases} T, & x = A; \\ G, & x = C; \\ C, & x = G; \\ A, & x = T. \end{cases} \quad (7)$$

where  $\{3\}$  represents an appropriate finite length sequence consisting of multiple copies of integer 3. The numerical calculation can then be conducted base by base. For example, the numerical sequence of a DNA segment  $X = AGGCAT$  is  $f(X) = f(AGGCAT) = 022103$ . The complementary segment of  $X$  can be easily obtained as  $\bar{X} = f^{-1}(\{3\} - f(X)) = f^{-1}(311230) = TCCGTA$ .

3) *Vector representation*: In vector space analysis, numerical value based DNA sequences can be expressed as rows of a matrix. Addition of such kinds of matrices can be regarded as DNA hybridization process. Scaler multiplication produces multiple copies of the sequences in a test tube. Consider the four DNA bases  $\{A, T, G, C\}^T$  as a vector, any DNA strand  $X = x_1x_2\dots x_n, n \in N$ , can then be expressed as a vector by a transfer matrix  $\Pi$  as

$$X = \Pi \begin{bmatrix} A \\ T \\ G \\ C \end{bmatrix} = \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ \vdots & \vdots & \vdots & \vdots \\ p_{n1} & p_{n2} & p_{n3} & p_{n4} \end{bmatrix} \begin{bmatrix} A \\ T \\ G \\ C \end{bmatrix} \quad (8)$$

where  $\sum_{j=1}^4 p_{ij} = 1, \forall i \in N$ . Specifically,

$$\begin{aligned} p_{i1} &= \begin{cases} 1, & x_i = A. \\ 0, & \text{otherwise.} \end{cases}, p_{i2} = \begin{cases} 1, & x_i = T. \\ 0, & \text{otherwise.} \end{cases} \\ p_{i3} &= \begin{cases} 1, & x_i = G. \\ 0, & \text{otherwise.} \end{cases}, p_{i4} = \begin{cases} 1, & x_i = C. \\ 0, & \text{otherwise.} \end{cases} \end{aligned} \quad (9)$$

In the above definition, each row of the matrix  $\Pi$  represents one DNA base. The complementary sequence of  $\bar{X}$  can then be obtained by simply swapping column one with column two, and column three with column four as follows

$$\bar{X} = \begin{bmatrix} p_{12} & p_{11} & p_{14} & p_{13} \\ p_{22} & p_{21} & p_{24} & p_{23} \\ \vdots & \vdots & \vdots & \vdots \\ p_{n2} & p_{n1} & p_{n4} & p_{n3} \end{bmatrix} \begin{bmatrix} A \\ T \\ G \\ C \end{bmatrix} \quad (10)$$

For example, the transfer matrix of a single-stranded DNA  $X = ACGTGGATCT$  is shown in  $\Pi_1$ . The complementary sequence of  $X$  is  $\bar{X} = TGCACCTAGA$ , whose transfer matrix is given in  $\Pi_2$ .

$$\Pi_1 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{bmatrix}, \Pi_2 = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} \quad (11)$$

The above definitions (9) and (10) make it possible to define a DNA strand as an  $n \times 4$  matrix.

The DNA nucleotides may also be defined as a vector directly. For example,  $A = [1 \ 0]^T$ ,  $T = [-1 \ 0]^T$ ,  $G = [0 \ 1]^T$ , and  $C = [0 \ -1]^T$ . Then a DNA sequence can be expressed as a  $2 \times n$  matrix, where  $n$  is the number of bases for the DNA sequence. For example, a single-stranded DNA sequence  $X = GATCCAGT$  can be expressed as

$$\begin{bmatrix} 0 & 1 & -1 & 0 & 0 & 1 & 0 & -1 \\ 1 & 0 & 0 & -1 & -1 & 0 & 1 & 0 \end{bmatrix} \quad (12)$$

In a biological process, mutations often occur [5]. A stochastic transfer matrix  $\Gamma$  can be defined as follows to reflect this phenomena based on the above definitions (9) and (10).

$$\Gamma = \begin{bmatrix} \rho_{aa} & \rho_{at} & \rho_{ag} & \rho_{ac} \\ \rho_{ta} & \rho_{tt} & \rho_{tg} & \rho_{tc} \\ \rho_{ga} & \rho_{gt} & \rho_{gg} & \rho_{gc} \\ \rho_{ca} & \rho_{ct} & \rho_{cg} & \rho_{cc} \end{bmatrix} \quad (13)$$

where  $\rho_{ij}$  represents the probability of transformation from DNA base  $i$  to  $j$ , where  $i, j \in \{a, t, g, c\}$ . Obviously,  $\rho_{ia} + \rho_{it} + \rho_{ig} + \rho_{ic} = 1$ ,  $\forall i \in \{a, t, g, c\}$ .  $\rho_{ii}$  is the probability for correct transformation.

The inner product (same definition as in linear algebra) of a DNA sequence  $X$  can then be expressed as  $X^T X$ , which is a diagonal  $2 \times 2$  matrix. The first and the last elements in the matrix represent the number of bases in  $X$  from the set  $\{A, T\}$  or  $\{G, C\}$ , respectively.

Once a DNA sequence is converted into a numerical sequence, many interesting properties can be investigated. Next, some theoretical results are presented.

### B. Some theoretical results

By definition (4) and in viewing a DNA sequence as a vector in the format of (12), the following results are obtained. Proof of these results are straightforward.

**Proposition 2.1:** If the base-by-base plus operation of two equal length numerical value based DNA sequences results in a zero vector, then the two DNA sequences are complementary to each other.

**Proposition 2.2:** If the base-by-base plus operation of all numerical value based DNA sequences in different test tubes

results in a zero vector, then the hybridization by mixing the test tubes should be complete. A complete hybridization means all single-stranded DNA sequences find their complementaries.

**Proposition 2.3:** Under the definition (4), if the inner product of two equal length sequences is not a real number, then the two sequences are not complementary to each other. It can be further claimed that they are not complementary in  $G$  and  $C$ .

**Proposition 2.4:** Under the definition (4),  $\forall X, Y \in R^n$  ( $n$  represents the number of DNA bases in the strands), if  $X$  and  $Y$  have a complete hybridization and  $X^T Y = 0$ , then  $X$  and  $Y$  have equal number of DNA bases from  $\{G, C\}$  and  $\{A, T\}$ . Similarly, if a complete hybridization occurs, but  $X^T Y > 0$ , it means that there are more bases from  $\{G, C\}$  set than from  $\{A, T\}$  set. Otherwise, if complete hybridization occurs, but  $X^T Y < 0$ , it means that there are more bases from  $\{A, T\}$  set.

**Proposition 2.5:** Under the definition (4),  $\forall X, Y \in R^n$  ( $n$  represents number of DNA bases in the strands), we have  $\|X^T Y\| < \|X\| \|Y\|$ , where  $\|X^T Y\|$  represents the length of the DNA strand after hybridization.  $\|X\|$  and  $\|Y\|$  represent the length of single-stranded DNA segments. If a complete hybridization occurs,  $\|X^T Y\| = n$ . If none of the bases is hybridized, then  $\|X^T Y\| = 2n$ . This relationship is similar to the well-known Cauchy-Schwartz inequality in linear algebra.

To investigate properties under the above formulation (12) in vector space, the following definitions are proposed.

**Definition:** Equivalent transfer matrices. Since each single strand of a double-stranded DNA uniquely determines the other strand, each single-stranded DNA can be alternatively used to describe the same DNA double strand. Transfer matrices of a single-stranded DNA and its complementary are regarded as equivalent to each other. For example, the above  $\Pi_1$  in (11) is an equivalent transfer matrix of the  $\Pi_2$  in (??) expressed as  $\Pi_1 \Leftrightarrow \Pi_2$ , and vice versa. Two DNA sequences are complementary to each other, if and only if their transfer matrices are equivalent.

Note, two DNA transfer matrixes are equivalent, if and only if one matrix is the result of swapping column one with column two, and column three with column four of the other matrix.

**Proposition 2.6:** If DNA transfer matrices  $A \Leftrightarrow B$  and  $B \Leftrightarrow C$ , then  $A$  is the same as  $C$ .

**Definition:** Similar DNA sequences. Two equal-lengthed DNA sequences that have less than certain percent (usually 10% in practice) different bases in order are regarded as similar sequences. The binding results for similar sequences may be hard to be distinguished using current molecular techniques. It is advised not to use similar sequences to code different words in DNA computation.

**Proposition 2.7: Necessary condition for similar sequences.** Under the formulation (12), if two DNA sequences are similar, then the sum of all columns of the transfer matrices has numerical value variations less than pre-defined percent of the length of a single DNA sequence.

#### IV. APPLICATIONS

The above mathematical formulation of DNA computation may be used in the following applications.

##### A. Word design for DNA computation

In DNA computation, to reliably store and retrieve information in synthetic DNA strands, DNA word design is very important. DNA word design is to design sets of equal-lengthed words over the alphabet  $\{A, T, G, C\}$  satisfying certain constraints. The primary constraint is  $\forall X, Y \in \{\text{DNA words}\}$  at least  $d$  mismatches between  $X$  and  $Y$ , and between  $X^R$  and  $\bar{Y}$ , where  $X^R$  represents the reverse of  $X$  [9].

Based on the above formulation, the word design problem is equivalent to the following mathematical problem: choosing  $X$  and  $Y$  sequences from characters  $\{A, T, G, C\}$ , so that the number of non-zero bases from the base-by-base operation of  $f(X)$  plus  $f(Y)$ , and  $f(X^R)$  plus  $f(\bar{Y})$  are greater or equal to  $d$ . Based on the above discussion, solvability and upper or lower bounds of the DNA word design can be solved [4].

##### B. Natural DNA processing

The above formulation may also be used to process natural DNA for sequencing, fingerprinting and mutation detection. The idea is to first convert character-based DNA sequences into numerical sequences, then apply numerical computation techniques.

An example is multiple DNA sequence alignment. The proposed mathematical formulation may save significant amount of machine time, if the sequences are expressed and compared in the numerical domain.

For example, to detect long DNA sequence mutations as shown in equation (14), the proposed method can be applied first to convert the character-based sequences into numerical sequences as shown in (15). A numerical minus operation can then be conducted. If the final result is non-zero, this implies that there is mutation. The comparison is efficient by avoiding tedious character-based side-by-side comparison. This is called re-coding DNA. As pointed in [10], the re-coding has great potential application for DNA engineering applications. The formulation proposed in this paper provides an ideal mechanism for the re-coding process.

$$\begin{aligned} &ATTCCAGA \cdots GACCTTGAGT \\ &ATTCCATA \cdots GACCTCGAGT \end{aligned} \quad (14)$$

$$\begin{aligned} &03311020 \cdots 2011332023 \\ &03311030 \cdots 2011312023 \end{aligned} \quad (15)$$

where “ $\cdots$ ” may represent thousands or millions of DNA bases.

##### C. Combinatorial chemistry

The mathematical formulation may also be used in combinatorial chemistry for pseudo-enzyme design [2]. The goal is to create molecules with desired properties, which may be difficult or expensive by direct experimental studies. With the above mathematical formulation, the molecule design process

can be done in a numerical simulation mode. The beauty of this approach is that the tedious and expensive pseudo-molecule experimental design process may be avoided.

It is also interesting to note that most advanced numerical operations for conventional computers require combinations of a number of basic silicon computer operations. However, they may be completed by much less operations using DNA computation as shown in the mathematical operations. The above symbolic operation may be used for algorithm implementation in DNA computation.

#### V. DISCUSSION AND CONCLUSIONS

The high error rate is a major concern for DNA computation. This paper proposes a genetic code based DNA computation method to reduce the error rate. The idea is inspired by the genetic code of biological systems, except that the codon sets have been reduced. Redundancy in the genetic codes plays an important role in reducing error rates. Different from many methods proposed in the open literature, the proposed method does not require any additional bio-molecular techniques or steps, and it is not limited to any specific type of errors, which is a great advantage compared with other methods in the open literature.

A mathematical formulation of DNA computation and some propositions have also been presented in this paper. The formulation can convert a character-based DNA computation problem into a numerical value based computing problem. This will allow researchers to build theoretical framework for DNA computation, and analyze algorithmic as well as computational efficiency of DNA computation. Some potential problems may be further studied within the proposed mathematical framework, such as how can we increase DNA computational efficiency? what is the ultimate limit of the error-rate for DNA computation using genetic coding? what are the solvability and the upper or lower bounds for the DNA word design? We hope more discussions can be inspired in these regard.

#### REFERENCES

- [1] L. Adleman, Molecular Computation of Solutions to Combinatorial Problems, *Science*, vol.266, pp. 1021-1024, Nov. 1994.
- [2] D. Bartel, and J. Szostak, Isolation of new ribozymes from a large pool of random sequences. *Science*, vol. 261, pp. 1411-1418, September 1991.
- [3] E. Baum, and D. Boneh, Running dynamic programming algorithms on a DNA computer. *the 2nd DIMACS workshop on DNA based computers*, pp. 141-147, 1996.
- [4] D. Boneh, and R. Lipton, Making DNA computers error resistant. *Technical report, Princeton University*, CS-TR-491-495, 1996.
- [5] R. Durrett, Probability models for DNA sequence evolution. *New York : Springer*, 2002.
- [6] T. Eng, and B. Serridge, A Surface-Based DNA Algorithm for Minimal Set Cover. *the 3rd DIMACS workshop on DNA based computers*, pp. 74-82, June 1997.
- [7] R. Lipton, Using DNA to solve NP-Complete Problem. *Science*, vol. 268, pp. 542-545, April 1995.
- [8] C. C. Maley, DNA computation: theory, practice, and prospects, *Evolutionary computation*, vol. 6, no. 3, pp. 201-229, 1998.
- [9] A. Marathe, A. E. Condon, and R. M. Corn, On combinatorial DNA word design, *J. Computational Biology*, vol. 8, no. 3, pp. 201-220, 2001.
- [10] J. H. Reif. Paradigms for Biomolecular Computation, *Unconventional Models of Computation*, edited by C. S. Calude, J. Casti, and M. J. Dinneen, Springer-Verlag, New York, January 1998, pp 72-93.