

## Mathematical Analysis of DNA Curves

**R. Sengothai**

*Ph.D. Research Scholar in Mathematics*

*University of Madras Chennai*

### **ABSTRACT**

*Effective representation of DNA sequences is one of the important tasks in the study of genome sequences. In this paper, we propose a graphical representation of DNA sequences based on nucleotide ring structure. In the proposed representation, we convert DNA sequences into 16 dinucleotides on the surface of the hexagon so that it can preserve nucleotide's chemical property and positional information. Our approach can provide capability of efficient similarity comparison between DNA sequences and also high comparison accuracy. Furthermore, our approach satisfies uniqueness and no degeneracy of DNA sequences. In the experimental study, we use phylogeny analysis for evolutionary relationship among different species. Extensive performance study shows that the proposed method can give better performance than existing methods in comparison with the degree of similarity. In this paper, we converted DNA sequences into DNA curves without any loss of information and degeneracy.*

**Keywords:  $\beta$ -globin gene, DNA curve, hexagon, ring structure**

### **I. INTRODUCTION**

The rapid growth of biological sequences, such as of DNA, RNA, and protein, has demanded effective analysis methods for large biological sequences. Additionally, the analysis results are very helpful to biological researchers for predicting genes' structure and function, as well as similarity comparison between genes and different species.

For biological sequence analysis, two approaches have been mainly used: (i) sequence alignment method and (ii) non-sequence alignment method. The first approach obtains a degree of similarity between DNA sequences by comparing alignment scores of two sequences. This approach suffers from expensive computational cost as the length of sequences grows exponentially. The second approach analyzes DNA sequences by establishing a statistical model or a graphical representation model, or some machine learning model of DNA sequences. Recently, this approach is popularly studied due to the fact that it can give better accuracy and low computational overhead.

In the case of non-sequence alignment method, effective DNA sequence representation or feature selection from DNA sequences is essential for DNA sequence analysis, in areas such as gene prediction, similarity comparison between genes of different species, and finding gene structure and function. For this purpose, several graphical representations have been proposed according to chemical structures of 4 nucleotides, reflecting their distribution with different chemical structure and allowing numerical characterization. As for feature selection,

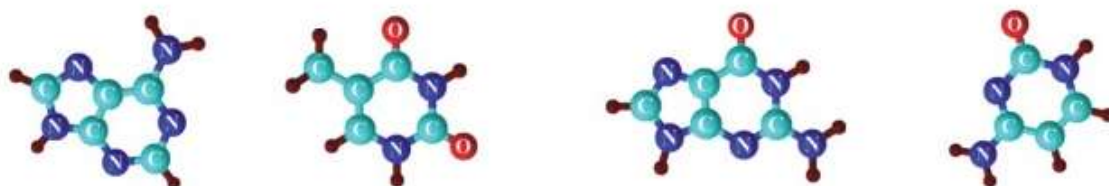
several machine learning techniques are effectively applied such as principal component analysis (PCA), neural network, and several classification models.

In this paper, we propose a graphical representation of DNA sequences based on nucleotide ring structure. In the proposed representation, we convert DNA sequences into 16 dinucleotides on the surface of the hexagon so that the nucleotide's chemical property and positional information is preserved. Our approach satisfies uniqueness and no degeneracy of DNA sequence is observed. It can also provide capability of efficient similarity comparison between DNA sequences in addition to high comparison accuracy. Extensive performance study shows that the proposed method can give better performance than existing methods in comparison with the degree of similarity.

## II. DNA SEQUENCE VISUALIZATION BY HEXAGONAL STRUCTURE

### 2.1 Chemical structure and classification of DNA bases

As stated previously, DNA sequences are the strings of four bases, that is, A, T, C and G. The core of these bases is heterocyclic organic compound, which forms ring in their chemical structure. Of them purines (A and G) have two rings while pyrimidines (C and T) have one. The chemical ring structures of those four bases are depicted in Figure 1.



a) Adenine (A)      b) Thymine (T)      c) Guanine (G)      d) Cytosine (C)

#### 1. Heterogenic cycle of four bases

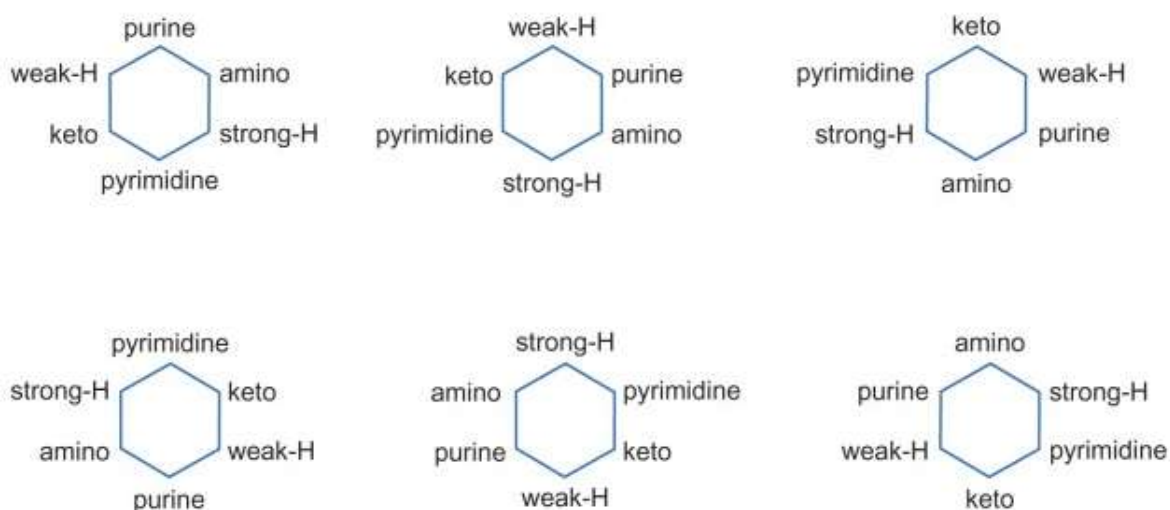
The element of these cycles are carbon and nitrogen, hydrogen and oxygen. In Figure 1, grey balls are carbon and blue balls are nitrogen. The hexagonal cycle has nitrogen in positions 1 and 3, and carbon in positions 2, 4, 5, and 6. Other than carbon and nitrogen, the bases have oxygen and hydrogen bonded with carbon and nitrogen in different number. Hence, the bases differ in molecular weight. The molecular weight of A, T, C and G are 135.13, 112.1, 111.1 and 151.13 respectively. Their ascending order in terms of molecular weight is C→T→A→G.

The bases also differ by heterogenic cycle, functionality, and their bonding with hydrogen. A and C fall into the amino category while G and T are in the keto group, based on their functionality. A and T are bonded by three hydrogen bonds, and hence are in strong-H group while G and C are in weak-H group as they are bonded by only two hydrogen bonds.

## 2.2 Proposed DNA encoding

The proposed encoding of dinucleotides for DNA sequence visualization is solely based on ring structure of DNA bases and their molecular weight. The bases are paired to make dinucleotides in such a way that their ascending/descending order in terms of molecular weight remains intact. The dinucleotides are placed on the 6 end of the heterocyclic hexagon as well as at the midpoint of each arm of the hexagon. The six dinucleotides which are placed on the 6 ends of the hexagon are in ascending order. The midpoint dinucleotides are positioned by descending order of molecular weight.

We place six ordered dinucleotides on opposite ends of the heterogenic cycle. The opposite ends are 1–4, 2–5, and 3–6. Any class (purine, pyrimidine; amino, keto; strong-H, weak-H) can be positioned at either end of the hexagon. Therefore, there are six possible combinations, as shown in Figure 2. The names of these combinations are Cycle 1, Cycle 2, Cycle 3, Cycle 4, Cycle 5, and Cycle 6 respectively.

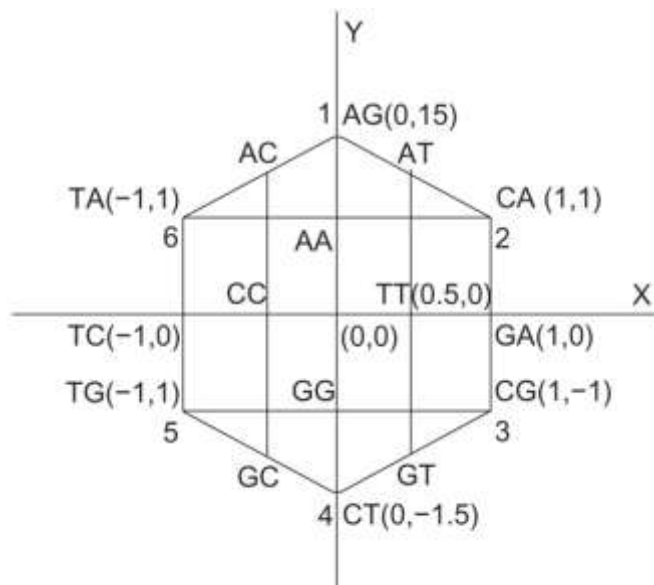


### 2. Six combinations of heterogenic cycle in 2D space

In Cycle 1, purines and pyrimidines are positioned at ends 1 and 4 of the hexagon respectively. A and G, the purines, form two dinucleotides: AG and GA. We keep AG on end 1 as it retains the sequence C→T→A→G. For the same reason, CT are placed on end 4, and CA (amino) and TG (keto) are placed on the 2 and 5 ends, respectively, and CG (strong-H) and TA (weak-H) are positioned on the 3 and 6 ends of the hexagon, respectively.

Conversely, midpoint of 2–3, 3–4 and 4–5 arms are determined by the following rule: take the uncommon nucleotides and form a dinucleotide with them such that descending order (G→A→T→C) of molecular weight prevails. As for example, the midpoint of 2–3 arms is GA because the commonality between CA and CG is C. So, G and A are uncommon. This rule is different for the midpoint of 5–6, 6–1 and 1–2 arms: take the common

nucleotide as well as the other which is not available on both ends. For example, the midpoint of 5–6 arms is TC because T is common between TA and TG, while C is neither in TA nor in TG. We follow these simple rules to position the 12 dinucleotides on the hexagon (six dinucleotide on six ends + six dinucleotide on midpoint of each arm of the hexagon). Based on the above discussions, Cycle 1 is drawn in the 2D Cartesian space, shown in Figure 3.



3. Cartesian coordinates of 16 dinucleotide in a hexagon

From Figure 3, we can derive the set of position coordinates of 16 dinucleotide:  $(0, 1.5) \rightarrow AG$ ,  $(0.5, 1.25) \rightarrow AT$ ,  $(1, 1) \rightarrow CA$ ,  $(1, 0) \rightarrow GA$ ,  $(1, -1) \rightarrow CG$ ,  $(0.5, -1.25) \rightarrow GT$ ,  $(0, -1.5) \rightarrow CT$ ,  $(-0.5, -1.25) \rightarrow GC$ ,  $(-1, -1) \rightarrow TG$ ,  $(-1, 0) \rightarrow TC$ ,  $(-1, 1) \rightarrow TA$ ,  $(-0.5, 1.25) \rightarrow AC$ ,  $(0, 1) \rightarrow AA$ ,  $(-0.5, 0) \rightarrow CC$ ,  $(-1, 0) \rightarrow GG$ ,  $(0.5, 0) \rightarrow TT$ .

Let  $S = \{s_1, s_2, \dots, s_N\}$  be a DNA sequence where  $s_i \in \Sigma = \{A, T, C, G\}$  and  $i = 1, 2, 3, \dots, N$ .  $S$  is mapped into a series of points  $P_1, P_2, \dots, P_{N-1}$ . We introduce a map function  $\phi$  such that  $S$  can be formulated as  $S = \phi(s_i s_{i+1}) \phi(s_{i+1} s_{i+2}) \dots \phi(s_{N-1} s_N)$  where,

$$P_i = \phi(s_i s_{i+1}) = \phi(x_i, y_i, i) = \phi(x_{s_i s_{i+1}}, y_{s_i s_{i+1}}, i), \quad i=1,2,3,\dots, N-1$$

$x_{s_i s_{i+1}}$ ,  $y_{s_i s_{i+1}}$  and  $i$  represent the x-coordinate, y-coordinate, and z-coordinate respectively.

Thus, we connect the  $N-1$  points from the first one and derive a 3D curve.

To locate the local and global features of the 3D curve as well as to visualize the 3D representation of this curve, we take another numerical representation. Let  $X_i = \sum_{k=1}^i x_k, Y_i = \sum_{k=1}^i y_k$ , we derive another mapping function for cumulative feature of the 3D curve such that

$$\lambda(S_i S_{i+1}) = (X_i, Y_i, i), \quad \text{where } i = 1, 2, 3, \dots, N - 1.$$

Connecting N-1 points from the first one, we get the proposed novel 3D zigzag curve.

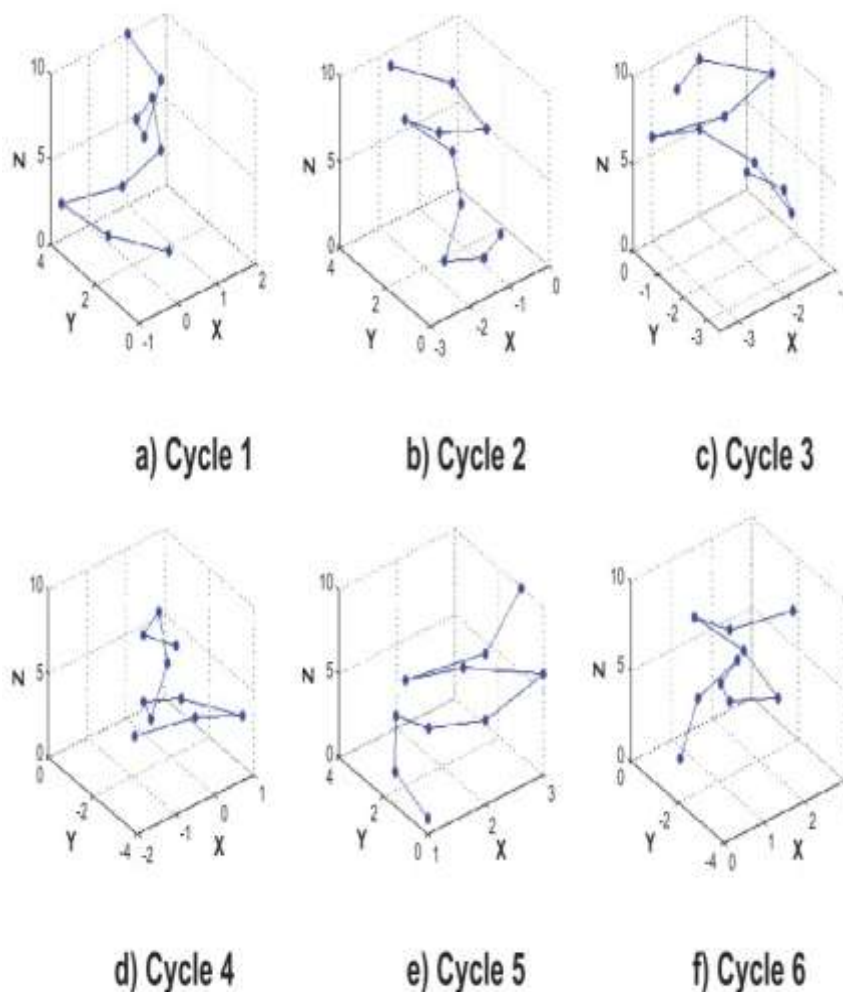
### 2.3 Example of the proposed method

The following example is used with the arbitrary DNA sequence ATACGATGCAG. The length of the string is 11, hence there are 10 di-nucleotide. The 3D coordinate for all cycles of the sequence is shown in Table 1.

Points Dinucleotide	Cycle 1		Cycle 2		Cycle 3		Cycle 4		Cycle 5			
	x	y	x	y	x	y	x	y	x	y		
P <sub>1</sub> AT	01	1.201	1	0	1	03	-1.201	-03	-1.201	-1	0	
P <sub>2</sub> TA	-01	1.201	1	10	2	13	-0.201	01	-1.201	-1	-11	
P <sub>3</sub> AC	-1.031	2	1.03	20	21	-0.201	1.0	-11	0	-1.5	-2.03	
P <sub>4</sub> CG	0	11	4	1.03	20.4	13	-1.20.4	0	-11	0	-1.5	-1.20.4
P <sub>5</sub> GA	1.0	11	0	2	0	1.0	-11	0	-1.0	-11	0	-2
P <sub>6</sub> AT	1.0	1.70.0	1	0	0	1.0	-1.70.0	-1.0	-1.70.0	-5	0	
P <sub>7</sub> TC	0.0	1.70.1	2	1	1	1.0	-1.20.1	-0.0	-1.70.1	-1	-1	
P <sub>8</sub> GC	0	1.3	0	1	1	0	-1.0	0	-1.3	0	-1	
P <sub>9</sub> CA	1.0	1.3	0	2	0	0	-1.0	0	-1	-1.3	0	
P <sub>10</sub> AG	1.0	1.3	1.0	1	1.0	-1.3	1.0	-1	-1.0	-1	1.0	

Table 1. 3D coordinates of ATACGATGCAG based on the proposed method

As for graphical representation, the 10 points P<sub>1</sub>, P<sub>2</sub>, ..., P<sub>10</sub> are plotted in 3D space for the example sequence ATACGATGCAG. The six possible DNA curves for the example sequence are shown in Figure 4.



**Figure 4. The graphical representation of the proposed model for the example sequence ATACGATGCAG**

In this way, each DNA sequence is converted into a series of points. Then DNA curves are drawn from those points. Connecting N-1 points from the first one, we get the proposed novel 3D zigzag curve in the 3D space. The DNA curve is helpful to easily distinguish among different species.

It can easily be seen that the example graphical representation does not hold any overlapping or loop. This property will be retained for any DNA sequence because the value of “i” in the proposed method is unique in every point.

### III. GRAPHICAL REPRESENTATION OF THE PROPOSED METHOD

The proposed model is useful to show the hidden properties of long DNA sequences which are not seen from the sequence. The pictorial presentation of the proposed method proves that it is very useful to understand the evolutionary similarity/dissimilarity of different species. Figure 5 shows the 3D zigzag curve based on Cycle 1 of first exon of  $\beta$ -globin for 11 different species. The graphical representation clearly shows that:

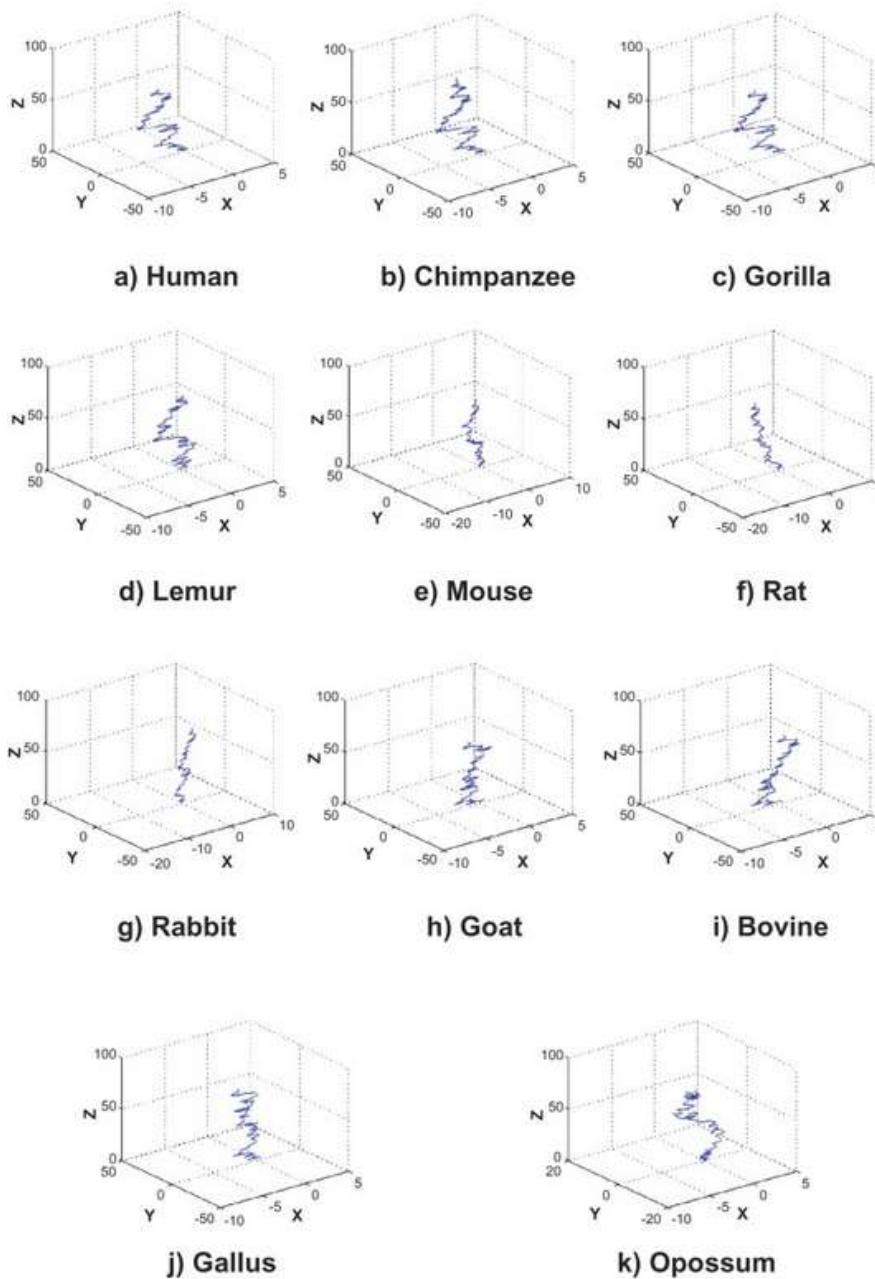


Figure 5. DNA curves of 11 different species

- i. DNA curves of human, gorilla, chimpanzee and lemur are closely similar;
- ii. Mouse and rat have also same DNA curves, so as rabbit's DNA curve;
- iii. Goat and bovine are similar; and that
- iv. Gallus and opossum seem to be outliers.

**IV. EXPERIMENTAL ANALYSIS**

**4.1 Performance metric, dataset and experimental environment**

To evaluate the performance, we illustrate the use of the proposed method with an examination of similarities/ dissimilarities among the  $\beta$ -globin gene of 11 different species, listed in Table 2, which were also previously studied. The table shows the different important characteristics of the dataset. First, we show the overall performance of the proposed method. To do this, two features are extracted from the DNA curves: (i) geometric center and (ii) mathematical descriptor. Each DNA sequence is finally represented by their mathematical descriptors. These descriptors form six dimensional feature vectors. After that, the Euclidian distance is calculated among feature vectors of the DNA sequences. Secondly, we draw the phylogenic tree from similarity/dissimilarity matrix using UPGMA method in PHYLIP package. Finally, we compare the proposed method with the already mentioned research works to show its superiority to others.

**Table 2**  
The first exon of  $\beta$ -globin gene of 11 different species.

Species	Accession	Database	Length
Human	U05333	NCBI	92
Chimpanzee	U02285	NCBI	105
Goat	U01109	NCBI	93
Leopard	U05728	NCBI	92
Rat	U04751	NCBI	92
Mouse	U01722	NCBI	93
Kabbit	U03952	NCBI	92
Owl	U01387	NCBI	88
Worm	U02115	NCBI	88
Spearmint	U03641	NCBI	92
Goat	U01109	NCBI	93

**Table 2. The first exon of  $\beta$ -globin gene of 11 different species**

**4.2 Numerical analysis of the proposed method**

As stated earlier, features from DNA curves are extracted two ways. Firstly, the geometric centers of the curves are calculated using the following equations. Table 3 shows the geometric center of 11 DNA curves.

**Table 3**  
Geometrical center of 11 different species.

Species	Cycle 1				Cycle 2				Cycle 3				G
	$x_1$	$y_1$	$x_2$	$y_2$	$x_3$	$y_3$	$x_4$	$y_4$	$x_5$	$y_5$	$x_6$	$y_6$	
Human	-4.405	-15.8275 86	-4.81	-3.75 86	-2.76725 867852 86	-0.4705							
Chimpanzee	-4.7481	-17.5624 52 5	-10.85	-4.88 52 5	-7.38125 5 847125 52 5	-0.7021							
Goat	-4.5161	-15.225 46 5	-4.9	-3.85 46 5	-6.8246 5 842025 46 5	-0.4728							
Leopard	-2.3132	-13.8181 86	-4.71	-2.23 86	-3.8936 6 03114 86	2.2876							
Rat	-0.7882	-11.3889 86	-4.78	1.51 86	-3.74176 12 74176 86	3.8258							
Mouse	-4.6882	-18.2281 47	-2.5	-1.85 47	-7.3177 12.25	47 2.8488							
Kabbit	-2.3716	-13.8481 86	-1.31	-4.88 86	-5.7811 2.23881 86	-2.8947							
Owl	-2.7824	-16.1881 43	-1.104	-4.88 43	-6.88275 3.79412 43	1.88827							
Worm	-4.7824	-16.8786 43	-1.347	-4.78 43	-6.88475 4.447059 43	0.11764							
Spearmint	-1.8771	-7.8176 86	-1.85	3.788 86	0.488811 7.7781 86	1.88881							

**Table 3. Geometrical center of 11 different species.**



$$u_x = \frac{1}{N} \sum_{i=1}^N x_i, \quad u_y = \frac{1}{N} \sum_{i=1}^N y_i, \quad u_z = \frac{1}{N} \sum_{i=1}^N z_i$$

The significance of the geometric center is that it shows the average value of x, y and z coordinate. Normally, if the geometric centers are plotted, then similar species fall into same cluster. So, the geometric center is an important feature of the DNA curves for the analysis of evolutionary relationship among different species.

To make our system unbiased, we take all the possible rotations (such as Cycle 1, Cycle 2, ..., Cycle 6) of hexagonal ring structure and extract the geometric center for each combination.

Secondly, mathematical descriptors are obtained from the first feature, the geometric center, using the following equation. Table 4. shows the mathematical descriptor of the 11 curves.

Species	Cycle 1	Cycle 2	Cycle 3	Cycle 4	Cycle 5	Cycle 6
Human	46.50017	47.12057	46.87512	46.86203	46.59829	47.59873
Chimpanzee	55.54423	53.83806	53.37854	53.51123	53.13654	54.30421
Gorilla	45.13718	47.88782	47.37582	47.46703	47.56275	46.07589
Lemur	47.79529	46.54827	46.53795	46.88248	46.79461	47.33921
Rat	47.58295	47.05385	47.28875	47.01807	46.5441	46.01626
Mouse	52.73956	46.76265	46.45373	46.95977	46.54146	46.14646
Rabbit	46.23836	46.03255	46.43998	46.6477	46.33272	46.93223
Goat	46.03942	46.23482	43.80219	44.05377	43.81217	46.79453
Worms	45.46871	43.05081	43.05269	44.18473	43.05765	44.31434
Opossum	46.13535	46.04480	46.06408	46.03293	46.24424	46.30385
Dolphin	47.48423	46.07130	46.21359	46.12039	46.27057	46.98216

Table 4. Mathematical descriptor of 11 different species.

We use the Euclidian distance for similarity measurement. Let two different species be i and j. The mathematical descriptor of i are p<sub>1i</sub>, p<sub>2i</sub>, p<sub>3i</sub>, p<sub>4i</sub>, p<sub>5i</sub> and p<sub>6i</sub>. The same descriptors for species j are p<sub>1j</sub>, p<sub>2j</sub>, p<sub>3j</sub>, p<sub>4j</sub>, p<sub>5j</sub> and p<sub>6j</sub>. The Euclidian distance of i and j are then calculated using the following equation:

The similarity/dissimilarity matrix found from the above Euclidian distance metric is shown in Table 5.

Species	Chimpanzee	Gorilla	Lemur	Rat	Mouse	Rabbit	Goat	Worms	Opossum
Human	0.6626	0.2602	0.0124	0.0109	0.0156	0.0114	0.0236	0.0156	0.0345
Chimpanzee		0.4058	0.0138	0.0135	0.0190	0.0135	0.0224	0.0137	0.0398
Gorilla			0.0126	0.0105	0.0190	0.0118	0.0225	0.0139	0.0344
Lemur				0.0134	0.0420	0.0080	0.0090	0.0156	0.0235
Rat					0.0255	0.0174	0.0223	0.0124	0.0342
Mouse						0.0447	0.0280	0.0190	0.0444
Rabbit							0.0311	0.0184	0.0235
Goat								0.0221	0.0494
Worms									0.0351
Opossum									

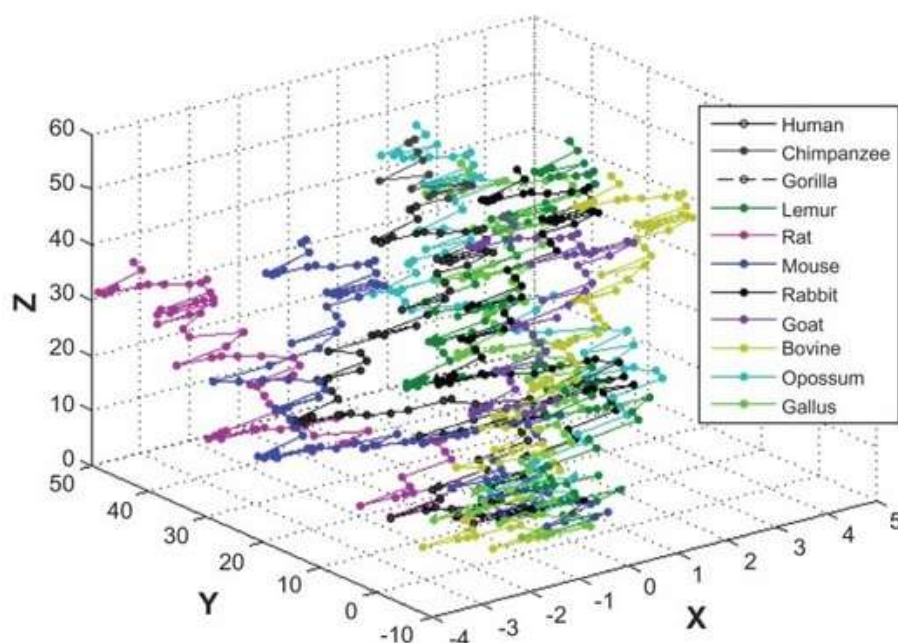
Table 5. Euclidian distance among 11 different species

Some observations are vividly depicted from Table 5 which are also consistent with the graphical representation portrayed in Section 3.4. They are as follows:

- i. The smallest entry is 0.0002 for the pair (human, gorilla), showing that human and gorilla are almost same in terms of evolutionary characteristics. The same is applied for the pair (human, chimpanzee) = 0.0020. Therefore, human, chimpanzee and gorilla are similar species;
- ii. The pair (goat, bovine) has the small entry 0.0131 which indicates the evolutionary similarity between goat and bovine. The biological taxonomy of bovine and goat proves that both of them are even-toed ungulates and belong to the family of "Bovidae";<sup>16</sup>
- iii. Rat and mouse also show a small entry which indicates their evolutionary closeness;
- iv. The remote mammalian opossum has the largest entry to all other mammals.

## V.CONCLUSION

A graphical method based on dinucleotides and their positional information is proposed in this research work. Graphical as well as numeric analyses of the model show that the proposed novel method is compatible with the natural consistency in terms of evolutionary relationship of 11 different species. In this paper, DNA sequences are transformed into 3D DNA curves, and features from those curves are then extracted. DNA curves are represented by their feature vector. Subsequently, Euclidian distance is applied to those feature vectors to deduce the evolutionary relationship among 11 different species. Tri-nucleotide based DNA sequence analysis using the proposed method would be one recommended future work.



DNA curves of 11 different species in one graph

## REFERENCES

- [1.] Hamori E, Ruskin J. H curves: a novel method of representation of nucleotide series especially suited for long DNA sequences. *J Biol Chem.* 1983;258(2):1318–27.
- [2.] Li Y, Huang G, Liao B, Liu Z. H-L curve: a novel 2-D graphical representation of protein sequences. *MATCH Commun Math Comput Chem.* 2009;61(2):519–32.
- [3.] Guo X, Randić M, Basak SC. A novel 2-D graphical representation of DNA sequences of low degeneracy. *Chem Phys Letter.* 2001;350(1–2):106–12.
- [4.] Jafarzadeh N, Iranmanesh A. A novel graphical and numerical representation for analyzing DNA sequences based on codons. *MATCH Commun Math Comput Chem.* 2012;68:611–20.
- [5.] Yu JF, Wang JH, Sun X. Analysis of similarities/dissimilarities of DNA sequences based on a novel graphical representation. *MATCH Commun Math Comput Chem.* 2010;63:493–512.
- [6.] Liao B, Zhu W, Liu Y. 3D graphical representation of DNA sequence without degeneracy and its applications in constructing phylogenetic tree. *MATCH Commun Math Comput Chem.* 2006;56(1):209–16.
- [7.] Cao Z, Liao B, Li R. A group of 3D graphical representation of DNA sequences based on dual nucleotides. *Internat J Quant Chem.* 2008;108(9):1485–90.
- [8.] Chi R, Ding K. Novel 4D numerical representation of DNA sequences. *Chemical Physics Letters.* 2005;407(1–3):63–7.
- [9.] Liao B, Li R, Zhu W, Xiang X. On the similarity of DNA primary sequences based on 5-D representation. *J Math Chem.* 2007;42(1):47–57.
- [10.] Liao B, Wang T. Analysis of similarity/dissimilarity of DNA sequences based on nonoverlapping triplets of nucleotide bases. *J Chem Inf Comput Sci.* 2004;44(5):1666–70.
- [11.] Wu R, Hu Q, Li R, Yue G. A novel composition coding method of DNA sequence and its application. *MATCH Commun Math Comput Chem.* 2012;67:269–76.
- [12.] Qi X, Wu Q, Zhang Y, Fuller E, Zhang CQ. A novel model for DNA sequence similarity analysis based on graph theory. *Evol Bioinform Online.* 2011;7:149–58.
- [13.] Ewens J, Grant G. *Statistical Methods in Bioinformatics: An Introduction.* 2nd ed. New York: Springer Science; 2005.
- [14.] Randić M, Vračko M, Lerš N, Plavšić D. Novel 2-D graphical representation of DNA sequences and their numerical characterization. *Chem Phys Lett.* 2003;368(1–2):1–6.
- [15.] Randić M, Vračko M, Lerš N, Plavšić D. Analysis of similarity/dissimilarity of DNA sequences based on novel 2-D graphical representation. *Chem Phys Lett.* 2003;371(1–2):202–7.