

Characterization of Pathogenic Genes through Condensed Matrix Method, Case Study through Bacterial Zeta Toxin

Uttam K. Mondal¹, Arnab Sen² and Asim K. Bothra^{1*}

¹*Cheminformatics Bioinformatics Laboratory,
Department of Chemistry,
Raiganj College (University College), Raiganj-733134,
Uttar Dinajpur, West Bengal, India.*

²*DBT Bioinformatics Facility, Department of Botany,
University of North Bengal, Siliguri-734013, West Bengal, India.*

**Corresponding Author E-mail: asimbothra@gmail.com*

Abstract

The Zeta protein causes Gram-negative bacteria (*Escherichia coli*) to cease growing and form long cells with many chromosomes, clearly unable to divide. Human cancer cells also die from the toxic effect of the Zeta poison. The molecular mechanism by which the Zeta toxin operates has not yet been discovered. Sequence based phylogeny has some limitation due to very low sequence similarity amongst the different zeta toxins and structure based phylogeny has also its limitation. In this study, zeta toxin nucleotide sequences of some pathogenic and non-pathogenic Bacteria were used for phylogenetic analysis. The uniqueness of this method is that it does not employ sequence alignment of complete nucleotide sequence of the corresponding gene.

Keywords: condensed matrix method, zeta toxin, eigenvalue and phylogram.

Introduction

The Zeta protein causes Gram-negative bacteria (*Escherichia coli*) to cease growing and form long cells with many chromosomes, clearly unable to divide. Similarly, when the toxin is introduced into yeast cells it changes their morphology and halts budding, and in large quantities it causes death. Human cancer cells also die from the toxic effect of the Zeta poison. The poison protein Zeta is extraordinarily large

compared other toxins (287 vs. around 100 amino acids) and with the exception of its frequently encountered nucleotide binding motif, it does not show any similarities to any known proteins. The molecular mechanism by which the Zeta toxin operates has not yet been discovered (1). Where alignment based and structure based phylogeny fails, nucleotide triplet based method may give light towards molecular phylogeny.

Availability of nucleotide sequences of zeta toxin motivated us to construct phylogram using their nucleotide sequence, which will complement the phylogram obtained by sequence similarity. In this work we have done the molecular phylogeny of zeta toxin using their nucleotide sequence and without making any sequence alignment. It is based on a method developed by basak et al associating DNA sequences with a set of sequence invariant. In this work we have quantified the string, which favors the direct comparison of the sequences. A sequence invariant, as considered as a number independent of the labels A, C, G, T standing for adenine (A), cytosine (C), guanine (G), and thiamine (T). We have form the matrix associated with each sequence and calculated the leading Eigen value of the matrices to see the variation of leading Eigen values associated with the string and the relationship between the enzymes. We have also build a phylogram using the Eigen values of the characteristics matrices of zeta toxin.

Our results complement the observation with the earlier studies based on multiple sequence alignment and structural alignment. The uniqueness of this method is that it does not employ sequence alignment of complete nucleotide sequence of the corresponding gene.

Material and Method

The nucleotide sequences of zeta toxin of some pathogenic and non pathogenic bacteria were obtained from www.img.jgi.doe.gov. In a DNA sequence of four letters, there are 64 possible triplets (subsequence of length3) that can occur, starting from AAA, AAT, AAG, AAC, ATA, ATT, ATG, ATC, AGA, AGT, AGG, AGC, ACA, ACT, ACG, ACC, etc. A $4 \times 4 \times 4$ cubic matrix with 64 entries that denote the frequencies of occurrence of all the 64 triplets in a DNA sequence are introduced. For the cubic matrix, three groups of $4 \times 4 \times 4$ matrices, $\{M_1, M_2, M_3, M_4\}$, $\{M_5, M_6, M_7, M_8\}$, $\{M_9, M_{10}, M_{11}, M_{12}\}$, can be obtained, each group of which contain all entries of the cubic (see Table I). Usually the group of 4×4 matrices $\{M_1, M_2, M_3, M_4\}$ as the representative of the cubic matrix. The four matrices contain not only the information about frequencies of occurrence of all triplets of a DNA sequence but also the information about the frequencies of occurrence of pairs and every letter in a DNA sequence. For example, the number of all TG-pair in a DNA sequence is equal to the row sum of the third row in M_2 plus ∂ , where $\partial = 0$ if the last two letters of the DNA sequence are not TG and $\partial = 1$ otherwise. The frequency of occurrence of any pair in a DNA sequence can obtain by the above method. In addition, the frequencies of occurrence of four letters A, T, G, C are, respectively, equal to the sum of all entries of M_1, M_2, M_3, M_4 plus ∂ , where ∂ are, respectively, equal to the number of

A, T, G, C in the last two letters of the DNA sequence. The column sums of M1, M2, M3 and M4 just denote the number of pairs of distance two in a DNA sequence (2, 3).

Table I: Three Groups of Four 4×4 matrices, $\{M_1, M_2, M_3, M_4\}, \{M_5, M_6, M_7, M_8\}$, and $\{M_9, M_{10}, M_{11}, M_{12}\}$ Listing All 64 possible XYZ Entries, Where X, Y, Z = A, C, G, T.

M ₁	M ₂	M ₃	M ₄
AAA AAT AAG AAC	TAA TAT TAG TAC	GAA GAT GAG GAC	CAA CAT CAG CAC
ATA ATT ATG ATC	TTA TTT TTG TTC	GTA GTT GTG GTC	CTA CTT CTG CTC
AGA AGT AGG AGC	TGA TGT TGG TGC	GGA GGT GGG GGC	CGA CGT CGG CGC
AGA ACT ACG ACC	TCA TCT TCG TCC	GCA GCT GCG GCC	CCA CCT CCG CCC
M ₅	M ₆	M ₇	M ₈
AAA AAT AAG AAC	ATA ATT ATG ATC	AGA AGT AGG AGC	ACA ACT ACG ACC
TAA TAT TAG TAC	TTA TTT TTG TTC	TGA TGT TGG TGC	TCA TCT TCG TCC
GAA GAT GAG GAC	GTA GTT GTG GTC	GGA GGT GGG GGC	GCA GCT GCG GCC
CAA CAT CAG CAC	CTA CTT CTG CTC	CGA CGT CGG CGC	CCA CCT CCG CCC
M ₉	M ₁₀	M ₁₁	M ₁₂
AAA TAA AAG CAA	AAT TAT GAT CAG	AAG TAG GAG CAG	AAC TAC GAC CAC
ATA TTA GTA TAC	ATT TTT GTT CTG	TAG TTG GTG CTG	ATC CTT TCG CTC
AGA TGA GGA CGA	AGT TGT GGT CGG	AGG TGG GGG CGG	AGC GGC GCG CGC
ACA TCA GCA ACC	ACT TCT GCT CCG	ACG TCG GCG CCG	ACC CCT GCC CCC

We developed our own program in C++ to count all the possible triplets of t-RNA synthetase and formed the matrices by using all the possible triplets. Also we have calculated the leading Eigen values of the matrices by using MATHLAB (Version 4) (4) software. We have constructed a distance matrix of the synthetases by summing the square of the difference of eigen values. A phylogram of the synthetases are constructed by the cluster analysis of the similarity matrix using phylip (5).

Result

The lengths of the zeta toxin of some pathogenic and non-pathogenic bacteria are given in Table 2. It is clear that the enzymes differ considerably in length. Firstly, we took the nucleotide sequence of zeta toxin of pathogenic and non pathogenic bacteria listed in Table 3 and counted the frequencies of occurrence of all the 64 triplets then the group of 4×4 matrices $\{M_1, M_2, M_3, M_4\}$ as the representative of the cubic matrix are constructed. The leading Eigen values of each matrix are evaluated. The leading Eigen values of each matrix of those bacteria are represented in Table 4. The distance matrices of the synthetases are constructed by summing up the square of the difference of eigen values. The distance matrix for bacteria mentioned in Tale 1 is given in Table 4. Using the distance matrices phylograms are constructed, which are represented in Figures 1. Figure 2 represented the Phylogram of zeta toxins based on ClustalW (6).

Table 2: The lengths of the zeta toxin of some pathogenic and non-pathogenic bacteria.

Bacteria Name	Short name	Length	Nature
<i>Pseudomonas fluorescens</i>	PSE	759	Non pathogenic
<i>Frankia sp. CcI3</i>	FR1	1053	Non pathogenic
<i>Frankia sp. CcI3</i>	FR2	1353	Non pathogenic
<i>Mesorhizobium sp. BNC1 plasmid 1</i>	MES	1761	Non pathogenic
<i>Alteromonas macleodii 'Deep ecotype</i>	ALT	720	Non pathogenic
<i>Streptococcus pneumoniae ATCC 700669</i>	STR	759	Pathogenic
<i>Neisseria cinerea ATCC 14685</i>	NEI	720	Pathogenic
<i>Enterococcus faecalis TX0104</i>	ENT	441	Pathogenic
<i>Oribacterium sinus F0268</i>	ORI	783	Pathogenic
<i>O.algarvensis Gamma1</i>	OLA	306	Non pathogenic
<i>Crenothrix polyspora</i>	CRE	2154	Non pathogenic

Table 3: The leading Eigen values of each matrix of bacteria are represented in Table II.

Name of Bacteria	Short Name	M1	M2	M3	M4
<i>Pseudomonas fluorescens</i>	PSE	63.9742	59.4986	66.1942	61.1813
<i>Frankia sp. CcI3</i>	FR1	40.9835	45.5497	83.7531	1.0101
<i>Frankia sp. CcI3</i>	FR2	51.1907	38.5775	82.8222	1.001
<i>Mesorhizobium sp. BNC1 plasmid 1</i>	MES	39.4057	47.476	87.3356	86.5084
<i>Alteromonas macleodii 'Deep ecotype</i>	ALT	80.975	64.4988	59.6351	51.9077
<i>Streptococcus pneumoniae ATCC 700669</i>	STR	101.2985	68.0832	59.42	49.6625
<i>Neisseria cinerea ATCC 14685</i>	NEI	92.048	69.8199	70.3941	41.2137
<i>Enterococcus faecalis TX0104</i>	ENT	93.5377	88.5706	54.7326	37.3998
<i>Oribacterium sinus F0268</i>	ORI	108.5584	71.2367	65.8962	40.403
<i>Olavius algarvensis Gamma1</i>	OLA	71.8841	53.1558	57.2465	69.9674
<i>Crenothrix polyspora</i>	CRE	72.5779	69.704	60.1925	54.6687

Table 4: The distance matrix of zeta toxins of some pathogenic and non-pathogenic bacteria mentioned in Table 1.

11	PSE	FR1	FR2	MES	ALT	STR	NEI	ENT	ORI	OLA	CRE
PSE	0	4652.03	4499.27	1836.57	443.05	1645.37	1311.01	2416.11	2557.36	260.05	256.61
FR1	4652.03	0	153.67	7328.99	5130.63	7104.81	4991.42	6779.14	7096.86	6470.41	5015.98
FR2	4499.27	153.67	0	7549.96	4688.15	6296.98	4416.93	6406.48	6196.68	6051.22	4818.6
MES	1836.57	7328.99	7549.96	0	3982.31	6392.28	5609.09	8093.65	7932.02	2266.07	3344.99
ALT	443.05	5130.63	4688.15	3982.31	0	430.98	381.04	971.79	977.8	543.17	105.54
STR	1645.37	7104.81	6296.98	6392.28	430.98	0	280.4	652.31	190.33	1505.05	853.16

NEI	1311.01	4991.42	4416.93	5609.09	381.04	280.4	0	613.64	295.49	1683.91	664.21
ENT	2416.11	6779.14	6406.48	8093.65	971.79	652.31	613.64	0	659.73	2790.05	1123.29
ORI	2557.36	7096.86	6196.68	7932.02	977.8	190.33	295.49	659.73	0	2620.79	1532.99
OLA	260.05	6470.41	6051.22	2266.07	543.17	1505.05	1683.91	2790.05	2620.79	0	517.05
CRE	256.61	5015.98	4818.6	3344.99	105.54	853.16	664.21	1123.29	1532.99	517.05	0

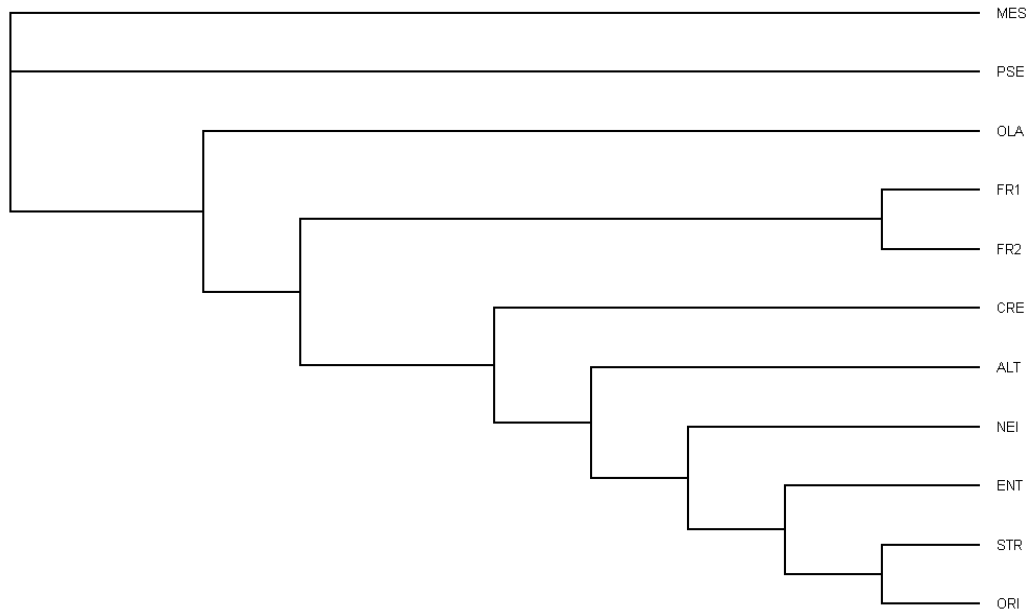


Figure 1: Phylogram of zeta toxins based on condensed matrix method.

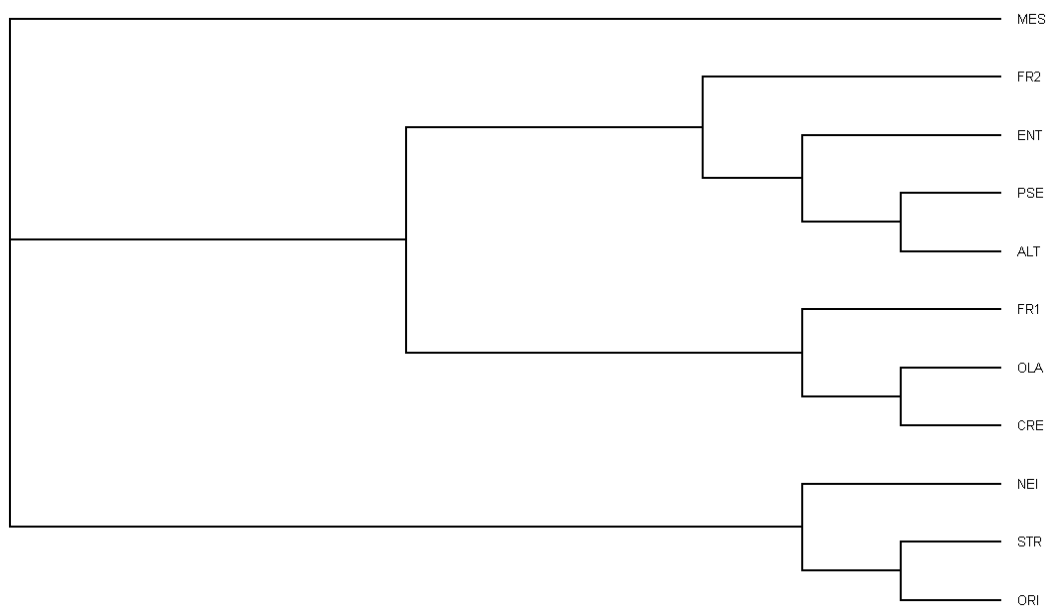


Figure 2: Phylogram of zeta toxins based on ClustalW.

Discussion

From Figure 1 is seen that, zeta toxin of pathogenic and non-pathogenic bacteria forms different cluster. It is also seen that two zeta toxins of Frankia sp. Cc13 are in the same clade. From Figure 2 it is seen that separation are not clear between pathogenic and non-pathogenic bacteria. It is also seen that two zeta toxins of Frankia sp. Cc13 are in the different clade which indicates superiority of condensed matrix method.

Sequence comparison quickly becomes unreliable at this and lower levels of sequence identity. In this regime of similarity, it becomes difficult to distinguish between correctly aligned homologous sequences and unrelated sequences or random alignments. Structure based phylogeny has limited scope because adequate number of structures are not yet solved to draw any general conclusion. From the separation of zeta toxin by condensed matrix method may help to identify pathogenic or non-pathogenic strain/species.

References

- [1] Zielenkiewicz, U. and Dmowski, M., 2009, "Systems ensuring plasmid inheritance Epsilon vs. Zeta," *Research in Progress Microbiology*, 1 (21), pp. 38-39.
- [2] Randic, M., Guo, X., and Basak, S. C., 2001, "On the characterization of DNA primary sequences by triplet of nucleic Acid Based," *J Chem Inf Comput Sci.*, 41(3), pp. 619-626.
- [3] Randic, M., and Basak, S. C., 2001, "Characterization of DNA Primary Sequences Based on the Average Distances between Bases," *J Chem Inf Comput Sci.*, 41(3), pp. 561-568.
- [4] Toh, K.C., Todd, M.J. and Tutuncu, R.H., 1999, "SDPT3 --- a Matlab software package for semidefinite programming, *Optimization Methods and Software*," 11(12), pp. 545--581.
- [5] Felsenstein, J., 1989 "PHYLIP - Phylogeny Inference Package (Version 3.2)," *Cladistics*, 5 (2), pp. 164-166.
- [6] Thompson J.D., Higgins D.G. and Gibson T.J., 1994, "CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice." *Nucleic Acids Res.* 22(22), pp. 4673-4680.