

**DEVELOPMENT OF AN *AB INITIO* PHYSICO-CHEMICAL MODEL FOR  
ANALYZING PROKARYOTIC GENOMES**

*by*

**POONAM SINGHAL**

**Department of Chemistry**

**THESIS SUBMITTED IN FULFILMENT  
OF THE REQUIREMENTS FOR  
THE DEGREE OF  
DOCTOR OF PHILOSOPHY**

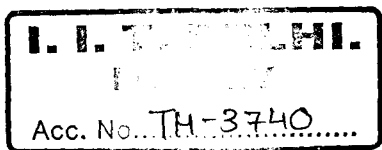
*to the*



**INDIAN INSTITUTE OF TECHNOLOGY, DELHI**

**HAUZ KHAS, NEW DELHI, INDIA**

**May, 2008**



7 H

579.202

SIN-D

*Dedicated to my parents*

## *Certificate*

This is to certify that the thesis entitled “Development of an *ab initio* Physico-Chemical Model for Analyzing Prokaryotic Genomes” being submitted by Ms. Poonam Singhal to the Indian Institute of Technology, Delhi for the award of the degree of Doctor of Philosophy in Chemistry is a record of bonafide research work carried out by her. Ms. Poonam Singhal has worked under my guidance and supervision, and has fulfilled the requirements for the submission of this thesis, which to my knowledge, has reached the requisite standard.

The results contained in this dissertation have not been submitted in part or full to other University or Institute for the award of any degree or diploma.



Prof. B. JAYARAM  
Head, Department of Chemistry,  
Indian Institute of Technology, Delhi

## *Acknowledgements*

*Guru Bhramaha, Guru Vishnu, Guru Devo Maheshwara,*

*Guru Sakshat Param Bhramaha, Tasmai Siri Gurve Namah.*

*The old saying is as true as it was in the ancient time. The credit of this project goes to my **Guru**, who made me to learn new things and taught me to take up new challenges in studies as well as in life. I am deeply obliged to my research supervisor and Head, Prof. B. Jayaram, Department of Chemistry, Indian Institute of Technology New Delhi, for giving me an opportunity to pursue my research interests in his group, his valuable suggestions, constant encouragement, keen interest and affectionate attitude during the course of this work. He has been a pillar of support to me during my research period and has always been available to give advice and to guide me. He is a great human being, who cares for his students.*

*I would like to acknowledge the financial support from Department of Biotechnology, Government of India and IIT for attending the international conference at Albany, USA. I am highly grateful to Prof. D. L. Beveridge and Dr. Surjit B. Dixit for their valuable scientific help and making me feel like home during my stay in USA.*

*I am grateful to CSIR, New Delhi for providing me with the fellowship during the tenure of my research.*

*Thanks are due to all the teachers who taught me at any point of time and have laid the foundation and helped me in growing both at academic as well as personal level. I am grateful to all my lab members of the Supercomputing Facility for Bioinformatics and Computational Biology and the Chemistry Department Lab, for their help and cooperation received in completing my research work and associated activities.*

*I thank all the faculty and staff of the Chemistry Department, IIT Delhi for their help and support received during this project.*

*I owe a lot to my friends. I feel privileged to have friends like Kumkum, Arti, Teena, and Garima.*

*I wish to express my deepest sense of gratitude and reverence to my loving and respectable parents for their unconditional love, encouragement and patience to pursue my interests. Without their initial momentum and continuous support, it would have been impossible for me to arrive where I am today. I thank them for providing me the best of education and making me realize the virtues of life. I am also indebted to my sisters Anchal, Jyoti and brother Sachin for their love, invaluable support and care that they have always shown for me. My mausi, mausaji and cousins Amit, Ankit and Nitin have also contributed much to this work in visible and invisible ways. This journey would have been very difficult or rather impossible without them.*

*Lastly I would like to thank Almighty God for giving me a direction and a sense of purpose in my life. Without His blessings, the work described herein would not have been possible.*

*Peonam*

## *Abstract*

The genome sequence represents the book of life. Buried in this large volume are genes, which are scattered as small DNA fragments throughout the genome and comprise a small percentage of the total text. Finding these indistinct 'needles' in a vast genomic 'haystack' can be extremely challenging.

The computational gene identification problem is an endeavor to interpret nucleotide sequences, in order to provide an annotation on the location and hopefully the functional class of protein-coding genes (mRNA) as well as genes for other RNAs (rRNA, tRNA etc.). This problem is of self-evident importance and remains far from full resolution. With the advent of whole-genome sequencing projects, there is considerable use for algorithms which scan genomic DNA sequences to find genes, particularly those that encode for proteins. Even though there is no substitute to experimentation in determining the exact locations of genes in a genome sequence, a prior knowledge of the approximate location of genes will hasten the process to a great extent apart from saving huge amount of laboratory time and resources. Several algorithms have been reported in the past for analysis of prokaryotic genomes based on sophisticated, statistical and mathematical methods. A physico-chemical understanding of the characteristics of the genes is lacking. The present thesis attempts to address this issue.

This thesis work presents an *ab initio* model for gene prediction in prokaryotic genomes based on physico-chemical characteristics of double helical trinucleotide sequences corresponding to the 64 codons calculated from molecular dynamics (MD) simulations. The methodology combines filters based on stereochemical properties of protein sequences and frequencies of occurrences of codons via their corresponding amino acids from 175000 Swissprot proteins to reduce false positives. The methodology has been validated on 372 prokaryotic genomes with sensitivity, specificity and correlation coefficients averaged over 356208 genes and an equal number of frame-shifted genes (non-genes) reaching 97.50%, 97.20% and 94.25% respectively. The protocol has been automated in the form of a web server (<http://www.scfbio-iitd.res.in/chemgenome2>) and is made freely accessible for usage by the scientific community.

This thesis is divided into eight chapters. Chapter 1 gives an introduction of the gene prediction methodologies. An overview of the computational methodologies (extrinsic, intrinsic, hybrid and comparative genomics approaches) is also presented. Chapter 2 describes the methodology developed as a part of this thesis for the evaluation of DNA sequences and prediction of protein-coding genes from the whole genome sequence of the organisms. The modular form of the methodology has been explained which allows for the automation of the methodology as a web server. Chapter 3 describes the results obtained from the gene evaluation

methodology explained in Ch. 2 on 372 prokaryotic genomes. An evaluation of the methodology with the virus genomes is also presented. Chapter 4 describes the results obtained from the gene prediction methodology explained in Ch. 2 on 372 prokaryotic genomes. Chapter 5 describes the performance appraisal of the gene evaluation methodology for eukaryotic genomes. Chapter 6 explains the protocol devised for extracting the physico-chemical information hidden in DNA sequences. The methodology explained in Ch. 2 has been automated in the form of web servers which are described in Chapter 7. Finally in Chapter 8, summary and some perspectives emerging from this thesis work on genome analysis are provided.

# *Contents*

<u><i>Certificate</i></u>	I
<u><i>Acknowledgements</i></u>	II-IV
<u><i>Abstract</i></u>	V-VII
<u><i>List of Figures</i></u>	VIII-X
<u><i>List of Tables</i></u>	XI-XII
<u><b>Chapter I: <i>Introduction</i></b></u>	<b>1-31</b>
1.1 Genome organization	2
1.2 Genome analysis	7
1.3 Applications of genome analysis	8
1.4 Computational methods of gene prediction	10
1.4.1 Extrinsic or similarity based methods	12
1.4.2 Intrinsic or ab-initio methods	18
1.4.3 Hybrid methods or consensus methods	24
1.4.4 Comparative genomics approaches	26
1.5 Scope of this thesis work	30
<u><b>Chapter II: <i>Methodology</i></b></u>	<b>32-59</b>

2.1	Introduction	33
2.2	A computational protocol	34
2.2.1	Development of the model parameters	37
2.2.2	Orthogonalization of the model parameters	52
2.2.3	Finding the best plane dividing genes from non-genes	52
2.2.4	<i>ChemGenome 1.1</i> : Gene Evaluator	53
2.2.5	<i>ChemGenome 2.0</i> : Gene Predictor	53
2.2.6	Whole genome to open reading frames	56
2.2.7	Selection of putative genes from open reading frames	56
2.2.8	Second stage selection based on stereochemical properties of proteins	57
2.2.9	Final selection of predicted genes	58
<b><u>Chapter III: Gene evaluation results on 372 prokaryotic genomes</u></b>		<b>60-79</b>
3.1	Introduction	61
3.2	Methodology	62
3.2.1	Analysis of DNA sequences based on physico-chemical properties of DNA	62
3.2.2	Description of dataset	63
3.2	Results and discussions	63
3.2.1	Assessment parameters of the results / Performance evaluation	63

3.2.2	Results obtained on a dataset of 372 test genomes	67
3.2.3	Results obtained on virus genomes	77
3.3	Conclusions	79
<b><u>Chapter IV: Gene prediction results on 372 prokaryotic genomes</u></b>		<b>80-95</b>
4.1	Introduction	81
4.2	Methodology	82
4.2.1	Gene prediction based on physico-chemical characteristics of codons	82
4.2.2	Description of dataset	82
4.3	Results and discussions	84
4.3.1	Results obtained on a dataset of 372 test genomes	84
4.3.2	Comparison with other gene prediction programs	92
4.4	Conclusions	93
<b><u>Chapter V: Extension of the methodology to eukaryotic genomes</u></b>		<b>96-107</b>
5.1	Introduction	97
5.2	Theory and methodology	100
5.3	Results and discussions	100
5.4	Conclusions	107

<b><u>Chapter VI: Octant analysis of DNA sequences</u></b>	<b>108-122</b>
6.1 Introduction	109
6.2 Methodology	111
6.2.1 Computational protocol	111
6.2.2 Description of dataset	113
6.3 Results and discussions	114
6.4 Conclusions	122
<b><u>Chapter VII: Webservers for (a) gene evaluation and (b) gene prediction</u></b>	<b>123-137</b>
7.1 Introduction	124
7.2 ChemGenome 1.1 Gene Evaluator web server	129
7.3 ChemGenome 2.0 Gene Predictor web server	132
7.4 Conclusions	137
<b><u>Chapter VIII: Summary and Perspectives</u></b>	<b>138-141</b>
<b><u>References</u></b>	<b>142-163</b>
<b><u>Appendices</u></b>	<b>164-225</b>
<b><u>Brief Bio-data</u></b>	<b>226-229</b>