

0052750

Reports of International Conference
Novosibirsk, August 28–September, 1990

MODELLING AND COMPUTER METHODS IN MOLECULAR BIOLOGY AND GENETICS

Edited by
V.A. Ratner and N.A. Kolchanov

NOVA SCIENCE PUBLISHERS, INC.

Nova Science Publishers, Inc.
6080 Jericho Turnpike, Suite 207
Commack, New York 11725

Book Production Manager: Janet Glanzman White
Graphics: Elenor Kallberg, Christopher Concannon,
and Michael A. Masotti

Library of Congress Cataloging-in-Publication Data
available upon request

ISBN 1-56072-077-8

© 1992 Nova Science Publishers, Inc.

All rights reserved. No part of this book may be reproduced, stored in a retrieval system or transmitted in any form or by any means: electronic, electrostatic, magnetic, tape, mechanical, photocopying, recording or otherwise without permission from the publishers.

Printed in the United States of America

**MODELLING AND
COMPUTER METHODS
IN MOLECULAR BIOLOGY
AND GENETICS**

PREFACE

An International Conference "Modelling and Computer Methods in Molecular Biology & Genetics" was organized by the Theoretical Department of the Institute of Gytology & Genetics, SB AS USSR, together with the All-Union scientific and technological program "GENINFORM". It was held on August 28 - September 1, 1990 in Akademgorodok, Novosibirsk Scientific Center of Siberian Branch of the USSR Academy of Sciences, at Maly Zal (the Small Hall) of the House of Scientists [1].

Being the 1-st International Conference on these topics, it was actually the 4-th All-Union Conference, because the soviet participants meet once in two years since 1984 in Akademgorodok under the patronage of Siberian Branch of the USSR Academy of Sciences and the "GENINFORM" program. Here we had found a compact range of cooperated institutes and persons, discussed many common works, methods, results and technical projects. The subject of all these conferences was the usage of all available software, data-banks, methods for solution of some problems of the Theory of Molecular Genetic Regulatory Systems (MGRS), analysis of primary and spacial structure of macromolecules, and development of the Theory of Molecular Evolution (TME) (see [2-5]).

To continue this tradition the present Conference had the meetings and poster sessions on the following subjects:

1. Computer analysis of polynucleotide sequences;
2. Data Banks and program packages;
3. Computer analysis and Modelling of Protein Structure;
4. Molecular Genetic Sistem Modelling;
5. Theory of Molecular Evolution.

Moreover, there was the evening Round Table for Discussions devoted to the two main problems:

- a). Genetical Language as a linguistic system;
- b). Unified Theory of Molecular Evolution.

We decided not to present here the text of discussion, but our final paper [6] which contains the list of the problems for discussion and some comments.

In general, we are satisfied with the results of the Conference. We hope, the opinions of our guests are the same. It seems to me that they were pleased to find in so a distant place as Novosibirsk such an active community of mathematical and computer geneticists. And the weather was also fine.

In my final words I want to welcome all the authors of this Issue, all the participants of our Conference and interested readers. My best regards to our Publisher, Mr. Frank Columbus (NOVA SCIENCE PUBLISHERS), who gave his "yes" just after the first words about this Conference were said. The emblem of the Conference was designed by Dr. Andrey Rzhetsky, the participant of the Conference. I must also thank all our colleagues who helped us to organize this meeting.

REFERENCES

- [1] Abstr. of Intern. Conf. "Modelling and Computer Methods in Molecular Biology & Genetics", Eds. V.A.Ramer, N.A.Kolchanov, ICG-Press, Novosibirsk, (1990), v. I,2.
- [2] Proc. of I-st All-Union Workshop "Theoretical Investigations and Data Banks in Molecular Biology & Genetics", Ed. V.A. Ramer, ICG-Press, Novosibirsk, (1986) (Russ.).
- [3] Abstr. of 3-rd All-Union Workshop "Theoretical Investigations and Data Banks in Molecular Biology & Genetics", Ed. V.A.Ratner, ICG-Press, Novosibirsk, (1988) (Russ.).
- [4] Computer Analysis of Structure, Function and Evolution of Genetic Macromolecules, Ed. N.A. Kolchanov, ICG-Press, Novosibirsk, (1989) (Russ.).
- [5] Computer Analysis of Genetic Texts. Ed. M.D. Frank-Kamenetzky, Nauka, Moscow (1990) (Russ.).
- [6] V.A. Ramer, this volume,

Prof. Vadim A. Ramer
Head of Theoretical Department,
Vice-President of Organizing
Committee of the Conference

November, 1990
Novosibirsk

Contents

Preface	xi
Current Problems of Computer Analysis of Genetic Texts <i>N.A. Kolchanov</i>	1
Complexity Profiles of Genetic Texts <i>V.D. Gusev, O.M. Chupakhina, and V.A. Kulichkov</i>	23
Base-Pair Selection in Specific Protein Binding Sites on DNA; A Statistical-Mechanical Model <i>O.G. Berg</i>	31
Prediction of Protein-Coding Regions Interrupted by Introns <i>M.S. Gelfand</i>	43
Computer Investigation of Structural Organization and Evolution of Functional Sites in Polynucleotide Sequences <i>A.E. Kel', N.A. Kolchanov, V.V. Solovyev, M.P. Ponomarenko, I.V. Ishchenko, Y.L. Orlov, and V.V. Kapitonov</i>	49
Intelligent System of Mutational Analysis <i>I.B. Rogozin, N.E. Sredneva, and N.A. Kolchanov</i>	63
Revealing of Complicated Structural-Functional Organization of Genetic Signals <i>V.V. Solovyov, A.K. Salikhova, and I.A. Seledtsov</i>	73
Similarity Search in Biological Sequences <i>M.A. Roytberg</i>	81
Local Parallel Analysis of Nucleotide Sequences <i>A.I. Adamatzky, S.F. Ivanov, and V.A. Bronnikov</i>	87
Determining Global Helical Parameters in Duplexed DNA <i>D.W. Deerfield II, D. York, and T. Darden</i>	95
Electrostatics of DNA. New Dipole Model <i>R.V. Polozov, D.A. Kuznetsov, V.G. Tumanyan, and N.G. Esipova</i>	97

Statistical Analysis of Dinucleotide Compositions and Its Application to the Identification of Coding Regions <i>J. Kleffe and M. Borodovsky</i>	103
The Increasing of Informational Content of the Protein Database <i>C. Sander</i>	119
Molecular-Genetic Database and Software Development in the Soviet Union. The Problem of Database Integration <i>A.A. Alexandrov</i>	121
"Context" Package of Computer Programs for Analysis of DNA, RNA and Protein Sequences. Searching for Sequence Similarities and Functional Sites <i>V.V. Solov'ov, I.B. Rogozin, A.K. Salikhova, I.A. Seledtsov, A.A. Salamov, A.E. Kel'</i>	125
"DNA-SUN"—A Powerful Computer Tool for DNA and Protein Analysis <i>A.A. Mironov, L.V. Lunovskaya-Gurova, N.Yu. Bogodanova, N.N. Alexandrov, A.V. Grigorjev, V. Lebedev, P.A. Pevzner, and M.E. Trukhan</i>	131
The <i>Escherichia Coli</i> K12 Genome: Database and Data Analysis <i>T. Kunisawa and A. Tsugita</i>	133
Genbee: A Package of Computer Programs for Biopolymer Sequence Analysis <i>L.I. Brodsky, N.B. Vasiliev, A.E. Gorbalenya, A.P. Gulyaev, D.R. Davydov, A.L. Drachev, E.V. Koonin, A.M. Leontovich, R.I. Tatuzov, and K.M. Chumakov</i>	139
Flexis: Systems for Management of Integrated and Inquiry-Education Databases <i>Eu.I. Golovanov and G.M. Subotch</i>	141
Graphical Representation of DNA and Amino Acid Sequences <i>M.T. Carroll, G. Varga, E. Hamori, and H.A. Lim</i>	143
DNA Physical Mapping and Discrete Optimization <i>P.A. Pevzner and A.A. Mironov</i>	149

A Software System for the Assembly, Manipulation, and Visualization of Digitally Simulated Chromosomes	151
<i>William D. Shoaff and John C. Hozier</i>	
Computer Support of Electrophoretic Experiment	163
<i>V.G. Babski and M.Yu. Zhukov</i>	
Application of Conformational Search to Antibody Modeling	169
<i>R.E. Brucoleri, E. Haber, and J. Novotny</i>	
From the Comparative Analysis of Proteins to Similarity-Based Modelling	191
<i>Mark S. Johnson, J. Overington, A. Sali and T.L. Blundell</i>	
Conformation of Triple Helix of Collagen as a Function of Primary Structure	197
<i>V.G. Tumanyan, V.N. Rogulenkova, and N.G. Esipova</i>	
Protein Design and Secondary Structure Prediction	203
<i>J. Garnier</i>	
Computer Analysis of Protein Tertiary Structure: Residue-Residue Approach	217
<i>S.G. Galaktionov and M.A. Rodionov</i>	
Interaction of Peptide and Protein Molecules with Lipophilic Environment: Computer Modelling	225
<i>V.M. Tseytin, I.A. Vakser, and S.G. Galaktionov</i>	
Prediction of Protein Tertiary Structure by Analogy—Modelling of Loop Regions	231
<i>F. Eisenmenger and H. Sklenar</i>	
Biological Adequacy of Protein Sequence Alignments	233
<i>E.J. Demchuk and V.G. Tumanyan</i>	
Distance Geometry and the Calculation of Protein Conformation	241
<i>G.M. Crippen</i>	
Computer System of Protein Structure and Function Predictions	251
<i>M.P. Ponomarenko, D.N. Benjikh, A.A. Salamov, V.V. Solovyov, I.N. Shindyalov, V.B. Streletz, Yu.L. Orlov, and N.A. Kolchanov</i>	

Protein Secondary Structure Calculation on the Base of Discriminant Analysis with the Use of Information on Homologous Protein Structure <i>V.V. Solov'ov and A.A. Salamov</i>	265
A Novel Approach for Computing the Global Minimum of Proteins III. Comparison with other Global Minimization Methods <i>R.L. Somorjai</i>	275
Property Patterns in Protein Sequences: Automatic Generation and Applications <i>P. Bork, K. Rohde and J. Reich</i>	283
Mathematical Modelling of Molecular Genetic System Regulation of the Interferon Induction and Antiviral Action <i>V.A. Likhoshvai, S.I. Bazhan, and O.E. Belova</i>	293
AIDs and AIDs Vaccine. Analysis of Vaccination Perspective for Struggle Against Infection Induced by Human Immunodeficiency Virus (HIV): Mathematical Model <i>V.V. Chuykov, S.I. Bazhan, and V.A. Kulichkov</i>	301
Optimization of Gene Expression for Granulozyte-Monozyte Colony Stimulating Factor (GM-CSF) by Gene Design and RNA Secondary Structure Prediction <i>K. Weller, I. Petry, and M. Hartmann</i>	309
Flux Stoichiometric Models of Cell Metabolism <i>L.N. Drozdov-Tikhomirov, G.I. Scurida, and V.V. Serganova</i>	329
Limiting Factors and Evolution of Molecular Genetic Regulatory Systems (MGRS) <i>V.A. Ratner</i>	335
"Parasitic" DNA and Genome: Some Evolutionary and Coevolutionary Aspects <i>S.N. Rodin and J.S. Krushkal</i>	351
On Molecular Recapitulation in Globin Genes System <i>S.N. Rodin, A.Yu. Rzhetsky, and A.A. Zharkikh</i>	357
The Intrinsic Causes of the Evolution of Biosystems <i>Yu.G. Matushkin</i>	363

A Darwinian Theory of Molecular Evolution <i>Francesco M. Scudo</i>	369
Method of Linear Invariants for Phylogenetic Reconstruction <i>W.-H. Li and Y.-X. Fu</i>	379
Evaluation of DNA or Protein Sequence Similarity by L-Tuple Frequencies <i>A.A. Zharkikh and A.Yu. Rzhetsky</i>	391
Phylogenetic Trees for Several Kinds of Ribosomal RNA Shed Light Upon Early Stages of Seed Plant Evolution <i>A.V. Troitsky, Yu.F. Melekhovets, G.M. Rakhimova, V.K. Bobrova, K.M. Valiejo-Roman and A.S. Antonov</i>	399
Codon Usage and the Silent Molecular Clock: Variation Among Genes and Among Organisms <i>P.M. Sharp</i>	407
Investigation of Regularities of Molecular Evolution on the Base of Comparative Analysis of Homologous Sequences <i>I.N. Shindyalov and N.A. Kolchanov</i>	409
Computer Analysis of Mobile Genetic Elements <i>V.V. Kapitonov and N.A. Kolchanov</i>	421
Non-Random Sequence Changes Typify Both Spontaneous Mutations and Evolutionary Substitutions <i>B.W. Glickman and G.B. Golding</i>	431
Evolutionary Analysis of Eight Tomato <i>U1</i> RNA Gene Candidates <i>H. Hegyi and F. Solymosy</i>	445
Translation Framing Pattern in mRNA—Compensation Effects <i>J. Lagunez-Otero and E.N. Trifonov</i>	451
The Grammatical Rule for DNA Language: Messages Written in Palindromic Verses <i>Susumu Ohno</i>	461
Rare Codons: Fortuity or Regularity? <i>V.A. Likhoshvai</i>	463

Codon Mutational Variability and Evolutionary Strategies of Genetic Texts	471
<i>V.V. Solovyov and A.A. Salamov</i>	
Nonrandomness of Genetic Code Symmetry	477
<i>V.V. Solovyov</i>	
Some Basic Notions and Problems of the Theory of Molecular Genetic Regulatory System (MGRS)	481
<i>V.A. Ratner</i>	
Subject Index	501

CURRENT PROBLEMS OF COMPUTER ANALYSIS OF GENETIC TEXTS

N.A.Kolchanov

Institute of Cytology and Genetics of the USSR Academy of Sciences,
Novosibirsk, 630090, USSR

The development of highly efficient methods of DNA sequencing stimulated working out methods of genetic text computer analysis. A number of such program packages are designed up to date [1,2]. However, they cannot provide the rate of analysis compared with that of DNA sequencing. In order to improve the situation, more efficient technique is to be developed. The most complicated problems of genetic text analysis are associated with the "Human genome" project. In this case even the primary data analysis requires highly productive computer systems, able to work in an automatic mode.

It should be emphasized that some fundamental problems are also associated with the genetic text analysis. They are:

1) The theory of structural & functional organization of genomes, which involves:

- a) Theory of genome regulatory regions and functional sites.
- b) Theory of mobile genetic elements.
- c) Theory of gene structural & functional organization.

2) Theory of mutational and recombinational processes in genomes.

3) Theory of genetic text evolution.

4) Theory of protein structure & function prediction on the base of their amino acid sequences.

The most efficient way for these problems to be solved is the use of a new informational technique, which is oriented to the creation of artificial intelligence systems and knowledge bases in the corresponding fields of molecular biology and genetics.

This paper is devoted to new approaches in the genetic text analysis at the example of the results obtained in the Laboratory of Theoretical Molecular Genetics of the Institute of Cytology and Genetics of the USSR Academy of Sciences.

TYPE	ELEMENT
OBLIGATORY ELEMENTS	TATA-BOX
	-35-BOX
MODULATORS	COMPLEMENTARY PALINDROMS [3]
	REGIONS WITH HIGH AND LOW ENERGY OF DNA MELTING [3]
	TANDEM NONPERFECT REPEATS [4]

Fig. 1. Examples of obligatory elements and modulators in prokaryotic promoters

Let us consider the problem of functional site computer analysis of polynucleotide sequences. Developing of artificial intelligence systems begins from the creation of conceptual models of phenomena in study. As for functional sites, the below concept is represented

We suppose that in case of a site to function normally, the two types of elements are essential:

1) Obligatory elements, which are vitally important for the site normal activity. As a rule (but not always), various types of consensus or some conservative blocks of consensus in functional sites are considered as obligatory elements. In most cases functional sites can be characterized by the more or less exact location of some obligatory elements.

2) Modulators - second type of elements, which are presented in functional sites with variable location and number.

It is supposed, that the location and number of modulators is the essential factor of site functioning. We have found various types of such modulators in functional sites (Fig. 1). Solovyev et al. [3] have shown that prokaryotic promoters contain such type of modulators as complementary palindromes. The greater the number of these palindromes, the more is the strength of the promoter. Modulators are regions of high and low energy of DNA melting and their location within promoters is different [3]. Kel' et al. [4] have shown that

promoter strength depends on saturation with tandem nonperfect repeats. So, obligatory elements and modulators are actual and their origin is different.

It should be noted, that the mechanisms of site functioning are unknown in most cases. Apparently it is difficult to suppose a certain type of structural elements to define functioning of a definite type of sites. Consequently, it is necessary to generate a wide range of hypotheses upon the nature of these elements while studying functional sites. This approach corresponds to one of the main principles of artificial intelligence - the principle of impartiality [5], according to which the hypothesis should never be neglected without serious reasons. This very approach was realized in a computer system designed by Kel' et al. [6,7].

This system has a number of standard blocks (see fig.4 in paper [7]). The first block is an intellectual interface, which was designed for input and output of information and also for the system operating.

The second block is a database, containing information about functional sites (type of functional site, its polynucleotide sequence, location, etc.).

The third block is the system of hypothesis generation about the role of various structural elements in site functioning. Here this block generates a hypothesis of the role of the two types of elements: 1) oligonucleotides; 2) various types of repeats and symmetrical sequences.

The fourth block is used in statistical examination of hypotheses. The hypotheses, whose justice is stated as being proved are considered as

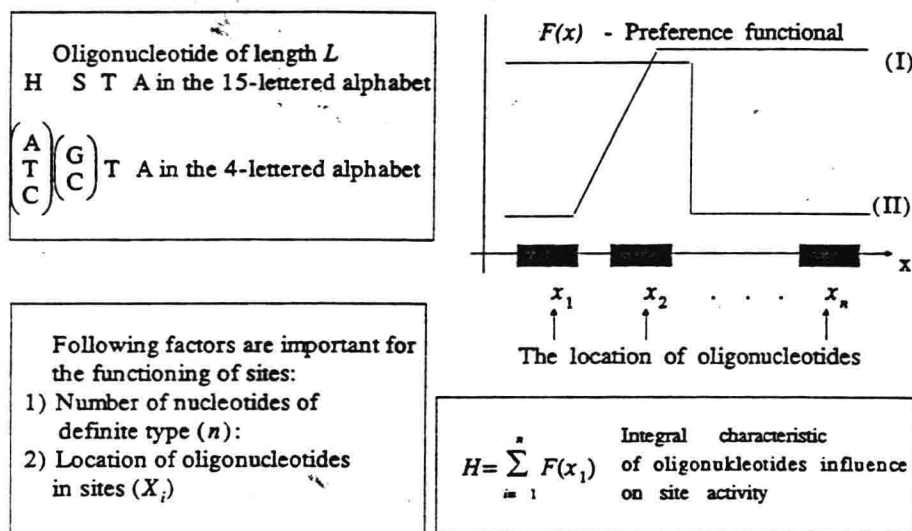


Fig. 2. Formal description of oligonucleotide influence on site activity

knowledges. Therefore, the above blocks are included in the system of knowledge production.

It is the knowledge base that distinguishes the artificial intelligence system from ordinary packages of programs. The knowledge obtained in the course of the system work is stored in the knowledge base.

In this system the two types of knowledge are stored in the knowledge base: 1) knowledge about structural elements which provide site functioning; 2) knowledge about the methods of those sites recognition. The methods are produced by the block of recognition. This block is also the part of the knowledge production system, as well as the intellectual interface controlling the process of their construction.

Let us consider the principles of system work at the example of oligonucleotides as structural elements (Fig. 2). Oligonucleotide is a fragment of DNA of length L and it has certain nucleotide sequence. It can be presented in the 15-lettered code or in 4-lettered one, as it is shown in the upper left part of Fig. 2.

Let us suppose that functioning of every type of sites depends on a specific set of oligonucleotides. In this case, for the site functioning the following parameters are important:

- 1) Number of oligonucleotides of a certain type in the site (N).
- 2) Their localization in the site (X_i).

It is supposed that the influence of the certain type of oligonucleotides on the site activity depends upon their position in the site.

It is the preference functional that is used in our system for the qualitative description of this influence. It is schematically represented in the right upper part of Fig. 2. For example, functional I reflects the increasing of oligonucleotide influence on the site activity from the 5'- to the 3'- direction. Functional II reflects the decreasing of the oligonucleotide influence on the site activity from its 5'- to the 3'-direction.

It is supposed that the contribution to the total site activity carried out by the oligonucleotide at the X_i -position is proportional to the quantity of the preference function $F(X_i)$. So, value H is the integral characteristic of a certain type of oligonucleotides to influence the site activity (Fig. 2).

Thus, the hypothesis about oligonucleotide influence on the site functioning can be described by the following parameters:

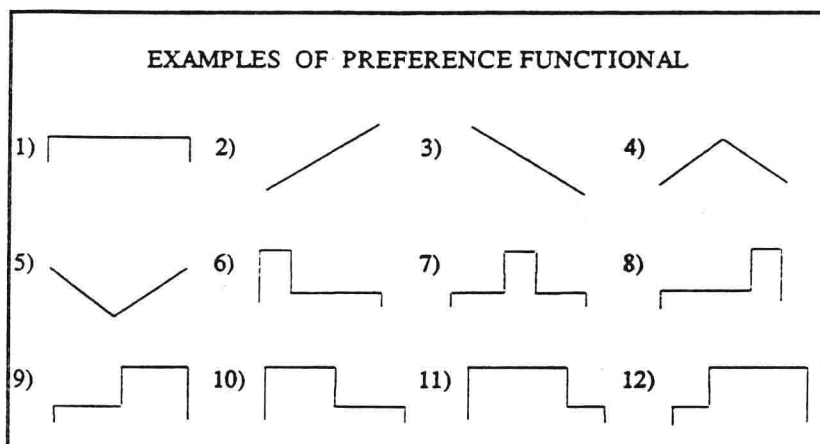
- 1) Oligonucleotide length L ;
- 2) Oligonucleotide type;
- 3) Function of preference.

The given set of parameters describes the model of oligonucleotide influence on the site functioning.

25 various variants of the preference function are examined in the system (Fig. 3). One could note that when the oligonucleotide length varies within the 1

PARAMETERS OF THE MODEL

- | | |
|--|--------------|
| 1) Length of oligonucleotides | $L = 4$ |
| 2) Variants of oligonucleotides | $V_o = 15^L$ |
| 3) Variants of the preference functional | $V_f = 25$ |



Total number of models: $N = V_o \times V_f = 15^4 \times 25 = 1.434.375$

Fig. 3. Variants of the preference function.

to L range and when 25 preference functionals are under consideration, the total number of model variants accounts to billions.

Let us fix a variant of the oligonucleotide, the variant of a preference function and consider a sample of sites (Fig. 4) to calculate H for each site. Now let us do the same for a sample of "non-sites", i.e. for sequences not having this activity. We calculate H for each "non-site".

These procedures are accomplished with the block of hypothesis generation. Then starts the block of hypothesis examination which compare the distributions of H and H^* as it is shown in Fig. 4. The more the difference between the distributions, the more significance the given type of oligonucleotides has for the site activity. Such oligonucleotides are more useful for constructing the method of this site recognition.

In our system the difference between H and H^* distributions is integrally estimated by value U which is called "Utility". The utility is calculated on the base of a special expert system designed in the Laboratory [8]. It contains knowledges on the theory of pattern recognition, theory of classification, methods of cluster analysis and the knowledge on mathematical statistics.