

## Table of Contents

<b>Introduction. About GeneNetWorks™.....</b>	<b>4</b>
<b>Part I. DNA Integration Level .....</b>	<b>6</b>
Chapter 1. Transcription Regulatory Regions Database (TRRD) .....	6
Example 1. How to view the map of regulatory regions of a gene? .....	12
Example 2. Access to TRRD database.....	13
Example 3. Access to functionally important gene systems (TRRD sections).....	13
Example 4. Searching for the genes of a particular species.....	14
Example 5. Searching for the genes by name .....	16
Example 6. Searching for the genes by name and species name .....	18
Example 7. Searching for the genes with particular chromosomal localisation .....	20
Example 8. Searching for the genes by KeyWords.....	21
Example 9. Search for the liver-specific hypersensitive to DNase I sites .....	23
Example 10. Search for the transcription factors binding sites by names .....	25
Example 11. How to search for a site by its sequence? .....	27
Example 12. How to search for sites by the factor name and the factor influence? .....	28
Example 13. How to search genes by inducer or repressor name? .....	29
Example 14. How to search expression patterns from genes that are expressed in particular organ?.....	29
Chapter 2. Site Recognition Module .....	31
1. B-DNA Site Video.....	31
1.1. B-DNA Site Video Databases .....	31
Example of SRS queries to the B-DNA Site Video Databases:.....	34
1.2. Software.....	35
1.2.1. bDNA Profiles DNA. Feature-Recognition Tools .....	35
Release 2003 .....	35
Example.....	37
1.2.2. DNA Property Plot.....	38
Example.....	39
Chapter 3. SELEX System .....	41
1. SELEX Knowledge Base.....	41
SELEX_DB .....	42
SELEX_BIB.....	42
2. SELEX Profiles Program.....	43
Example.....	45
Chapter 4. RegScan Module For DNA Functional Site Recognition.....	46
1. RegScan Databases .....	46
Example of query: .....	49
2. Software For Transcription Factor Binding Site Recognition .....	51
2.1. Program BinomSite .....	51
Example.....	52
2.2. Program MMSite .....	53
Example: .....	55
2.3. Programs for DNA Functional Sites Recognition .....	57
Example: Recognising of USF binding site in the sequence of interest .....	58
2.4. RGSiteScan Program.....	60
Example.....	61

Example .....	61
2.5. RecGroup program .....	62
Example .....	63
2.6. ARGO system .....	63
Example .....	66
3. Programs for Promoter Recognition .....	68
3.1. Program to recognize eukaryotic promoters.....	68
Example .....	69
3.2. ARGO-Viewer .....	70
3.3. POLIIISCAN .....	74
4. Other Programs .....	76
4.1. NASCA program .....	76
Example .....	78
4.2. Program for estimation of stochastic complexity of genetic texts.....	80
<u>Example .....</u>	83
Chapter 5. ACTIVITY System .....	85
1. ACTIVITY Database .....	85
Examples of SRS queries to the ACTIVITY database.....	88
2. Predicting Activities of Functional Sites in DNA/RNA .....	91
Example .....	92
Chapter 6. DNA Nucleosomal Organisation.....	94
1. PROFILES Databases .....	94
Example of SRS queries to the PROFILE database.....	95
2. Software .....	99
2.1. Nucleosome binding site recognition .....	99
Example .....	100
2.2. Recognition tools .....	101
Example .....	102
<b>Part II. RNA Integration Level.....</b>	<b>104</b>
Chapter 1. Leader RNA .....	104
1. LEADER_RNA Knowledge Base .....	104
2. Software .....	107
2.1. Programs for mRNA Translatability Prediction.....	107
Examples .....	109
2.2. MatrixSS: Building of E-score plot for RNA sequence .....	113
Example .....	114
2.3. GArna Program.....	115
Example .....	116
<b>Part III. Protein Integration Level.....</b>	<b>118</b>
Chapter 1. Databases.....	118
1. EnPDB .....	118
Example 1 .....	121
Example 2 .....	124
2. PDBSite Database .....	127
Example 1 .....	129
Example 2 .....	131
Example 3 .....	133

Example 4.....	134
3. PDBSiteScan.....	136
Example.....	137
4. Artificial Selected Peptides/Proteins Database (ASPD).....	138
5. DCS - Database on residue coordination spheres.....	140
Example 1.....	142
Example 2.....	145
6. Database on local conformations of protein chains -"conformons" (ConfDB) ...	148
Chapter 2. Software .....	151
1. CRASP .....	151
1.1. Analysis of pairwise positional correlations.....	151
Example: .....	155
1.2. Analysis of protein integral physico-chemical characteristics .....	155
Example: .....	160
<b>Part IV. Gene Networks Integration Level.....</b>	<b>161</b>
Chapter 1. GeneNet database .....	161
Example 1 .....	165
Example 2 .....	168
Chapter 2. GeneNet Software .....	170
1. GeneNet Viewer. ....	170
Example.....	171
2. GeneNet Modelling.....	175
Example 1 .....	177
Example 2 .....	182
Example 3.....	183
3. Mathematical models of gene networks in SBML format.....	184
3.1 About mathematical models of gene networks in SBML format.....	184
3.2 List of mathematical models in SBML format of gene networks represented in <a href="#">GeneNet</a> database .....	184

# INTRODUCTION.

## ABOUT GENENETWORKS™

Integrated system **GeneNetWorks™** is designed for accumulation of experimental data, data navigation, data analysis, and analysis of dependencies in the field of gene expression regulation. It integrates the databases and programs for processing the data about structure and function of DNA, RNA, and proteins, together with the other information resources important for gene expression description. The unique property of above described system is that all the resources within the system **GeneNetWorks™** are divided according to the natural hierarchy of molecular genetic systems and has the following levels:

- 1) DNA;
- 2) RNA;
- 3) proteins; and
- 4) gene networks.

Each module contains:

- 1) experimental data represented as a database or some sample;
- 2) program for data analysis;
- 3) results of an automated data processing;
- 4) tools for the graphical representation of these data and the results of the data analyses.

The database system **GeneNetWorks™** is the product developed at the Laboratory of Theoretical Genetics of the Institute of Cytology and Genetics of Siberian Branch of Russian Academy of Science in Novosibirsk (Russia). Institute of Cytology and Genetics since early 80's has been developing the database system, which trade name now is **GeneNetWorks™**. The using of the Institute name in advertising and commercial purpose is allowed only after written agreement with the Institute.

Institute of Cytology and Genetics was founded in 1957 and has large experience of scientific study in fields of genetic, cytology, microbiology etc., is involved in couple of international research projects. Since early nineties, the Institute is active in commercial projects, had organized few joint ventures with commercial structures. Now it has organized internal structures for performing of commercial projects, which includes leaders of the groups, project managers and executors of project. At present time, the experience is enough to be involved into international cooperation and performance of international commercial projects.

### **The tasks that could be solved by the **GeneNetWorks™** system (GNW):**

<b>Task</b>	<b>Tool</b>	<b>Section</b>
What is known about the regulation of the gene XX?	<b>TRRD, TRRD-Viewer</b>	Part I. Chapter 1.
Which known genes are expressed in the tissue NN, or only at embryonic stage of development, or only during the S-phase of the cell cycle?	<b>TRRD</b>	Part I. Chapter 1
Which structural and functional characteristics are typical for transcription factor XX binding site?	<b>TRRD, ACTIVITY</b>	Part I. Chapters 1, 5
Which transcription factors might bind to the sequence XX?	<b>TRRD, B-DNA Site Video, RegScan Programs</b>	Part I. Chapters 1, 2.1, 4.2

Which transcription factors might regulate gene XX (where no information is available in TRRD)?	<b>B-DNA Site Video, RegScan Programs</b>	Part I. Chapters 2.1, 4.2
Which conformational or physico-chemical properties are characteristic for the sequence XX and which functional DNA sites have the same properties?	<b>B-DNA Site Video, RegScan Programs, NASCA</b>	Part I. Chapters 2.1, 4.2, 4.4.1
Which nucleotides in the sequence of the functional DNA or RNA site are important for binding of a protein?	<b>TRRD, SELEX</b>	Part I. Chapters 1, 3
Whether the sequence XX might be a promoter?	<b>RegScan Programs</b>	Part I. Chapter 4.2
Which regions of the sequence XX might be bound to a nucleosome?	<b>Nucleosome Site Recognition</b>	Part I. Chapter 6.2.1
Genes X, Y, and Z are coregulated in an experiment E (but the genes A, B, and C are not). Which sequences might be involved in the regulation?	<b>RegScan Programs</b>	Part I. Chapter 4.2
Homology to which known extended regulatory region is typical for the sequence XX?	<b>TRRD, RegScan Programs</b>	Part I. Chapters 1, 4.2
How strongly might the transcription factors bind the sequence XX?	<b>ACTIVITY</b>	Part I. Chapter 5
How to evaluate translation efficiency of RNA sequence XX?	<b>LEADER_RNA</b>	Part II. Chapter 1
Is RNA XX highly expressed?	<b>LEADER_RNA, GARna</b>	Part II. Chapters 1, 2.3
How can I use GNW to predict the secondary structure of RNA?	<b>GARna, MatrixSS Programs</b>	Part II. Chapters 2.3, 2.2
From published Phage display experiments what are the interacting residues?	<b>ASPD</b>	Part III. Chapter 1.3
How do I use GNW to find sequences of proteins and peptides selected by phage display?	<b>ASPD</b>	Part III. Chapter 1.3
How can I use GNW to search for proteins according the data on protein function and structural characteristics of protein active sites?	<b>PDBSite, EnPDB</b>	Part III. Chapters 1.2, 1.1
How can I use GNW to search for proteins based on the data about sites subjected to biochemical modifications?	<b>PDBSite</b>	Part III. Chapter 1.2
Which physico-chemical properties of the protein XX are conserved and which are variable?	<b>CRASP</b>	Part III. Chapter 2.1
How the substitution of amino acid in a protein sequence might influence chemical and physical protein characteristics?	<b>CRASP</b>	Part III. Chapter 2.1
To which extent are independent point mutations that cause amino acid substitutions in a protein sequence?	<b>CRASP</b>	Part III. Chapter 2.1
Which proteins, genes and non-proteinaceous substances support the functioning of the gene network XX?	<b>GeneNet Database</b>	Part IV. Chapter 1
In which processes participates the protein XX?	<b>GeneNet Database, GeneNet Viewer</b>	Part IV. Chapters 1, 2.1
How to display the general visualisation of structure-functional organisation of a gene network described in GNW?	<b>GeneNet Viewer</b>	Part IV. Chapter 2.1
Which models are available in GNW? How can models be changed?	<b>GeneNet Modelling</b>	Part IV. Chapter 2.2

For questions, please contact scientific supervisor of the project: Prof. N. A. Kolchanov, telephone: +7-3832-333468; e-mail: [kol@bionet.nsc.ru](mailto:kol@bionet.nsc.ru)

# PART I. DNA INTEGRATION LEVEL

## CHAPTER 1. TRANSCRIPTION REGULATORY REGIONS DATABASE (TRRD)

Release 2003

### Database description:

Transcription Regulatory Regions Database (**TRRD**) is a unique information resource, accumulating information on structural and functional organisation of transcription regulatory regions of eukaryotic genes. Only experimental information is included into **TRRD**.

### Access to **TRRD**:

<http://www.domain.com/mgs/gnw/trrd/>

### Database content:

SRS table	Description	Number of entries
<b>TRRDGENES</b>	General information about genes	1862
<b>TRRDUNITS4</b>	Descriptions of regulatory units (promoters, enhancers, silencers)	2832
<b>TRRDEXPR</b>	Gene expression patterns	10575
<b>TRRDSITES4</b>	Descriptions of transcription factor binding sites	8326
<b>TRRDFACTORS4</b>	Descriptions of transcription factors	6706
<b>TRRDLCR</b>	Descriptions of Locus Control Regions	14
<b>TRRDBIB4</b>	References	6233

### List of biological tasks that could be solved by using **TRRD**:

- to extract the list of genes by gene name and/or by species name;
- to view the map of transcription regulatory regions reconstructed by **TRRD-Viewer**
- to obtain patterns of gene expression regulation;
- to obtain data on transcription regulatory regions and transcription regulatory units involved in regulation of expression of a particular gene;
- to obtain data on structural and functional characteristics of transcription factor binding sites;
- to obtain descriptions of transcription factors used in the experiments on either binding capacities or functional activities of the corresponding sites;
- to extract references to original publications containing data listed above;
- to compare novel DNA sequences with regulatory units stored in **TRRD** using BLAST tool and find sequence alignments of sequences.

SRS tables format:

#### **TRRDGENES**

Line code	Field name	Field description
AC	GeneAC	TRRD gene accession number
ID	GeneID	TRRD identifier
DT	Updated	Date of the last update
GN	TransfacLink	Link to the TRANSFAC database
OS	Species	English and Latin names of species
NA	GeneName	Full, short gene names and synonyms
SN	GeneName_Brief	Short gene name
NG	GeneName_Full	Full gene name
SY	GeneSynonym	Synonyms of short or full gene names
EC	EnzymeClass	Enzyme classification
KW	KeyWords	Key words
CH	Chromosome	Chromosomal localisation
RG	RegRegion	Regulatory region
AP	RegUnitAC	Name and localisation of the regulatory unit; site links
PR	RegUnit	Name and localisation of the regulatory unit; site links
AG	ExperimentCodes	Cells, assay codes, reference to the paper
CE	CompElement	Composite element description
AL	Alignment	Alignment of extended regulatory regions
MP	StartPoints	Distance between start points
BI	DNABankLink	Link to EMBL/GenBank
CC	Comments	Comments
HN	HSS_AC	TRRD accession number of the hypersensitive site
HS	HSS_Position	Name of the hypersensitive site and its location
HD	HSS_Descript	Functional characteristics of the hypersensitive site
HR	HSS_Reference	Bibliographic reference
HG	HSS_Genes	Genes regulated by the LCR
HC	HSS_Comments	Comments on hypersensitive site
DR	BankLink	Links to other databases

### TRRDUNITS

Line code	Field name	Field description
ID	RegUnitAC	Accession number of a regulatory unit
GN	GeneID	Identifier of an entry in the TRRD database
RG	RegRegion	Transcription regulatory region
PR	RegUnit	The description of promoter, enhancer or silencer (name and location; start point name; the list of all site accession numbers located in this regulatory unit)
AQ	DNA_BankLink	Link to EMBL/GenBank, and the first and last nucleotide positions of the sequence according to the EMBL/GenBank
LQ	LeftTrunc	The left part of the sequences is truncated (nucleotide number)
RQ	RightTrunc	The right part of the sequences is truncated (nucleotide number)
SL	SeqLength	The length of the sequence indicated in the field SE
SE	Sequence	The nucleotide sequence of the unit
PT	PromotTisSp	The information about promoter tissue-specific characteristics
PI	PromotInd	The information about promoter induction pattern
AG	ExperimentCodes	The field includes name of cells (tissue or organ) under experiment, codes of experiments described in the articles, and reference to the original paper
SE	Sequence	DNA sequence

### TRRDEXP

Line code	Field name	Field description
RE	ExpressionPatternAC	Pattern identifier
ID	GeneID	Identifier of an entry (gene) in the TRRDGENES4
RT	ExpressionDetectionDevice	Molecular product used to estimate the expression level (protein or mRNA)
RY	CellCycleStage	Cell cycle stage (from G0 to M)
RD	StageOrgDev	Developmental stage (embryo, fetus, etc.)
RO	Organ	Organ
RU	Tissue	Tissue
RN	Cells	Cell type
RF	StageCellDiff	Cell differentiation stage
RL	ExpressionLevel	Gene expression level
RI	IndReprName	Inductor or repressor name
RH	InductionTime	Duration of the inductor's or repressor's effect
FF	Influence	The effect produced by external signal
RP	RegUnitLink	Accession number(s) of the regulatory unit(s) involved in regulation
RS	SiteLink	Accession number(s) of the site(s) involved in regulation
RC	Comments	Comments on the expression pattern
RR	Reference	Reference to the paper
RM	ExpComparison	Comparison of the expression levels from different expression patterns
CC	TextComments	Comments
RX	Sex	Sex of an organism

#### TRRDSITES4

Line code	Field name	Field description
AN	SiteAC	Site accession number
NM	SiteName	Site name
GID	GeneID	Identifier of an entry (gene) in the TRRDGENES4
NP	PreferredName	Preferred site name
AP	RegUnitAC	Accession number of a regulatory unit which includes this site
AT	FactorInfluence	Alteration of transcription level caused by binding of the factor to the site
NY	SiteNameSynonym	Site name synonym
NS	TransFacSiteReference	TRANSFAC link
WW	PredictionProgramLink	The reference to the Internet-accessible program for recognition of the site
DR	DatabaseReference	Link to the external database
TF	FactorName	A name of the protein or protein complex interacting with the site
SQ	Sequence	Site sequence
NI	SiteIndex	TRRD site index
PQ	SequencePosition	Positions of the site
PF	FootprintSequencePosition	Footprint positions
SC	SeqContradiction	In the case when the annotator finds a discrepancy in the site sequence between the paper annotated and the corresponding data from EMBL/GenBank, the sequence from the paper annotated are given in the field SeqContradiction (SC), while the site sequence corresponding to EMBL/GenBank are indicated in the fields Sequence (SQ).
PC	PosContradiction	In the case when the annotator finds a discrepancy in the site positions between the paper annotated and the corresponding data from EMBL/GenBank, the positions from the paper annotated are given in the field PosContradiction (PC), while the site positions corresponding to EMBL/GenBank are indicated in the fields SequencePosition (PQ) of the database TRRDSITES.
IP	ImportantPos	The nucleotides within transcription factor binding sites that are important for interactions with the corresponding proteins.
KK	TextComments	Comments
EF	ComparativeFactorAffinity	Comments on comparative affinity of binding sites
HM	SiteHomology	Comments on homologous sites in other genes
BF	DNA_BankLink	Position of the site sequence first nucleotide according to the EMBL/GenBank
AG	ExperimentCodes	Cells, assay codes, reference to the paper

#### TRRDFACTORS

Line code	Field name	Field description
ID	Identifier	Identifier
GI	GeneID	Identifier of an entry (gene) in the TRRDGENES4
AN	SiteAC	Site accession number
TF	FactorName	The abbreviated and full name of the factor

FS	FactorSubunitName	Factor subunit name (for heteromeric proteins)
TY	FactorNameSynonym	Synonymous names
TS	FactorOrigin	Species specificity
NF	TRANSFAC_link	TRANSFAC link
TO	FactorSource	The source of the factor (in vitro synthesised, recombinant, etc.)
TG	Organ	Organs used for isolation of the factor
TE	StageOrgDev	Stage of an organism development
TU	Tissue	Tissues used for isolation of the factor
TC	Cells	Cells used for isolation of the factor
TD	IndReprName	External signal
TR	Reference	Bibliographic reference
CC	Comments	Comments on the protein

#### TRRDLCR

Line code	Field name	Field description
AC	gene_locus_AC	gene cluster accession number
ID	gene_locus_ID	gene cluster identifier
OS	species	species
OC	Biological_classification	Biological classification (taxonomy)
LO	chromosomal_cytological_location	chromosomal and/or cytological location of the gene cluster
LM	locus_map	The map of the gene cluster
GI	gene_name	Gene name
DR	external_database	Link to SWISS-PROT, TRRD
AL	Locus_Control_Region_AC	Accession number of the Locus Control Region in TRRDLCR
LI	Locus_Control_Region_ID	Locus Control Region identifier in TRRDLCR
TL	LCR_tissue_specificity	Tissue specificity of the LCR
EL	LCR_functional_group	The functional group of LCR elements
BI	EMBL_GenBank_ID_AC	EMBL/GenBank identifier and accession number
AR	functional_group_AC	Accession number of the functional group of LCR elements in LCRTT RD
RF	function	The function of the functional group of LCR elements
NM	Name_of_RE	Name of the regulatory element
NY	Synonym_of_RE_name	Synonym of the regulatory element
RT	funct_group_tissue_specificity	Tissue specificity of the functional group of LCR elements
RE	functional_group_elements	Functional group elements (hypersensitiv sites or regulatory elements)
CC	comments	Comments
HI	HS_RE_ID	Identifier of the hypersensitiv site or regulatory element in TRRDLCR
AP	TRRD_accession_number_of_RE	Accession number of the hypersensitiv site or regulatory element in TRRDLCR
DE	HS_RE_name	Name of the hypersensitiv site or regulatory element
HT	HS_RE_tissue_specificity	Tissue specificity of the hypersensitiv site or regulatory element
DP	zero_point_of_HS_SEQUENCE	“Zero” point of the hypersensitiv site or regulatory element sequence
PT	map_localization_of_HS_RE	Approximate map localization of the hypersensitiv site or regulatory element
PS	The_1st_relatively_pos	The 1st position of the hypersensitiv site or regulatory

	sition_of_HS	element relatively to “zero” point
PE	last_position_of_HS_RE	The last position of the hypersensitiv site or regulatory element relatively to “zero” point
CD	HS_characteristic	The hypersensitiv site characteristic (strong, weak )
SQ	HS_RE_sequence	Nucleotide sequence of hypersensitiv site or regulatory element
PH	position_of_HS_RE_in_EMBL	Position of hypersensitiv site or regulatory element in EMBL/GenBank
PA	positive_negative_arguments	Positive or negative arguments
OA	organ	Organ
TA	tissue	Tissue
CA	cell_type	Cell type
AG	arguments	Experiments for hypersensitiv site or regulatory element revealing and description
AU	authors	Authors
TI	title	Title
SO	source	Source
YR	year	Year
VL	volume	Volume
IS	issue	Issue
PGF	page_from	The first page
PGL	page_to	The lat page
ML	MEDLINE_index	MEDLINE index

### TRRDBIB

Line code	Field name	Field description
ID	ReferenceId	Identifier of the reference
GID	GeneID	Identifier of an entry (gene) in the TRRDGENES4
AU	Authors	Authors of the paper
AI	ArticleInd	TRRD paper index
AD	ArticleDescr	Type of the annotated information
TI	Title	Title of the paper
SO	Journal	Journal name
VL	Volume	Volume number
YR	Year	Year
PG	Pages	Pages
ML	MEDLINE_link	MEDLINE link

## TRRD Usage Examples:

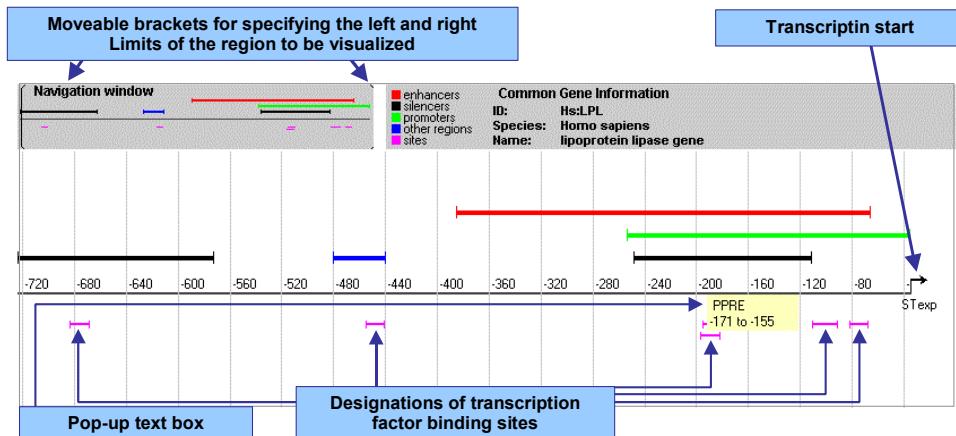
### Example 1. How to view the map of regulatory regions of a gene?

This entry is from: <b>TRRDGENES4</b>	<b>TRRDGENES4:A00410</b>
<a href="#">Save</a>	<b>GeneID</b> Hs:CYP7 ( <a href="#">TRRD Viewer</a> )
<a href="#">Link</a>	<b>Links:</b> <a href="#">Binding sites</a> <a href="#">Transcription factors</a> <a href="#">Gene expression regulation</a> <a href="#">Bibliography</a>
<a href="#">Printer Friendly</a>	<b>Updated</b> 09/10/98
	<b>GeneAC</b> A00410
	<b>TransfacLink</b> G001193
	<b>Annotators</b> Ignatieva E.V.
	<b>Species</b> human, Homo sapiens

The option to access the map of regulatory regions of a gene in question is provided by the table **TRRDGENES**. Each entry of the table has a special link [TRRD Viewer](#). By clicking this hyperlink, you access to the map of regulatory regions of corresponding gene:

**TRRD Viewer** is Java applet and is developed using JDK 1.1.8. For the Viewer to work properly, either Internet Explorer 4.0 (or further releases) or Netscape Communicator 4.61 (or further releases) are required.

Visualisation of a gene regulatory map is exemplified below. Three windows are provided: (1) navigation window, (2) text box with the relevant information and designations, and (3) window with the map of gene regulatory regions. In this figure, visualisation of the map of gene regulatory regions by **TRRD Viewer** exemplified with human cholesterol 7-alpha-hydroxylase gene (CYP7) is displayed. A hooked black arrow indicates the transcription start of a gene. Shown on the axis is the distance from the reference point (here, the transcription start). A pop-up text box (yellow rectangle) presents the information on SRE binding site.



### Options of TRRD Viewer:

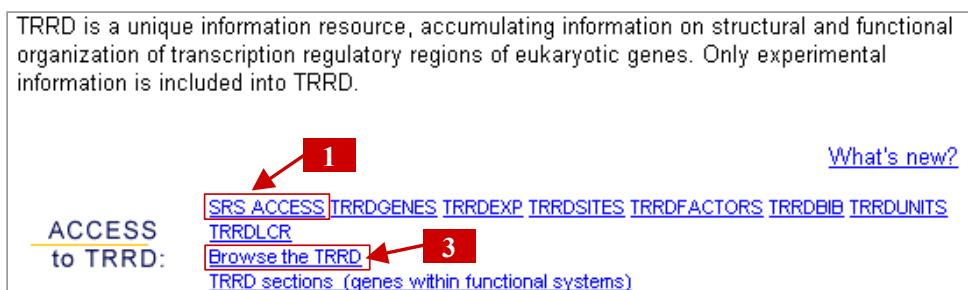
1. Movable brackets for specifying the left and right limits of the fragment to be visualised in the window below;

## I. DNA. Chapter 1. Transcription Regulatory Regions Database (TRRD)

2. Indicating a transcription factor binding site with the mouse cursor pops up the prompt with the site name and positions;
3. Indicating a regulatory unit (promoter, enhancer, or silencer) with the mouse cursor pops up the prompt with its name, reference point, and positions;
4. Clicking a regulatory unit (promoter, enhancer, or silencer) allows the binding sites contained in it to be shown, as upon clicking, all the corresponding sites will take the colour of this particular regulatory unit (green, blue, red, or black).

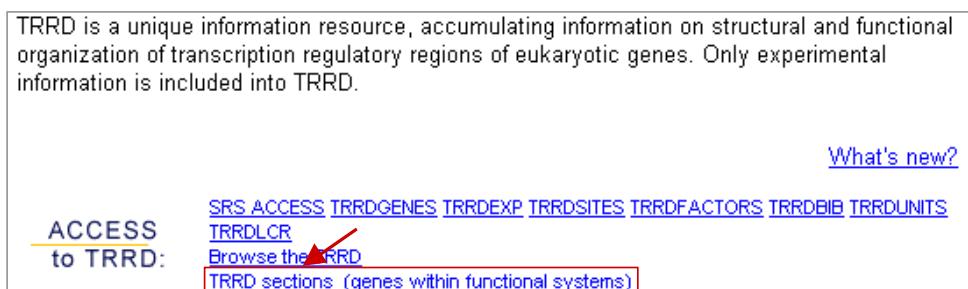
### Example 2. Access to TRRD database

Access to **TRRD** by SRS (1), BLAST (2), Browser by species and Browser by gene name (3) is available at the **TRRD home page** via corresponding links:



### Example 3. Access to functionally important gene systems (TRRD sections)

To get access to **TRRD sections**, click the link TRRD sections in the top menu of **TRRD homepage**:



This will bring up a table with section names and hyperlinks. To get access to the section you are interested in, click the corresponding link (in the figure, **ESRG-TRRD** is chosen as an example):

## I. DNA. Chapter 1. Transcription Regulatory Regions Database (TRRD)

TRRD Section	Short name and link	Compiler
Heat Shock-Induced Genes	<a href="#">HS-TRRD</a>	<a href="#">Stepanenko I.L.</a>
Interferon-Inducible Genes	<a href="#">IIG-TRRD</a>	<a href="#">Ananko E.A.</a>
Erythroid-Specific Regulated Genes	<a href="#">ESRG-TRRD</a>	<a href="#">Podkolodnaya O.A.</a>
Genes of Lipid Metabolism	<a href="#">LM-TRRD</a>	<a href="#">Ignatieva E.V.</a>
Endocrine System Genes	<a href="#">ES-TRRD</a>	<a href="#">Ignatieva E.V.</a>
Glucocorticoid-Regulated Genes	<a href="#">GR-TRRD</a>	<a href="#">Merkulova T.I.</a>
Plant Genes	<a href="#">PLANT-TRRD</a>	<a href="#">Goryachkovsky T.N.</a>
Cell Cycle Genes	<a href="#">CCG-TRRD</a>	<a href="#">Turnaev I.I.</a>
Redox-Sensitive Genes	<a href="#">ROS-TRRD</a>	<a href="#">Stepanenko I.L.</a>
Genes Expressed in Endocrine Pancreas	<a href="#">EP-TRRD</a>	<a href="#">Ignatieva E.V.</a>
Macrophage-Expressed Genes	<a href="#">MG-TRRD</a>	<a href="#">Ananko E.A.</a>
Genes, controlling blood coagulation and fibrinolysis	<a href="#">BCF-TRRD</a>	<a href="#">Khlebodarova T.M., Podkolodnaya O.A.</a>
Apoptosis Genes	<a href="#">Apoptosis-TRRD</a>	<a href="#">Stepanenko I.L.</a>
Genes, controlling circadian rhythm, and genes with circadian expression	<a href="#">CLOCK-TRRD</a>	<a href="#">Khlebodarova T.M.</a>
Genes encoding proteins involved in the Fe metabolism	<a href="#">FM-TRRD</a>	<a href="#">Mischenko E.L., Podkolodnaya O.A.</a>

### Example 4. Searching for the genes of a particular species

To display the list of genes of a particular species, use the standard (simple) Query Form of **TRRDGENES** table. Starting from **TRRD home page**, follow the hyperlink [TRRDGENES](#):

The screenshot shows the TRRD home page. On the left, there's a logo for 'TRANSCRIPTION REGULATORY REGIONS DATABASE' with a stylized 'TRRD' logo. Below it are links for 'General information' and 'How to cite TRRD?'. On the right, there's a section titled 'ACCESS to TRRD:' with several hyperlinks: 'SRS ACCESS', 'TRRDGENES' (which has a red arrow pointing to it), 'TRRDEXP', 'TRRDSITES', 'TRRDFACTORS', 'TRRDDB', 'TRRDUNITS', 'TRDLCR', 'Browse the TRRD', and 'TRRD sections (genes within functional systems)'. A 'What's new?' link is also present.

On the 'Library Information' page of **TRRDGENES** then displayed, click 'Search' button:

The screenshot shows the 'TRRDGENES' search page. At the top, there's a navigation bar with links: 'TOP PAGE', 'QUERY', 'RESULTS', 'SESSIONS', 'VIEWS', 'DATABANKS', and 'HELP'. Below the navigation bar, there's a search form with a 'Name' input field containing 'TRRDGENES4'. To the right of the input field is a 'Search' button. A red arrow points to the 'Search' button.

You will get the standard Query Form. Choose the data field 'Species' (arrow 2) by selecting the appropriate item in the drop-down list box 'GeneAC' (arrow 1 in the figure below):

## I. DNA. Chapter 1. Transcription Regulatory Regions Database (TRRD)

The screenshot shows the TRRD Query interface. At the top, there are tabs for TOP PAGE, QUERY (which is selected), RESULTS, SESSIONS, VIEWS, DATABANKS, and HELP. Below the tabs, there is a search bar with the text "search TRRDGENES4" and an "Info" button. A dropdown menu titled "about field GeneAC" is open, showing options like GeneAC, GeneID, TransfacLink, Species, GeneName, GeneName\_Brief, and GeneName\_Full. A red arrow labeled "1" points to the "Species" option in the dropdown. Another red arrow labeled "2" points to the "Species" dropdown in the main query form. To the right of the dropdowns, there is a text input field "retrieve entries of type [Entry]" with a dropdown arrow.

You may enter into the query text box the species name (or its part followed by wildcard symbol asterisk) in English or in Latin. Note that if you input a common species name in English, you may get the list of genes referring to several related species. For example, by entering 'mouse', you will get the list of genes of four murine species (*Mus caroli*, *Mus domesticus*, *Mus hortulans*, *Mus musculus*).

Example: search for the genes of grey (*Rattus norvegicus*), but not black (*Rattus rattus*) rat by the species name in English.

Enter the terms 'rat' (see arrow 1 in the figure below) and 'black rat' (arrow 2) into the appropriate consecutive 'Species' fields (unnumbered arrows). In the 'combine searches with' drop-down list select the 'BUTNOT' item (arrow 3), and then click the 'Submit Query' button (arrow 4):

This screenshot shows the TRRD Query interface with the following configuration:

- Search Bar:** search TRRDGENES4
- Dropdowns:** The "Species" dropdown has "Species" selected. The second "Species" dropdown has "rat" entered. The third "Species" dropdown has "Species" selected. The fourth "Species" dropdown has "black rat" entered.
- Buttons:** The "Submit Query" button is highlighted with a red arrow labeled "4". The "BUTNOT" button in the sidebar is highlighted with a red arrow labeled "3".
- Sidebar Options:** "append wildcards to words" is checked. "combine searches with BUTNOT" is selected. "Number of entries to display per page" is set to 30.
- Bottom Options:** "Use predefined view" and "Create your own view" buttons, and a dropdown menu set to "\* Names only \*".

The query brings up the resulting list of genes containing 291 entries. Use the vertical scrolling bar to fit the content of interest in the window. Click 'next' button to view the next portion of resulting data:

The screenshot shows a web-based application with a navigation bar at the top: TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, and DATABASES. Below the navigation bar, there is a search query: "Query "[trrdgenes4-Species: rat\*] ! [trrdgenes4-Species: black\* & rat\* | black rat\*]" found 291 entries". A red box highlights the number "291". To the right of the number is a "next" button, also highlighted with a red box and an arrow pointing to it.

**Perform operation**

- on all but selected
- on selected

**Link**

**Save**

**View**

\* Names only \*

Number of entries to display per page **30**

**Printer Friendly**

A list of gene identifiers follows:

- [TRRDGENES4:A00356](#)
- [TRRDGENES4:A00118](#)
- [TRRDGENES4:A00090](#)
- [TRRDGENES4:A00311](#)
- [TRRDGENES4:A00394](#)
- [TRRDGENES4:A00296](#)
- [TRRDGENES4:A00160](#)
- [TRRDGENES4:A00416](#)
- [TRRDGENES4:A00662](#)
- [TRRDGENES4:A00848](#)
- [TRRDGENES4:A00873](#)
- [TRRDGENES4:A00914](#)
- [TRRDGENES4:A00940](#)
- [TRRDGENES4:A00010](#)
- [TRRDGENES4:A00631](#)

### Example 5. Searching for the genes by name

To display the genes with the similar name, use again the standard Query Form for the **TRRDGENES** table.

Use the following data fields (by selecting appropriate items in the drop-down list boxes): 'GeneName\_Full' and/or 'GeneSynonym', in case you know the complete gene name, or 'GeneName\_Brief' and/or 'GeneSynonym' if you know the brief gene name.

Example: querying the genes encoding interleukines.

Enter the term 'interleukin\*' (see arrow 1 in the figure below) into the data field 'GeneName\_Full' (arrow 2). Check the drop-down list box 'Use predefined view', the default option \*Names only\* should be selected, if not, select it (arrow 3). Click 'Submit Query' button (4):

## I. DNA. Chapter 1. Transcription Regulatory Regions Database (TRRD)

The screenshot shows the TRRD Query interface. The search bar contains "search TRRDGENES4". Below it, the query input field contains "GeneName\_Full interleukin\*". To the right of the input fields is a dropdown menu set to "GeneAC". A red arrow labeled "1" points to the first entry in the list. A red arrow labeled "2" points to the "Submit Query" button. A red arrow labeled "3" points to the "Use predefined view" dropdown set to "\* Names only\*". A red arrow labeled "4" points to the "append wildcards to words" checkbox.

This will bring up the list of entries accumulated in **TRRD**. In total, 18 genes will be displayed (see arrow 1):

To receive more details click any hyperlink denoted in accordance with the gene index in **TRRD**, for example the first one on the top of the list (arrow 2).

The screenshot shows the TRRD Results interface. At the top, a message says "Query '[trrdgenes4-GeneName\_Full: interleukin\*]' found 18 entries". The main area displays a list of 18 gene entries, each preceded by a checkbox. The first entry, "TRRDGENES4:A00476", is highlighted with a red box and a red arrow labeled "1". On the left, there is a sidebar with options for "Perform operation" (radio buttons for "on all but selected" and "on selected"), "Link", "Save", "View", and a dropdown set to "\* Names only\*". Below the sidebar is a "Number of entries to display per page" dropdown set to "30". A "Printer Friendly" button is at the bottom.

The resulting detailed view is shown below. Its content exceeds the window dimensions; so use the vertical scrolling bar marked by arrows to fit it in:

This entry is from: [TRRDGENES4:A00476](#)

**GeneID**: Hs:IL5 ([TRRD Viewer](#))

**Links:**

- [Binding sites](#)
- [Transcription factors](#)
- [Gene expression regulation](#)
- [Bibliography](#)

**Updated:** 04/10/99

**GeneAC:** A00476

**TransfacLink:** G000315

**Annotators:** Kel O.

**Species:** human, Homo sapiens

**GeneName\_Brief:** IL-5

**GeneName\_Full:** interleukin 5

**DNABankLink:** EMBL: [HSIL5](#); [J03478](#); ST:509

**DataBankLink:** SWISS-PROT; [IL5\\_HUMAN](#); [P05113](#) (Expasy server)

**EPD Class:** 6.1.5.9.

#### Example 6. Searching for the genes by name and species name

Use the standard Query Form for the **TRRDGENES** table in this and all below mentioned cases unless otherwise specified.

Example. Search for the gene encoding rat serine protease inhibitor 2.3.

Enter the term 'rat' into the field 'Species' (steps 1 and 2 in the figures). Compose the query by one of the following three options:

- a) You may use the complete gene name. In the data field 'GeneName\_Full', enter the term 'serine protease inhibitor gene 2.3' (see steps 3 and 4 in the figure below):

search [TRRDGENES4](#)

**Info** about field GeneAC

**Submit Query** (Step 6)

separate multiple values by & (and), | (or), ! (and not)

Species	rat (Step 2)	3
GeneName_Full	serine protease inhibitor gene 2.3 (Step 4)	4
GeneAC		5
GeneAC		6

append wildcards to words

combine searches with AND (Step 5)

retrieve entries of type Entry

## I. DNA. Chapter 1. Transcription Regulatory Regions Database (TRRD)

b) If you know the synonym, e.g., in our case 'kallikrein-binding protein', enter this term into the data field 'GeneSynonym':

The screenshot shows the TRRDGENES4 search interface. The search bar contains 'search TRRDGENES4'. The 'QUERY' tab is selected. The search parameters are as follows:

- Species:** rat (highlighted by red arrow 2)
- GeneSynonym:** kallikrein-binding protein (highlighted by red arrow 4)
- GeneAC:** (empty)
- GeneAC:** (empty)

On the left, there are configuration options:

- append wildcards to words:** checked (highlighted by red arrow 3)
- combine searches with:** AND (highlighted by red arrow 5)

**Submit Query** button (highlighted by red arrow 6) and **Info about field GeneAC** link.

c) If you know the standard gene abbreviation, for example, 'Spi 2.3', enter it into the data field 'GeneName\_Brief'.

Combine the searches with the data field 'Species' and one of the fields mentioned above by 'AND', then submit the query (steps 5 and 6 in the figures):

The screenshot shows the TRRDGENES4 search interface. The search bar contains 'search TRRDGENES4'. The 'QUERY' tab is selected. The search parameters are as follows:

- Species:** rat (highlighted by red arrow 2)
- GeneName\_Brief:** Spi 2.3 (highlighted by red arrow 4)
- GeneAC:** (empty)
- GeneAC:** (empty)

On the left, there are configuration options:

- append wildcards to words:** checked (highlighted by red arrow 3)
- combine searches with:** AND (highlighted by red arrow 5)

**Submit Query** button (highlighted by red arrow 6) and **Info about field GeneAC** link.

The query brings up the following result. Data items of interest are underlined by red:

## I. DNA. Chapter 1. Transcription Regulatory Regions Database (TRRD)

This entry is from:  
[TRRDGENES4](#)

[Save](#)  
[Link](#)  
[Printer Friendly](#)

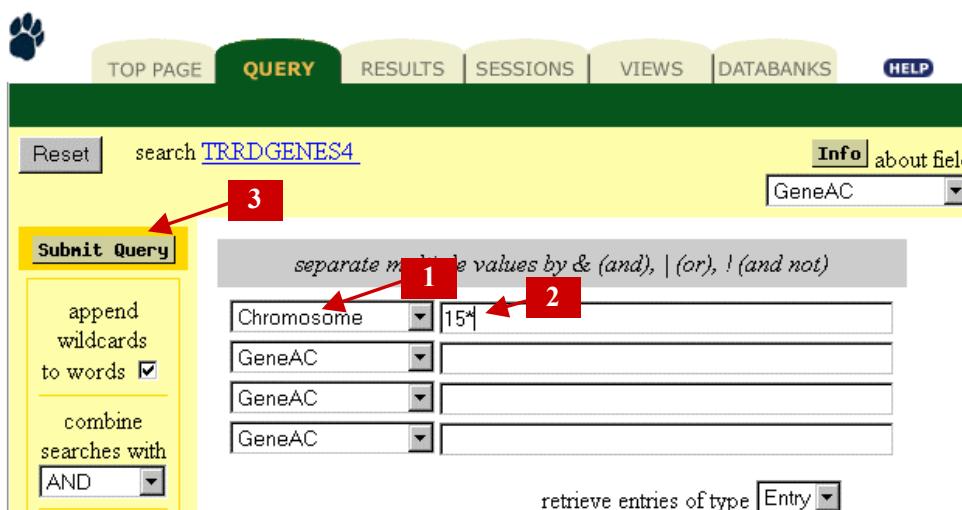
**GeneID**  
Rn:SPI23 ([TRRD Viewer](#))  
**Links:**  
[Binding sites](#)  
[Transcription factors](#)  
[Gene expression regulation](#)  
[Bibliography](#)  
**Updated**  
12/09/98  
**GeneAC**  
A00106  
**TransfacLink**  
G000793  
**Annotators**  
Kochetov A.V., Ananko E.  
**Species**  
rat, Rattus norvegicus  
**GeneName\_Brief**  
Spi 2.3  
**GeneName\_Full**  
serine protease inhibitor gene 2.3  
**GeneSynonym**  
SPI  
serpin  
gene 2.3  
contrapsin-like protease inhibitor precursor  
kallikrein-binding protein  
GHR-p63



### Example 7. Searching for the genes with particular chromosomal localisation

To display the list of genes localised at particular chromosome, select the data field 'Chromosome' and enter the chromosome number, for example '14\*' or '15\*'. Note that the asterisk should supplement the number, because the database may contain the record '14' as well as the record '14q32.32'.

Example: querying the genes localised at chromosome 15:



The screenshot shows the TRRD search interface with the following details:

- Top Navigation Bar:** Includes links for TOP PAGE, QUERY (highlighted in green), RESULTS, SESSIONS, VIEWS, DATABASES, and HELP.
- Search Bar:** Contains a "Reset" button, the search term "search TRRDGENES4", and an "Info about field" dropdown set to "GeneAC".
- Query Form:**
  - Submit Query** button (highlighted with a red box labeled 3).
  - Search Fields:** A list of fields with dropdown menus and input boxes. The "Chromosome" field is highlighted with a red box labeled 1 and contains the value "15\*". The "GeneAC" field next to it is highlighted with a red box labeled 2.
  - Instructions:** A note above the fields says "separate multiple values by & (and), | (or), ! (and not)".
  - Advanced Options:** On the left, there are checkboxes for "append wildcards to words" (checked) and "combine searches with" (radio buttons for AND, OR, NOT). The "AND" option is selected.
  - Buttons at the bottom:** "retrieve entries of type" dropdown set to "Entry" and a "Submit" button.

The result of querying contains 8 entries:

The screenshot shows a web-based search interface. At the top, there are tabs: TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, and DATABANK. Below these, a search bar displays the query "trrdgenes4-Chromosome: 15\*" and indicates "15 entries" found. To the left, a sidebar titled "Perform operation" has radio buttons for "on all but selected" (selected) and "on selected". It also contains buttons for "Link", "Save", and "View". The main area lists 15 entries, each preceded by a checkbox and a blue hyperlink: TRRDGENES4:A00951, TRRDGENES4:A00091, TRRDGENES4:A01184, TRRDGENES4:A00283, TRRDGENES4:A00498, TRRDGENES4:A01636, TRRDGENES4:A01626, and TRRDGENES4:A00152.

Click the hyperlink pointed to by the arrow 1 in the figure above to receive the complete record data of the first found gene (depicted below):

This screenshot shows the detailed record for the gene TRRDGENES4:A00036. The left sidebar contains buttons for "Save", "Link", and "Printer Friendly". The main content area includes the following fields:

- GeneID:** Hs: AAC ([TRRD Viewer](#))
- Links:**
  - [Binding sites](#)
  - [Transcription factors](#)
  - [Gene expression regulation](#)
  - [Bibliography](#)
- Updated:** 04/04/00
- GeneAC:** A00036
- TransfacLink:** G000193
- Annotators:** Kel O., Sidorenko
- Species:** human, Homo sapiens
- GeneName\_Brief:** CA-ACT
- GeneName\_Full:** cardiac alpha-actin
- DNABankLink:** GenBank; [HUMACTCA](#); M13483; ST:486
- EPD\_Class:** 6.1.2.5.1.
- KeyWords:** structural protein, contractile protein
- Chromosome:** [15q11-qter](#)
- RegRegion:** 5' region

#### Example 8. Searching for the genes by KeyWords

Enter some known to you characteristic properties of the gene (or its product) to be found into the field 'KeyWords'. By querying the keywords, it is possible to extract the following information:

## I. DNA. Chapter 1. Transcription Regulatory Regions Database (TRRD)

- the list of genes with characteristic structural features, for example, by querying TRRD for keywords 'TATA-less promoter', 'multiple transcription initiation sites', 'GC-rich promoter', 'alternative exon', 'alternative promoters', 'Alu repeat';
- the list of genes with typical patterns of functioning, by ordering the query by the keywords 'alternative splicing', 'allergic triggering', 'autoregulation', 'DNA damage-induced';
- the list of genes by their products, by means of searching for keywords 'glycoprotein', 'hormone precursor', 'transcription factor', 'growth factor', 'enzyme';
- the list of genes according to peculiarities of their products' functioning by querying the keywords 'DNA binding protein', 'G-protein coupled receptor', 'mitogen', 'monooxygenase';
- the list of genes participating in functioning of this or that organism system, or the organism as a whole, by querying the keywords 'glyoxylate cycle', 'lignin synthesis', 'iron homeostasis', 'intracellular protein', 'housekeeping gene', 'homeotic gene', 'heme biosynthesis pathway', 'growth regulation', 'embryo-specific gene';
- the list of genes acting in formation of this or that organism structure, by ordering the query by keywords 'exocuticle', 'eye lens';
- the list of genes according to evolutionary homology, for example, by querying the keywords 'growth hormone family', 'heat shock-related').

Example: to search for the genes by combination of keywords 'TATA-less promoter', 'protooncogene', 'autoregulation'. For this purpose, it is necessary to select the data field 'KeyWords', to input these keywords, to select combination 'AND' and to submit the query:

TOP PAGE    QUERY    RESULTS    SESSIONS    VIEWS    DATABANKS    HELP

Reset search TRRDGENES4 Info about field GeneAC

Submit Query

append wildcards to words

combine searches with  5

separate multiple values by & (and), | (or), ! (and not)

KeyWords 1 TATA-less promoter 2  
KeyWords 3 protooncogene 3  
KeyWords 4 autoregulation 4

GeneAC 5

retrieve entries of type Entry

This will bring up the resulting list of **TRRD** entries. In our example, 4 entries were found (see arrows 1):

TOP PAGE    QUERY    RESULTS    SESSIONS    VIEWS    DATABANKS    HELP

Reset

Query "[trrdgenes4-KeyWords: TATA-less\* & promoter\* | TATA-less promoter\*] & [trrdgenes4-KeyWords: protooncogene\*] & [trrdgenes4-KeyWords: autoregulation\*]" [found 4 entries]

Perform operation 2

on all but selected  
 on selected 4

TRRDGENES4:A00348  
 TRRDGENES4:A00290  
 TRRDGENES4:A00347  
 TRRDGENES4:A00308

## I. DNA. Chapter 1. Transcription Regulatory Regions Database (TRRD)

By clicking the above marked hyperlink [A00348](#) (arrow 2), you may receive the complete record data related to this gene (shown in the figure below). The keywords due to which the entry of the gene A00348 was extracted are underlined by red and pointed to by arrows:

The screenshot displays a gene record for Hs:E2F1 (GeneID A00348). The 'KeyWords' section is highlighted with red arrows pointing to the following terms: DNA binding protein, transcription factor, nuclear protein, autoregulation, cell cycle regulator, regulation of cell proliferation, multigene family, TATA-less promoter, and protooncogene.

### Example 9. Search for the liver-specific hypersensitive to DNase I sites

Use the field 'HSS\_Descript'. Input into this field the name of an organ with an asterisk, 'liver\*' in our case. You may choose the \*Complete entries\* option in the 'Use predefined view' drop-down list box for convenience. Submit the query:

## I. DNA. Chapter 1. Transcription Regulatory Regions Database (TRRD)

The screenshot shows the TRRD Query interface. The search term 'search TRRDGENES4' is entered in the search bar. The query parameters are set as follows:

- Field:** GeneAC (highlighted with red box 1)
- Value:** liver\* (highlighted with red box 2)
- Search Type:** AND (highlighted with red box 3)
- View:** \*Complete entries\* (highlighted with red box 4)

Other settings include 'append wildcards to words' checked and 'Number of entries to display per page' set to 30.

A fragment of one of the TRRD entries found is shown in the figure:

The screenshot shows the TRRD Results interface for entry A00481. The entry details are as follows:

- GeneID:** Mm:HNF3G (TRRD Viewer) (highlighted with red arrow)
- Links:**
  - Binding sites
  - Transcription factors
  - Gene expression regulation
  - Bibliography
- Updated:** 20/07/00
- GeneAC:** A00481
- TransfacLink:**
- Annotators:** Merkulov V.M.
- Species:** mouse, Mus musculus
- GeneName\_Brief:** HNF3gamma
- GeneName\_Full:** hepatocyte nuclear factor 3 gamma
- DNABankLink:** EMBL; MMHNF3GEN; Y12559; ST:
- DataBankLink:**

The left sidebar includes options for 'Perform operation' (radio buttons for 'on all but selected' or 'on selected'), 'Link', 'Save', and 'View'. The 'View' button is highlighted with a red arrow. Other settings include 'Number of entries to display per page' set to 30 and a 'Printer Friendly' link.

### Example 10. Search for the transcription factors binding sites by names

To display the list of transcription factors binding sites, use the standard Query Form of **TRRDSITES** table. Starting from **TRRD home page**, follow the hyperlink [TRRDSITES](#):

The main content area displays a brief description of TRRD: "TRRD is a unique information resource, accumulating information on structural and functional organization of transcription regulatory regions of eukaryotic genes. Only experimental information is included into TRRD." Below this is a section titled "ACCESS to TRRD:" which includes links for "SRS ACCESS TRRDGENES TRRDEXP", "TRRDSITES" (which is highlighted with a red arrow), "TRRFACATORS", "TRRDDBB", and "TRRDUNITS". There is also a link "What's new?" and a link "TRRDLCR" with the subtext "Browse the TRRD sections (genes within functional systems)".

On **TRRDSITES** 'Library Information' page then displayed, click the 'Search' button:

This will bring up the standard Query Form (see the figure below).

Example: Let us search for insulin responsive element and insulin responsive sequence by transcription factor binding site names and their synonyms. Since abbreviations of these regulatory sequences, IRE and IRS, coincide with abbreviations of interferone-regulated sites, it is recommended to set the query by complete name. Enter the term 'insulin\*' (arrows 1) into the data fields 'SiteName' and 'SiteNameSynonym' (arrows 2) and combine the searches with 'OR' (arrow 3):

After clicking the screen button 'Submit Query' (arrow 4) you may look through the following resulting list of transcription factor binding sites:

## I. DNA. Chapter 1. Transcription Regulatory Regions Database (TRRD)

Query "[trrdsites4-SiteName: insulin\*] | [trrdsites4-SiteNameSynonym: insulin\*]" found 21 entries

**Perform operation**

on all but selected  
 on selected

**Link**  
**Save**  
**View**  
 \* Names only \*

Number of entries to display per page **30**

**Printer Friendly**

- [TRRDSITES4:S4095](#)
- [TRRDSITES4:S6395](#)
- [TRRDSITES4:S6784](#)
- [TRRDSITES4:S4642](#)
- [TRRDSITES4:S7319](#)
- [TRRDSITES4:S7136](#)
- [TRRDSITES4:S7489](#)
- [TRRDSITES4:S6590](#)
- [TRRDSITES4:S6592](#)
- [TRRDSITES4:S6595](#)
- [TRRDSITES4:S6595](#)
- [TRRDSITES4:S6600](#)
- [TRRDSITES4:S7094](#)
- [TRRDSITES4:S7095](#)
- [TRRDSITES4:S7098](#)
- [TRRDSITES4:S4277](#)
- [TRRDSITES4:S9047](#)

Click the first, for example, item in the list (pointed to by red arrow in the figure above) and by means of hyperlink you will receive the complete entry data of appropriate site:

This entry is from: [TRRDSITES4:S4277](#)

**SiteAC**  
 S4277  
**GeneID**  
 Gene: Hs:MTP  
**RegUnitAC**  
 REGULATORY UNIT: P01192  
**SiteName**  
 IRE; insulin response element  
**SiteIndex**  
 2  
**FactorInfluence**  
 decrease  
**Sequence**  
 gcAGCCCCACCTACGtt  
**SequencePosition**  
 -125 to -110  
**DNA\_BankLink**  
 X83013:620  
**ExperimentCodes**  
 HepG2: 6.2, 6.5 (insulin) [[Hagan](#) D.L. et al., 1994]  
 7.1 [[Hagan](#) D.L. et al., 1994]

### Example 11. How to search for a site by its sequence?

Use again the standard Query Form for the **TRRDSITES** table. Select the data field 'Sequence' and enter or copy from clipboard buffer the sequence of transcription factor binding site to be found. Note that the sequence to be entered should not contain any symbols except nucleotides (the only exception is for asterisks that may border the sequence, see the figure). In case of searching for long sequences, the sequences with deletions, insertions, duplications, etc., you should better use the option 'Blast search TRRD database'.

Let us search for transcription factor binding sites containing the sequence 'tcaaggcag'. Enter the sequence to be found **\*tcaaggcag\*** limited by asterisks into the data field 'Sequence':

Submit the query and view the resulting list of transcription factor binding sites:

Click any site name (which is the hyperlink at the same time) to receive the complete entry content:

### Example 12. How to search for sites by the factor name and the factor influence?

Start from the standard Query Form for the TRRDSITES table, as in previous cases. In order to search for the sites binding to particular transcription factors, which influence gene transcription in a particular manner, it is necessary to make combined query by two fields, 'SiteName' and 'FactorInfluence'. For this purpose, it is sufficed to enter into the data field 'FactorInfluence' one of the terms 'increase' or 'decrease'.

Let us search for transcription factor binding sites Pit-1, interaction of which with the factor of the same name activates transcription. Into the data field 'SiteName', enter the site name or its part (in this example, 'Pit-1\*'). In the field 'FactorInfluence' enter 'increase'. Combine the search by 'AND', submit query and view the resulting list.

The screenshot shows the 'QUERY' tab selected in the top navigation bar. The search term 'search TRRDSITES4' is entered. The 'Info' button shows 'about field SiteAC'. The 'Submit Query' button is highlighted with a red arrow. On the left, there are checkboxes for 'append wildcards to words' and 'combine searches with', and a dropdown menu set to 'AND' with a red arrow pointing to it. The search criteria are listed in a table:

separate multiple values by & (and),   (or), ! (and not)	
SiteName	Pit-1*
FactorInfluence	increase
SiteAC	
SiteAC	

The screenshot shows the 'RESULTS' tab selected. The query string is displayed as 'Query "[trrdsites4-SiteName: Pit-1\*] & [trrdsites4-FactorInfluence: increase\*]" found 12 entries'. The results table lists 12 entries, each preceded by a checkbox:

<input type="checkbox"/> TRRDSITES4.S3951
<input type="checkbox"/> TRRDSITES4.S3953
<input type="checkbox"/> TRRDSITES4.S4332
<input type="checkbox"/> TRRDSITES4.S4333
<input type="checkbox"/> TRRDSITES4.S4334
<input type="checkbox"/> TRRDSITES4.S214
<input type="checkbox"/> TRRDSITES4.S215
<input type="checkbox"/> TRRDSITES4.S4983
<input type="checkbox"/> TRRDSITES4.S4984
<input type="checkbox"/> TRRDSITES4.S4971
<input type="checkbox"/> TRRDSITES4.S4972
<input type="checkbox"/> TRRDSITES4.S4991

On the left, there is a 'Perform operation' section with radio buttons for 'on all but selected' (selected) and 'on selected', and buttons for 'Link', 'Save', and 'View'. A dropdown menu 'Names only' is also present. Below it, a dropdown menu 'Number of entries to display per page' is set to 30.

### Example 13. How to search genes by inducer or repressor name?

Use the standard Query Form for the **TRRDEXP** table. In order to extract the list of genes, about which there is information that their expression is modified by a particular factor, it is necessary to query the data field 'IndReprName'. Under the term 'factor' we understand any impact, for example, hormone, transcription factor, diet, pathology, surgery, or emotional action.

Example: Search for the genes that contain information about the influence of TSH (thyroid stimulating hormone) on their expression. Into the field 'IndReprName', enter 'TSH\*':

The screenshot shows the 'Submit Query' interface for the TRRDEXP table. On the left, there are checkboxes for 'append wildcards to words' and 'combine searches with AND'. The main search area has a dropdown for 'IndReprName' set to 'TSH\*' and a dropdown for 'ExpressionPatternAC'. Below these are two scrollable lists of gene records. The first list contains genes starting with 'TRRDEXP4:A00979...' up to 'TRRDEXP4:A00962.015'. The second list shows 'TRRDEXP4:A00962.003' in detail. To the right of this gene record, arrows point from the 'IndReprName' field in the search form to the 'IndReprName' field in the gene details, and from the 'TSH\*' entry in the search form to the 'TSH (thyroid stimulating hormone)' entry in the gene details. The gene details also show other fields like 'GeneID', 'Rn:TSHR', 'mRNA', 'Cells', 'ExpressionLevel', 'IndReprName', 'Influence', and 'InductionTime'.

### Example 14. How to search expression patterns from genes that are expressed in particular organ?

Use again the standard Query Form for the **TRRDEXP** table. To search for expression patterns related to the genes expressed in a particular organ, it is necessary to make query by two fields, 'Organ' and 'ExpressionLevel'. It is necessary to account that among the terms in the field 'ExpressionLevel' may be the following: 'exclusive expression', 'high', 'low', 'maximal level', 'med', 'minimal level', 'none', 'present', 'undetectable', 'very high', 'very low'. To search for the gene patterns that are expressed in particular organ, it is necessary to except from the query the patterns of expression that contain the records 'none' and 'undetectable' in the field 'ExpressionLevel'.

Example: Search for expression patterns related to genes expressed in liver. Choose for the querying data field 'Organ', enter the term 'liver'. Choose the querying field 'ExpressionLevel', enter the terms 'none' and 'undetectable', combine them by the sign '&'. Combine the search by 'BUTNOT'. Submit the query:

## I. DNA. Chapter 1. Transcription Regulatory Regions Database (TRRD)

The screenshot shows the TRRD interface with the following details:

**Top Navigation Bar:** TOP PAGE, QUERY (highlighted in yellow), RESULTS, SESSIONS, VIEWS, DATABASES, HELP.

**Search Bar:** Reset, search [TRRDEXP4](#), Info about field ExpressionPatternAC.

**Query Form:** Submit Query, append wildcards to words (checkbox checked), combine searches with BUTNOT (dropdown menu).

**Search Criteria:** Organ: liver, ExpressionLevel: none & undetectable, ExpressionPatternAC: (empty).

**Result Summary:** Query "[trrdexp4-Organ: liver\*] ! [trrdexp4-ExpressionLevel: none\*]" found 851 entries.

**Perform operation:** on all but selected (radio button selected), Link, Save, View, \*Names only\* dropdown, Number of entries to display per page: 30.

**Result List:** A list of 851 entries, starting with TRRDEXP4:A00374.002, each with a checkbox. The first entry is highlighted.

**Result Details (for TRRDEXP4:A00150.006):**

- ExpressionPatternAC: A00150.006
- GeneID: Hs:APOD
- ExpressionDetectionDevice: mRNA
- Organ: liver
- ExpressionLevel: present
- Reference: [Drayna D. et al., 1986]

## CHAPTER 2. SITE RECOGNITION MODULE

### 1. B-DNA Site Video

#### 1.1. B-DNA Site Video Databases

Release 2003

##### Databases description:

These databases are designed to study the sets of various transcription factor binding sites, providing evidence that transcription factor binding sites are characterized by specific sets of significant conformational and physico-chemical DNA properties.

##### Access to B-DNA Site Video Databases:

<http://www.domain.com/mgs/gnw/bdna/>

##### Databases content:

SRS table	Description	Number of entries
FEATURES	Knowledge base on significant B-DNA properties of sites	51
PROPERTY	Database on sequence-dependent conformational and physico-chemical B-DNA properties	38
SAMPLES	Database on functional site sequences	77
PROFILE_LIST	This SRS table accumulates sets of significant physico-chemical properties profiles of the sites.	6

##### List of biological tasks that could be solved by using the B-DNA Site Video Databases:

- to extract specific sets of significant conformational and physico-chemical DNA properties of various transcription factor binding sites;
- to browse information about all conformational and physico-chemical DNA properties that differ significantly for the sequences of transcription factor binding sites and random sequences (the sequences of transcription factor binding sites occurring in nature are stored in the SAMPLES database);
- for a site given, by using the B-DNA features selected for recognition of this particular site, to generate the C-program recognising this site on the basis of the site properties stored in the database B-DNA-VIDEO;

**SRS tables format:**

<b>FEATURES</b>		
Line code	Field name	Field description
MI	Entry ID	This field indicates an identifier.
MN	Site Name	This field contains the name of transcription factor binding site.
HN	SCIENTIST	This field contains the link to SCIENTIST database and indicates a contributor of the entry.
DR	SRS-links	This field contains the links to the supplementary databases installed under SRS (SAMPLES, etc.).
WW	Web-link to Recognition Tools	This field contains the link to the Web-based tools implementing each C program documented within the entry and aimed to recognition of transcription factor binding site by its conformational and physico-chemical features within an arbitrary DNA sequence.
DP	Link to PROPERTY database	This field contains the link to the PROPERTY database installed under SRS.
PV	Property Name	This field contains the short name of the property investigated.
HL	Feature Indicator	This field indicates high or low deviation of the conformational or physico-chemical feature of the site studied from the corresponding value for random sequences.
AB	Analyzed Region	This field indicates the DNA region, for which the conformational or physico-chemical feature differing significantly the sites studied from the random sequences was revealed. Positions are given in bp relative to the site center.
UT	Utility	This field indicates the utility of the conformational or physico-chemical feature for discrimination of the sites studied from the random sequences.
ST	Means, Standard Deviation, False Negatives for Control Sequences	This field provides the mean value of the property at [a; b] region (first estimate), the standard deviation of this mean value (second estimate), and the I type error rate (false negatives), which have been determined for the independent control sequences of the site by the program given below in the field C-CODE.
NT	Means, Standard Deviation, False Positives for Random Sequences	This field provides the mean value of the property at [a; b] region (first estimate), the standard deviation of this mean value (second estimate), and the II type error rate (false positives), which have been determined for random sequences by the program given below in the field C-CODE.
FG	Graphical Representation of Test Results	This field represents the links to the graphical representation of the results obtained by the program, given below in the field C-CODE, over the independent control sequences of the site.
C-CODE		This field contains the recognition program by the sequence-dependent conformational or physico-chemical feature or Mean recognition in the 'C' language of the ANSI standard. Each C program documented within the B-DNA-Video entry has a check-box in the MENU window of the Recognition Tools (see field WW).

## PROPERTY

See Chapter 5.1. ACTIVITY Databases.

### SAMPLES

Line code	Field name	Field description
FI	SampleID	brief name of the sample
NM	SampleName	name of the sample and explanation of its biological meaning
OR	Organization	title of the organization
AU	Author	the first and last name of the author
DA	Date	date of creation
LU	LustUpdate	the date and the author of the last update/modification (Last update); in case of no modification, the date and the authors name are from the DA and AU fields
FV	FormatVersion	Number of the Format Version
ST	SiteTypeDescription	ST {X,Y} [Left_border,Right_border] Point; Description; Factor_name. - (Structure) description of the site (not a particular site, but type of the sites); several ST fields are allowed; X=0....9 - the number of site batch (sites are joined into batches if they in the aggregate they represent an integral structure-function unit); Y=0....9 – the number of the site in a site batch; [Left_border,Right_border] – The position of the site in case this particular site type possesses a fixed location relative to a specific point, for example, relative to transcription start. If no, put nothing; Point – the point relative to which the site location is fixed; If the site has no fixed position, put nothing Description - conventional name of the site (abbreviated name if available) Factor_name - a name of site-binding factor
WA		reference to 'Aligned' database
WF		reference to 'Features' database
WW		hyperlink to x-ray structure of a binding complex
ID	CardID	identifier of a card
AC	CardAC	accession number of a card
OS	Specia	organism specia
OC	Taxon	organism taxon
DR	SourceDatabase	reference to the source database
CC	Comments	comments
FT	FeatureTable	the field contains positions of site described in ST field. The format is as follows: FT {X,Y} [left;right]; Method X = 0...9 → number of the site batch Y = 0...9 → number of the site in the site batch [left;right] → positions from the start of the sequence Method = EXP (experimental), GBS (Gibbs Sampler alignment), RCG (Recognition Group alignment);
SQ	Sequence	nucleotide sequence

## PROFILE\_LIST

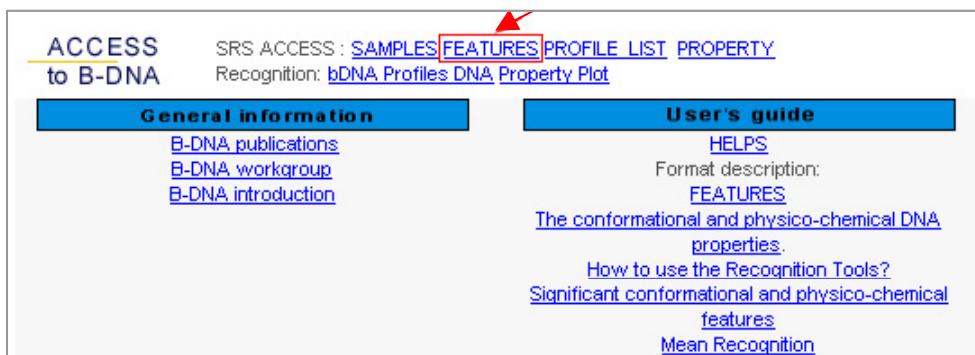
See Chapter 6.1. PROFILES Databases.

### Example of SRS queries to the B-DNA Site Video Databases:

What are conformational and physico-chemical DNA properties differing significantly the transcription factor binding site AP-1 sequences from the random sequences?

To make such a query, you should perform the following:

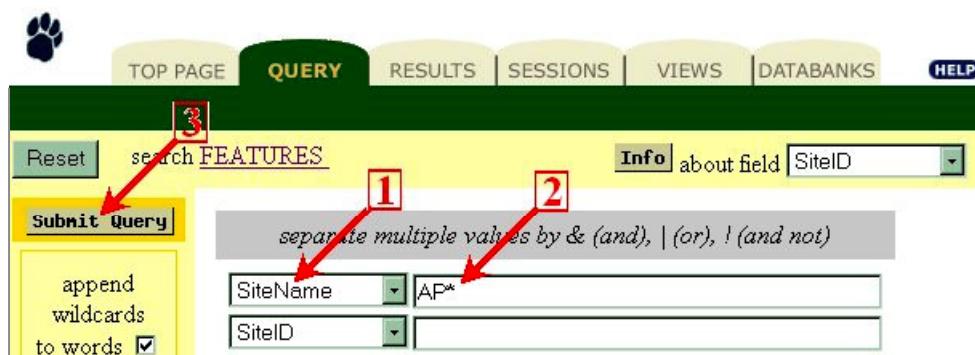
1. Choose 'FEATURES' SRS table on the page 'access to B-DNA':



2. This will bring up the home page of the chosen 'FEATURES' SRS table. Click the button 'Search':



3. Select the fields to be searched for from the list. You will need the field 'SiteName' (1), Type terms to be searched in the text window: 'AP\*' (2) (An asterisk marks "any symbol"). Then click the button 'Submit Query' (3):



4. This will bring up the resulting window with the list of matching entries (FEATURES AP-1). Click the link:



5. The query result will be displayed as the complete list of entries containing significant B-DNA features of the AP-1 site. Red arrows denote links to the recognition tools for this site (1) and graphical representation of the results for chosen property (2):

FEATURES:AP-1

```
MI AP-1
MN AP-1 transcription factor binding DNA-region
YY
HN SCI00002
YY
DR SAMPLES: AP-1;
YY
WW GALLERY, http://www.sgi.sscc.ru/Programs/bdna/gallery/AP1\_bGal.html
YY
WW PROGRAM, http://www.sgi.sscc.ru/Programs/bdna/api\_hdna.html
YY
RM Here, the ST-lines are containing the control test results that
RM has been obtained with using the control 50%-subsets of this site
RM sequences randomly selected, the NT-lines for 1000 random sequences.
XX
CF SEQUENCE-DEPENDENT CONFORMATIONAL FEATURE
CT PROPERTY AVERAGED FOR REGION [A:B]
DP P0000001
PV Twist
HL Highest
AB -18 4
UT 0.650
ST 36.443 (0.584) 36.5%
NT 36.105 (0.645) 39.0%
FG DIAGRAM, http://www.sgi.sscc.ru/Programs/bdna/images/AP1\_b00.html
XX
```

Red arrows point to the links: (1) [AP-1](#) and (2) [http://www.sgi.sscc.ru/Programs/bdna/gallery/AP1\\_bGal.html](http://www.sgi.sscc.ru/Programs/bdna/gallery/AP1_bGal.html).

Comments and questions are welcome to Mikhail P. Ponomarenko ([pon@bionet.nsc.ru](mailto:pon@bionet.nsc.ru)).

## 1.2. Software

### 1.2.1. bDNA Profiles DNA. Feature-Recognition Tools

Release 2003

**Program description:**

The transcription factor binding sites are characterized by specific sets of significant conformational and physico-chemical DNA properties. The site recognition programs are generated by the sets of significant conformational and physico-chemical DNA properties (DNA features).

### Access to Recognition Tools:

<http://www.domain.com/mgs/gnw/bdna/>; then choose the link [bDNA Profiles DNA](#):

**The biological task that could be solved by using the program:** DNA functional site recognition by most significant DNA conformational and physico-chemical properties.

#### Data input:

Input the DNA sequence of interest into the text window (1). Sequence should be in a plain text format (a, t, g, c in upper- or lower-cases, spaces or tabulation are accepted). The sequence to be input for analysis should be of maximal length of 32 kbp and minimal length of 100 bp.

### Program options:

Select the necessary recognition program by clicking one of the radio buttons (2). Each one refers to a significant conformational or physico-chemical DNA property. Generalized (mean) recognition profile is calculated using all known significant conformational and physico-chemical features of the site, contained in **B-DNA features database**.

### Program execution

Start the tools processing by clicking the button 'Execute' (arrow 3 in the figure above).

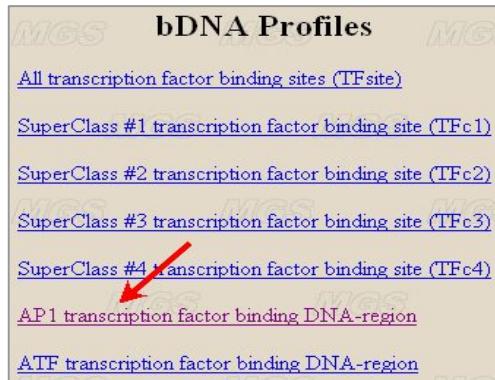
### Data output

The tools output represents the profile of the score value. The positive peaks of this profile pinpoint to potential site recognized.

### Example

Search for potential binding sites of the AP-1 transcription factor in promoter region of the human c-ets-1 protooncogene.

- 1) On the **B-DNA home page**, click the link [bDNA Profiles DNA](#) (see above, "2) Access to Recognition Tools" and illustrating figure).
- 2) On the page 'bDNA Profiles', click the link [AP1 transcription factor binding DNA-region](#):



- 3) Input into the text-box (arrow 1) the following sequence:

```
acgcctgac tcaagatccg gctggagtcc aataactccta  
aaggccttg aggacacggg ctcacgaatc ccctgcgcct  
gcctgcacgc tcgcttcatc cacatgcctc acgtcctgtg  
tgtcagtctt tgtggaatga atgatgtaca cgcaacttgg  
aaactatgct gctactggg gggggcgaga gcgggtgacc  
aaggcctcaa gaatgcgtgg agaatcagac ggactttccc  
gaaacggtgg aggccggctg tgcacccagc ctgcacaccc  
gctccggcc cttccggccc ctgcctggc tccgaggcccc  
ggggctccac gcactgctcc tccgcggctcg gcccggcccc  
gctgcgtccca gccccttctt tcgccttggg ccggggcgga  
gattggccgc gtgctggcc cggccggccgc tccccggccct  
gccccgacgc cccgccccctc gct
```

(human c-ets-1 protooncogene sequence from -461 to +1 bp relatively transcription start site, GenBank; X65469) Use program options set by default. Click the button 'Execute' (arrow 2):

**AP-1 transcription factor binding DNA-region**

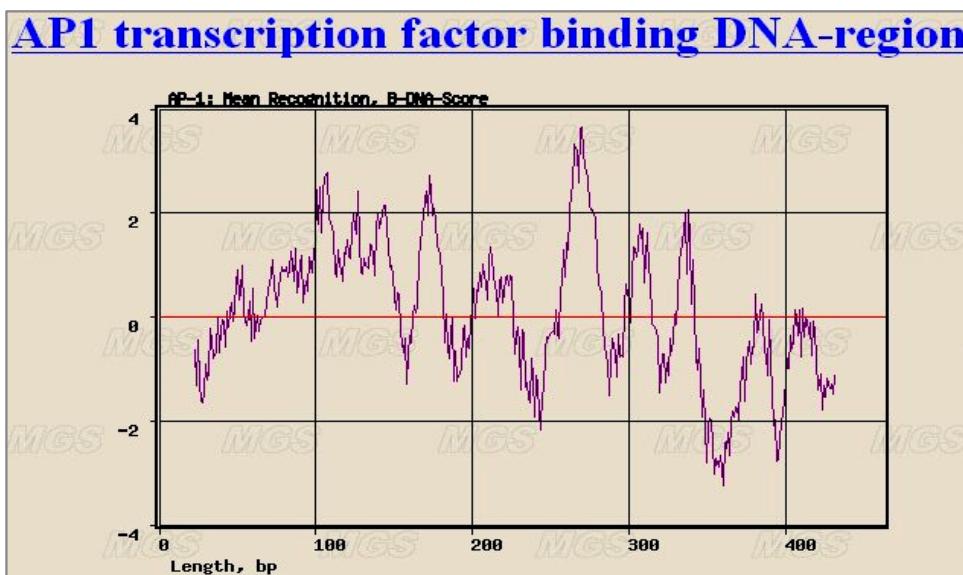
Input DNA Sequence :  from Screen: **1**

acgcctgac	tcaagatccg	gctggaggcc	aatactccca
aagccctttg	aggacacggg	ctcacgaatc	ccctgcccct

from DB: **2**  Bases Available:  
SRS5 from Heidelberg (EMBL) by ID

from File:   [File formats here.](#)

4) Data output window:



Comments and questions are welcome to Mikhail P. Ponomarenko ([pon@bionet.nsc.ru](mailto:pon@bionet.nsc.ru)).

### 1.2.2. DNA Property Plot

Release 2003

#### Program description:

This program is designed to construct the profile of DNA conformational or physico-chemical properties.

#### Access to Property Plot:

<http://www.domain.com/mgs/gnw/bdna/>; then choose the link [Property Plot](#).

**The biological task that could be solved by using the program:** investigation of DNA conformational or physico-chemical properties.

### Data input

Input the DNA sequence of interest into the text-box 'Input DNA Sequence' (arrow 1 in the figure below). Sequence should be in plain textual format. Only the symbols a, g, c, t or A, T, G, C are permissible, line feeds and blanks are ignored. Sequence length should be not less than 10 bp and not longer then 1000 bp in order to view the results in details.

### Program options:

- Select the necessary DNA conformational physico-chemical property from the list by clicking the field 'Property' to choose the bDNA physico-chemical or conformational property of interest from the Property database (arrow 2).
- Set the values in the fields 'Start Position' and 'End Position' to set the scope of the sequence region of interest (arrow 3).
- Select the averaging window size value in the field 'Window Size' (arrow 4).

### Program execution

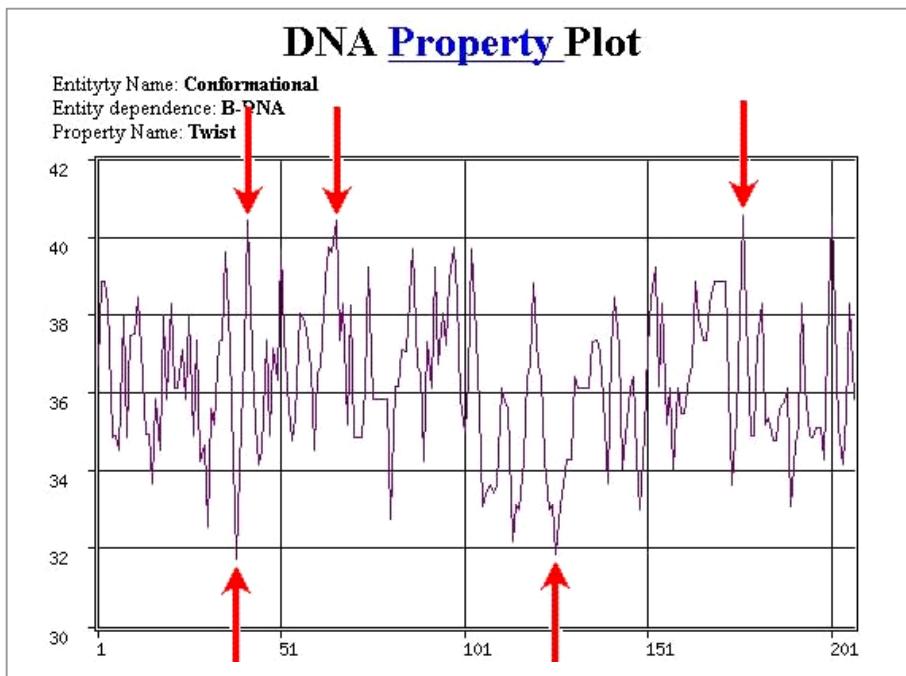
Start the program execution by clicking the button 'Execute' (5).

### Data output

The output represents the profile of the selected physico-chemical property. The Y value on the plot corresponds to selected property magnitude at positions numbered on the X-axis.

### Example

Data output for random sequence calculated for values of 'Twist' conformational property. Sequence length = 207, Window size = 3. Arrows mark regions displaying maximal and minimal twist values:



Comments and questions are welcome to Mikhail P. Ponomarenko ([pon@bionet.nsc.ru](mailto:pon@bionet.nsc.ru)).

## CHAPTER 3. SELEX SYSTEM

### 1. SELEX Knowledge Base

Release 2003

#### Description:

**SELEX Knowledge Base** is designed for accumulation of experimental data on selected affinity-enriched sequences from different combinatorial libraries.

#### Access to SELEX Knowledge Base:

<http://www.domain.com/mgs/gnw/selex/>

#### SELEX Knowledge Base content:

SRS table	Description	Number of entries
SELEX_DB	a database on selected randomized DNA/RNA sequences	116
SELEX_BIB	a database on annotated in SELEX_DB papers	70
SELEX_TOOLS	a database on computer programs for recognition of sites accumulated in SELEX_DB	32

#### List of biological tasks that could be solved by using the SELEX Knowledge Base:

- to extract the sequences binding to a particular protein, which DNA or RNA binding site is annotated in **SELEX**;
- to extract sequence positions that are most important for binding of a particular DNA/RNA site;
- to recognise the site documented in **SELEX** within an arbitrary sequence in on-line mode;
- for a site accumulated in **SELEX**, to generate the C-program recognising this site according to appropriate weight matrix calculated by data stored in the database **SELEX**;
- planning of novel experiments applying **SELEX** technology.

**SELEX Knowledge Base SRS tables formats:**

**SELEX DB**

Line code	Field name	Field description
ID	CardID	Identifier
AC	Accession	Accession number
CR	Referee	Reference on SELEX_DB database
OS	Specia	Organism
OC	Taxon	Taxon
NF	LigandName	Name of a ligand
KW	Key Words	Keywords
DA		Date of creation
DT		Date of the last update
FV		Release number
MN		Name of an entry
RF		Reference to the literature source
TE		Templates for amplification
EX		Type of an experiment
EC		Experimental conditions
NS		Number of sequences
AA		Aligned sequences as they are represented in the original paper
WA		Weight impact of the A nt at functionally important positions
WT		Weight impact of the T nt at functionally important positions
WG		Weight impact of the G nt at functionally important positions
WC		Weight impact of the C nt at functionally important positions
CN		Consensus sequence
DR		Links to the other databases
WW		Link to recognition tools
NM		Number of sequences in the set
SQ		Sequences
CC		Comments of an annotator concerning the functional role of a factor or peculiarities of consensus evaluation

**SELEX BIB**

Line code	Field name	Field description
Id	CardID	Reference identifier
AU	Authors	Reference authors
TI	PaperTitle	Reference heading
SO	Source	Reference journal
VL	VolumeNumber	Reference journal volume
IS	Issue	Reference journal issue
YR	Year	Reference year

**SELEX TOOLS**

Line code	Field name	Field description
Id	SiteID	Site identifier
MN	SiteName	Name of RNA/DNA site binding the ligand
SC	SC	Homology Score
HN		Author of the entry
DR SELEX DB		Link to the SELEX_DB SRS Table
WW		Link to Recognition Tools
RM		Comments
AB		Sequence region
ST		Means, standard deviation, false negatives
NT		Means, standard deviation, false positives
C-CODE		C-code procedure for calculation of the data for an arbitrary DNA

Comments and questions are welcome to Mikhail Ponomarenko ([pon@bionet.nsc.ru](mailto:pon@bionet.nsc.ru))

## 2. SELEX Profiles Program

Release 2003

### Program description:

The aim of the analysis is to determine positions of DNA in an arbitrary sequence binding to particular proteins, information about which is stored in the **SELEX\_DB** database. The recognition procedure is based on 15 different variants of calculation, including homology score, matrix similarity, weighting consensus match scores, etc.

### The biological task that could be solved by using the SELEX Profiles program:

Searching for DNA sites binding to some proteins or transcription factors (in total 23 DNA binding sites) in an arbitrary sequence entered by a user.

### Access to the SELEX Profiles program:

<http://www.domain.com/mgs/gnw/selex/>; link [Selex Profiles](#).

### Data input

Input the data into the input textbox (1) shown in the figure below. The sequences should be entered in upper- or lower-case letters; line feeds and blanks are ignored.

SELEX\_DB: DNA binding the transcription factor YY1, core CCAT

Input DNA Sequence :  1

from Screen:

from DE: 3 2

Bases Available: SRS5 from Heidelberg (EMBL) by ID

from File:  Browse... File formats here.

Execute Reset form

**Program options:**

By clicking the hyperlink, choose the protein or transcription factor from the list, to which you would compare your sequence with:

**Selex Profiles**

These modules are a part of [GeneExpress](#).  
[SELEX](#) database is available through [SRS5](#).

[Retinoid receptor-related Testis-assotiated Receptor, RTR](#)

[DNA specificities of the AHR, ARNT, and SIM proteins](#)

[DNA specificities of the AHR-ARNT heterodimer \(hAHR\)](#)

[DNA sequence binding bHLH-zip domain of N-Myc protein](#)

[DNA specificities of the ARNT homorodimer](#)

[DNA specificities of the SIM-\(ARNT homorodimer\) complex](#)

[Downstream DNA sites binding the hXBP protein](#)

[Upstream DNA sites binding the hXBP protein](#)

[DNA binding the transcription factor YY1, core CCAT](#)



The sequence to be analysed could be entered in different ways: from the screen in the textbox (arrow 1 in the previous figure), by typing in from the keyboard or by cut & paste operation; or from file in 'FASTA' format: in this case, click the 'BROWSE' button and select the source file (arrow 2).

Choose one recognition model from the list by clicking the appropriate radio-button:

Select one of 20 YY1-related recognition model (the site core CCAT) listed below:

(A) Fifteen YY1(CCAT)-site recognitions optimized over the YY1 SELEXed data containing the YY1 core CCAT:

YY1(CCAT): Frequency matrix for Homology Score  
 YY1(CCAT): Frequency matrix for Information content  
 YY1(CCAT): Frequency matrix for Matrix similarity  
 YY1(CCAT): Frequency matrix for Sensibility score  
 YY1(CCAT): Frequency matrix for Fisher's Linear Discriminant  
 YY1(CCAT): Frequency matrix for Free energy change  
 YY1(CCAT): Frequency matrix for Chance amount  
 YY1(CCAT): Consensus with Homology Score  
 YY1(CCAT): Consensus with Information content



**Program execution:**

Click the button 'Execute' to execute the program (see arrow3 in the figure).

**Data output:**

The program output represents the recognition score profile.

**Example**

Let us analyse the fragment of the Moloney murine leukemia virus complete genome (EMBL:J02255) containing YY1 transcription factor.

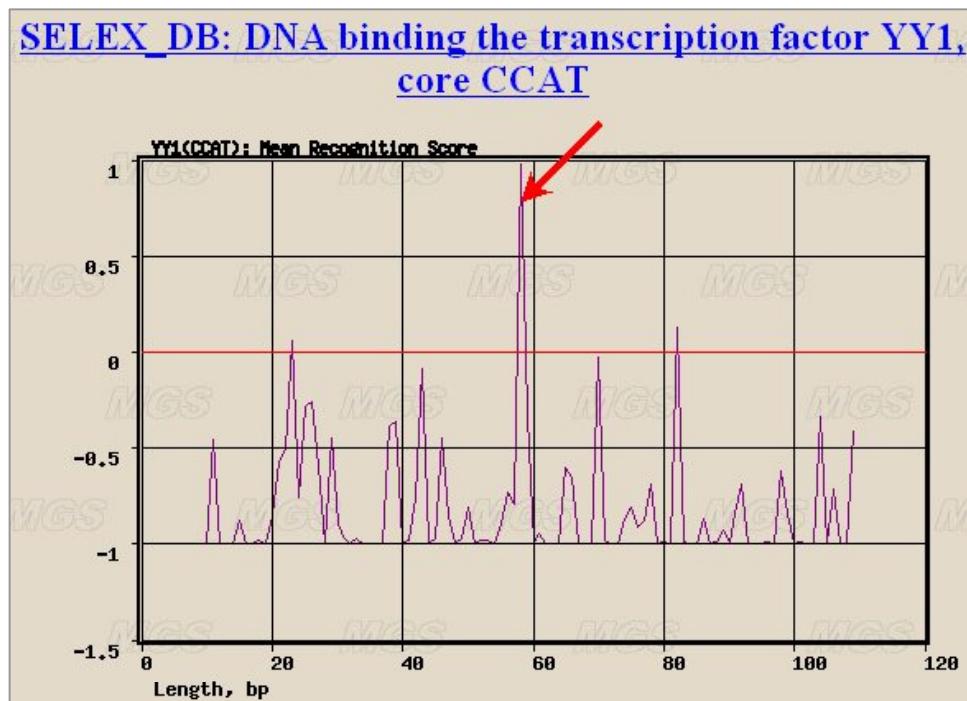
Select recognition model 'SELEX\_DB: DNA binding the transcription factor YY1, core CCAT'.

For analysis, insert the fragment of this sequence in-between positions 7805-7924 into the text-box by choosing the option 'from Screen'.

Choose 'Mean Recognition Score' procedure.

Click the button 'Execute'.

The output window displays the YY-1 recognition score profile with the pick marked by arrow, which corresponds to the experimentally identified YY-1 transcription factor binding site (positions 7860-7868).



Comments and questions are welcome to Mikhail Ponomarenko ([pon@bionet.nsc.ru](mailto:pon@bionet.nsc.ru))

## CHAPTER 4. REGSCAN MODULE FOR DNA FUNCTIONAL SITE RECOGNITION

### 1. RegScan Databases

Release 2003

#### Databases description:

**MATRIX** Database accumulates oligonucleotide frequency matrices of transcription factor binding sites, the natural sequences of which are contained in the **SAMPLES** Database (see **Chapter 2.1.1. B-DNA Site Video Databases**).

#### Access to RegScan Databases:

<http://www.domain.com/mgs/gnw/regscan/>

link [MATRIX](#) for **MATRIX** database

link [ALIGNED](#) for **ALIGNED** database

link [Property](#) for **PROPERTY** database

#### Databases content:

SRS table	Description	Number of entries
MATRIX	Database of oligonucleotide frequency matrices of transcription factor binding sites	42
ALIGNED	A database of regulatory active genomic regions aligned by Gibbs Sampler. Investigation of transcriptional factors binding sites of eukaryotes.	43
PROPERTY	A distributed and intelligent database for the activities of the functional sites in DNA and RNA.	38

#### The biological task that could be solved by using RegScan Databases:

Recognition of functional sites by their specific oligonucleotide content.

**SRS tables format:**

**MATRIX**

Line code	Field name	Field description
MI	SiteID	This field indicates an identifier.
MN	SiteName	This field contains name of the transcription factor binding site.
HN	SCIENTIST	This field contains link to SCIENTIST database and indicates a contributor of the entry.
DR	SRS-link	This field contains the links to the supplementary databases installed under SRS (SAMPLES, ALIGNED, etc.).
WW	Web-link to Recognition Tools	This field contains the link to the Web-based tools implementing each C program documented within the entry to recognize the transcription factor binding site by its oligonucleotide frequency matrices within an arbitrary DNA sequence.
DP	Oligonucleotide Alphabet Length	This field indicates the length of <a href="#">oligonucleotide alphabet</a> used by the recognition program given below in the field ' <a href="#">C-CODE</a> '.
PV	Oligonucleotide Alphabet	This field indicates the <a href="#">oligonucleotide alphabet</a> used by the recognition program given below in the field ' <a href="#">C-CODE</a> '.
AB	Analyzed Region	This field indicates the DNA region, for which the weight matrix was constructed. Positions are given in bp relative to the first position of the site core multiply aligned by the Gibbs potential function (Lawrence C., 1994, Comput. Chem., 18, 255-258) (see <a href="#">ALIGNED</a> database).
ST	Means, Standard Deviation, False Negatives for Control Sequences	This field provides the mean value of the <a href="#">Score</a> (first estimate), the standard deviation of this mean value (second estimate), and the I type error rate (false negatives), which have been determined on the independent control sequences of the site by the program given below in the field ' <a href="#">C-CODE</a> '.
NT	Means, Standard Deviation, False Positives for Random Sequences	This field provides the mean value of the <a href="#">Score</a> (first estimate), the standard deviation of this mean value (second estimate), and the II type error rate (false positives), which have been determined on the random sequences by the program given below in the field ' <a href="#">C-CODE</a> '.
FG	Graphical Representation of Test Results	This field represents the links to the graphical representation of the results obtained by the program, given below in the field ' <a href="#">C-CODE</a> ', over the independent control sequences of the site.
C-CODE	C-CODE	Recognition program by the <a href="#">Partial recognition Score</a> or <a href="#">Mean recognition Score</a> in the C language of the ANSI standard. Each C program documented within the MATRIX entry has a check-box in the MENU window of the Recognition Tools (field <a href="#">Web-link to Recognition Tools</a> )

**ALIGNED**

Line code	Field name	Field description
FI	SampleID	brief name of the sample
NM	SampleName	name of the sample and explanation of its biological meaning
OR	Organization	title of the organization
AU	Author	the first and last name of the author
DA	Date	date of creation
LU	LustUpdate	the date and the author of the last update/modification (Last update); in case of no modification, the date and the authors name are from the DA and AU fields
FV	FormatVersion	Number of the Format Version
ST	SiteTypeDescription	<p>ST {X,Y} [Left_border,Right_border] Point; Description; Factor_name. - (Structure) description of the site (not a particular site, but type of the sites); several ST fields are allowed;</p> <p>X=0....9 - the number of site batch (sites are joined into batches if they in the aggregate they represent an integral structure-function unit);</p> <p>Y=0....9 – the number of the site in a site batch;</p> <p>[Left_border,Right_border] – The position of the site in case this particular site type possesses a fixed location relative to a specific point, for example, relative to transcription start. If no, put nothing;</p> <p>Point – the point relative to which the site location is fixed; If the site has no fixed position, put nothing</p> <p>Description - conventional name of the site (abbreviated name if available)</p> <p>Factor_name - a name of the site-binding factor</p>
WS		reference to 'SAMPLES' database
WC		reference to 'CONSENSUS' database
WM		reference to 'MATRIX' database
ID	CardID	identifier of an entry
AC	CardAC	accession number of an entry
OS	Specia	organism specia
OC	Taxon	organism taxon
DR	SourceDatabase	reference to the source database
FT	FeatureTable	<p>This field contains positions of site described in ST field. The format is as follows:</p> <p>FT {X,Y} [left;right]; Method</p> <p>X = 0...9 → number of the site batch</p> <p>Y = 0...9 → number of the site in the site batch</p> <p>[left;right] → positions from the start of the sequence</p> <p>Method = EXP (experimental), GBS (Gibbs Sampler alignment), RCG (Recognition Group alignment);</p>
SQ	Sequence	nucleotide sequence

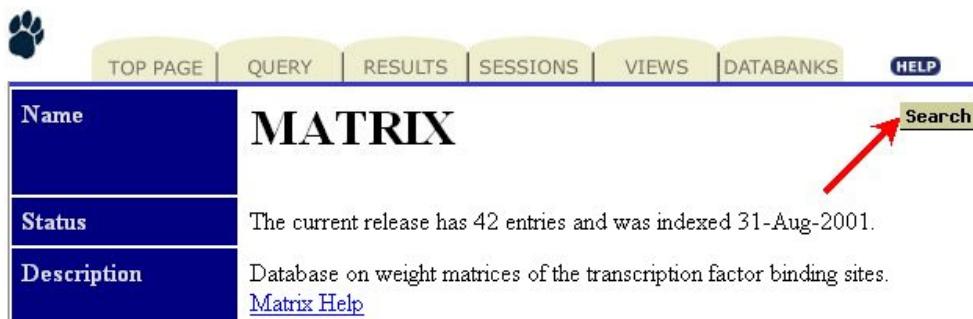
## PROPERTY

See **Chapter 5. ACTIVITY Databases.**

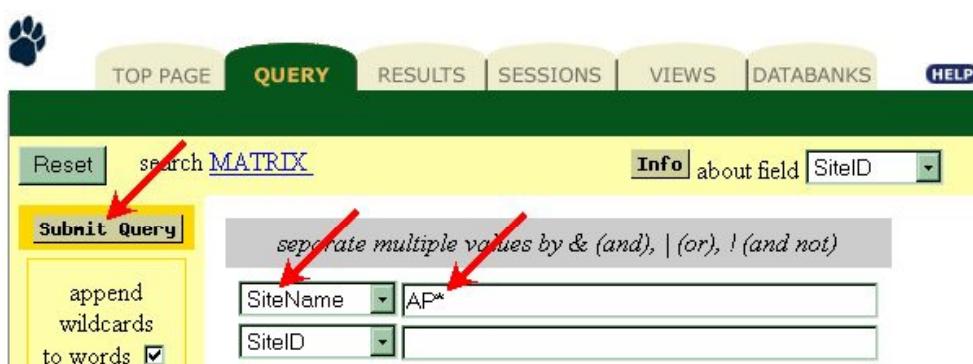
**Example of query:**

To extract C-code for the AP-1 transcription factor recognition from the **MATRIX Database**, or apply WWW-accessible **AP-1 Recognition Tool** to an arbitrary sequence of interest.

1. Click the button 'Search':



2. Select the fields to be searched for from the list. You will need the field 'SiteName'. Type term to be searched in the text window: 'AP\*'; and click the button 'Submit Query':



3. This will bring up the resulting window with the list of matching entries (MATRIX: AP-1):

```
MATRIX:AP-1

MI AP-1
MN AP-1 transcription factor binding region
YY
HN SCI00001
YY
DR SAMPLES: AP-1;
DR ALIGNED: AP-1;
YY
WW http://www.sqi.sscc.ru/Programs/Freq/ap1\_freq.html
YY
RM Here, the ST-lines are containing the control test results that
RM has been obtained with using the control 50%-subsets of this site
RM sequences randomly selected and the 1000 random DNA sequences.
XX
CF SEQUENCE-DEPENDENT RECOGNIZING PROCEDURE
CT FREQUENCY OF REGION [A;B]
DP lbp
PV (A,T,G,C)
AB -3 12
ST 0.762 (0.384) 2.9%
NT -0.853 (0.462) 5.3%
FG http://www.sqi.sscc.ru/Programs/Freq/images/AP1\_alf\_.html
XX
C-CODE
/*
 * (00) AP-1: lbp-FREQUENCY of the alphabet {A,T,G,C} */
double AP1_lbp_Freq_ATGC (char *seq){
    double A[ 15]={
        0.263, 0.237, 0.263, 0.026, 0.026, 0.868, 0.079, 0.053, 0.211, 0.763,
        0.158, 0.237, 0.132, 0.211, 0.263};
    double T[ 15]={
        0.368, 0.211, 0.053, 0.895, 0.053, 0.053, 0.053, 0.684, 0.079, 0.079,
        0.263, 0.316, 0.132, 0.184, 0.237};
    double G[ 15]={
        0.132, 0.395, 0.342, 0.026, 0.868, 0.053, 0.026, 0.184, 0.026, 0.053,
        0.342, 0.184, 0.395, 0.263, 0.263};
    double C[ 15]={
        0.237, 0.158, 0.342, 0.053, 0.053, 0.026, 0.842, 0.079, 0.684, 0.105,
        0.237, 0.263, 0.342, 0.342, 0.237};
    double X;int AP1Len=15, AP1Start=3, len, i;char nuc, *s;
    X=0.;len=strlen(seq);if(len < AP1Start) return(-9999.);s=&seq[-AP1Start];
    for(i=0;i < AP1Len;i++){nuc=s[i];switch(nuc){
        case'A': X+=A[i];break;case'T': X+=T[i];break;
        case'G': X+=G[i];break;case'C': X+=C[i];break;}}
    return ((X- 5.3409)/ 1.9207);}
XX
```

This entry displays the frequency matrices for AP-1 site, including C-code for recognition program (marked with red rectangle) and recognition tools for this site. The link to the entry is marked by arrow:

Comments and questions are welcome to Mikhail P. Ponomarenko ([pon@bionet.nsc.ru](mailto:pon@bionet.nsc.ru)).

## 2. Software For Transcription Factor Binding Site Recognition

### 2.1. Program BinomSite

Release 2003

#### Program description:

This software was developed to search for potential transcription factor binding sites by using binomial criterion for estimation of similarity score between the regions of the sequence analysed and the sequences of transcription factor binding sites.

#### Access to BinomSite:

<http://www.domain.com/mgs/gnw/regscan/> link [BinomSite](#)

#### Biological task that could be solved by using BinomSite:

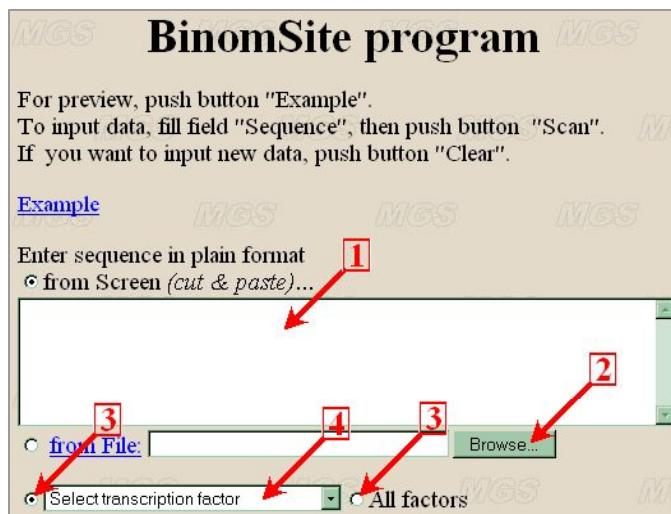
This program searches for potential transcription factor binding sites in the arbitrary sequence entered by a user.

#### Data input

To search for transcription factor binding sites, do the following:

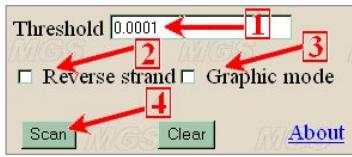
1) Enter the nucleotide sequence into the text-box (red arrow 1 in the figure below). It is possible to enter the sequence from clipboard or from file by using the standard dialog window, which pops up after clicking the screen button 'Browse' (arrow 2). Note that the sequence should contain only the symbols a, c, t, g, A, C, T, G and the length should range from 30 to 100 000 bp. If the sequence is entered from file, it should be of FASTA format.

2) By using the radio buttons (arrows 3) and drop-down list (arrow 4) located below the text-box 'Sequence', select the mode of searching for sites. It is possible to search for a single transcription factor binding site or for all of them contained in the list. If a single transcription factor is searched for, select its name in the list:



### Program options:

Enter the significance level value into the text box 'Threshold' (red arrow 1 in the figure. The threshold value  $10^{-4} = 0.0001$  is set there). It is recommended to set the threshold ranging from  $10^{-12}$  to  $10^{-4}$ . The more is the threshold value, the more regions of the sequence analysed will be annotated by the program as true transcription factor binding sites. In case the threshold value is set near by 1 (or more), each position will be annotated as a site. In case the threshold value is equalling or less than zero, then the program will not annotate any site at all.



You may if necessary select for processing the complementary DNA sequence strand by clicking the check-box 'Reverse strand' (arrow 2). If this box is not checked, the search will be done in the sequence entered, "as is".

Choose the mode of the data output. For displaying the results graphically, click the check-box 'Graphic mode' (arrow 3). If this option is off, then the result will be output in a textual form.

### Program execution

Start the program by clicking the button 'Scan' (see arrow number 4 in figure given above).

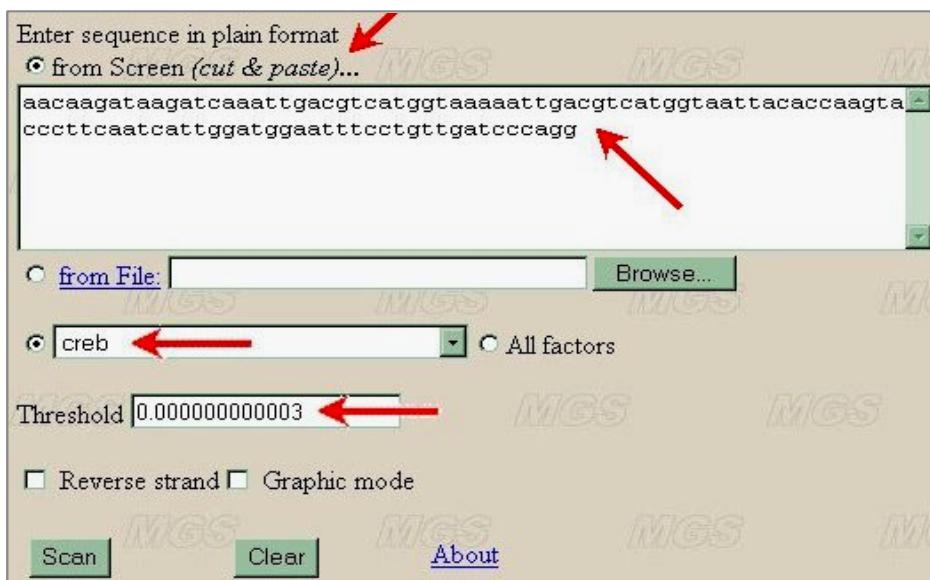
### Data output

See the example below.

### Example

As an example, which is invoked by clicking the appropriate hyperlink in the BinomSite starting page, the sequence of steroid-responsive region annotated in TRRD with the number 'P00082' and located in-between positions (-160; -60) relatively transcription start site was analysed. The data input is illustrated for searching for a single transcription factor binding site in direct strand. The example sequence has the CREB transcription factor binding site ('caaataTGACGTCatgg'; TRRD number S173) at position -147 relatively transcription start site. This position corresponds to 14-th position in the regulatory region sequence, which was entered for searching in it potential CREB binding sites.

In the figure, the data input is illustrated for searching for a single transcription factor binding site, with the threshold  $3 \cdot 10^{-12}$ , in direct strand, with the textual data output:

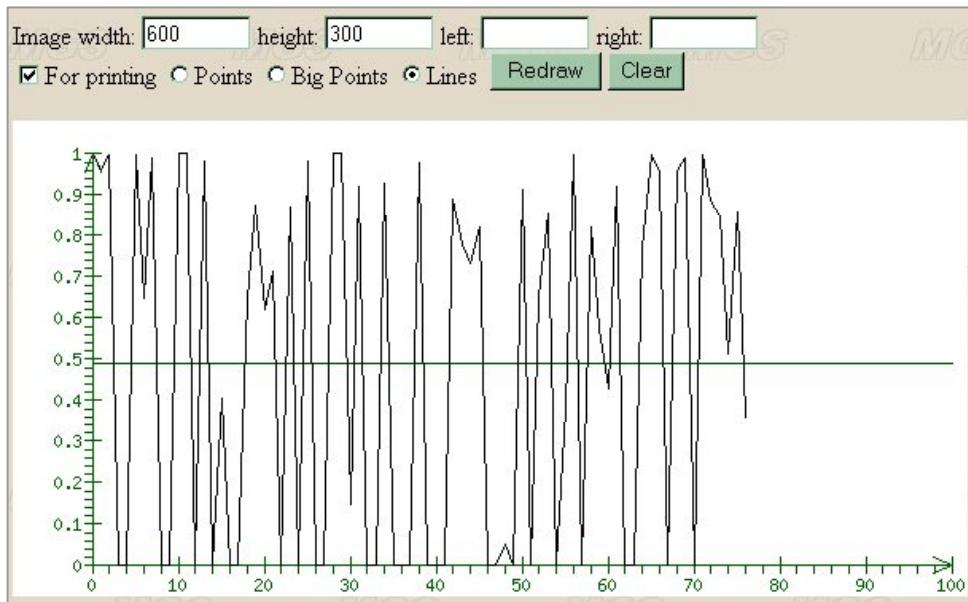


The textual data output is illustrated in the next figure. The program BinomSite displays a list of positions with potential transcription factor binding sites for a particular transcription factor, the name of this transcription factor, DNA strand ('+' or '−', that is, direct or complementary DNA chain), as well as the fragments of the sequence ordered that correspond to potential transcription factor binding sites in positions indicated:

Name	Position	Strand	Site
creb	12	+	ataagatcaaattgacgtcatggtaaaaa
creb	30	+	catgtaaaaattgacgtcatggtaattac

As seen, the program BinomSite has detected a potential site with the sequence 'ACGTCA' at position 12, with the length of 29 nt and the sequence 'ataagatcaaattgACGTCAttggtaaaaa'. This sequence includes the real CREB binding site annotated in TRRD. In the sequence analysed, the BinomSite program has detected one more potential CREB binding site at position 30.

The graphic data output is illustrated in the next figure. The program BinomSite displays the plot of similarity score at each position of the sequence analysed with the CREB transcription factor binding sites. The green horizontal line marks the average similarity score estimated for all positions of a sequence:



Comments and questions are welcome to Mikhail Pozdniakov ([mike@bionet.nsc.ru](mailto:mike@bionet.nsc.ru))

## 2.2. Program MMSite

Release 2003

### Program description:

This program realizes a novel approach for searching for potential transcription factor binding sites in the arbitrary sequence entered by a user.

**Access to MMSite:**

<http://www.domain.com/mgs/gnw/regscan/> link [MMSite](#)

**Biological task that could be solved by using MMSite:**

This software was developed for searching for potential transcription factor binding sites in the arbitrary sequence entered by a user. User may enter one of possible four methods of transcription factor binding site recognition. This enables to choose appropriate recognition method for each type of transcription factor binding sites.

**Data input:**

To search for transcription factor binding sites, do the following:

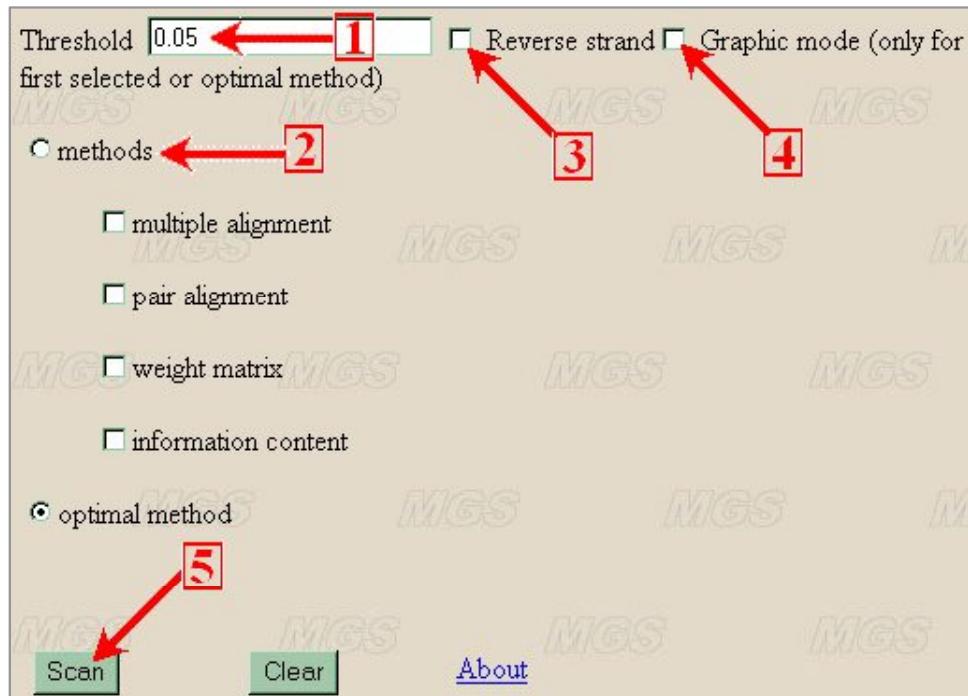
- Enter the nucleotide sequence into the text-box. It is possible to enter the sequence from clipboard or from file by using the dialog window 'Browse'. Note that the sequence should contain only the symbols a, c, t, g, A, C, T, G and the length should range from 30 to 3000 bp. If the sequence is entered from file, it should be of 'FASTA' format. The text-box for data input is marked by red arrow 1 in the figure.
- Choose the transcription factor that you want to recognise from the drop-down list 'Select transcription factor'. This list for choosing the name of transcription factor is also marked by red arrow in the figure.



**Program options:**

Select significance level value in the drop-down list 'Threshold' (arrow 1). It is recommended to set the threshold ranging from  $10^{-5}$  to  $10^{-2}$ . The more is the threshold value, the more regions of the sequence analysed will be annotated by the program as true transcription factor binding sites. In case the threshold value is set near by unit (or more), each position will be annotated as a site. In case the threshold value is equalling to zero or less than zero, then the program will not annotate any site at all. Click the check-box 'Methods' (arrow 2) for choosing the method or a combination of methods that you want to apply or the option 'Optimal method' that enables automated choosing of the best method for particular site of your interest. If you click the check-box 'Optimal method',

the program automatically selects one of four possible methods that will at best discriminate particular transcription factor binding sites indicated out of random sequences.



#### Program execution:

Start the program by clicking the button 'Scan' (see arrow number 5 in figure given above).

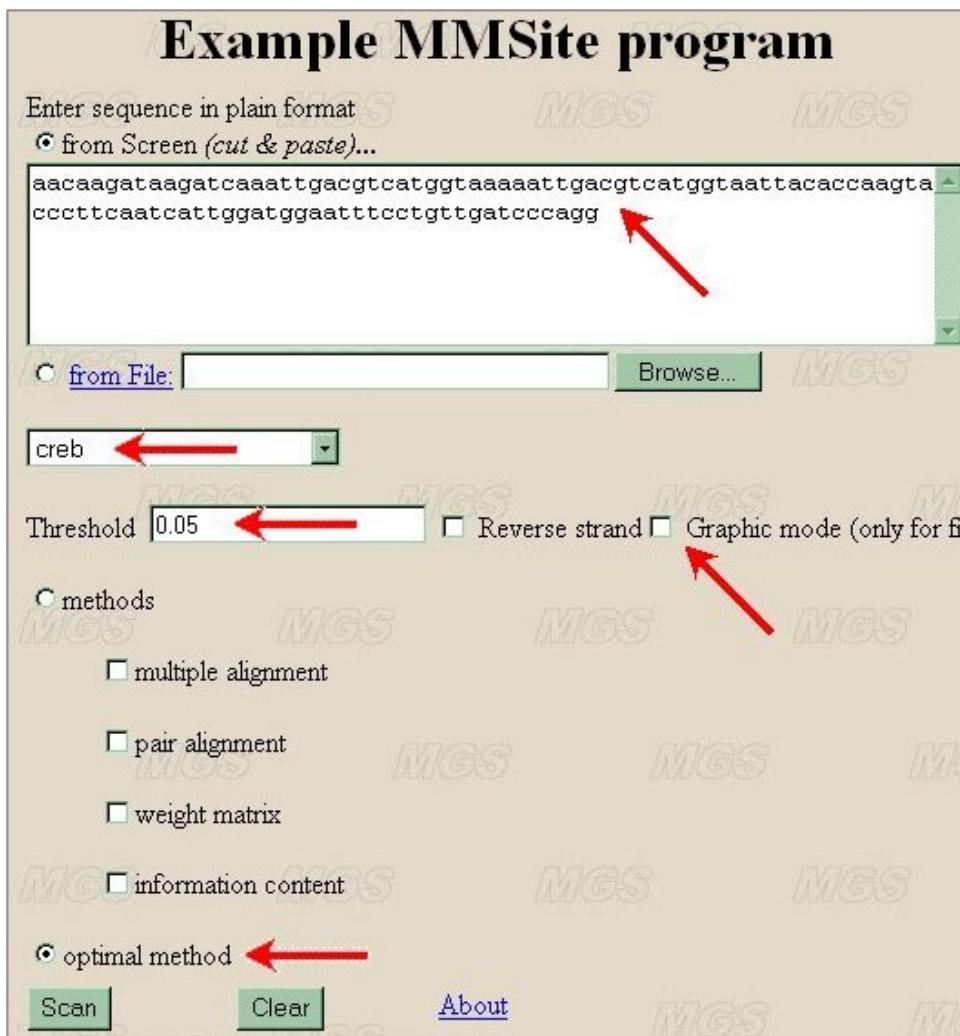
#### Data output:

- 1) Select the strand of DNA sequence to be analysed by clicking the check-box 'Reverse strand' (arrow 3). If this box is not checked, the search will be done in the sequence entered, otherwise, in complementary, reverse DNA strand.
- 2) Choose the mode of the data output. For displaying the results graphically, click the button 'Graphic mode' (arrow 4).

#### Example:

An example of data input, selecting the site type and recognition method of the program **MMSite**. As the example, the sequence of steroid-responsive region annotated in TRRD with the number 'P00082' and located in-between positions (-160; -60) relatively transcription start site.

In the figure, the data input is illustrated for searching for a single transcription factor binding site, with the threshold 0.05, in direct strand, with the textual data output:



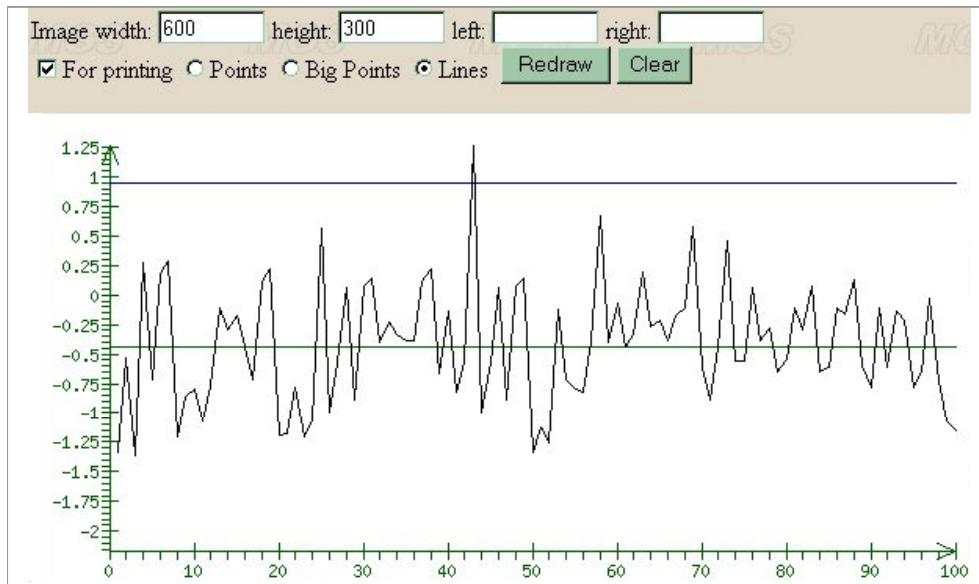
Below, the textual data output of the program 'MMSite' is illustrated. It displays the list of positions where potential transcription factor binding sites are found, the name of potential transcription factor binding sites, DNA strand ('+' or '-', that is, direct or inverse DNA), as well as the regions in particular sequence corresponding to potential binding sites in positions indicated. These regions are aligned with the sample of the sites of particular type, so, the symbol of deletion (-) may occur in the sequences of potential transcription factor binding sites '-'.

The data output by the program 'MMSite' searching for potential transcription factor binding sites in the sequence shown in figure below. The program 'MMSite' has found two potential 'CREB' transcription factor binding sites. The first of these sites is annotated in TRRD database.

As an example, we have performed an analysis of steroid-responsive region annotated in TRRD with the number 'P00082' and with experimentally detected 'CREB' transcription factor binding site ('S173', 'caaattGACGTCA~~t~~gg') at position -147 relatively transcription start site. Note that this position corresponds to position 14 of the sequence entered in above figure. The program 'MMSite' has detected potential 'CREB' binding site at position 21, which is a part of the sequence of real 'CREB' binding site annotated in TRRD. Also, the program has recognised three more potential 'CREB' binding site: at position 39, 53 and 64, which are not annotated in TRRD.

Name	Position	Strand	Site	Method
creb	22	+	-cggtca	multiple alignment
creb	39	+	acgtca	multiple alignment
creb	53	+	acacca	multiple alignment
creb	64	+	ccttca	multiple alignment

Graphic mode of resulting data output. In the Figure, the plot illustrates the results of analysis by the program 'MMSite' of the example sequence shown in Figure above. The program displays similarity Score at each position of the sequence of interest with the 'CREB' transcription factor binding sites. The blue horizontal line marks the threshold set by user. The green horizontal line marks the average similarity Score estimated for all positions of a sequence. If the region displays similarity Score peak exceeding the blue line, this region is considered as potential binding site of a certain transcription factor.



Comments and questions are welcome to Mikhail Pozdniakov ([mike@bionet.nsc.ru](mailto:mike@bionet.nsc.ru))

## 2.3. Programs for DNA Functional Sites Recognition

Release 2003

**Programs description:**

**RegScan module** contains programs for transcription factor binding sites recognition.

**Access to the programs for DNA regulatory regions and functional sites recognition:**

<http://www.domain.com/mgs/gnw/regscan/>

- Click the link [FreqMatrix](#) for searching for the sites by the method of frequency matrices.
- Click the link [Activity](#) for searching for the sites by the method of binding activity of the site.
- Click the link [bDNA](#) for searching for the sites by the method of physico-chemical and conformational properties.

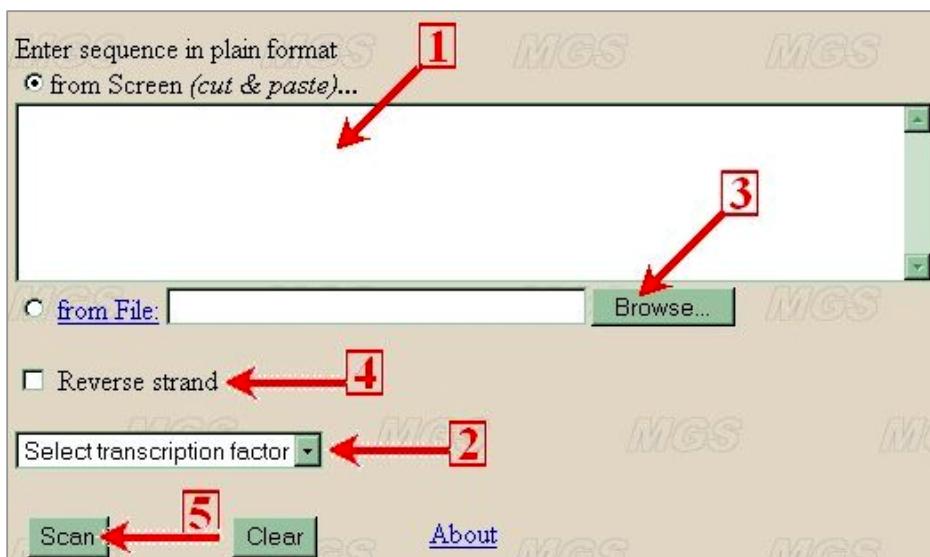
**Biological task that could be solved by using the programs:**

Recognition of transcription factor binding sites in DNA sequences by using 3 different methods, based on data on the frequency matrices, physicochemical and conformational properties or binding activity of a site.

**Data input:**

Input window of **RegScan Recognition Module** contains input textbox for sequence (1) of interest and the list of available regulatory regions or functional sites to recognise (2). It is also possible to input sequence from a text file by clicking the 'BROWSE' button (3).

Essential notes: sequences should be entered in FASTA format using upper- or lower-case letters; line feeds and blanks are ignored; Sequence length should not exceed 32000 bp and should not be less than 50 bp.



**Programs options:**

The search for transcription factor binding sites could be provided by the reverse strand, in this case, click the check box 'Reverse strand' (4).

**Programs execution:**

Click the button 'Scan' (5).

**Data output:**

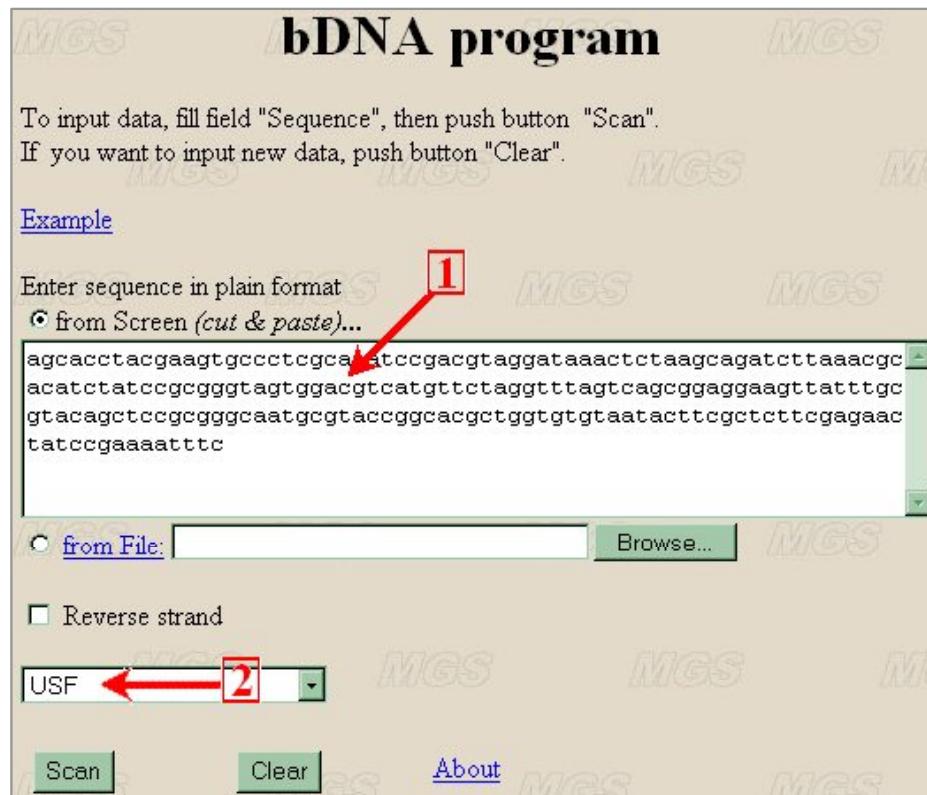
The tools output represents the profile of the Score value. The positive peaks of this profile pinpoint to the potential site recognized.

**Example: Recognising of USF binding site in the sequence of interest**

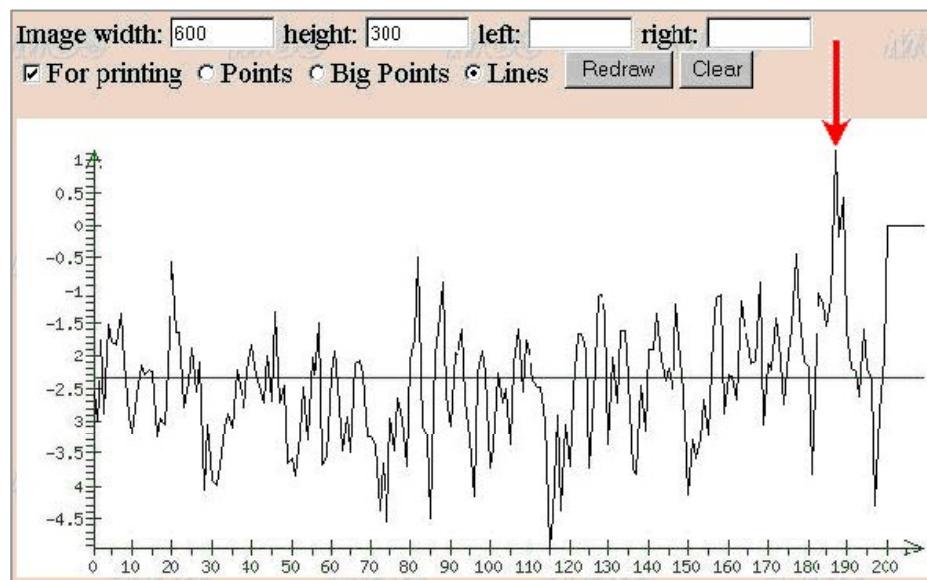
An example of recognising USF binding site in the sequence of interest by using one of the recognition program, **bDNA**.

Into the text-box, input the random sequence with the length 200 bp (1).

Select from the list (2) transcription factor USF:



After the program execution (by clicking the button 'Scan') view the graphic mode of resulting data output:



In the figure, the plot illustrates the results of analysis by the program **bDNA** of the example sequence shown in figure above. The program displays similarity Score at each position of the sequence of interest with the UTF binding sites. The horizontal line marks the average similarity Score estimated for all positions of this sequence. If the region displays similarity Score peak, this region is considered as potential binding site of a certain transcription factor (marked by arrow).

**Comments and questions are welcome to Mikhail P. Ponomarenko, ([pon@bionet.nsc.ru](mailto:pon@bionet.nsc.ru))**

## 2.4. RGSiteScan Program

Release 2003

### Program description:

This program is designed to predict the potential binding sites of transcription factors in the target DNA sequence.

### Access to RGSiteScan program:

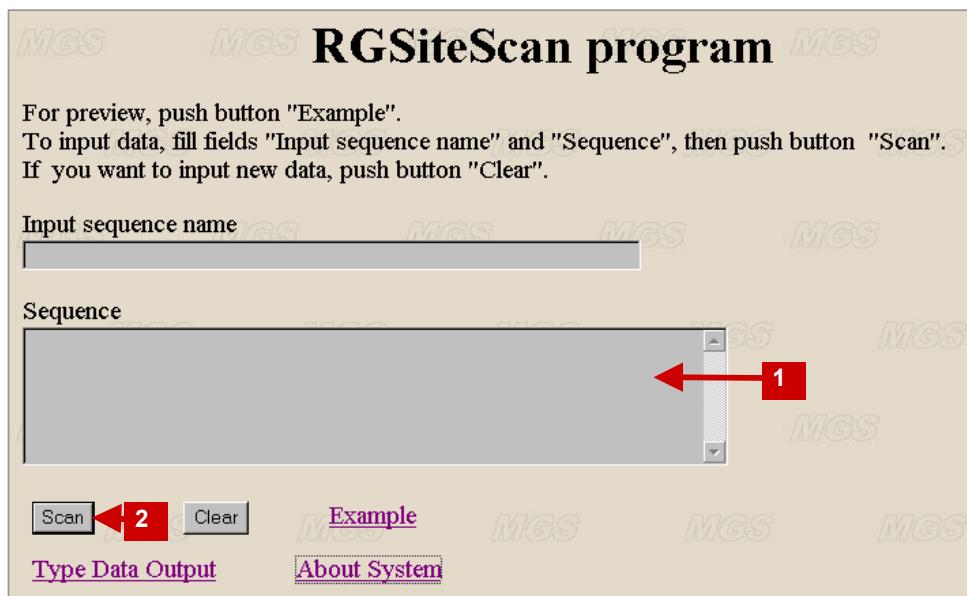
<http://www.domain.com/mgs/gnw/regscan/> link RGSiteScan

### List of biological tasks that could be solved by using the RGSiteScan program:

- ◆ Recognition of the potential binding sites of transcription factors in the target DNA sequence.

### Data input

The sequence to be analysed should have the plain textual format (a, t, g, c, printed in lower case, without gaps) and the sequence length should be less than 32000 bp.  
Enter the sequence to be analysed into the text-box 'Sequence' (1).



### Program options:

It is not necessary to install parameters.

### Program execution

The program is executed by clicking the button 'Scan' (2).

### Data output

The results of the program execution are presented as the oligonucleotides that correspond to potential transcription factor binding sites. The names of these factors, as well as positions of respective oligonucleotides are also displayed.

### Example

The task is to predict the potential binding sites of transcription factors in GGMYHE standard; DNA; VRT; (promoter region:1-2211 BP)

```
*****
* Prediction of potential binding sites *
* of transcription factors in given *
* DNA-sequence on the base of recognition *
* groups. *
*****
The name of sequence: GGMYHE      standard; DNA; VRT; (promoter region:1-2211
The number of all predicted binding sites =      49

      Name      Position      Site
      AR        17          ttttct
      AR        92          ttttct
      COUP      98          tgacct
      AR        166         ttttct
      ICSBP     185         agtttc
      APF       196         gttaag
      AR        200         agaaaa
      AR        207         tgctct
      AR        217         agcaca
      AR        231         aaaaca
      AR        297         agaaaa
      AR        300         aaaaca
      AR        380         tgtcct
      AR        416         ttttct
      AR        427         aaaaca
      APF       450         attacc
      Pit-1     465         tatacat
      AP-1      473         tgactga
      AR        479         agaaaa
```

As a result, positions of 49 potential binding sites of transcription factors were predicted.

**Comments and questions are welcome to  
Yury Kondrakhin ([kondrat@bionet.nsc.ru](mailto:kondrat@bionet.nsc.ru)).**

### Example

The task is to predict the potential binding sites of transcription factors in GGMYHE standard; DNA; VRT; (promoter region:1-2211 BP)

```
*****
* Prediction of potential binding sites *
* of transcription factors in given *
* DNA-sequence on the base of recognition *
* groups. *
*****
The name of sequence: GGMYHE      standard; DNA; VRT; (promoter region:1-2211
The number of all predicted binding sites =      49

      Name      Position      Site
      AR        17          ttttct
      AR        92          ttttct
      COUP      98          tgacct
      AR        166         ttttct
      ICSBP     185         agtttc
      APF       196         gttaag
      AR        200         agaaaa
      AR        207         tgctct
      AR        217         agcaca
      AR        231         aaaaca
      AR        297         agaaaa
      AR        300         aaaaca
      AR        380         tgtcct
      AR        416         ttttct
      AR        427         aaaaca
      APF       450         attacc
      Pit-1     465         tatacat
      AP-1      473         tgactga
      AR        479         agaaaa
```

As a result, positions of 49 potential binding sites of transcription factors were predicted.

**Comments and questions are welcome to  
Yury Kondrakhin ([kondrat@bionet.nsc.ru](mailto:kondrat@bionet.nsc.ru)).**

## 2.5. RecGroup program

Release 2003

### Program description:

The program is designed for searching for recognition groups of oligonucleotides of fixed length that are specific for transcription factor binding sites.

### Access to RecGroup program:

<http://www.domain.com/mgs/gnw/regscan/> link RecGroup

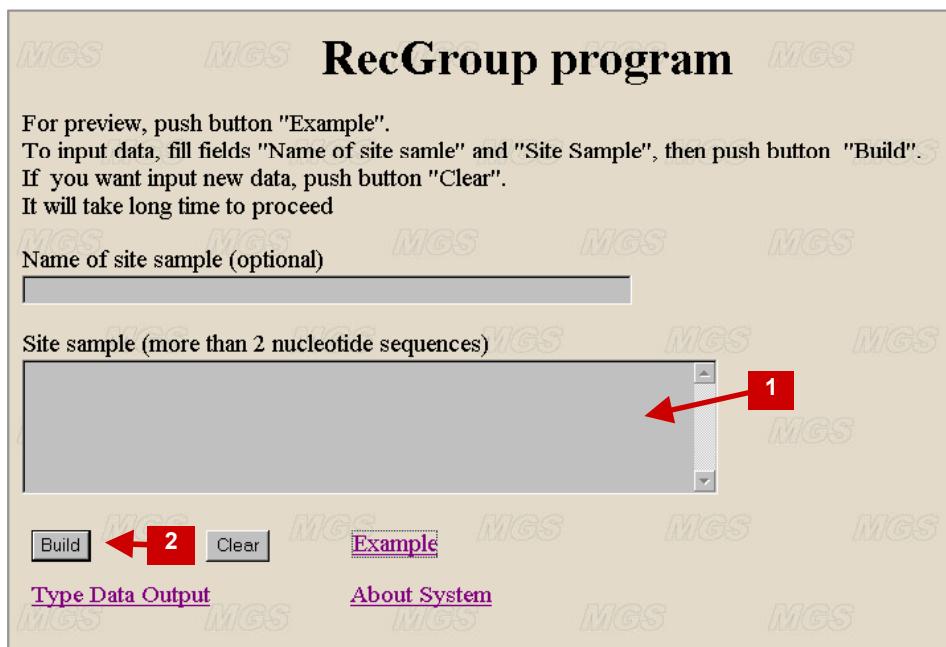
### List of biological tasks that could be solved by using the RecGroup program:

- ♦ Recognition Group retains most of the information and is, thus, well suited to describe the transcription factor binding sites.

### Data input

The sequences to be analysed should have the plain textual format (a, t, g, c in lower case, without spaces). Note that the number of sequences should be more than 2.

Enter the sequence to be analysed into the text-box 'Site sample' (1).



### Program options

For the program execution, it is not necessary to set the parameters.

### Program execution

The program is executed by clicking the button 'Build' (2).

### Data output

The program outputs the results as a column of the weighted oligonucleotides.

### Example

The goal is to build recognition group for a set of sequences containing NF-Y binding site. As a result, we get the list of weighted oligonucleotides, which could be used for recognition of transcription factor NF-Y binding site.

```
It will take long time to proceed. Please wait ...
.....
*****
* Calculation of optimal parameters *
* l (site length) and ANM (the *
* admissible number of mismatches) *
* and construction of recognition *
* group. *
*****
The name of transcription factor: Sequences Containing NF-Y
The binding site length = 7
Admissible number of mismatches = 2

Recognition group:
agccaat    7
aaccaat    3
agccact    1
gaccaaat   4
agctcat    1
agcctct    1
aggcaag    1
atgcaat    1
agccatg    1
ggcaaat    1
```

**Comments and questions are welcome to  
Yury Kondrakhin (kondrat@bionet.nsc.ru).**

## 2.6. ARGO system

**Release 2003**

### Program description:

The program is designed for searching for degenerated oligonucleotide motifs that are specific for gene regulatory regions.

### Access to ARGO system:

<http://www.domain.com/mgs/gnw/regscan/> link 'Search for degenerate oligonucleotide motifs'

### List of biological tasks that could be solved by using the ARGO system:

- ◆ Oligonucleotide analysis of promoter regions,
- ◆ Oligonucleotide analysis of transcription factor binding sites,
- ◆ Oligonucleotide analysis of enhancers and silencers, as well as other regulatory elements.
- ◆ Oligonucleotide analysis of nucleosome binding sites.

## Data input

The sequence to be analysed should have the plain text format (a, t, g, c in lower case, without spaces). Note that the number of sequences should be less than 100. Length of sequences should be less than 100 bp.

Enter the sequence to be analysed into the text-box 'Sequence'.

## Program options

'Length of oligs' (1) – Length of oligonucleotides to be considered.

The length could be varied in dependence of the length of the signals to be searched for. The length value could vary from 6 to 12. The recommended length is 8.

'Hamming's distance value' (2).

In our case, the Hamming distance between two oligonucleotides is a number of non coinciding positions in these oligonucleotides. For example, for the pair of oligonucleotides AATGCGT and TAGGCTT, Hamming distance is equal to 4 (\*A\*\*G\*\*T).

If Hamming distance  $R$  between oligonucleotides from different sequences is lower than the threshold value  $r_o$ , they are united into one class. We recommend to use  $r_o > (1/2)$  of the oligonucleotides length.

### 'Minimal share of sequences with olig' (3)

This parameter denotes the minimal percent of sequences in the set analysed, in which every found motif should present. This parameter ranges from 1 to 100. Recommended value equals to 80.

#### 'Minimal binomial probability' (4)

This parameter means the decimal logarithm of the significance level.

The oligonucleotide motif is considered significant, if it meets the following conditions:  
 (1) the fraction f of the RGS containing the motif is higher than a certain level fo and  
 (2) the binomial probability P(n,N) to observe this motif by accident in n and more number of RGS from N RGS considered, is lower than a significance level a.

The value of this parameter may be varied from -1 to -999.

The more this value is close to -1, the more probable to find false motifs in the set of motifs recognised. If this value is close to -999, then the probability to loose from consideration the true signals that are present in a set.

We recommend to use lg(a)<-8.

#### 'Alphabet size' (5)

This parameter has two values:

4 – to be used for analysing the motifs with canonical letters A, T, G,C.

14 – to be used for analysis of motifs written by extended alphabet.

#### 'Analysis of complementary chains' (6)

This parameter enables a user to analyse complementary chains of the sequences to be analysed. That is, the motifs could be searched for both in '+' and in '-' directions.

#### 'Reset' (7)

This motif is needed to clear up the text-box 'Sequence'.

**Search for degenerate oligonucleotide motifs**

Input the sequence sample to be analysed  
Note that the number of sequences should be less than 100. Length of sequences should be less than 100.

Sequence

```
caggccccgtaccatgttatataaaggagacactgggacaaggcaccat
ggctggggccacgtccctgttatataaaggggacctggggctgagcacta
ggctggctaggatagaagaataaaaggaaaggcacccctcagcgttccaca
cagaatcccgccatggatagaataaaggccatggccatggcggcggcggc
gagccaggccatggatataaaggtaggttaggttagtgcgttcttacatt
aggggccaggggctggccataaaagtcaaggccatggcgttatttttttttt
gaggggcccccggggaggcgataaaatgtggggacacagacggccggctacc
gcgggtgtacaggagatataaggatggtcggccctgcaggctccatca
```

Length of oligs 8      1  
 Hamming's distance 4      2  
 Minimal share of sequences with olig 80 (%)      3  
 Minimal binomial probability -10      4  
 Alphabet size 4      5  
 Analysis of complementary chains YES      6

8 → Submit      Reset      Help Example 7  
 The process may take few minutes...

## Program execution

The program is executed by clicking the button 'Submit' (8).

## Data output

The program outputs the results as a column of the motifs found.

The screenshot shows the RegScan software interface. At the top, there are dropdown menus for 'Alphabet size' (set to 14) and 'Analysis of complimentary chains' (set to YES). Below these are buttons for 'Submit', 'Reset', and 'Help'. A message says 'The process may take few minutes...'. In the center, there is a logo of the Institute of Cytology and Genetics SB RAS. To the right, a vertical scrollable list displays a large number of motifs found in the sequences. A red arrow points from the text 'The process may take few minutes...' towards the scrollable list of motifs.

NNGATTAG
GNATNNGC
NNNAGCGC
NNNGCGCT
NNGCGCTN
NGCGCTNN
GCGCTNNN
CGCTNNNG
GNNAGNGC
NNAGNGCG
NNAAGCGC
NGNGCGCT
<b>GNGCGCTG</b>
NGCGCTGN
GCGCTGNN
CGCTGNNN
NNTGNGCG
NNAGCGCG
NAGCCCGN
AGCGCGNN
GCGCGNTN
CGCGNNCN
NTNCGACA
ANNTGCNG
CATNCAGC
NNGCAGCT
TNCAGCNG
CNGCTGNA
AGCTGCNN
GCTGNATG
TGTCGNAN
NGNNCGCG
NNACCGCG
NNCCGGCT

## Example

The goal is to make oligonucleotide analysis of a set of sequences containing the regions of erythroid-specific promoters in-between -50 - +1 bp relatively transcription start site in order to find in these sequences the motifs that possibly correspond to the known transcription factor binding sites.

The parameters are set as follows:

Length of oligs - 8

Hamming's distance - 4

Minimal share of sequences with olig - 50%

Minimal binomial probability - -8

Alphabet size -14

Analysis of complementary chains – NO

### Search for degenerate oligonucleotide motifs

Input the sequence sample to be analysed  
Note that the number of sequences should be less than 100. Length of sequences should be less than 100.

Sequence

```
gctgtccccccgcgcggcaaggataaaaccctggcgctcgccggccggca
gaggcggcttatttcggcgccgcacacggccggcgtaacggccgg
ggcctggaaagataacagctagcaggtaaggctcagacactgacatttgc
gggcgcgcggggaggggcgagggtcaggggctggggacgcgcgtgggg
gacgctgcagcggttttcgaaagggtttccgcctgtcgacgtcg
ataagaccagcagttaggttacacttcccccagtcgcgtttgc
ccaatcagatgtggcagacaggagccctccaagaaaacttctagcctc
```

Length of oligs   
 Hamming's distance   
 Minimal share of sequences with olig  (%)  
 Minimal binomial probability   
 Alphabet size   
 Analysis of complementary chains

[Submit](#) [Reset](#) [Help](#) [Example](#)

As a result, we get the list of motifs, some of them are similar by nucleotide content, for example, to consensus transcription factor TBP binding sites (they are marked by red rectangles at the Figure).

### Search for degenerate oligonucleotide motifs

Input the sequence sample to be analysed  
Note that the number of sequences should be less than 100. Length of sequences should be less than 100.

Sequence

```
gctgtccccccgcgcggcaaggataaaaccctggcgctcgccggccggca
gaggcggcttatttcggcgccgcacacggccggcgtaacggccgg
ggcctggaaagataacagctagcaggtaaggctcagacactgacatttgc
gggcgcgcggggaggggcgagggtcaggggctggggacgcgcgtgggg
gacgctgcagcggttttcgaaagggtttccgcctgtcgacgtcg
ataagaccagcagttaggttacacttcccccagtcgcgtttgc
ccaatcagatgtggcagacaggagccctccaagaaaacttctagcctc
```

Length of oligs   
 Hamming's distance   
 Minimal share of sequences with olig  (%)  
 Minimal binomial probability   
 Alphabet size   
 Analysis of complementary chains

[Submit](#) [Reset](#) [Help](#) [Example](#)

SMGSNCNC
DGNRYAHA
DGNABAHA
VAGGNSVB
GRSVCNGS
RSNCTGVV
SMVGSDSC
SNDGCHSC
GCHSNBCM
<b>RNVNATAA</b>
NCVGCNSY
CVGCHSYN
BDSKGCMs
SKGCWGNN
GCWGDNHV
SSVGSBCY
GCNGBBCN
SNBGGCHS
GCWGVSNV
VMDGNCDG
<b>NDNATAAV</b>
DNATAAVN
SNGBGCWB
ASHCWGNV
NCVGCVSY
GGBCMSDS
BCHSBBCM
VCMCTBVV
VMDGYSDG
VDGNSKGC
GNSDGCHS

Comments and questions are welcome to  
Oleg V. Vishnevsky (oleg@bionet.nsc.ru).

## 3. Programs for Promoter Recognition

### 3.1. Program to recognize eukaryotic promoters

Release 2003

#### Program description

The program is designed for recognition of promoters in an arbitrary nucleotide sequence.

#### Access to the program

<http://www.domain.com/mgs/gnw/regscan/> link 'Recog'

#### List of biological tasks that could be solved by using the program

- Recognition of promoters in DNA sequences.

#### Data input

Input DNA sequence to be analysed into the field 'Input DNA Sequence' (1). Sequence length should be at least 400 bp and at most 32 kbp. (a, t, g, c in upper or lower case, line feeds and blanks are ignored).

#### Program options

The analysis of complementary DNA strand is executed by clicking the check-box 'Reverse strand' (2). To select the data output mode, choose the option by clicking the check-box 'Graphic mode' (3). Choose the promoter type: by species (Human or Drosophila) or by cell type (erythroid), type of an inducer (interferon, heat shock), functional system (lipid or endocrine metabolism). In all these cases, only one of TATA-containing or TATA-less promoter types should be selected, i.e., in our example, positive or negative regulation by glucocorticoid inducer (4). Choose the confidence level in the check-box 'Confidence Level' (default value 0.95). By varying the Confidence Level and 'CF' value, it is possible to find a compromise between sensitivity and specificity. By default, the 'CF' value is set at 0.95; whereas the reasonable values are within the range from 0.5 to 0.99. The 'CF' value denotes the portion of correctly recognised promoters in the sequences of the training set. By increasing the 'CF' value, it is possible to get more hits (5).

**Promoter recognition program**

For preview, click the button "Example".  
 Enter your sequence into the text-box "Sequence".  
 Set program parameters and data output mode (see [About](#)).  
 Click the button "Scan".  
 To reload the data, click the button "Clear".

**Example**

Enter sequence in plain format  
 from Screen (*cut & paste*)...  
 from File:

**Sequence:**  
 gtgaactccctgtacctttgtggactgacagttttacagt cgtgacacagtcaaaaca  
 ttaacttgggtatcgattttggccatatatataatataatataaagttagggaggggcgaa  
 cctctggcaggagcaaggccatggctgtggagtccacggccgacccacttgtctcg  
 ggccgtgtgtgtgtgtgtgtgtgtgtgtcccatgtggaaagatactgtgtatccc  
 agtgatggcaggccactggctgagcatgttggccatccagcagctgcagcagaggggac  
 atgaatatagtgtcttagcacctgacgcctgttt

**Graphic mode**  **Forward strand**  Reverse strand

**Confidence Level** (The higher the value allocated, the more hits.) **Interval [0 .. 1]** **5**

**Promoter types** (TATA+ means TATA-containing promoter, TATA- means TATA-less promoter)

By species	By tissue type, functional system or inducer type		
<input checked="" type="radio"/> Human TATA+	<input type="radio"/> Lipid metabolism system	<input type="radio"/> Erythroid-specific TATA-	<input type="radio"/> Interferon-induced TATA-
<input type="radio"/> Human TATA-	<input type="radio"/> Lipid metabolism system	<input type="radio"/> Erythroid-specific TATA+	<input type="radio"/> Interferon-induced TATA+
<input type="radio"/> D.melanogaster TATA+	<input type="radio"/> Glucocorticoid, negative regulated	<input type="radio"/> Endocrine system TATA-	<input type="radio"/> Heat shock-induced
<input type="radio"/> D.melanogaster TATA-	<input type="radio"/> Glucocorticoid, positive regulated	<input type="radio"/> Endocrine system TATA+	

**Buttons:** **Scan** **6** **Clear** **About**

## Program execution

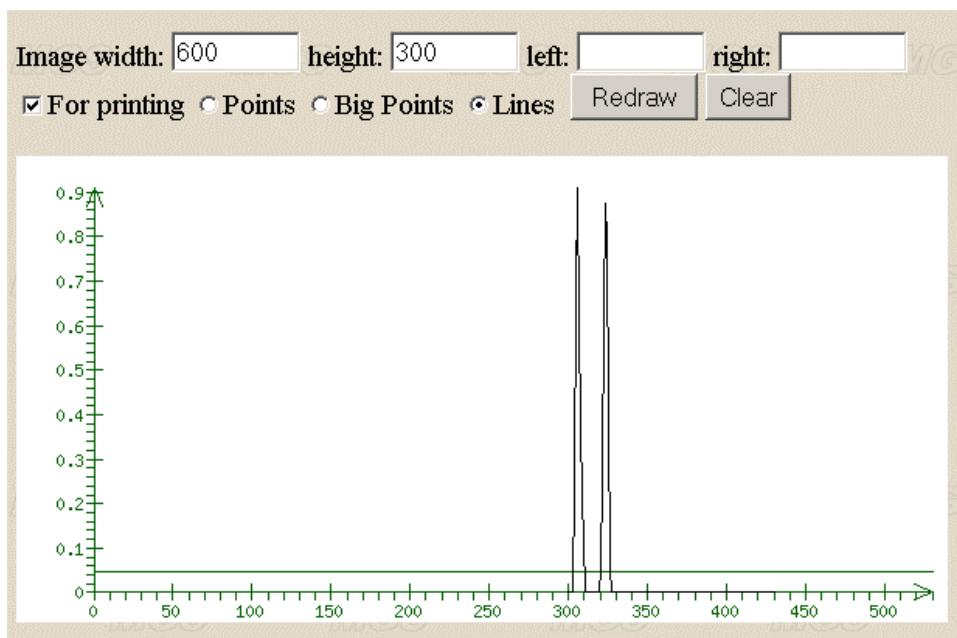
Click the button 'Scan' and wait for the program execution (6).

## Output data

Output data could be presented in two modes: graphical representation and numerical delivery.

### Example

The output gives the scores within the range in-between 0 and 1. The best hit equals to 1. For example, enter the sequence of the human bilirubin UDP-glucuronosyltransferase (1-1, UGT1A1) gene, AC AF352795, with the TATA-box sequence (279...295) and with transcription start at 314 nt. This gene is considered as TATA-containing human gene. Output data could be displayed as the graphical representation or as the numerical delivery. The graphical representation is shown below.



The numerical delivery is displayed as follows:

pos	value	method
305	0.245268	recon2
306	0.670882	recon2
307	0.908197	recon2
308	0.605870	recon2
309	0.375586	recon2
310	0.210228	recon2
311	0.112294	recon2
321	0.051683	recon2
322	0.262601	recon2
323	0.626194	recon2
324	0.876315	recon2
325	0.792281	recon2
326	0.400228	recon2
327	0.053702	recon2

**Comments and questions are welcome to  
Victor Levitsky ([levitsky@bionet.nsc.ru](mailto:levitsky@bionet.nsc.ru))**

### 3.2. ARGO-Viewer

Release 2003

#### 1. Program description:

A software program ARGO-Viewer was designed for recognition of regulatory gene regions in arbitrary extended nucleotide sequences. Promoters are recognised in extended genome regions by estimating the set R of region-specific motifs found by the program ARGO.

#### 2. Access to ARGO-Viewer system

<http://www.domain.com/mgs/gnw/regscan/> link ARGO-Viewer

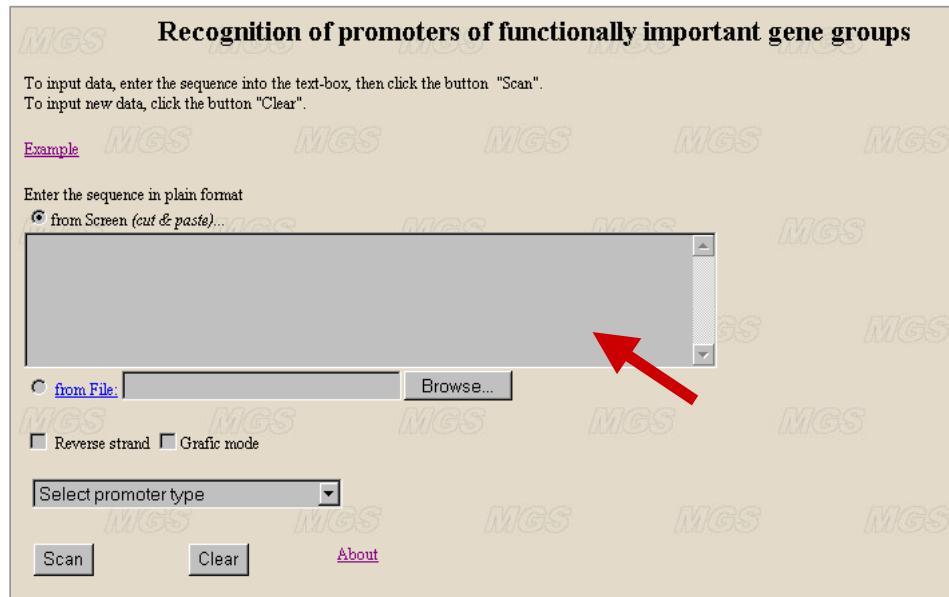
### 3. List of biological tasks that could be solved by using the ARGO system.

Recognition of promoter regions in tissue-specific genes in eukaryotes.

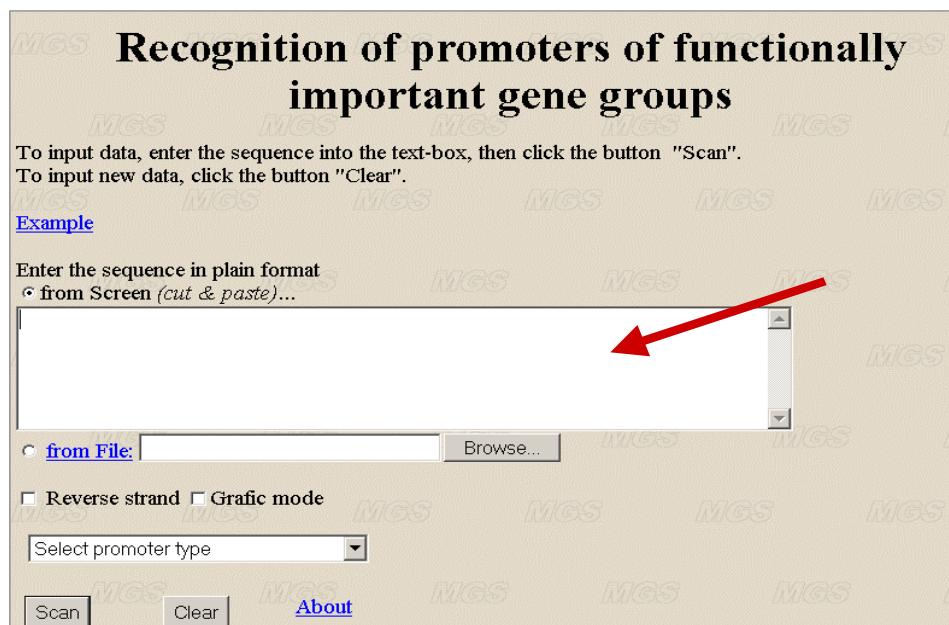
#### 4. Data input

The sequence to be analysed should have the plain text format (a, t, g, c, printed in lower case, without gaps) and the length should range from 400 to 80000 bp.

Enter the sequence to be analysed into the text-box.



If you input the sequence from your file, use the dialog menu by clicking the button 'Browse'.



## 5. Program options

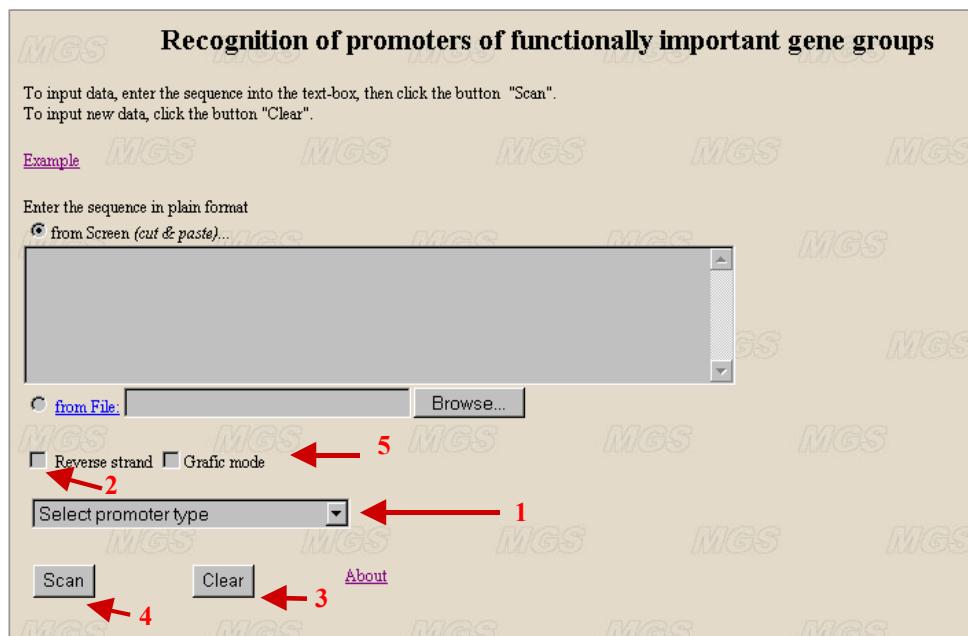
'Select promoter type' (1) is an option for choosing the type of tissue-specific promoters for recognition.

'Reverse strand' (2) is a parameter enabling perform recognition not only in direct strand, but in complementary strand too.

By clicking the button 'Clear' (3), it is possible to clear up the text box for data input.

## 6. Program execution

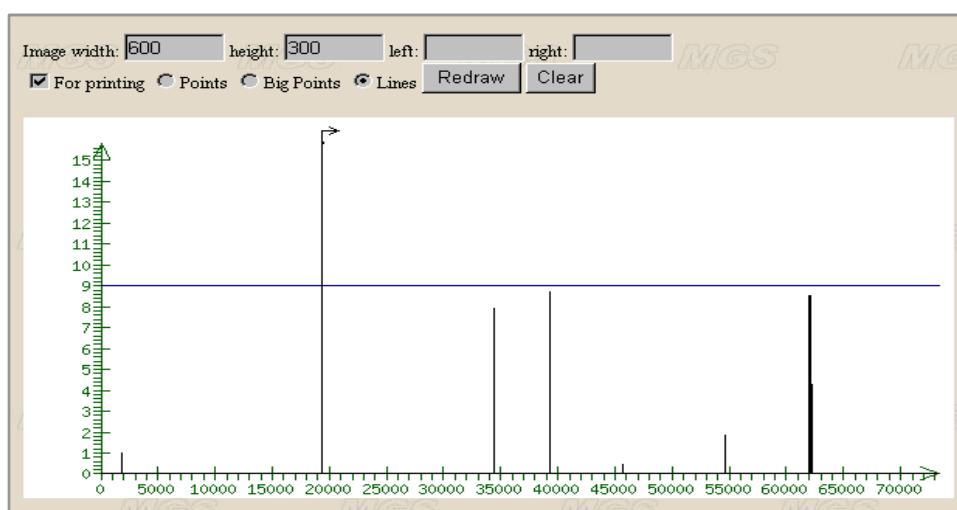
The program is executed by clicking the button 'Scan' (4).



## 7. Data output

The option 'Grafic mode' (5) is designed for graphical visualisation of recognition function profile.

The results are output as a graphical profile of the sequence analysed. The peaks at the profile correspond to potential promoter regions.



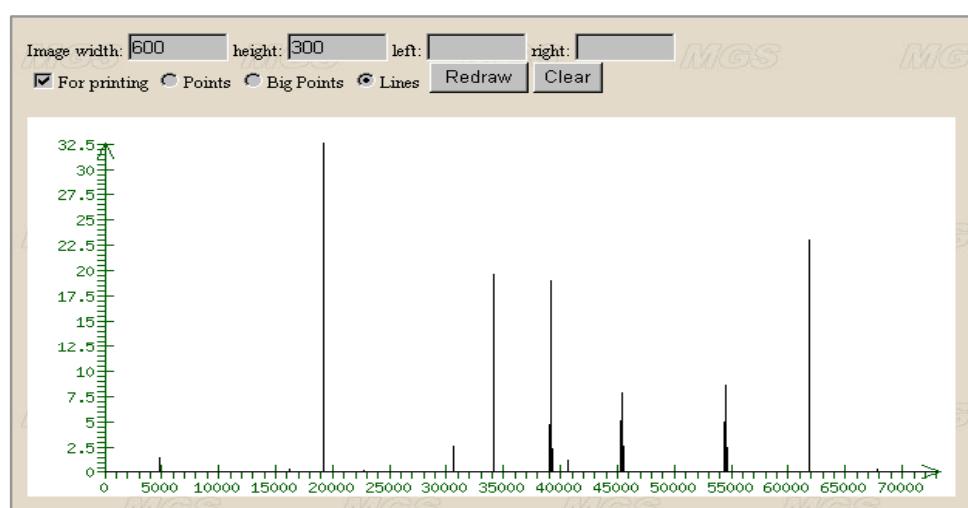
### Example

The task is to search for potential erythroid-specific promoter regions in Human beta globin region (ID HSHBB), of 73308 bp in length, mapped on the chromosome 11.

Choose the option 'Grafic mode' and select 'Erythroid-Specific Regulated genes' in the pull-down menu for choosing the type of promoters (these settings are marked by red rectangles in the Figure).



The resulting profile of the recognition function clearly illustrates 6 potential promoter regions, which correspond to experimentally detected transcription start sites at positions 19289, 34478, 39414, 54740, 62137 of five real promoters and to one promoter region of a pseudogene in vicinity of position 45557.



**Comments and questions are welcome to**  
Oleg V. Vishnevsky (oleg@bionet.nsc.ru).

### 3.3. POLIIISCAN

Release 2003

#### 1. Program description:

The program is based on a new method for recognition of type 2 promoters in the eukaryotic tRNA genes, which takes into consideration not only nucleotide distribution at individual box positions, but also the multiple interactions between the different parts of the boxes. The recognition system is determined by the module organization of the A and B boxes within intragenic promoters of the tRNA genes.

#### 2. Access to POLIIISCAN

<http://www.domain.com/mgs/gnw/regscan/> link 'PolIIIScan'

#### 3. List of biological tasks that could be solved by using the POLIIISCAN system.

- ♦ the program is designed for the recognition of type 2 promoters in the eukaryotic tRNA genes

#### 4. Data input

The sequence to be analysed should have the plain text format (a, t, g, c in lower case, without spaces). Enter the sequence to be analysed into the text-box (1).

Recognition procedure was constructed on the base of analysis of the A/B box samples. These samples were composed by applying the developed alignment algorithm to the set of sequences in [Compilation of tRNA sequences](#).

Two samples are available: [MGS](#)   [MGS](#)   [MGS](#)   [MGS](#)

[The A box](#)  
[The B box](#)

Select an subject of recognition:  2

Select 1-st threshold for the A-box recognition:  3

Select 2-nd threshold for the A-box recognition:  4

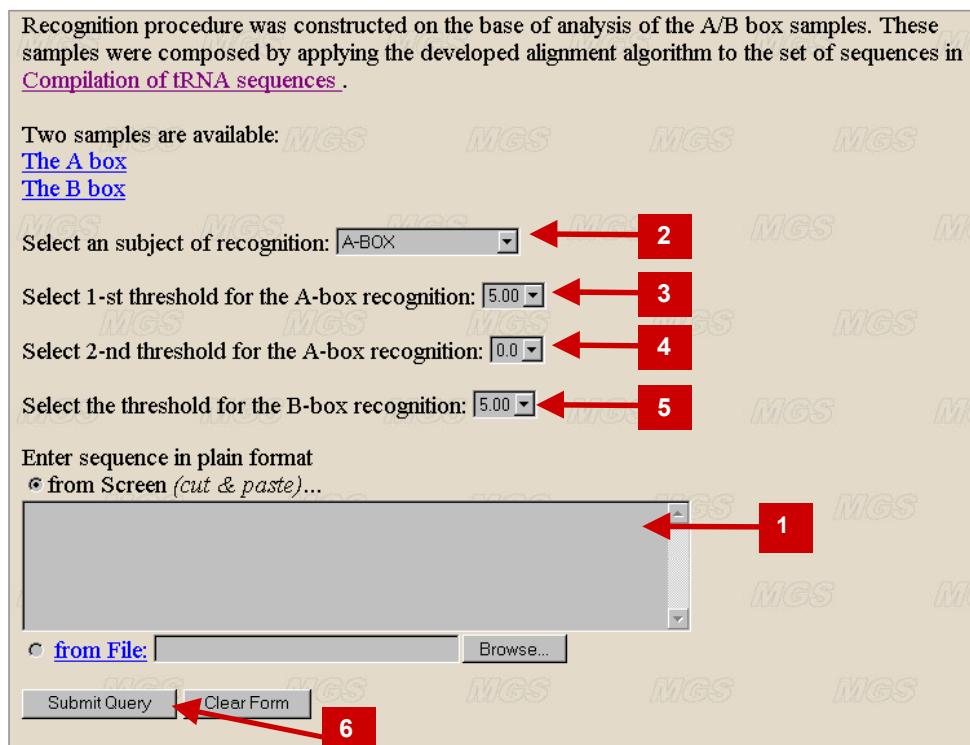
Select the threshold for the B-box recognition:  5

Enter sequence in plain format  
• from Screen (cut & paste)... 1

• from File:  Browse...

Submit Query 6

Clear Form



#### 5. Program options

Option 'Subject of recognition' (2) determines the type of regulatory region for recognition.

Option '1-st threshold for the A-box recognition' (3) orders the first limiting value for recognition of the A-box. It is recommended to set the values within the range from 5.00 to 7.70.

Option '2-nd threshold for the A-box recognition' (4) sets the second limiting value for recognition of the A-box. This value varies in-between 0.00 - 2.00.

Option 'threshold for the B-box recognition' (5) sets the first threshold value for recognition of the A-box. This value varies in-between 5.00 - 10.00.

## 6. Program execution

The program is executed by clicking the button 'Submit Query' (6).

## 7. Data output

The program outputs the results as a column of the regulatory regions found.

Prediction of the promoters									
1) A-Box:	[ 30,	40], tcacccatgg F=5.627000,	f=1.112957; B-Box:	[ 66,	76], gtcttagtatcc F=5.506000,				
2) A-Box:	[ 39,	49], tggatatataaa F=2.684000,	f=0.089701; B-Box:	[ 1,	11], gaattttaatc F=6.413000,				
3) A-Box:	[ 39,	49], tggatatataaa F=5.092000,	f=0.455149; B-Box:	[ 110,	120], gtctgcaagcc F=5.173000,				
4) A-Box:	[ 44,	54], tataaaagatgt F=1.206000,	f=0.584718; B-Box:	[ 1,	11], gaattttaatc F=6.413000,				
5) A-Box:	[ 44,	54], tataaaagatgt F=5.301000,	f=0.156145; B-Box:	[ 110,	120], gtctgcaagcc F=5.173000,				
6) A-Box:	[ 49,	59], aagatgtgtttgt F=1.972000,	f=0.245848; B-Box:	[ 1,	11], gaattttaatc F=6.413000,				
7) A-Box:	[ 49,	59], aagatgtgtttgt F=5.129000,	f=0.720929; B-Box:	[ 110,	120], gtctgcaagcc F=5.173000,				
8) A-Box:	[ 57,	67], ttgtctactgtc F=3.135000,	f=0.345514; B-Box:	[ 1,	11], gaattttaatc F=6.413000,				
9) A-Box:	[ 57,	67], ttgtctactgtc F=5.365000,	f=0.850499; B-Box:	[ 110,	120], gtctgcaagcc F=5.173000,				
10) A-Box:	[ 73,	83], atcccctaagta F=5.334000,	f=1.272424; B-Box:	[ 1,	11], gaattttaatc F=6.413000,				
11) A-Box:	[ 106,	116], aatagtctgcaa F=5.151000,	f=0.345514; B-Box:	[ 65,	75], tgctctagtatc F=5.050000,				
12) A-Box:	[ 118,	128], gccaggagggt F=3.141000,	f=0.451828; B-Box:	[ 65,	75], tgctctagtatc F=5.050000,				
13) A-Box:	[ 118,	128], gccaggagggt F=5.273000,	f=1.106312; B-Box:	[ 169,	179], gggaggactgc F=5.102000,				
14) A-Box:	[ 118,	128], gccaggagggt F=5.273000,	f=1.106312; B-Box:	[ 177,	187], tgcttgagctc F=5.020000,				
15) A-Box:	[ 121,	131], aggagtgggtggc F=3.273000,	f=0.372094; B-Box:	[ 65,	75], tgctctagtatc F=5.050000,				
16) A-Box:	[ 121,	131], aggagtgggtggc F=5.426000,	f=0.664452; B-Box:	[ 169,	179], gggaggactgc F=5.102000,				
17) A-Box:	[ 121,	131], aggagtgggtggc F=5.426000,	f=0.664452; B-Box:	[ 177,	187], tgcttgagctc F=5.020000,				
18) A-Box:	[ 121,	131], aggagtgggtggc F=5.426000,	f=0.664452; B-Box:	[ 191,	201], agttttagatt F=5.372000,				
19) A-Box:	[ 129,	139], tggctcatgtct F=2.157000,	f=0.146180; B-Box:	[ 65,	75], tgctctagtatc F=5.050000,				
20) A-Box:	[ 129,	139], tggctcatgtct F=2.157000,	f=0.146180; B-Box:	[ 91,	101], ggaatttagtca F=5.216000,				
21) A-Box:	[ 129,	139], tggctcatgtct F=5.807000,	f=1.617940; B-Box:	[ 169,	179], gggaggactgc F=5.102000,				
22) A-Box:	[ 129,	139], tggctcatgtct F=5.807000,	f=1.617940; B-Box:	[ 177,	187], tgcttgagctc F=5.020000,				
23) A-Box:	[ 129,	139], tggctcatgtct F=5.807000,	f=1.617940; B-Box:	[ 191,	201], agttttagatt F=5.372000,				
24) A-Box:	[ 129,	139], tggctcatgtct F=5.807000,	f=1.617940; B-Box:	[ 195,	205], tgatattatcc F=5.026000,				
25) A-Box:	[ 146,	156], tccagcactgtca F=4.226000,	f=0.634552; B-Box:	[ 91,	101], ggaatttagtca F=5.216000,				
26) A-Box:	[ 146,	156], tccagcactgtca F=5.386000,	f=0.724253; B-Box:	[ 191,	201], agtctgcaagc F=5.144000,				
27) A-Box:	[ 146,	156], tccagcactgtca F=5.386000,	f=0.724253; B-Box:	[ 195,	205], tgatattatcc F=5.026000,				
28) A-Box:	[ 146,	156], tccagcactgtca F=5.386000,	f=0.724253; B-Box:	[ 195,	205], tgatattatcc F=5.026000,				
29) A-Box:	[ 151,	161], cactggagagggt F=2.086000,	f=0.827244; B-Box:	[ 91,	101], ggaatttagtca F=5.216000,				
30) A-Box:	[ 151,	161], cactggagagggt F=2.086000,	f=0.827244; B-Box:	[ 109,	119], agtctgcaagc F=5.144000,				
31) A-Box:	[ 151,	161], cactggagagggt F=5.268000,	f=1.305647; B-Box:	[ 191,	201], agttttagatt F=5.372000,				
32) A-Box:	[ 151,	161], cactggagagggt F=5.268000,	f=1.305647; B-Box:	[ 195,	205], tgatattatcc F=5.026000,				
33) A-Box:	[ 154,	164], tggagaggtaga F=2.795000,	f=0.149502; B-Box:	[ 91,	101], ggaatttagtca F=5.216000,				
34) A-Box:	[ 154,	164], tggagaggtaga F=2.795000,	f=0.149502; B-Box:	[ 109,	119], agtctgcaagc F=5.144000,				
35) A-Box:	[ 154,	164], tggagaggtaga F=2.795000,	f=0.149502; B-Box:	[ 122,	132], ggagggtggc F=5.353000,				
36) A-Box:	[ 154,	164], tggagaggtaga F=5.461000,	f=0.345515; B-Box:	[ 191,	201], agttttagatt F=5.372000,				

Comments and questions are welcome to Yury Kondrakhin (kondrat@bionet.nsc.ru).

## 4. Other Programs

### 4.1. NASCA program

Release 2003

#### 1. Program description:

The aim of the analysis is to determine the positions or regions of aligned sequences where chosen physico-chemical or conformational feature variations depend on each other. Such dependence could indicate to possible importance of these regions for functioning of the sites and could point to possible molecular mechanisms of these regulatory sequences functioning.

#### 2. List of biological tasks that could be solved by using the Prediction programs:

- Searching for conformationally or physico-chemically interdependent regions in DNA sequence alignment of functional sites or regulatory regions.
- Searching for conformational or physicochemical properties of bDNA double helix which indicate to possible importance of regions for functioning of regulatory sequences and point to possible molecular mechanisms of functioning of these sequences.

#### 3. Access to the NASCA program

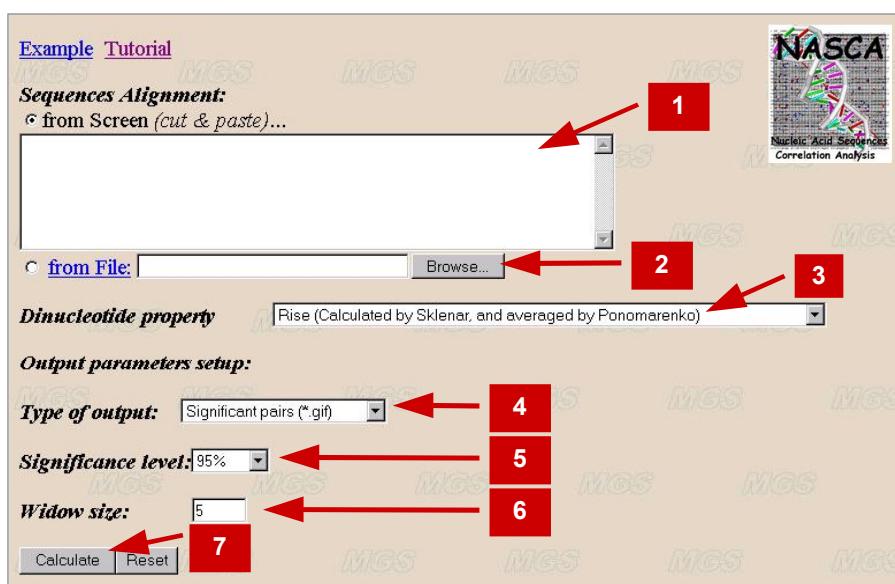
<http://www.domain.com/mgs/gnw/regscan/> link 'NASCA'

#### 4. Data input

Input the data into the input textbox (1) shown in the figure below for sequence alignment of interest.

##### Essential notes:

sequences should be entered in uper- or lower-case letters; line feeds and blanks are ignored; The sequence number should be at least 10. Note that accuracy of the method increases with sequence number. The sequences should be of equal length that should not exceed 130 bp and should not be less than 10. Gaps are forbidden, so positions of alignment that contain gaps will be eliminated from analysis.



## 5. Program options:

Input sequence alignment by different ways:

- from screen in the textbox (1) (by typing in from the keyboard or by cut & paste operation)
- or from file in 'FASTA' format: in this case, click the 'BROWSE' button (2) and select the source file.

Choose the property of interest in the list of available conformational and physico-chemical properties for calculation (3).

Choose the form of the output data display (4) in the 'Output parameters setup' menus. Output data may be represented by matrix colour picture with marked significantly correlated pairs (recommended), colour matrix diagram, HTML table of correlation coefficients matrix, and ASCII table of correlation coefficients matrix.

Choose significance level (5) in the 'Output parameters setup' menu to calculate the critical value for the correlation coefficients, which is calculated in accordance with the Student criterion at a chosen significance level.

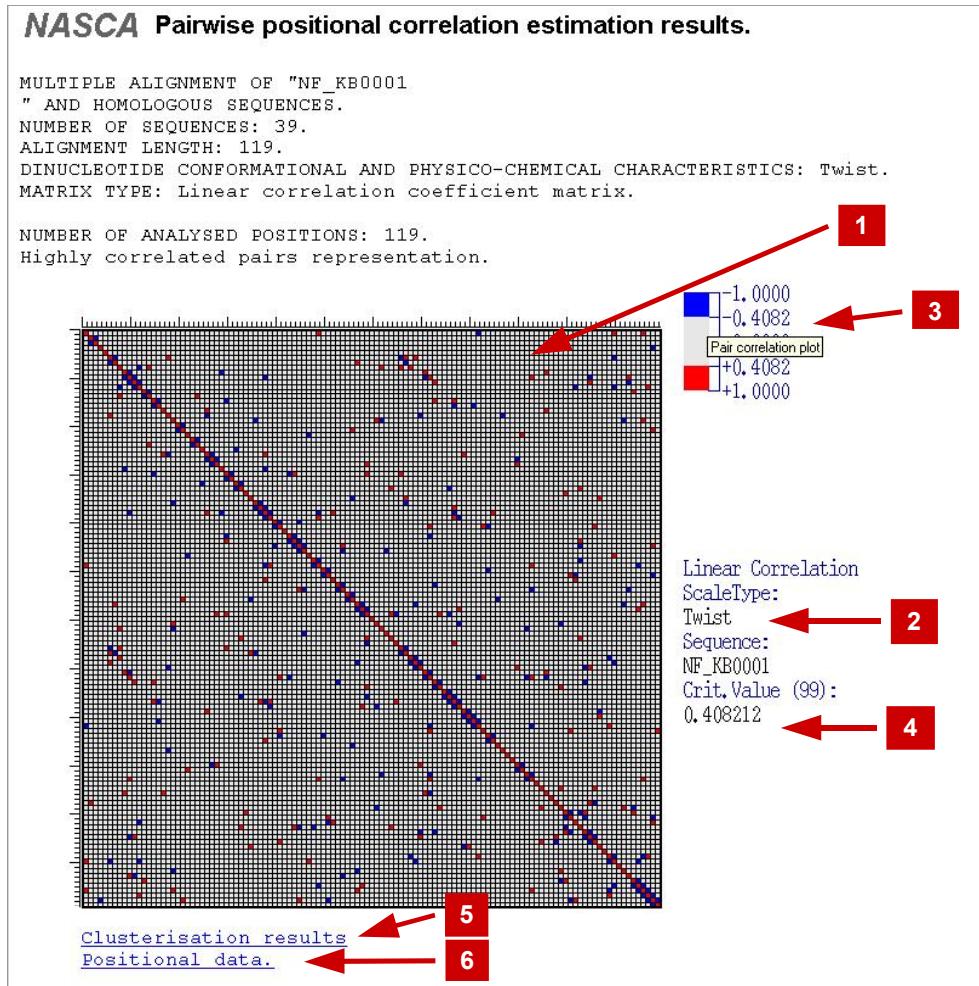
Type the value in the range 1-20 in the text-box 'Window size' (6) to choose the size of the window for clusterisation process. Values between 3 and 10 are recommended.

## 6. Program execution

Click the button 'Calculate' to execute the program (7).

## 7. Data output

The program output represents the matrix of correlation coefficients (1) of the chosen property (2) values for a given alignment. Insignificant correlations are marked grey, significant are coloured. The value scale is shown in the top right corner (3). The critical value (4) is calculated in accordance with the Student criterion at a chosen significance level.



Link to 'Clusterisation results' (5) will lead to representation of the matrix of positioning of significant blocks in correlation coefficients matrix calculated for chosen property values for a given alignment. By X- and Y-axes, positions of clusters are marked. The critical value is calculated in accordance with the Student criterion at a chosen significance level. Clusterisation was made for the window with the chosen size and chosen significance level calculated according to binomial distribution. The blocks are marked by colour. Red colour corresponds to the centre of significant block, where the number of positively and significantly correlated pairs exceeds the number of negatively and significantly correlated pairs. In other case, the centre position of significant block is marked by blue.

Link to 'Positional data' (6) leads to the table where results for significantly correlated pairs are represented in numerical values.

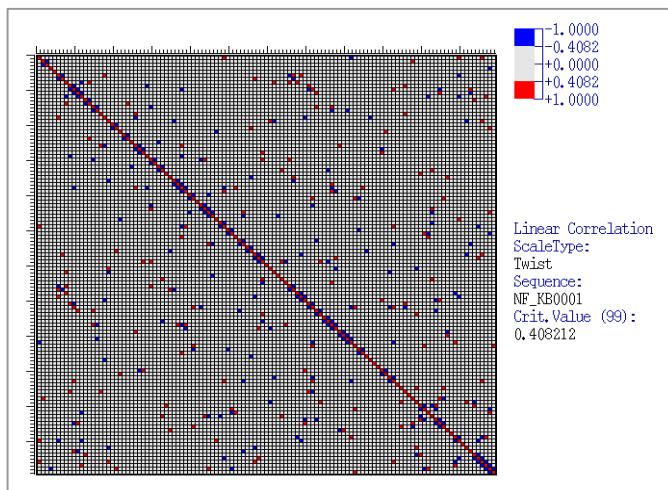
### Example

Example of the program execution for alignment of 39 sequences of NF-kappaB transcription factor binding sites (available in the SAMPLES database), with the sample length of 120 nucleotides and chosen property 'Twist(Averaged for X-rays)' .

### Results

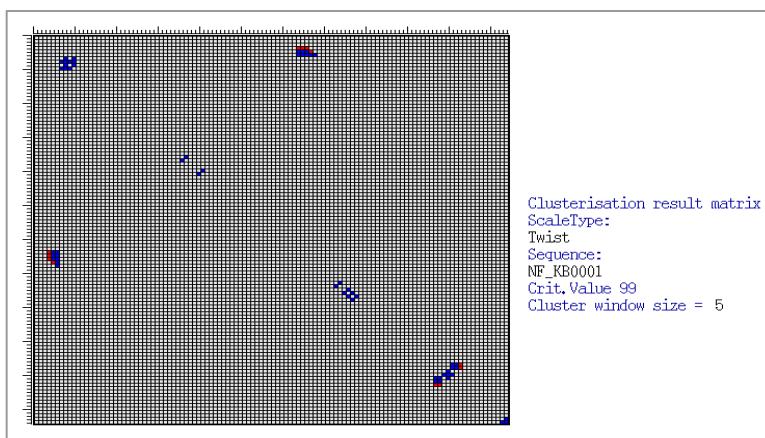
This figure exemplifies the matrix of correlation coefficients of the Twist (Averaged for X-rays) values for a sample of 39 sequences of NF-kappaB sites. Insignificant correlations are marked grey,

significant are coloured. The value scale is shown in the top right corner. The critical value is calculated in accordance with the Student criterion at a 99% significance level:



Matrix has both isolated elements corresponding to the pairs of significantly correlating positions and the clusters of such elements. To reveal such clusters (blocks), the clusterisation method could be applied. By clicking the hyperlink 'Clusterisation results', the results of the program execution are displayed.

This Figure illustrates positioning of significant blocks in correlation coefficients matrix calculated for Twist (Averaged for X-rays) property values for NF-kB sites. By X- and Y-axes, positions of clusters are marked. The value of correlation coefficient significance is estimated by the formula (2) under the confidence level of 99%. Clusterisation was made for the window with the size 5x5 and confidence level of 99%, calculated by the formula (3). The blocks are marked by colour. Red colour corresponds to the centre of significant block, where the number of positively and significantly correlated pairs exceeds the number of negatively and significantly correlated pairs. In other case, the centre position of significant block is marked by blue.



Thus, it is possible to detect the sequence regions that could be important for functioning of regulatory regions as well as relationships of these regions to each other. For example, as seen from the Figure, the cluster located in the upper central part indicates that the left flank of the site (positions 4-6) is related to the site centre (positions 74-79) according to the Twist property value. Thus, for the proper site functioning, there are functional restrictions for nucleotide content of the sequence in this regions.

**Comments and questions are welcome to Dmitry Y Oshchepkov ([diman@bionet.nsc.ru](mailto:diman@bionet.nsc.ru))**

## 4.2. Program for estimation of stochastic complexity of genetic texts

Release 2003

**Program description:** The program is designed for estimation of symbolic data compression.

**Access to the program:**

<http://www.domain.com/mgs/gnw/regscan/> link Complexity

**List of biological tasks that could be solved by using the program for estimation of stochastic complexity:**

- Estimation of general regularities in the context structure in a set of nucleotide sequences with low homology.
- Estimation of general regularities in the context structure in a set of amino acid sequences with low homology.
- Evaluation of homogeneity in a group of sequences.
- Detection of statistically significant short contexts in a sequence or in a group of sequences.

**Data input.**

The sequence to be analysed should be entered into the text-box in the FASTA-format. The divisor between separate sequences is the line with the first symbol '>'. The sequence can be entered from file of the user's computer by clicking the option 'From file'.

**Estimation of text complexity. Construction of context tree**

**DNA sequences :**

Standard alphabet {A,T,G,C}  GC context {A/T,G/C}

User-defined nucleotide alphabet in brackets  [at][gc]

**Amino acid sequences:**

Hydrophobic-hydrophilic residues (2-letter: (-)hydrophobic, (+)hydrophilic)

Residue charge (3-letter: (+)base, (0)neutral, (-)acid)

Residue surface (3-letter: (+)outer, (0)ambivalent, (-)inner)

User-defined amino acid alphabet in brackets  [ailmfpwv][rndcqeghksty]

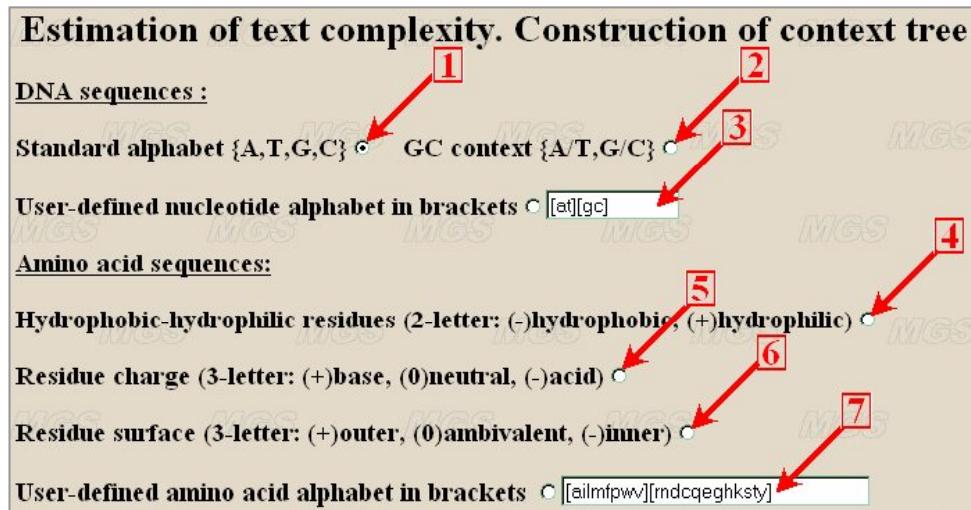
Enter sequence in plain format  
 from Screen (cut & paste)... 

from File: 

(The sequences should be in FASTA format, may be only one sequence without comment line)

If there is only a single sequence to be analysed, you may input this sequence in a plain text format without the comment line. The program has no formal limitations for a sequence length. It is recommended to analyse sequence not less than 500 bp.

### Program options.



Let us determine the alphabet. If it is necessary to analyse DNA sequence, choose the alphabet for analysis of the text. It is possible to use the nucleotide symbols {A, T, G, C} (the standard alphabet, by default) (number 1 in the figure), or their combinations (for example, A/T, G/C) for studying GC-content (number 2). A user may use another combination of letters by choosing the option 'User defined alphabet'. Then a user may set in a text-box the alphabet (number 3), for example, in a form [AT][GC] or [TC][AG], or A[CGT]. The symbols in square brackets are interpreted as a single symbol. The symbols that are not indicated in square brackets are ignored.

To analyse amino acid sequences, choose the alphabet for analysis of the protein sequence. In this case, the parameters of nucleotide alphabet are ignored. The program searches for the contexts in the alphabet with the length up to 5 symbols. That is why 20 symbols denoting amino acid residues should be united in the groups. We suggest the following variants of grouping of 20 residues: by hydrophobicity-hydrophilicity (number 4 in the figure), by charge (number 5), by location (inner positioning or the surface one; number 6). A user may order his own variant of partitioning by ordering in appropriate window the line indicating how to group the symbols (number 7). By ordering by a user of his own alphabet, the residues that are not indicated will be ignored. In case the alphabet of more than 5 symbols is ordered, the estimation of complexity will be incorrect.

Only a single type of the alphabet should be chosen (see numbers 1-7).

The program parameters are ordered as indicated by arrows in the figure given below.

Context length (2<n<12)  1  
 Text report. Output all found context 2  
 Graphic output 3  
 Tree types - Standard tree or  Round tree 4  
 Graphic output parameters: Letters in image Yes 5  
 Max. width of picture (in pixels) 6  
 Max. height of picture (in pixels) 7  
 8  
 9

(1) The context length could vary from 2 to 12. The default value is 3. This is the length of the maximal context significant for generation of the symbol (dimension of the Markov model). It is recommended to input the maximal value of a parameter equalling to 12.

(2) By clicking the check box 'Text report' the program displays the complete statistical data in a text form. By default, the program will display these data. Since these data may have a large volume (it exceeds  $4^d$  lines, where 'd' is the context length), a user may discard this option.

(3) To display the graphical data output, select the appropriate check box by default. Select the form of graphical representation of the prefix context tree.

(4) By default, the tree with suspended vertexes (leaves) is ordered ('Standard tree'). The contexts are read from top to the bottom or from leaves to the root. This corresponds to disposition of contexts in the sequence from the left to the right (for DNA, from 5' to 3' end; for proteins, from N- to C-end). A user may choose the tree representation in a form of the radial tree. This representation is more compact. In this case, the contexts are read from outer surface to the centre (root of a tree).

(5) Additional parameters for the tree representation set the data output without denotation of vertexes (without letters) for simplicity. By default, the graphical output is displayed with the letters from the alphabet ordered by a user. If the alphabet is not standard, the symbols are replaced by numbers.

(6) A user may also set the options for the scale of graphical representation or to make the picture less for more convenient copying from screen and further usage. The absolute size of the picture by height and width in pixels are set in appropriate text-boxes. The maximal size is set by default. It is recommended to set the size of at least 200 pixels.

## Program execution.

The program is executed by clicking the button 'Execute' (see number 7 in figure given above).

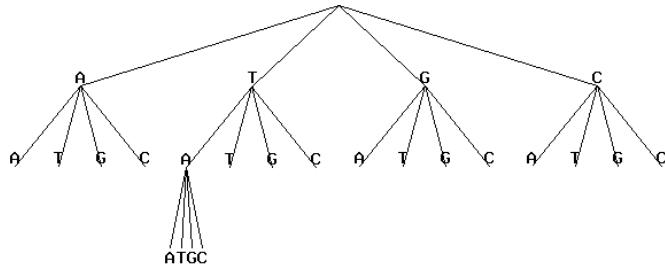
## Data output.

The program outputs the results in a textual or graphical format as was indicated by a user, by means of clicking the check boxes 2 or 3 shown in the figure above.

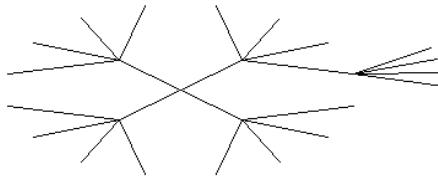
The program displays in the text format the value of complexity for a sequence entered and a set of selected significant contexts.

Example of the data output:

In a graphical representation, the visualisation of the context tree in a standard or radial form is displayed as a GIF-file. An example of representation of the context tree in a standard form is given below:



The same context tree could be represented in a radial form:



### Example.

Let us consider the task of searching for context dependencies in the set of AP-1 binding sites (32 kB, Eukaryotic Transcription Factor Binding Sites Compilation <http://wwwmgs.bionet.nsc.ru/mgs/dbases/nsamples/>).

Enter this set in a FASTA-format into the input text box (step 1).  
Set the standard alphabet {A, T, G, C} (step 2).

**Estimation of text complexity. Construction of context tree**

DNA sequences : 2

Standard alphabet {A,T,G,C}  GC context {A/T,G/C} 1

User-defined nucleotide alphabet in brackets  [at] [gc]

Amino acid sequences:

Hydrophobic-hydrophilic residues (2-letter: (-)hydrophobic, (+)hydrophilic)

Residue charge (3-letter: (+)base, (0)neutral, (-)acid)

Residue surface (3-letter: (+)outer, (0)ambivalent, (-)inner)

User-defined amino acid alphabet in brackets  [ailmfpvw] [rndcqeghksty]

Enter sequence in plain format  
 from Screen (cut & paste)... 1

> AP\_10001  
ggaactgggc ggagtttaggg gccccatggg cggagttagg ggcggggacta tggttgc  
ctaattgaga tgcattgttt gcatacttct gcctgtggg gagcctgggg actttcc  
> AP\_10002  
catgctttgc atacttctgc ctgtgggg aacctggggac ttccacacc tggttgc  
ctaattgaga tgcattgttt gcatacttct gcctgtggg gagcctgggg actttcc

from File:  Browse...  
*(The sequences should be in FASTA format, may be only one sequence without comment line)*

Set all the parameters by default.

Set the context length equal to 3 (step 3 in the figure below).

Set the textual output of contexts (step 4).

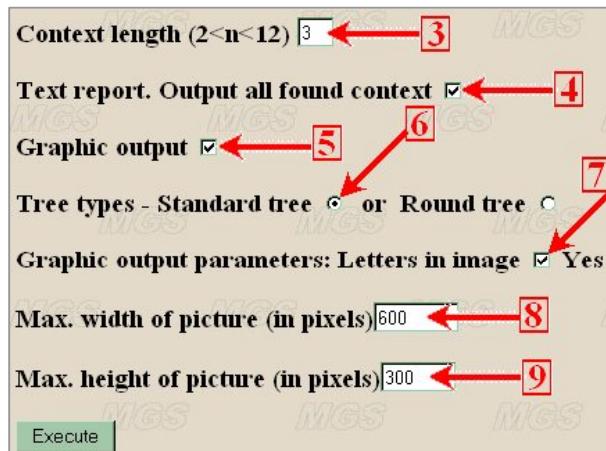
Set the graphic mode of the output data (step 5).

Set tree type as a standard tree (step 6).

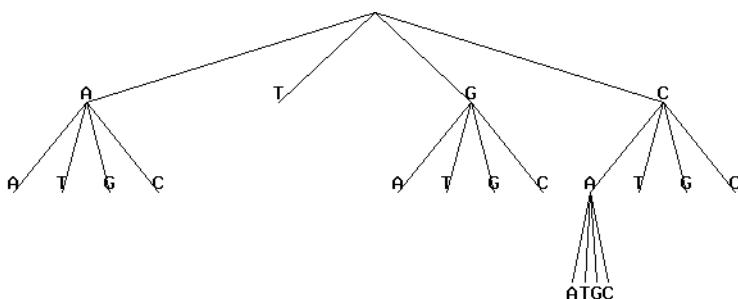
Set parameter 'letters in image' equal to 'Yes' (step 7).

Set maximal width of the picture equal to 600 pixels (step 8).

Set maximal height of the picture equal to 300 pixels (step 9).



The results of the program execution will be displayed as the figure shown:



As seen from the figure, for the sequence considered, the statistically significant are the dinucleotides of the form NA, NG, NC, where N=A/T/G/C. This fact should be taken into account for constructing recognition programs. As illustrated additionally, the contexts of the form NAC are significant too.

**Comments and questions are welcome to Orlov Yu.L. ([orlov@bionet.nsc.ru](mailto:orlov@bionet.nsc.ru))**

# CHAPTER 5. ACTIVITY SYSTEM

## 1. ACTIVITY Database

Release 2003

### 1. Database description:

We suggest a new approach to predict the activity of DNA functional sites that is focused on the perspicuity of the prediction in terms of "a probable molecular mechanism of the site functioning". The biological novelty of the method is in the involvement of physico-chemical and conformational DNA properties to provide clear interpretation of the obtained activity predictions in terms of a probable molecular mechanism of the site functioning.

### 2. Access to ACTIVITY database:

<http://www.domain.com/mgs/gnw/activity/>

### 3. Database content

SRS table	Description	Number of entries
ACTIVITY	A database for accumulation of activity values of DNA/RNA functional sites. Each entry in the ACTIVITY database describes a set of "sequence-activity" data measured in a fixed experimental system. The entry presents the "sequence-activity" data in a computable format.	554
KNOWLEDGE	It contains significant features of DNA/RNA site sequences determined experimentally (i.e. sequence-activity relationships under fixed experimental conditions), and the programs for predicting activity of these sites.	22
REFERENCE	Database on references to original publications describing the experimental data on DNA/RNA site sequences with known activity magnitudes.	265
PROPERTY	Database on sequence-dependent conformational and physico-chemical B-DNA parameters	38
WEIGHT	Database on weight	33
SCIENTIST	Database on scientists annotating the literature data	29

### 4. List of biological tasks that could be solved by using ACTIVITY database

- To browse all available experimental data concerning exploration of DNA functional sites activity.
- To explore data about the significant features of the DNA/RNA site sequences determined experimentally (i.e. sequence-activity relationships under fixed experimental conditions). Information is stored within KNOWLEDGE supplementary database.
- To predict possible type of DNA/RNA functional sites activity in the sequence of interest. For a fixed type of activity, the predicting C-program is stored in the KNOWLEDGE database.
- ACTIVITY is useful for molecular biology, pharmacogenetics, metabolic engineering, drug-design, and biotechnology.

### 5. SRS table format

### ACTIVITY

Line code	Field name	Field description
OC	Taxon Specificity	This field describes the taxon specificity of the gene containing the site analyzed, if the sequences are not synthetic.
FF	Site Name	This field describes the name of the site analyzed and the site localization in the gene, if the sequences are not synthetic.
AN	Type of Activity's Measurement	This field describes the magnitude measured in an experiment.
AU	Measurement Units	This field describes the units of a magnitude.
PN	Sequence Phasing Point	This field describes the site sequence phasing point.
SC	Site's Variant	This field contains the name of the site's variant analyzed.
SQ	Site Sequence	This field contains the sequence analyzed ("+"-chain, -5'-3' direction).
SA	Activity Magnitude	This field contains the value of a magnitude measured in an experiment for the site's variant with the sequence given in the field SQ.
SD	Standard Deviation	This field contains the standard deviation value of a magnitude measured in an experiment for the site's variant with the sequence given in the field SQ.
PA	Position of the Phasing Point in the Sequence	This field indicates position of the phasing point (see field PN) relative to the start of the sequence given in the field SQ.

### KNOWLEDGE

Line code	Field name	Field description
MI	Entry ID	This field indicates an identifier.
MN	Entry Name	This field contains name of the entry.
HN	Link to SCIENTIST database	This field contains link to the SCIENTIST database and indicates a contributor of the entry.
DA	Link to Activity database	This field contains link to the ACTIVITY database.
DR	SRS-link	This field contains the links to other databases installed under SRS (ACTIVITY, SCIENTIST, SYSTEM, PROPERTY, etc.).
WW	Web-link to Activity Tools	This field contains the link to the Web-based tools implementing each C program documented within the entry to recognize the site by its significant feature or predicting activity value within an arbitrary DNA sequence.
GF	Mathematical Model	This field describes the mathematical model
CT	Computational Method	This field describes the computational method.
DW	Link to WEIGHT database	This field contains the link to the WEIGHT database installed under SRS.
DP	Link to PROPERTY database	This field contains the link to the PROPERTY database installed under SRS.
PV	Property Name	This field contains short name of the property investigated
HL	Indicates Feature Deviation	This field indicates high or low deviation of the conformational or physico-chemical feature of the site studied from the random sequences.

AB	Analyzed Region	This field indicates the DNA region, for which the conformational or physico-chemical feature differing significantly the sites studied from the random sequences was revealed. Positions are given in bp relative to the first position of the site sequences.
UT	Utility	This field indicates the utility of the context-dependent feature for discrimination of the sites studied from the random sequences
LC	Linear Correlation Coefficient	This field describes the linear correlation coefficient.
AL	Significance Level	This field describes the significance level.
FG	Graphical Representation of Test Results	This field represents the links to the graphical representation of the results obtained by the program, given below in the field C-CODE, over the independent control sequences of the site.
C-CODE		This field contains the program calculating the profile of the sequence-dependent feature or activity predicted values along arbitrary sequence in the 'C' language of the ANSI standard. Each C program documented within the KNOWLEDGE entry has a check-box in the MENU window of the Activity Tools (see field WW).

#### REFERENCE

Line code	Field name	Field description
RN	Entry ID	This field contains an identifier.
RA	Authors	This field indicates the authors of the article
RT	Title	This field contains an article title.
RJ	Journal	This field contains the journal name
RV	Volume	This field contains the number of the journal volume.
RP	Pages	This field indicates the numbers of pages.
RY	Year	This field indicates the year of publication.
RR	Abstract	This field contains the abstract of an article.
RS	Data Source	This field indicates the number of a figure or a table, which contain the data used in the ACTIVITY database.
DR	SRS-link	This field contains the links to the supplementary databases (ACTIVITY, SYSTEM, etc.).

#### PROPERTY

Line code	Field name	Field description
MI	Entry ID	This field indicates an identifier.
MN	EntityName	Each PROPERTY entry describes the only B-DNA property, either "conformational" or "physico-chemical" one.
MD	EntityDependence	Entity Dependence
ML	Step	Step
RN	Reference	The field RN links to the database REFERENCE citing the literature sources, out of which the property values were taken
PN	PropertyName	This property is entitled by the field PN
PM	Method_for_PropertyDefinition	Method for Property Definition
PV	PropertyID	
PU	Unit_for_PropertyExpression	contains the property unit
DINUCLEOTIDE		quantitative value assigned to each dsDNA dinucleotide step

## I. DNA. Chapter 5. ACTIVITY System

### SCIENTIST

Line code	Field name	Field description
ID	Identifier	Identifier
HH	Authors	Authors name
HP	Tel	Telephone
HF	FAX	FAX
HE	Email	E-mail
HL	Laboratory	Laboratory name
HI	Institute	Institute name
HD	Department	Department name
HM	Minister	Ministry name
HC	Country	Country name
HA	Address	Home address

### WEIGHT

Line code	Field name	Field description
ID	EntityID	This field indicates an identifier.
MN	EntityName	This field contains name of the entry.
MD	EntityDependence	Entity Dependence
WP	WeightName	Weight Name
WL	WeightPicking	Weight Picking
WW	WeightWidth	Weight Width
WI	WeightID	Weight ID

### Examples of SRS queries to the ACTIVITY database

The query is:

What experimental data is available concerning relative binding activity of different variants AP-1 site sequences?

To make such a query, you should perform the following:

1. Choose ACTIVITY SRS table on the page 'access to ACTIVITY'.

## I. DNA. Chapter 5. ACTIVITY System

2. This will bring up the home page of the chosen ACTIVITY SRS table. Click the button 'Search'.

The screenshot shows the ACTIVITY database homepage. At the top, there's a navigation bar with links for TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, DATABANKS, and HELP. Below the navigation bar, the word "ACTIVITY" is prominently displayed in a large, bold font. To the right of "ACTIVITY" is a "Search" button with a red arrow pointing to it. On the left side of the main content area, there are four dark blue boxes labeled "Name", "Status", "Description", and "Literature". The "Literature" box contains a numbered list of references:

1. Ponomarenko M.P. et al., Bioinformatics, 1999, 15, 7/8, 687-703.
2. Kolchanov N.A. et al., Bioinformatics, 1999, 15, 7/8, 669-686.
3. Ponomarenko M.P. et al., Database on functional DNA and RNA site activity, ACTIVITY. Russian Patent Agency (RossPatent) Certificate No. 980059, Sept. 2, 1998, Moscow, Russia (Russ)
4. Ponomarenko J.V. et al., Proceedings of BGRS'98, Ed. N.Kolchanov et al., ICG, Novosibirsk, Russia, P. 62-65. Full Text ([html](#))
5. Kolchanov N. A. et al., Proceedings of ISMB-98. Ed.: Glasgow J. et al., AAAI Press, P. 95-104. Full Text ([html](#))
6. Kolchanov N.A. et al., Mol. Biol. (Mosk), 1998, 32, 2, 255-267. (Russ).

In the center of the page, there's a message: "The current release has 554 entries and was indexed 24-Jun-2001. A distributed and intelligent database for the activities of the functional sites in DNA and RNA."

3. Select the fields to be searched for from the list. You will need the field 'SiteName' (1), Type terms to be searched in the text window: 'AP\*' (2) (An asterisk means "any symbol") and click the button 'Submit Query' (3).

The screenshot shows the ACTIVITY query interface. At the top, there's a navigation bar with links for TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, DATABANKS, and HELP. Below the navigation bar, there's a search bar with the text "search ACTIVITY" and an "Info about field EntityID" dropdown. On the left, there's a sidebar with buttons for "Reset", "Submit Query" (highlighted with a red arrow), "append wildcards to words" (with a checked checkbox), "combine searches with AND", "Number of entries to display per page" (set to 30), and "Extended query form". The main search area has a dropdown menu for "Info about field EntityID". Below this, there's a text input field containing "AP\*" with a red arrow pointing to it. To the left of the input field is a dropdown menu for "SiteName" (highlighted with a red arrow). To the right of the input field is a dropdown menu for "retrieve entries of type Entry" and a checkbox for "\*Names only\*". At the bottom, there's a section for "Select fields to display:" with a list of fields: EntityID, EntityName, GenomeRegionName, Specie, Taxon, SiteName, and ActivityName.

## I. DNA. Chapter 5. ACTIVITY System

4. This will bring up the resulting window with the list of matching entries/  
To extract all information contained in the entry, click any link

The screenshot shows a web-based application with a navigation bar at the top containing links for TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, DATABANKS, and HELP. A logo of a paw print is on the left. Below the navigation bar, a yellow header bar displays the query "Query "[activity-SiteName: AP\*]" found 2 entries". On the left, a sidebar titled "Perform operation" has two radio button options: "on all but selected" (selected) and "on selected". It contains three buttons: Link, Save, and View. A dropdown menu is open, showing the option "\* Names only \*". Below the sidebar, there is a section for "Number of entries to display per page" with a dropdown set to 30, and a "Printer Friendly" link. The main content area lists two entries: "ACTIVITY:A00J0049" and "ACTIVITY:A00LK013". A red arrow points to the link for "ACTIVITY:A00J0049". At the bottom of the page, there is a footer with the text "SRS 6.0.7.3 | feedback".

5. The query result is the whole set of information about experiment and relevant binding activity of the AP-1 site sequence variants based on the features of the AP-1 site.

The screenshot shows a detailed view of the entry ACTIVITY:A00J0049. The page title is "ACTIVITY:A00J0049". The content includes experimental details: MI A00J0049, MN binding activity of Jun-Fos dimer to AP-1 and CRE sites in vitro, YY, HN Expert Database: SCI00002, YY, OG synthetic oligonucleotides, OS in vitro, OC EUKARYOTA, YY, FF AP-1 or CRE site, AN relative binding activity, AU percent, PN start of sequences, YY. The TD section lists numerous specific sites, each followed by "(TRRD)". Some examples include: site 1125 (TRRD), site 2281 (TRRD), site 2203 (TRRD), site 2284 (TRRD), site 2286 (TRRD), site 1383 (TRRD), site 1952 (TRRD), site 2374 (TRRD), site 1503 (TRRD), site 1892 (TRRD), site 1893 (TRRD), site 981 (TRRD), site 982 (TRRD), site 1903 (TRRD), site 1904 (TRRD), site 2074 (TRRD), site 2075 (TRRD), and site 2240 (TRRD).

Comments and questions are welcome to Mikhail P. Ponomarenko, (pon@bionet.nsc.ru).

## 2. Predicting Activities of Functional Sites in DNA/RNA

### 1. Program description

The identification of the sequence-dependent DNA features correlating with affinity magnitudes of DNA sites interacting with a protein could pinpoint to molecular event limiting this protein/DNA recognition machinery. This approach is realized in computer system ACTIVITY containing the databases on site activity and on conformational and physical-chemical DNA/RNA parameters. By using the system ACTIVITY, an analysis of some sites was provided and the methods for predicting the site activity were constructed.

### 2. Access to Predicting activities of functional sites in DNA/RNA

<http://www.domain.com/mgs/gnw/activity/> link 'Predicting activities of functional sites in DNA/RNA'

### 3. List of biological tasks that could be solved by using the program

- ◆ Prediction of site activity based on its primary nucleotide sequence.

### 4. Data input

Input the DNA sequence of interest into the field 'Input DNA Sequence' (1). Sequence should be in a plain text format (a, t, g, c in up- or low case, tabulation or spaces are accepted). The sequence to be entered should be of at most length of 32 kbp and of least length equaling to the site sequence length.

**Predicting activities of functional sites in DNA/RNA**

**CRP-activator binding site in E.coli operons**

**Input DNA Sequence :**

from Screen: agccgggtgc gccccccca gtgcgcgcgg ccgggtgtt cgctcgga cgcgaagtgac  
ctcgccccgg taccactggcgggtata tcagcgcggg gctgtgtcag gcagcggccc 1

from DB: Bases Available: SRS5 from Heidelberg (EMBL) ▼

from File: Browse... [File formats here.](#)

Execute Reset form [Example](#)

See 3

CRP-activator specificity 2  
 CRP-activator specificity prediction via Inclination

This resource has been developed in Institute of Cytology and Genetics. Novosibirsk, Russia

Authors: [Misha Ponomarenko, Anatoly Frolov](#)  
Contributors: [Julia Ponomarenko, N.L. Podkolodny](#)  
Leader: [Nikolay A. Kolchanov](#)

### 5. Program options

Specify feature of prediction program by clicking one of check-boxes (2).

### 6. Program execution

Start the tools processing by clicking the button 'Execute' (3).

## 7. Data output

The tools output represents the profile of the significant feature or predicting activity value.

### Example

To search for activity of the CRP-activator binding site in E.coli operons

1. Choose the link ‘predicting activities of functional sites in DNA/RNA’.

2. Choose what kind of activity should be predicted. Each C program documented within the KNOWLEDGE entry has a check-box in the MENU windows:

3. Input the sequence of the CRP-activator binding site into the text-box window (1)

```

ggaactgggc  ggagtttaggg  gcgggatggg  cggagtagg  ggcgggacta  tgggtgctga
ctaattgaga  tgcatacgctt  gcataacttct  gcctgctggg  gagcctgggg  actttccaca
ggaactgggc  ggagtttaggg  gcgggatggg  cggagtagg  ggcgggacta  tgggtgctga
ctaattgaga  tgcatacgctt  gcataacttct  gcctgctggg  gagcctgggg  actttccaca

```

Choose the option ‘CRP-activator specificity’ (2) and click the button ‘Execute’ (3).

**Predicting activities of functional sites in DNA/RNA**

**CRP-activator binding site in E.coli operons**

**Input DNA Sequence :**

from Screen:  
`agccgggtgc gccccggccca gtgcgcgcgg ccgggtgttt cgcttggagc aadgtgac  
ctcgccccgg taccactggcggtgata tcagcgccgg gctgtgtcag gacgcggcccg`

from DB:      **Bases Available:** SRS5 from Heidelberg (EMBL)

from File:      [Browse...](#) [File formats here.](#)

**Execute**    **Resetform**    [Example](#)

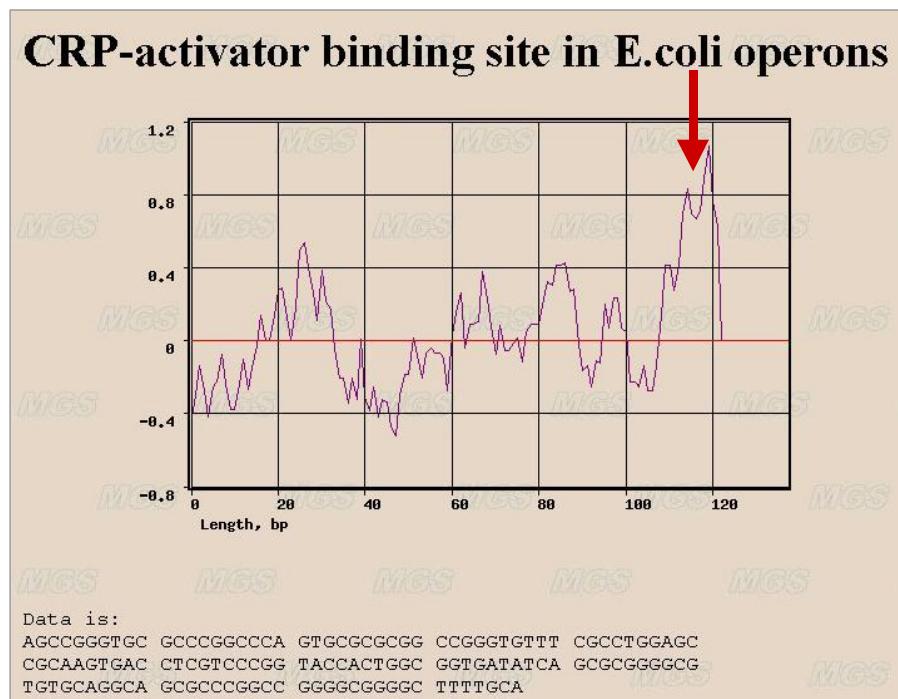
**See**

CRP-activator specificity      **2**

CRP-activator specificity prediction via Inclination

**1**

4. The output represents the profile of the Score value. The positive peaks of this profile point to the potential site recognized.



Comments and questions are welcome to Mikhail P. Ponomarenko, (pon@bionet.nsc.ru).

# CHAPTER 6. DNA NUCLEOSOMAL ORGANISATION

## 1. PROFILES Databases

Release 2003

### 1. Database description

PROFILES and PROFILE\_LIST databases accumulate an information about the profiles of conformational and physico-chemical DNA parameters, their significant extremum points and linear trends.

### 2. Access to PROFILES Databases

<http://www.domain.com/mgs/gnw/nucleosom/> links 'PROFILE' 'PROFILE\_LIST'

### 3. Database content

SRS table	Description	Number of entries
PROFILE_LIST	This SRS table accumulates sets of significant physico-chemical properties profiles of the nucleosome formation sites.	6
PROFILE	database on profiles of conformational and physico-chemical B-DNA properties	228

### 4. List of biological tasks that could be solved by using the PROFILES and PROFILE\_LIST databases

- ◆ Searching for and characteristics of extremal values and linear trends for the profiles of conformational and physico-chemical properties of nucleosomal DNA

### 5. SRS table format

#### PROFILE\_LIST

Line code	Field name	Field description
ID	ProfileListID	Identifier
AC	AccessNumber	Accession number
SD	SiteDescription	Site description

#### PROFILE

This SRS table accumulates significant physico-chemical properties profiles of the nucleosome formation sites.

#### SRS table format

Line code	Field name	Field description
ID	ProfileID	identifier of an entry
SD	SiteDescription	site description
PW	ProfileWindow	profile window size
PA	ProfileAverage	profile average value
PD	StandardDeviation	standard deviation
ET	ExtremumType	type of extremum: minimum or maximum
EV	ExtremumValue	profile value in the position of extremum
EL	ExtremumLevel	significance level of extremum
GC	SignCriteria	criteria for gradient significance evaluation

The database PROFILE\_LIST contains entry headers, which present site samples and the lists of DNA conformational and physico-chemical properties. Each property in the list refers to individual property entry, which is contained in the database PROFILE. The property entry presents the profile of DNA conformational physico-chemical property and description of its extremum points and linear trends. To construct a profile, various sliding window sizes were applied. The evaluation of the best suited window size was made by comparison of the sample profile with the profile constructed for random sequences with the same dinucleotide content as that of the site. The best selected profile is included into the database entry. Then the profiles were checked for linear trends and the resulted data are also compiled in the database.

### Example of SRS queries to the PROFILE database

The example of querying the database is as follows:

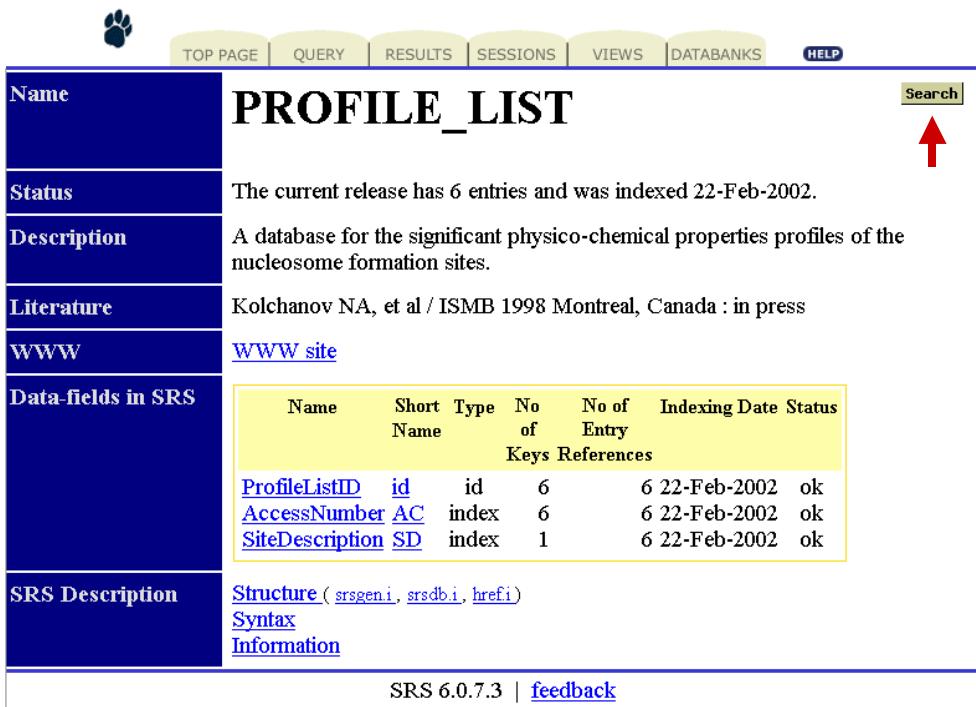
Which conformational and physico-chemical DNA properties profiles are significant for nucleosome binding DNA-region in 5' & 3' gene regions?

To make such a query, you should perform the following:

1. Choose 'PROFILE\_LIST' SRS table on the page 'access to DNA NUCLEOSOMAL ORGANIZATION'.

The screenshot shows the Genenetworks website interface. At the top, there is a navigation bar with links for HOME, DNA, RNA, PROTEIN, GENENETWORKS, and MAP. Below the navigation bar, there is a logo for TFBSR (Transcription Factor Binding Site Recognition) and a main content area. The main content area has several sections: 'General information' (with links to 'How to cite Nucleosomal organization?', 'Nucleosomal organization publications', 'The latest report on Nucleosomal organization', 'Nucleosomal organization Workgroup', 'User's guide', 'Module scheme', 'Help', and two links under 'Help'), 'ACCESS' (with links to 'SRS access: PROFILE\_LIST', 'Context properties: Nucleosome site recognition', and 'Conformational (physicochemical) properties: Nucleosomal organization profiles'), 'General information' (with links to 'How to cite Nucleosomal organization?', 'Nucleosomal organization publications', 'The latest report on Nucleosomal organization', and 'Nucleosomal organization Workgroup'), and 'User's guide' (with links to 'Module scheme', 'Help', and two links under 'Help').

2. This will bring up the home page of the chosen PROFILE\_LIST SRS table. Click the button 'Search'.



**PROFILE\_LIST**

The current release has 6 entries and was indexed 22-Feb-2002.

A database for the significant physico-chemical properties profiles of the nucleosome formation sites.

Kolchanov NA, et al / ISMB 1998 Montreal, Canada : in press

[WWW site](#)

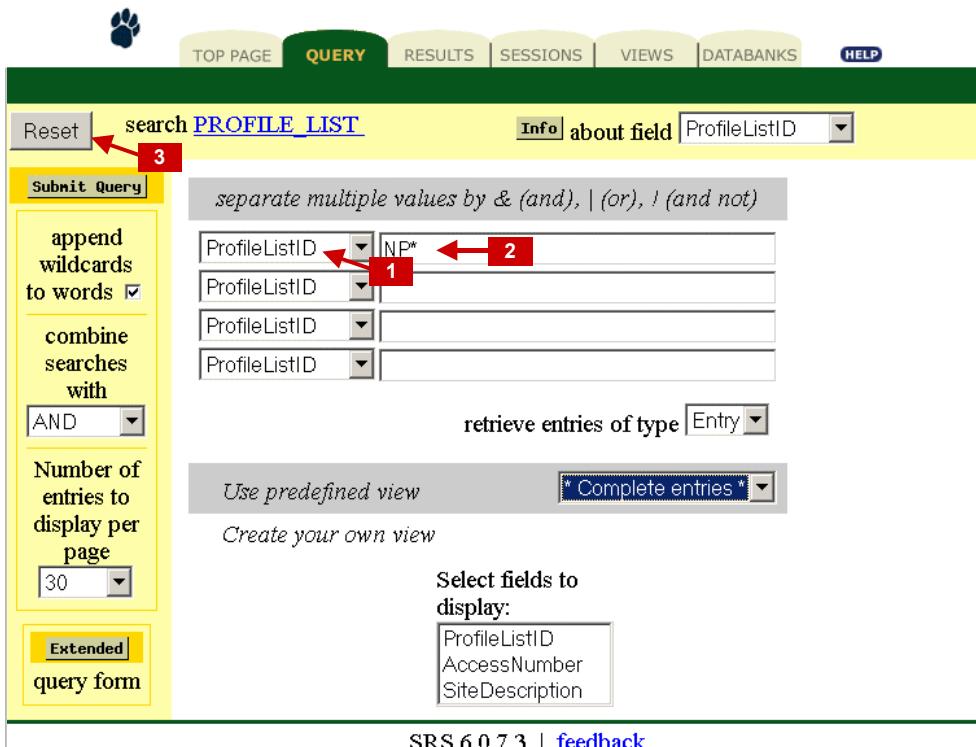
Name	Short Name	Type	No of Keys	No of References	Indexing Date	Status
<a href="#">ProfileListID</a>	<a href="#">id</a>	id	6		6 22-Feb-2002	ok
<a href="#">AccessNumber</a>	<a href="#">AC</a>	index	6		6 22-Feb-2002	ok
<a href="#">SiteDescription</a>	<a href="#">SD</a>	index	1		6 22-Feb-2002	ok

**SRS Description**

[Structure](#) ([srsgen.i](#), [srldb.i](#), [href.i](#))  
[Syntax](#)  
[Information](#)

SRS 6.0.7.3 | [feedback](#)

3. Select the fields to be searched for from the list. You will need the field 'ProfileListID' (1), Type terms to be searched in the text window: 'NP\*' (2) (An asterisk means "any symbol") and click the button 'Submit Query' (3).



search **PROFILE\_LIST**

**Info** about field **ProfileListID**

separate multiple values by & (and), | (or), ! (and not)

ProfileListID	NP*
ProfileListID	
ProfileListID	
ProfileListID	

retrieve entries of type **Entry**

**Use predefined view**      **\* Complete entries \***

**Create your own view**

Select fields to display:

ProfileListID
AccessNumber
SiteDescription

30

**Extended query form**

SRS 6.0.7.3 | [feedback](#)

The query result will be displayed as the complete list of entries containing significant B-DNA property profiles of nucleosome binding DNA in 5' & 3' gene regions.  
Click the link 'NP00001; Twist'.

The screenshot shows a web-based application for querying nucleosome profiles. At the top, there's a navigation bar with links for TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, DATABANKS, and HELP. A logo of a paw print is on the left.

The main area has a yellow header bar with the text "Query '[profile\_list-ProfileListID: NP\*]' found 1 entries". Below this, a form titled "Perform operation" contains two radio buttons: "on all but selected" (selected) and "on selected". It includes buttons for "Link", "Save", and "View". A dropdown menu labeled "\* Complete entries \*" is also present.

Below the form, there's a section for "Number of entries to display per page" with a dropdown set to "30". A "Printer Friendly" link is available.

The main content area displays a list of profile entries under the heading "PROFILE\_LIST:NP". Each entry consists of a code (e.g., ID, AC, CV, SD, LD, LP, SL, SZ, PI) followed by a description. The descriptions include terms like "Nucleosome sites from 5' & 3' gene region sequences", "center of the footprint region", and various DNA properties such as Twist, Rise, Bend, Tip, Inclination, Major groove width, Depth, Minor groove width, Depth, Free DNA Roll, Twist, Tilt, Slide, DNA/Protein-complex Roll, Twist, Tilt, Slide, B-DNA Twist, Wedge, Direction, Persistence length, and Melting temperature.

The result will be displayed as the complete entry containing the profile description of the nucleosome binding DNA in 5' & 3' gene regions of the conformational B-DNA property 'Twist'

## I. DNA. Chapter 6. DNA Nucleosomal Organisation

```
PROFILES NP00001

PI NP00001
SD Nucleosome sites from 5' & 3' gene region sequences
PR p0000001; Twist
IP http://www.sqi.sscc.ru/Programs/acts2/images/TWIST.html
YY
CR Student criterion
FG http://www.sqi.sscc.ru/Programs/NucSitRec/images/ptwist1\_1.gif
PW 9
FG http://www.sqi.sscc.ru/Programs/NucSitRec/images/ptwist1\_2.gif
PA 36.538
PD 0.123
PN 2 2
ET Minimum
EN 1
EV 36.230
EL 0.02
EP -20.5
ET Maximum
EN 2
EV 36.230
EL 0.02
EP 20.5
ET Maximum
EN 1
EV 36.830
EL 0.02
EP 41.5
ET Maximum
EN 2
EV 36.830
EL 0.02
EP -41.5
YY
GC Fisher-Snedecor criterion
FG http://www.sqi.sscc.ru/Programs/NucSitRec/images/ptwist1.gif
GW 29
GX [-44.5, -16.5]
GM -30.5
GA 36.544688
GL 0.050000
GK -0.026441
GI 0.021562
GF Y(X) = 36.544688 + (-0.026441 * (X + 30.5))
YY
GW 29
GX [16.5, 44.5]
```

**Comments and questions are welcome to Victor Levitsky (levitsky@bionet.nsc.ru)**

## 2. Software

### 2.1. Nucleosome binding site recognition

Release 2003

#### 1. Program description

Program generates a profile of the so-called nucleosome positioning potential.

#### 2. Access to the program:

<http://www.domain.com/mgs/gnw/nucleosom/> link ‘Nucleosome site recognition’.

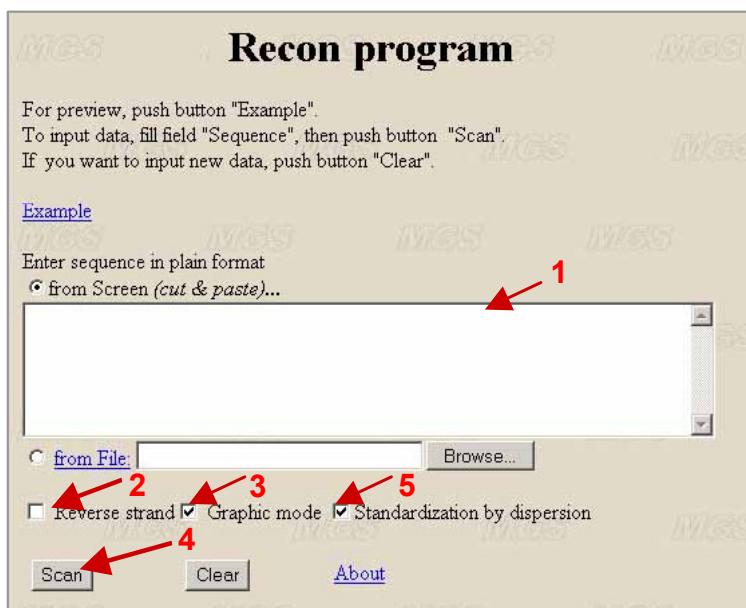
#### 3. Biological task that could be solved by using the program:

Recognition of nucleosome formation sites by oligonucleotide content.

#### 4. Input Data

Input the DNA sequence into the field ‘Input DNA Sequence’ (1). The sequence Length should be at least 160 bp and at most 32 kbp. The sequence should be in a plain textual format (a, t, g, c in upper- or lower case, line feeds and blanks are ignored). The sequence to be analysed should have the maximal length of 32 kb and the minimal length of 160 bp). The dinucleotide relative abundance distance was chosen as an additional restriction for input data to exclude the sequences with poor dinucleotide content. Positions of the sliding window that are ignored by the program are marked by colour in the graphical representation of the data output and by symbol \*, in numerical delivery of the program.

#### 5. Program options



The analysis of complementary DNA strand is executed by clicking the check-box ‘Reverse strand’ (2). To select the data output mode, choose the option by clicking the check-box ‘Graphic mode’ (3). To select the output profile transformation mode, choose the option by clicking the check-box ‘Standardization by dispersion’ (5).

## 6. Program execution

Click the button ‘SCAN’ and wait for the program execution (4).

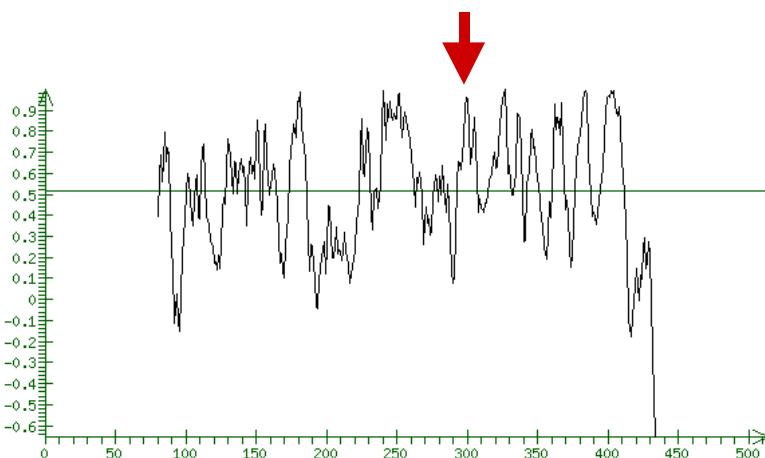
## 7. Output Data

Output data could be displayed in two modes: graphical representation and numerical delivery.

### Example

Nucleosome potential is constructed in a way so that its mean value over the training set of nucleosome site sequences equals +1; whereas over the set of non-site (random) sequences, -1. This means that the values close to +1 correspond to a higher probability of nucleosome positioning. The option 'Standardization by dispersion' transforms nucleosome potential values to the interval  $\leq 1$ , here value +1 corresponds to the best prediction, and interval [0; 1] corresponds to reliable nucleosome formation sites. For example, we consider here the analysis of *Xenopus laevis* TFIIIA gene, with the center of nucleosome site placed approximately at 290 position relatively the sequence start (marked in the Figure above). Click the button ‘Example’ to display the example of the program execution.

The graphical representation is illustrated in the Figure below.



This Figure illustrates numerical delivery.

81	0.396568
82	0.589130
83	0.686781
84	0.562456
85	0.657379
86	0.793726
87	0.653407
88	0.716841
89	0.688931
90	0.414675
91	0.231619
92	0.151196
93	-0.115555
94	-0.051597
95	0.026091
96	-0.114321
97	-0.150909
98	0.020630
99	0.209199
100	0.380815
101	0.510850
...	...

Comments and questions are welcome to Victor Levitsky  
(levitsky@bionet.nsc.ru)

## 2.2. Recognition tools

### 1. Program description

The nucleosome binding sites are characterized by specific sets of significant conformational and physico-chemical DNA properties. The site recognition programs are generated by the sets of significant conformational and physico-chemical DNA properties (DNA features).

### 2. Access to RECOGNITION TOOLS

<http://www.domain.com/mgs/gnw/nucleosom/> link ‘Nucleosomal organization profiles’.

### 3. List of biological tasks that could be solved by using the program

- ◆ nucleosome site recognition by the most significant DNA conformational and physico-chemical properties.

### 4. Data input

Input the DNA sequence of interest into the text window (1). Sequence should be in a plain textual format (a, t, g, c in up- or low cases, spaces or tabulation are accepted). The sequence to be input for analysis should be of maximal length of 32 kbp and minimal length of 160 bp.

**Nucleosome binding DNA-region (5' & 3' gene parts)**

Input DNA Sequence >200 bp:

from Screen: 1  
 from DB: Bases Available: SRS5 from Heidelberg (EMBL) by ID  
 from File: Browse... [File formats here.](#)

Execute Resetform

N Choose profile

1. 5' & 3' gene region nucleosome sites character is the Highest Twist   
2. 5' & 3' gene region nucleosome sites character is the Lowest Rise   
3. 2 5' & 3' gene region nucleosome sites character is the Highest Roll   
4. 5' & 3' gene region nucleosome sites character is the Highest Tilt   
5. 5' & 3' gene region nucleosome sites character is the Highest Nucleosomeness   
6. 5' & 3' gene region nucleosome sites character is the Lowest Propeller   
7. 5' & 3' gene region nucleosome sites character is the Lowest Clash   
8. 5' & 3' gene region nucleosome sites character is the Highest Entalpy   
9. 5' & 3' gene region nucleosome sites character is the Highest Entropy   
10. 5' & 3' gene region nucleosome sites character is the Highest Free Energy   
11. ALL above Mean-Likeness

### 5. Program options

Select the necessary recognition program by clicking one of the check-boxes (2). Each check-box refer to a significant conformational or physico-chemical DNA property. Generalized (mean) recognition profile is calculated using all known significant conformational and physico-chemical features, contained in B-DNA features database.

## 6. Program execution

Start the tools processing by clicking the button 'Execute' (3).

## 7. Data output.

The tools output represents the profile of the Score value. The positive peaks of this profile pinpoint to potential nucleosome site recognized.

### Example

Search for potential nucleosome sites of the *Xenopus laevis* TFIIIA gene 5' region and partial coding region, EMBL ID XLTf3A5.

1. On the DNA NUCLEOSOMAL ORGANIZATION home page, click the link 'Nucleosomal organization profiles'.

Nucleosomal organization. Nucleosomal DNA property database comprises sets of nucleosomal DNA sequences, most important conformational and physico-chemical properties of nucleosomal DNA, descriptions of conformational and physico-chemical properties profiles of nucleosomal DNA, programs for nucleosome site recognition based on the context or conformational and physico-chemical properties of nucleosomal DNA.

**ACCESS**

SRS access: [PROFILE](#) [PROFILE LIST](#) [Context properties: Nucleosome site recognition](#) [Conformational \(physicochemical\) properties: Nucleosomal organization profiles](#)

**General information**

[How to cite Nucleosomal organization?](#) [Nucleosomal organization publications](#) [The latest report on Nucleosomal organization](#) [Nucleosomal organization Workgroup](#)

**User's guide**

[Module scheme](#) [Help:](#)

- [program to recognize nucleosome binding sites](#) [- DNA Property plot](#)

2. On the page 'Nucleosomal organization profiles' click the link 'NP nucleosome binding DNA-region'

**Nucleosomal organization profiles**

[NA - Nucleosome binding DNA-region](#)  
[NP - Nucleosome binding DNA-region](#) [NR - Nucleosome binding DNA-region](#)  
[NS - Nucleosome binding DNA-region](#)  
[NuclHM - Nucleosome binding DNA-region](#)  
[NV - Nucleosome binding DNA-region](#)

**General information**

[How to cite Nucleosomal organization?](#) [Nucleosomal organization publications](#) [The latest report on Nucleosomal organization](#)

[Nucleosomal organization Workgroup](#)

**User's guide**

[Module scheme](#) [Help:](#)

- [program to recognize nucleosome binding sites](#) [- DNA Property plot](#)

3. Input into the text-box (1) the following sequence:

```
agatctattgagaaaggccttactgtgtgctgttaatttagatgtgttagttatcgcaactcctgtg
tggaccattgcattccatcacattcacaacagttacagttctccaacaccagcagctgctgcacac
cgtttcctcgcttcatgttattatcacgtgctccactaggactcaaccactaagaacgagggaggt
gtccagaaaacacccaacttgtgaaataaacatcgctgacataaaacacaggaatttaacatcctt
tttttaagttgcagcgcaattactgtgaaacttcccgatgtgcgataatggttgtcctagagctat
gccaatccttcagacatcgcaaaacttcccgatgtgcgataatggttgtcctagagctat
ataaacaggcacacatggcgctacagtgcgttctacaagttcagaggaagccgagggcagcttag
ttactgaaggagatgggagagaaggcgctgccgggtgtataagcggtac
```

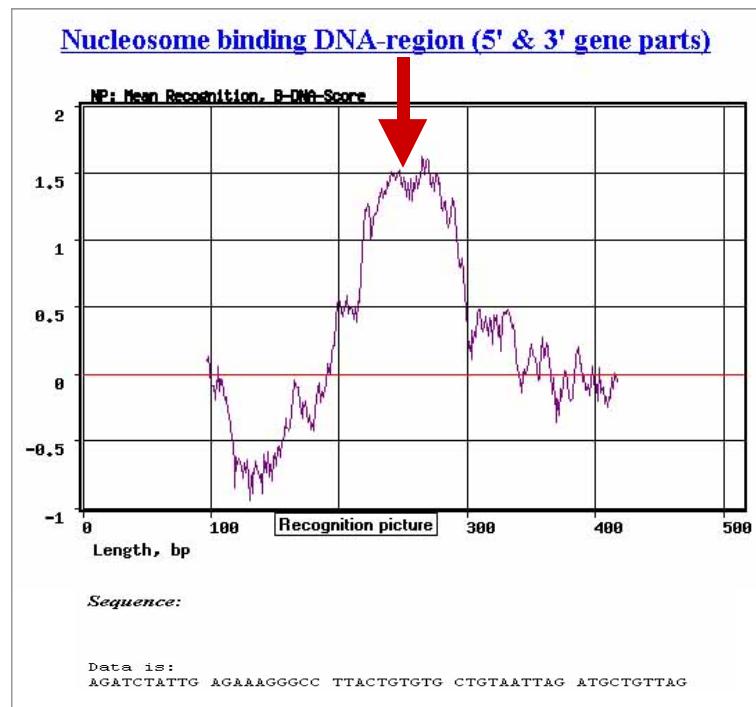
(*Xenopus laevis* TFIIIA gene 5' region and partial coding region, EMBL ID XLTf3A5).

Use program options set by default.

Click the button 'Execute' (2)

The screenshot shows the 'Input DNA Sequence' field containing the provided DNA sequence. A red arrow labeled '1' points to the 'Execute' button. Another red arrow labeled '2' points to the 'Execute' button again, indicating it twice. Below the sequence input, there are three radio button options: 'from Screen', 'from DB', and 'from File'. The 'from Screen' option is selected. To the right of the sequence input, there is a 'Bases Available' dropdown set to 'SRSS5 from Heidelberg (EMBL) by ID'. Below the input area, a message says 'Select one of the B-DNA Features listed below to analyse the Sequence inputted:'.

4. Data output window:



Comments and questions are welcome to Victor Levitsky (levitsky@bionet.nsc.ru)

## PART II. RNA INTEGRATION LEVEL.

### CHAPTER 1. LEADER RNA.

#### 1. LEADER\_RNA Knowledge Base.

##### LEADER\_RNA Knowledge base description.

The LEADER\_RNA was designed to evaluate translational activity of mRNAs of mammals, dicot plants and monocot plants.

##### Access to LEADER\_RNA prediction program:

<http://www.domain.com/mgs/gnw/leader/>

##### LEADER\_RNA content

Knowledge base LEADER\_RNA is linked to the prediction program allowing to evaluate similarity of user-defined mRNA to high and low expression ones. It consists of five data files installed at the SRS platform. Note: this knowledge base may not be used efficiently by itself but as a data source for both the prediction program and mRNA features description.

SRS table	Description	Number of entries
LEADER_SQ	sequences of 5'UTRs of high- and low-expression eukaryotic mRNAs	879
LEADER_REF	brief summaries of several papers concerning mRNA translational activity.	22
LEADER_WHY	contains brief description of mRNA contextual features influencing translatability	27
LEADER_KN	contextual features of mRNA 5'UTRs significantly different between high and low expression genes and C-codes of the prediction program	34
LEADER_SCI	reference authors	2

##### 4. List of biological tasks that could be solved by using the LEADER database

- evaluation of mRNA translational efficiency in higher plant or mammalian cells. It answers to the following question: to what extent the user-defined mRNA is similar to mRNAs of high or low expression genes? Consequently, similarity of this mRNA 5'UTR on this region of high expression genes can indicate high expression rate of its gene.
- gene engineering experiments (e.g., for prediction of translation properties of modified mRNAs in model experiments, as well as for evaluation of translation efficiency of transgene mRNA in new host organism).

##### LEADER\_RNA knowledge base: SRS table format

Note, this knowledge base may be used efficiently only in combination with the prediction program to get info on some analysis details: it was not designed for independent usage.

**Leader\_SQ**

Line code	Field name	Field description
MI	EntityID	Identifier of taxon-specific subdatabase
MN	EntityName	General information on the database
KN	CrossRef	Reference on LEADER_KN database
OG	GenomeRegionName	Description of stored genes
OS	Specia	Taxon info
OC	Taxon	Taxon info
FF	SiteName	Gene region
AN	ActivityName	Expression characteristics
AU	Unit of ActivityExpression	Estimation type, 1 or -1 for high and low expression mRNAs, respectively
PN	SequenceSuprempositionName	Functional site (numbered 1)
SC	SiteVariant	Identifier of corresponding EMBL entry
SQ	Sequence	5'UTR nucleotide sequence
SA	ActivityValue	sequence translational activity (1 or -1, see AU)
PA	SequenceSuprempositionPointer	Sequence length up to functional site (PA)

**Leader WHY**

Line code	Field name	Field description
MI	EntityID	mRNA 5'UTR feature name
MN	EntityName	mRNA 5'UTR feature type
MD	EntityDependence	Gene region
ML	Step	Sequence feature
RN	CrossRef	Reference on Leader_REF database
PN	PropertyName	Type of analysis (models and hypothesis)
PM	Method for property definition	Type of analysis (statistical, modelling, etc)
PV	Property ID	Details of analysis scheme (training test, GA, modelling, etc)
PU	Unit for Property Expression	significance criteria
REASONS	Description	Basis of this feature usage as a discriminative characteristic in the prediction process

**Leader REF**

Line code	Field name	Field description
RN	Identifier	Reference identifier
RA	Authors	Reference authors
RT	Title	Reference heading
RJ	Journal	Reference journal
RV	Vol	Reference journal volume
RP	PP	Reference journal pages
RY	Year	Reference year
RR	Abstract	Reference abstract

**Leader\_KN**

Line code	Field name	Field description
CF	KnowledgeName	Type of mRNA feature
CT	KnowledgeNature	influence on mRNA translational activity
DP	FeatureName	Feature name
PV	SequenceFeature	Types of a sequence
AB	Values	Mean values of the features for high and low expression mRNAs
LC	LinearCorrelation	Significance of statistical difference
ST	ContrHigh	Results of prediction of a high expression control set: <mean values> <standard errors> <the percentages of incorrect prediction>
NN	ContrLow	Results of prediction of a low expression control set: <mean values> <standard errors> <the percentages of incorrect prediction>
C-CODE	CODE	C-code of the prediction program

**Comments and questions are welcome to Alex Kochetov (ak@bionet.nsc.ru)**

## 2. Software.

### 2.1. Programs for mRNA Translatability Prediction

Release 2003

#### Program description:

LEADER\_RNA contains distinct programs for prediction of mRNA translation efficiency in organisms of each of these taxa: dicot plants, monocot plants and mammals.

#### List of biological tasks that could be solved by using Prediction programs:

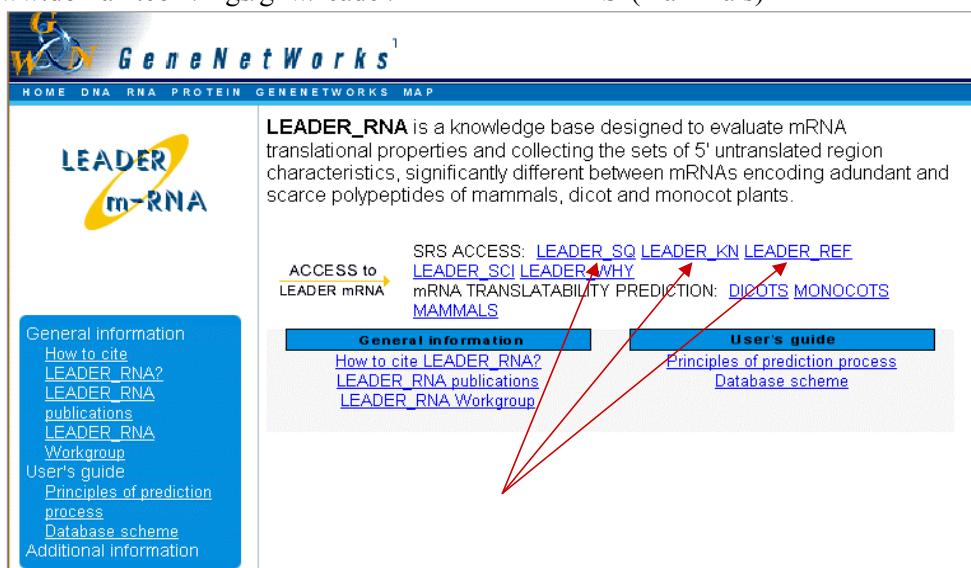
- evaluation of mRNA translational efficiency in higher plants or mammals. The program estimates to what extent the mRNA defined by a user is similar to mRNAs of high or low expression genes? Consequently, similarity of a particular mRNA 5'UTR to that of high expression genes possibly may indicate the fact that this gene is also expressed at high expression rate.
- gene engineering experiment design (e.g., for prediction of translation properties of modified mRNAs in model experiments, as well as for evaluation of translation efficiency of transgene mRNA in a novel host organism).

#### Access to Prediction programs:

<http://www.domain.com/mgs/gnw/leader/> link 'DICOTS' (dicot plants)

<http://www.domain.com/mgs/gnw/leader/> link 'MONOCOTS' (monocot plants)

<http://www.domain.com/mgs/gnw/leader/> link 'MAMMALS' (mammals)



#### Data input

Input window of LEADER\_RNA prediction program is presented below in the Figure. A user may select one of three available taxa (dicot plants, monocot plants, and mammals).

Next window contains input text-box for 5'UTR sequence and the list of discriminative features (determined by choosing of appropriate taxon at the previous page). It is possible to input 5'UTR sequence by several ways:

- from screen (by typing the sequence from the keyboard) (1).

- from file at your PC (by browsing it from the user's PC hard disk). For this operation, you should use the file with a single 5'UTR sequence in a FASTA format. Click the button 'BROWSE' (2) and choose the source file.
- from databases (currently, only EMBL and GeneBank are included in the database list; in these cases the entries should contain 5'UTR sequence only; in future, we plan to add the links to specialised mRNA 5'UTR databases) (3).

**Predicting High/Low mRNA expression of a monocot plant gene**

Input DNA Sequence :

from Screen:  
 from DB:  
 from File:

Bases Available:  
 SRS5 from Heidelberg (EMBL) by ID

Execute    Reset form    [Example](#) [Related Paper](#)

File formats here.

Expert Weights (0-10 are valid)

1.  5 Translation INCREASES with DECREASING the Leader length  
 2.  5 Translation INCREASES with DECREASING [T] content  
 3.  5 Translation INCREASES with DECREASING [AUG]-[AUG] disbalance  
 4.  5 Translation INCREASES with INCREASING [A]-[T] ratio

#### Essential notes:

- 5'UTR sequences should be entered in upper- or lower-case letters; both T and U letters are allowed (anyway, T is considered as U); line feeds and blanks are ignored;
- DO NOT include start codon of the ORF in 5'UTR sequence. INCLUDE 5'UTR up to -1 nucleotide (located just upstream AUG codon of the main ORF);
- Try to use COMPLETE 5'UTR sequence (if possible) to increase prediction accuracy.
- You may choose weights (see 5 in the Figure) of discriminative features (ranging from 0 to 10; 5 by defaults). This feature is likely to be helpful for users experienced in translation machinery organization in eukaryotes.

#### Program options

The list of discriminative features in the figure demonstrating the input window of the program is marked by number 5. It contains various contextual characteristics used in prediction process. By default, all these features are used with the same expert weights. Prediction program is based on the mean recognition approach, so the contribution of a single feature into the final evaluation is likely not so critical. The simplest way is to use the default options (each parameter equals to 5).

However, this program provides a user with opportunity to change the expert weights of discriminative features in order to optimise prediction process. Each 5'UTR feature listed is supplied with a brief description allowing to evaluate its significance and to assign the weight. In some cases, it may be useful to apply features, which significance was supported experimentally (e.g., weight of start codon context, the presence of leader AUG triplets and weights of their contexts).

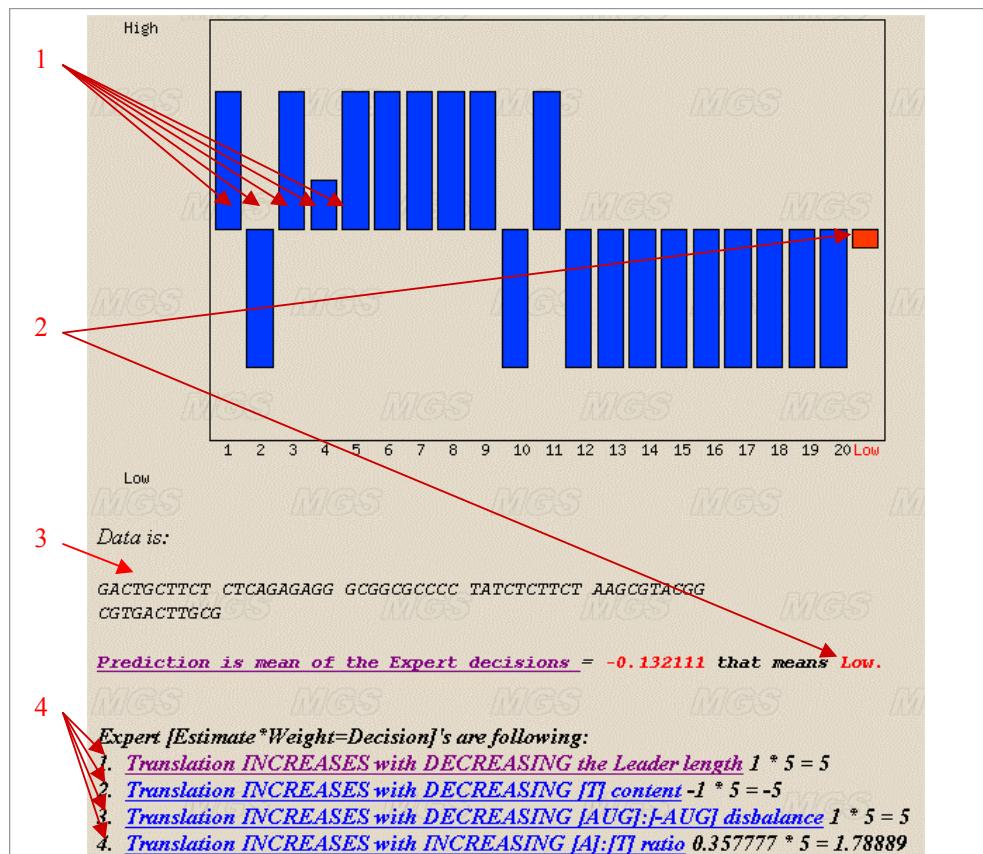
### Program execution

Click the button ‘Execute’ (4).

### Data output

This window presents the results of prediction (demonstrated both in the figure and table). Values of each of discriminative feature are shown (4). It is possible to evaluate similarity of 5'UTR analysed with those of high and low expression mRNAs of corresponding taxon (1).

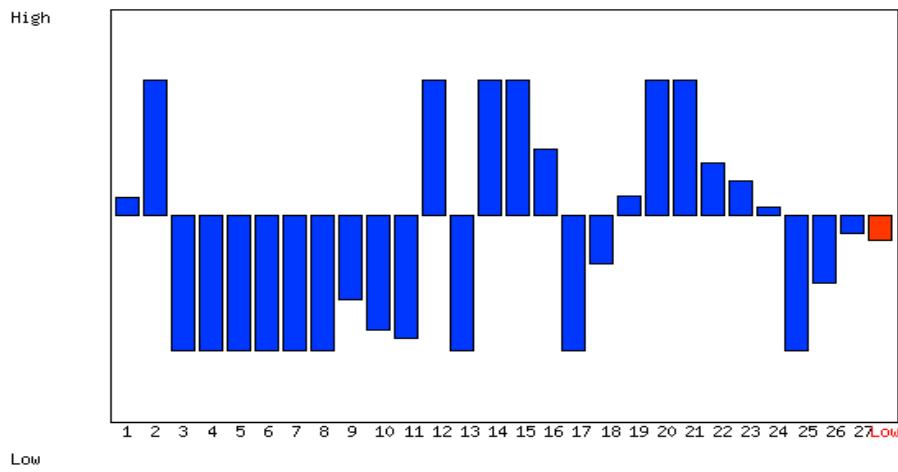
Each field in the table of discriminative features (4) contains links to LEADER\_RNA knowledge base allowing extract additional information. The input sequence is also shown (3). Integral prediction result is marked by 2.



### Examples

- (1) Random sequence with equal shares of nucleotides (default parameters were used for prediction).

## II. RNA. Chapter 1. Leader RNA.



Data is:

TGGCCGGTGG CTTATGGACT GTCACAGCTC GTGTCAGTGA TGAAGATAAG  
GGTATCCTCC CGCTGTGGGA CCATAAAGAA CCACACGATC ATTTTAATCG

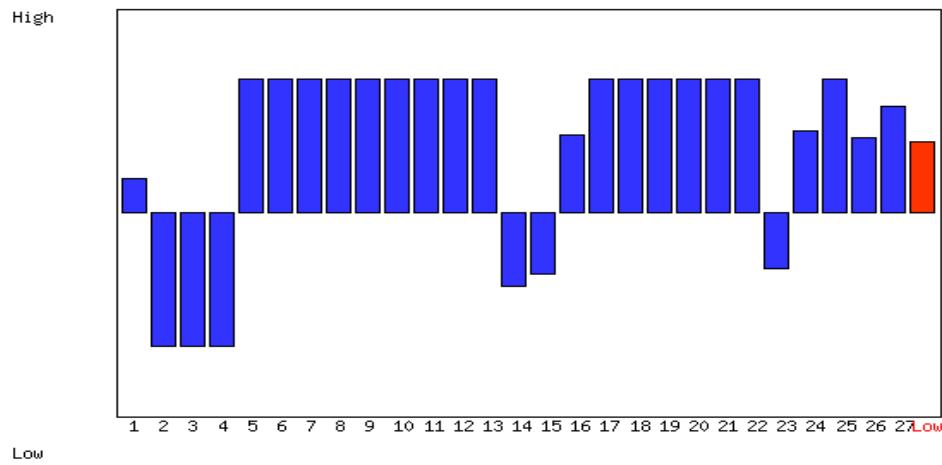
Prediction is mean of the Expert decisions = -0.181138 that means Low.

Expert [Estimate\*Weight=Decision]'s are following:

1. Translation INCREASES with DECREASING the Leader length  $0.137257 * 5 = 0.686285$
2. Translation INCREASES with DECREASING IT content  $1 * 5 = 5$
3. Translation INCREASES with DECREASING [AUG];[AUG] disbalance  $-1 * 5 = -5$
4. Translation INCREASES with INCREASING [A];[T] ratio  $-1 * 5 = -5$
5. Translation INCREASES with INCREASING [AUG];[AUG] ratio  $-1 * 5 = -5$
6. Translation INCREASES with DECREASING [A];[T] disbalance  $-1 * 5 = -5$
7. Translation INCREASES with DECREASING [AUG] content  $-1 * 5 = -5$
8. Translation INCREASES with DECREASING [AUG] framed  $-1 * 5 = -5$
9. Translation INCREASES with DECREASING [AUG] optimized  $-0.623815 * 5 = -3.11908$
10. Translation INCREASES depends on the "3 position" rule  $-0.840345 * 5 = -4.20173$
11. Translation INCREASES with DECREASING [AUG] "3"-ruled  $-0.903936 * 5 = -4.51968$
12. Translation INCREASES with DECREASING [K] content of f-17:-1]  $1 * 5 = 5$
13. Translation INCREASES with DECREASING [K] content of f-17:-1]  $-1 * 5 = -5$
14. Translation INCREASES with INCREASING High-consensus matches  $1 * 5 = 5$
15. Translation INCREASES with DECREASING Low-consensus matches  $1 * 5 = 5$
16. Translation INCREASES with INCREASING High/ShortFreqMatr  $0.483426 * 5 = 2.41713$
17. Translation INCREASES with INCREASING High/Low 1bp-FreqRatio  $-1 * 5 = -5$
18. Translation INCREASES with INCREASING High/Low 1bp(KM)-FreqRatio  $-0.360723 * 5 = -1.80362$
19. Translation INCREASES with INCREASING High/Low 2bp(KM)-FreqRatio  $0.146544 * 5 = 0.73272$
20. Translation INCREASES with INCREASING High/Low 3bp(KxM)-FreqRatio  $1 * 5 = 5$
21. Translation INCREASES with INCREASING High/Low 5bp(KM)-FreqRatio  $1 * 5 = 5$
22. Translation INCREASES with INCREASING High/Low 6bp(KM)-FreqRatio  $0.3873 * 5 = 1.9365$
23. Translation INCREASES with INCREASING High/Low 3bp(ATGCx)-FreqRatio  $0.250668 * 5 = 1.25334$
24. Translation INCREASES with INCREASING High/Low 5bp(ATGCx)-FreqRatio  $0.057873 * 5 = 0.289365$
25. Translation INCREASES with INCREASING High/Low 3bp(KxM)-FreqRatio  $-1 * 5 = -5$
26. Translation INCREASES with INCREASING High/Low 5bp(KxM)-FreqRatio  $-0.497674 * 5 = -2.48837$
27. Translation INCREASES with INCREASING High/Low 7bp(KxM)-FreqRatio  $-0.127296 * 5 = -0.63648$

(2) 5'UTR sequence of dicot plant gene with the high level of expression (EMBL accession number: AT30SRS13) (default parameters were used for this prediction).

## II. RNA. Chapter 1. Leader RNA.



Data is:

CGTTTGCCTT ATCCGTTTCAG CTCATCTTCT TCTTCTTCCT CGTCACCTCT  
GAATTAGTTT CCCAGAAATCC GAAATTCCCTA GGAAGAGAAC ATAACA

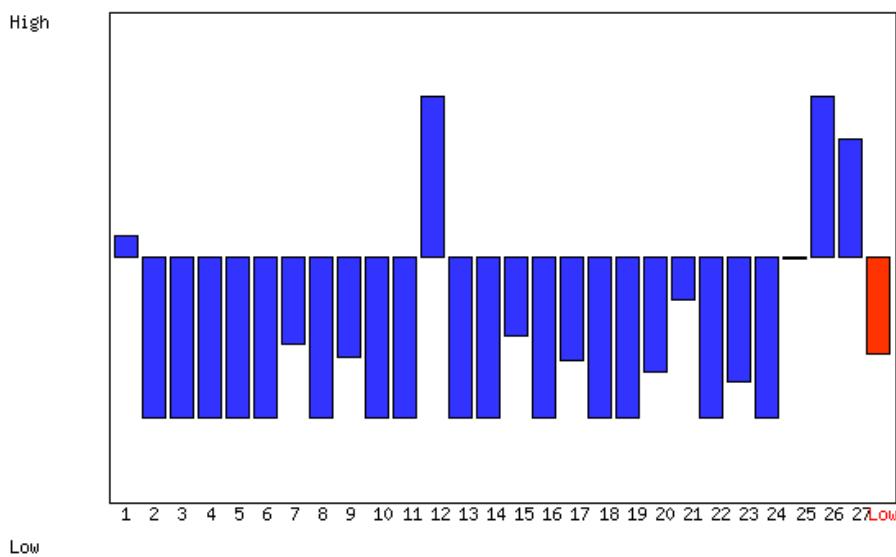
**Prediction is mean of the Expert decisions = 0.532137 that means High.**

Expert [Estimate\*Weight=Decision]'s are following:

1. Translation INCREASES with DECREASING the Leader length  $0.253257 * 5 = 1.26629$
2. Translation INCREASES with DECREASING fT content  $-1 * 5 = -5$
3. Translation INCREASES with DECREASING f[AUG];f[AUG] disbalance  $-1 * 5 = -5$
4. Translation INCREASES with INCREASING fA:fT ratio  $-1 * 5 = -5$
5. Translation INCREASES with INCREASING f[AUG];f[AUG] ratio  $1 * 5 = 5$
6. Translation INCREASES with DECREASING fA:fT disbalance  $1 * 5 = 5$
7. Translation INCREASES with DECREASING f[AUG] content  $1 * 5 = 5$
8. Translation INCREASES with DECREASING f[AUG] framed  $1 * 5 = 5$
9. Translation INCREASES with DECREASING f[AUG] optimized  $1 * 5 = 5$
10. Translation INCREASES depends on the "-3 position" rule  $1 * 5 = 5$
11. Translation INCREASES with DECREASING f[AUG] "-3"-ruled  $1 * 5 = 5$
12. Translation INCREASES with DECREASING fK content of f-17:-11  $1 * 5 = 5$
13. Translation INCREASES with DECREASING fK content of f-17:-11  $1 * 5 = 5$
14. Translation INCREASES with INCREASING High-consensus matches  $-0.54879 * 5 = -2.74395$
15. Translation INCREASES with DECREASING Low-consensus matches  $-0.455469 * 5 = -2.27735$
16. Translation INCREASES with INCREASING High-Short freqMatr  $0.583073 * 5 = 2.91537$
17. Translation INCREASES with INCREASING High/Low 1bp-FreqRatio  $1 * 5 = 5$
18. Translation INCREASES with INCREASING High/Low 1bp(KM)-FreqRatio  $1 * 5 = 5$
19. Translation INCREASES with INCREASING High/Low 2bp(KM)-FreqRatio  $1 * 5 = 5$
20. Translation INCREASES with INCREASING High/Low 3bp(KxM)-FreqRatio  $1 * 5 = 5$
21. Translation INCREASES with INCREASING High/Low 5bp(KM)-FreqRatio  $1 * 5 = 5$
22. Translation INCREASES with INCREASING High/Low 6bp(KM)-FreqRatio  $1 * 5 = 5$
23. Translation INCREASES with INCREASING High/Low 3bp(ATGCx)-FreqRatio  $-0.421578 * 5 = -2.10789$
24. Translation INCREASES with INCREASING High/Low 5bp(ATGCx)-FreqRatio  $0.608727 * 5 = 3.04363$
25. Translation INCREASES with INCREASING High/Low 3bp(KxM)-FreqRatio  $1 * 5 = 5$
26. Translation INCREASES with INCREASING High/Low 5bp(KxM)-FreqRatio  $0.557566 * 5 = 2.78783$
27. Translation INCREASES with INCREASING High/Low 7bp(KxM)-FreqRatio  $0.79092 * 5 = 3.9546$

- (3) 5'UTR sequence of dicot plant gene with low level of expression (EMBL accession number: AT31834) (default parameters were used for the prediction).

## II. RNA. Chapter 1. Leader RNA.



Data is:

TTCAACAACT GTGTTATCTC TATGAACAGT TCTTTGATG AAGAACACAA  
AGTGAAAGTT GCTGTCTTA TAACCAGGAT TTGGTAATTG CCATTGTTTC

**Prediction is mean of the Expert decisions = -0.598938 that means Low.**

Expert [Estimate\*Weight=Decision]'s are following:

1. Translation INCREASES with DECREASING the Leader length  $0.137257 * 5 = 0.686285$
2. Translation INCREASES with DECREASING [T] content  $-1 * 5 = -5$
3. Translation INCREASES with DECREASING [AUG];[f-AUG] disbalance  $-1 * 5 = -5$
4. Translation INCREASES with INCREASING [A];[T] ratio  $-1 * 5 = -5$
5. Translation INCREASES with INCREASING [AUG];[f-AUG] ratio  $-1 * 5 = -5$
6. Translation INCREASES with DECREASING [A];[T] disbalance  $-1 * 5 = -5$
7. Translation INCREASES with DECREASING [AUG] content  $-0.537042 * 5 = -2.68521$
8. Translation INCREASES with DECREASING [AUG] framed  $-1 * 5 = -5$
9. Translation INCREASES with DECREASING [AUG] optimized  $-0.623815 * 5 = -3.11908$
10. Translation INCREASES depends on the "-3 position" rule  $-1 * 5 = -5$
11. Translation INCREASES with DECREASING [AUG] "-3"-ruled  $-1 * 5 = -5$
12. Translation INCREASES with DECREASING [K] content of I-17:-11  $1 * 5 = 5$
13. Translation INCREASES with DECREASING [K] content of I-17:-11  $-1 * 5 = -5$
14. Translation INCREASES with INCREASING High-consensus matches  $-1 * 5 = -5$
15. Translation INCREASES with DECREASING Low-consensus matches  $-0.487738 * 5 = -2.43869$
16. Translation INCREASES with INCREASING High-ShortFreqMatr  $-1 * 5 = -5$
17. Translation INCREASES with INCREASING High/Low 1bp-FreqRatio  $-0.639225 * 5 = -3.19612$
18. Translation INCREASES with INCREASING High/Low 1bp(KM)-FreqRatio  $-1 * 5 = -5$
19. Translation INCREASES with INCREASING High/Low 2bp(KM)-FreqRatio  $-1 * 5 = -5$
20. Translation INCREASES with INCREASING High/Low 3bp(KxM)-FreqRatio  $-0.708711 * 5 = -3.54355$
21. Translation INCREASES with INCREASING High/Low 5bp(KM)-FreqRatio  $-0.264822 * 5 = -1.32411$
22. Translation INCREASES with INCREASING High/Low 6bp(KM)-FreqRatio  $-1 * 5 = -5$
23. Translation INCREASES with INCREASING High/Low 3bp(ATGCx)FreqRatio  $-0.776691 * 5 = -3.88346$
24. Translation INCREASES with INCREASING High/Low 5bp(ATGCx)FreqRatio  $-1 * 5 = -5$
25. Translation INCREASES with INCREASING High/Low 3bp(KxM)-FreqRatio  $-0.00717 * 5 = -0.03585$
26. Translation INCREASES with INCREASING High/Low 5bp(KxM)-FreqRatio  $1 * 5 = 5$
27. Translation INCREASES with INCREASING High/Low 7bp(KxM)-FreqRatio  $0.736632 * 5 = 3.68316$

Comments and questions are welcome to Alex Kochetov ([ak@bionet.nsc.ru](mailto:ak@bionet.nsc.ru)).

## 2.2. MatrixSS: Building of E-score plot for RNA sequence

Release 2003

**Program description:** 'MatrixSS' enables searching for the regions with the stable secondary RNA structure in extended (from 250 to 100000 nt) genome sequences.

**Access to 'MatrixSS':**

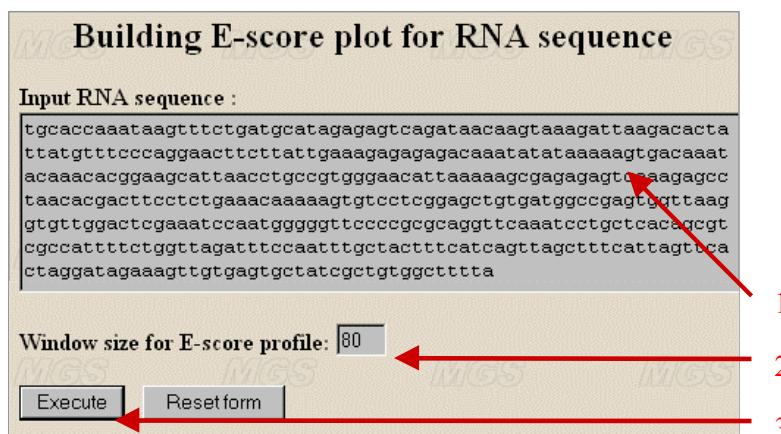
<http://www.domain.com/mgs/gnw/garna/> link 'Primary structure analysis (MatrixSS program)'

**Biological tasks that could be solved by using the 'MatrixSS':**

- ◆ The program MatrixSS is designed for searching in genome sequences for the genes encoding RNA with a potentially significant secondary structure (SS).

### Data input

Into the text-box (1), enter or insert from the clipboard the RNA sequence to be analysed. The sequence should be within the range of 250 – 100000 nt in length. Use the alphabets ATGC(atgc) or AUGC(augc) (any other symbol will give an error). Blanks or line folding are ignored.



### Program options

Set the window length for calculation of the E-score value (2). The window length should correspond to the length of the RNA molecule searched for (by default, the window length equals to 80 nt, which is the length of the tRNA molecule).

### Program execution

Click the button 'Execute' (3).

### Data output

As the data output, two graphical plots are presented. The first one shows the distribution of the E-score value along the sequence analysed. The dotted line marks the mean E-score value typical for the structural RNAs (tRNA, rRNA, etc.). The solid lines to the up and down of the dotted line mark the limits of E-score variation in the class of structural RNA. The regions of the resulting line, which lie between these limits, are marked by red. These are the regions that potentially encode structural RNA genes.

The second graphical plot is a coloured E-score matrix. This matrix describes cross-complementarity of the potential or real RNA regions to each other. The colour of the matrix

## II. RNA. Chapter 1. Leader RNA.

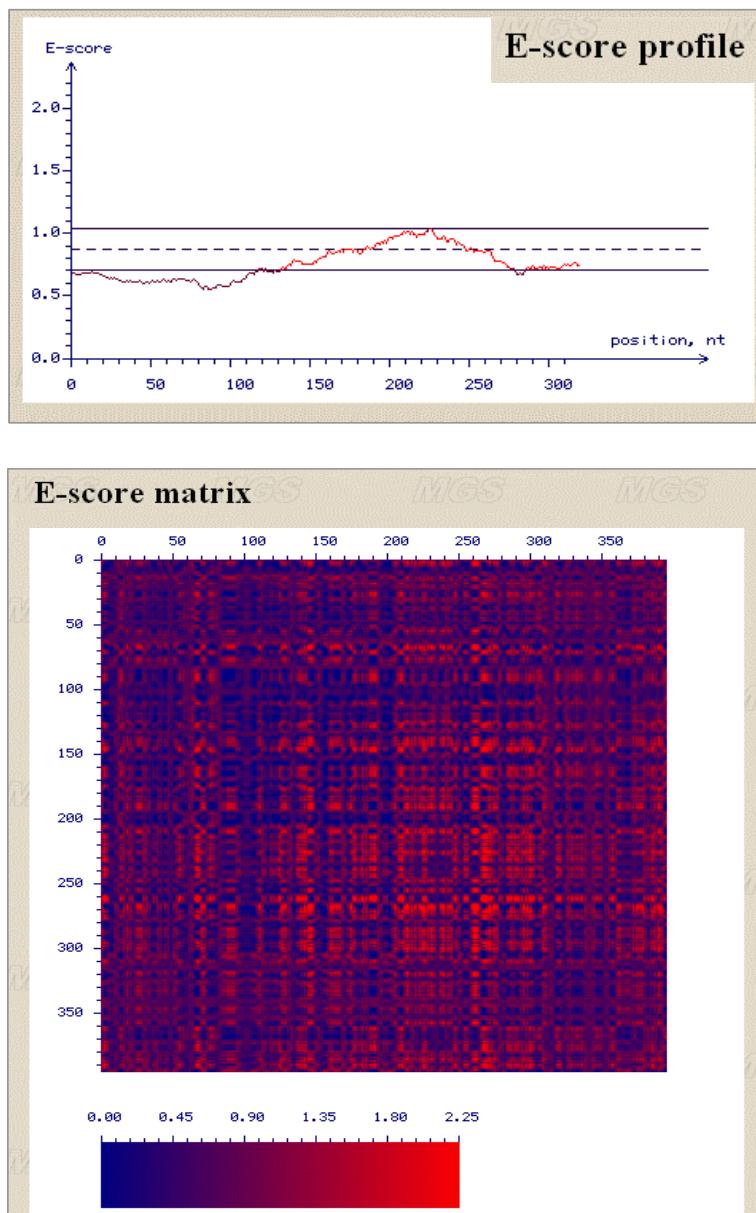
elements varies with respect to the expected stability of RNA SS, which is formed by two these RNA regions. Red colour corresponds to the higher values of expected stability.

### Example

As an example, we search for the tRNA Ser sequence (of 82 nt in length), which is present in the sequence extracted from the EMBL databank (ID AF184043). The total length of this sequence is 400 nt. The tRNA Ser gene is located in-between positions 218 - 299 nt.

In this search we use 'window size' equaling to 80.

The output is represented by the following plots:



As seen in the Figure, by the example of tRNA Ser, the region of predicted potential SS contains the real tRNA gene, that is, the program made correct prediction.

**Comments and questions are welcome to Igor Titov ([titov@bionet.nsc.ru](mailto:titov@bionet.nsc.ru))**

## 2.3. GArna Program

Release 2003

### 1. Program description:

GArna is a program for secondary structure prediction in short (up to 250 nt) RNA sequences and its visualization.

### 2. Access to 'GArna':

<http://www.domain.com/mgs/gnw/garna/> than follow link ' Secondary structure prediction (GArna program)'

### 3. Biological tasks that could be solved by using the GArna program:

- to calculate and visualise optimal secondary structure in short (at most 250 nt) RNA sequences
- to account available experimental information about single strand RNA regions
- to evaluate stability of the secondary RNA structure in comparison to the random sequences of the same length and nucleotide content (by calculating Z-score value of the secondary structure energy)
- to model rearrangements of the secondary RNA structure under the action of the anti-sense oligonucleotide.

### 4. Data input

Into the text-box, enter or insert from the clipboard the RNA sequence to be analysed. The sequence should be within the range of 5 – 250 nt in length. Use the alphabets ATGC or AUGC (any other symbol will give an error). Blanks or line feeds are ignored.

### 5. Program options

The program GArna has the following parameters:

- Stop criterion D (option 2 on the figure). This value varies within the range  $0 < D < 1$ . The less is D, the more quick is calculation and less the prediction accuracy (default value giving the most accurate prediction is 0.9). The low stop criterion values (about 0.5) may be applied for searching for alternative structures and/or kinetic intermediates for particular RNA.
- Minimal helix length (option 3 in the figure). This parameter determines the least stem length necessary to form SS. Default value is 2. Minimal value equals to 2, maximal value equals to 6. The more is 'Minimal helix length' value, the more quick is calculation and less the accuracy of the secondary structure prediction.
- Selection temperature (option 4 on the figure). It is recommended to use default value 4 kcal/mol. Any other value gives less precise prediction.
- Randomization parameter (r.p.; option 5 in the figure). Any integer value is allowed. If the particular r.p. is ordered (under constancy of the rest parameters), the program always outputs one and the same result. By varying r.p., different results could be obtained. This parameter may be used, for example, to determine the robustness of prediction.
- Positions covered by oligonucleotide (from X to Y) (option 6 in the figure). This parameter indicates the range of nucleotide positions, which are prohibited to form the duplexes with the other nucleotides in the RNA molecule studied. This parameter is used (a) for calculation if the anti-sense oligonucleotide is present, (b) for calculations using *a priori* information about single-stranded regions of SS. X and Y values could vary within the range from 1 to sequence\_length, X is less or equals to Y. By default, the values X=0, Y=0 are set. These values mean that all nucleotides in the sequence analysed may form duplexes (since numeration in the input sequence begins with 1). For calculations with non-zero X and Y values, z-score value is not specified.

**GArna: Predicting 2D structure of RNA by genetic algorithm**

If your sequence exceeds 250 nt please use [MatrixSS](#) program

Input RNA sequence :

```
GTAAATGTTAGCTTATAATAAGCAAAAGCACTGAAAATGCTTAGATGGATTCAAAAATCCCATAAACAA
```

1

Stop calculation when population degeneracy D exceeds  (0 < D < 1) 2

Minimal helix length:  nucleotides 3

Selection temperature:  Kcal/mol 4

Randomization parameter:  5

Positions covered by oligonucleotide, from:  to:  6

7 → Execute Reset form

## 6. Program execution

Click the button ‘Execute’ (7).

## 7. Data output

Program execution will bring up the resulting window with the data output. They include the following:

- the energy of the secondary structure calculated;
- deviation of stability of the secondary structure calculated from expected value for the random nucleotide sequences of the same length and content (z-score);
- graphical representation of the secondary structure calculated;
- representation of the calculated secondary structure in the format CT, which is convenient to use while exploiting the other programs (RNAstructure, GCG, etc.) for analysis of secondary RNA structure. [CT File Format description. A CT (Connectivity Table) file contains secondary structure information for a sequence. It may contain multiple structures for a single sequence. When entering a structure to calculate the free energy, the following format shown below should be used. The number at the beginning of the first line is the number of bases in the sequence. Next is the title of the structure. Each of the following lines provides information about a base in the sequence, each base being described in its own line. First is the base number, n, followed by the base denotation (A,C,G,U, or T). The third column is n-1, whereas the fourth column is n+1. The fifth column is the number of the base to which n is paired. The lack of pairing is indicated by 0 (zero). The last column is called the natural numbering.]

### Example

The calculation of rat Phe tRNA secondary structure (the entry DF5280 in the database accumulating tRNA sequences; <http://www.uni-bayreuth.de/departments/biochemie/sprinzl/trna/>). All options for this example were set to default values.

## II. RNA. Chapter 1. Leader RNA.

Your sequence has the optimal structure with  
Energy = -8.4 kcal/mol

with a z-score of -2.94, which means INCREASED stability compared to random sequences of the same length and nucleotide composition (expected Energy = -3.1 kcal/mol, standard deviation = 1.8 kcal/mol)

Secondary structure in MFOLD .ct file format:

67 ENERGY = -8.4 GA-generated
1 G 0 2 66 1
2 U 1 3 65 2
3 U 2 4 64 3
4 A 3 5 0 4
5 A 4 6 62 5

Move    Move    Zoom    Reset

Comments and questions are welcome to Igor Titov ([titov@bionet.nsc.ru](mailto:titov@bionet.nsc.ru))

# **PART III. PROTEIN INTEGRATION LEVEL.**

## **CHAPTER 1. DATABASES.**

### **1. EnPDB**

**Release 2003**

#### **1. Database description:**

ENPDB is search system for query implementation to the PDB (Protein Data Bank).

#### **2. Access to ENPDB database:**

<http://www.domain.com/mgs/gnw/enpdb/>

#### **3. Database content**

SRS table	Description	Number of entries
ENPDB	Search system for query implementation to the PDB	16777

#### **4. List of biological tasks that could be solved by using the ENPDB database**

- To provide extended querying of the PDB database accumulating spatial protein structures and their functional characteristics.

#### **5. SRS table format**

##### **ENPDB**

Line code	Field name	Field description
ID	ID	PDB ID code. This identifier is unique for PDB
HEADER	Header	PDB classification of the entry.
DATE	Date	Deposition date is the date when the coordinates were received by PDB. The date field contains mainly the date when the entry was created, always stored in the index as an eight-digit number of the format "yyyymmdd" (y = year, m = month, and d = day), e.g., "19940117". It is also possible to type the date to be searched in a different, more intuitive, format: "dd-mmm-yy" or "dd-mmm-yyyy", e.g., "1-jan-97" or "01-jan-1997".
TITLE	Title	This field contains the title for experiment or analysis described in the entry. The field content corresponds to that in PDB.

**III. Protein Integration Level.** Chapter 1. Databases.

Line code	Field name	Field description
COMPOUND	Compound	It describes the macromolecules contained in the entry. Each macromolecule of the entry is defined with a set of token: value pairs, and is referred to as a component of the COMPOUND field. The field content corresponds to that in PDB. For each macromolecular component, the molecule name, synonyms, number assigned by the Enzyme Commission (EC), and other relevant details are specified.
COMPOUND	Molecule	This specialised field is not present in the entry as a separate line. It contains names of macromolecules from the COMPND of PDB and is designed to search for entries by the names of macromolecules.
COMPOUND	Synonym	This specialised field is not present in the entry as a separate line. It contains synonymous names of macromolecules from the COMPND of PDB and is designed to search for entries by the synonymous names of macromolecules.
COMPOUND	EC	It contains the Enzyme Commission number associated with the molecule. If there is more than one EC number, they are presented as a comma-separated list.
COMPOUND	BioUnit	If a MOLECULE functions as a part of a larger biological unit, the entire functional unit may be described. The field content is selected from the COMPND of PDB using the token BIOLOGICAL UNIT.
SOURCE	Gene	It identifies the gene through the gene names taken from the SOURCE field of PDB.
MOL_SOURCE	MolSource	It specifies biological and/or chemical sources of all the biological molecules in the entry. There are three values: BIOLOGICAL, SYNTHETIC, and MIXED. SYNTHETIC means that all the molecules with the entry were chemically synthesised; BIOLOGICAL, all the molecules were not synthesized; and MIXED, some molecules were synthesized, whereas the rest are natural.
SOURCE	Source	This field specifies the biological source of each biological molecule in the entry. Sources are described by both the common and scientific names, e.g., genus and species. Strain and/or cell line for immortalized cells are given when they help to uniquely identify the biological entity studied. Note that the content of this field is not a replica of the SOURCE of PDB. The original PDB file is divided into two parts: all the information concerning the biological source is retained in this field, while all the data related to synthesis is comprised in the new field SYNTHESIS. We believe that this division allows user to specify the region to be searched for more precisely.
SOURCE	Synthesis	This field specifies the data on expression systems, e.g. strain, variant, cell line, etc. The content originates from the SOURCE of PDB. See also the description of Field Source.
KEYWORD	Keyword	It contains keywords describing the macromolecule. The content corresponds to that of KEYWDS of PDB.
TECHNIQUE	Technique	It identifies the experimental technique used. This may refer to the type of radiation and sample, or include the spectroscopic or simulation technique. The content originates from the EXPDTA of PDB.
AUTHOR	Author	It indicates the names of the experts responsible for the contents of the entry and corresponds to the AUTHOR field in PDB.
JRNL	Jrnl	It indicates the reference to original publication that describes the experiment and defines the coordinate set. Its content originates from the JRNL field of PDB.

**III. Protein Integration Level.** Chapter 1. Databases.

Line code	Field name	Field description
JRNL AUTH	JrnAuthor	It contains the list of authors of the paper cited or contribution to a larger work. Its content originates from the JRNL field of PDB.
JRNL TITL	JrnTitle	It specifies the title of the reference and is used for the title of a journal article, chapter, or part of a book. Its content originates from the JRNL field of PDB.
JRNL REF	JrnRef	It contains name of the publication. Its content originates from the JRNL field of PDB.
JRNL VOLUME	JrnVolume	It contains the volume of the publication. Its content originates from the JRNL field of PDB.
JRNL YEAR	JrnYear	It indicates the year of the publication. Its content originates from the JRNL field of PDB.
REMARK_1	Remark_1	It lists important publications related to the structure described in the entry. These citations are chosen by the depositor. The content originates from the REMARK 1 of PDB.
RESOLUTION	Resolution	It is derived from REMARK 2 in the PDB file. No resolution is given for NMR structures and models. The field indicates the highest resolution in Angstroms used in building the model.
CHAIN_AMOUNT	ChainAmount	It indicates the number of chains in the entry, calculated from the SEQRES field of PDB.
CHAIN_SIZES	ChainSizes	It specifies the lengths of the chains in the entry, calculated from the data contained in the SEQRES of PDB.
HELIX_AMOUNT	HelixAmount	It indicates the number of helices in the entry and is derived from the MASTER field of PDB.
SHEET_AMOUNT	SheetAmount	It indicates the number of beta-sheet structures in the entry and is derived from the MASTER field of PDB.
DNA_RNA_AMOUNT	DnaRnaAmount	It specifies the number of DNA/RNA strands in the entry, calculated from the SEQRES field of PDB.
PROTEIN_AMOUNT	ProteinAmount	It indicates the number of protein chains, calculated from the SEQRES field of PDB.
HET_AMOUNT	HetAmount	It indicates the number of unusual residues, such as prosthetic groups, inhibitors, solvent molecules, and ions, supplemented with their coordinates. The data are calculated from the HET field of PDB.
HETEROGEN	Heterogen	It gives the chemical name and the synonyms of unusual residues, such as prosthetic groups, inhibitors, solvent molecules, and ions, supplemented with their coordinates. The data are calculated from the HETNAM and HETSYN fields of PDB.
LINK_EMBL	LinkEmbl	It links to EMBL Data Bank through SWISS-PROT. For example, we find a SWISS-PROT entry with references to both PDB and EMBL entries. In this case, we consider that the PDB and EMBL entries are linked.
LINK_PIR	LinkPir	It links to PIR Data Bank through SWISS-PROT.
LINK_SWISS-PROT	LinkSwissProt	It links to SWISS-PROT Data Bank as its entries contain references to PDB.
LINK_TRANSFAC	LinkTransfac	It links to Transfac Data Bank through SWISS-PROT. For example, we find a SWISS-PROT entry with references to both PDB and TRANSFAC entries. In this case, we consider that the PDB and TRANSFAC entries are linked.
LINK_TRRD4	LinkTrrd4	It links to TRRD Data Bank.

### Examples of SRS queries to the ENPDB database

#### Example 1

The query is 'To search for receptors of TNF, with the known tertiary structures represented in PDB'.

To make such a query, you should perform the following:

1. Choose the option SRS ACCESS: ENPDB as shown in the Figure by red arrow.

**Gene Networks**<sup>1</sup>

HOME DNA RNA PROTEIN GENENETWORKS MAP

**EnPDB**

General information  
How to cite EnPDB  
EnPDB publications  
Contact us  
User's guide  
Brief manual on the database  
ENPDB  
SRS queries (examples)  
Current release  
Additional information  
Blast search ENPDB (PDB) database  
Links to other databases and programs

EnPDB database is made by reformatting PDB in a way allowing for extended search possibilities on PDB. To broaden the possibilities of search in the database additional fields are added that are absent in PDB. Coordinates of atoms are not included in the EnPDB. Database access is realised by means of SRS system (Sequence Retrieval System). EnPDB format implies full indexing of all the fields. EnPDB is integrated by means of hyperlinks with different databases (SWISS-PROT, PIR, TRRD, etc).

**ACCESS to EnPDB** → SRS ACCESS: EnPDB  
[Blast search ENPDB \(PDB\) database](#)

**General information**  
How to cite EnPDB?  
EnPDB publications  
Contact us

**User's guide**  
Brief manual on the database ENPDB  
SRS queries (examples)

**Current release**  
EnPDB is updated after each new PDB release.  
The current release has 16777 entries and was indexed 21-Mar-2003.

**Additional information**  
[Blast search ENPDB \(PDB\) database](#)  
[Links to other databases and programs](#)

**Gene Networks**<sup>1</sup>

HOME DNA RNA PROTEIN GENENETWORKS MAP

**EnPDB**

General information  
How to cite EnPDB  
EnPDB publications  
Contact us  
User's guide  
Brief manual on the database  
ENPDB  
SRS queries (examples)  
Current release  
Additional information  
Blast search ENPDB (PDB) database  
Links to other databases and programs

EnPDB database is made by reformatting PDB in a way allowing for extended search possibilities on PDB. To broaden the possibilities of search in the database additional fields are added that are absent in PDB. Coordinates of atoms are not included in the EnPDB. Database access is realised by means of SRS system (Sequence Retrieval System). EnPDB format implies full indexing of all the fields. EnPDB is integrated by means of hyperlinks with different databases (SWISS-PROT, PIR, TRRD, etc).

**ACCESS to EnPDB** → SRS ACCESS: EnPDB  
[Blast search ENPDB \(PDB\) database](#)

**General information**  
How to cite EnPDB?  
EnPDB publications  
Contact us

**User's guide**  
Brief manual on the database ENPDB  
SRS queries (examples)

**Current release**  
EnPDB is updated after each new PDB release.  
The current release has 16777 entries and was indexed 21-Mar-2003.

**Additional information**  
[Blast search ENPDB \(PDB\) database](#)  
[Links to other databases and programs](#)

### III. Protein Integration Level. Chapter 1. Databases.

2. To load the extended search form, click the button 'Extended'.

The screenshot shows the ENPDB query interface with the following details:

- Top Navigation:** TOP PAGE, QUERY (highlighted), RESULTS, SESSIONS, VIEWS, DATABASES, HELP.
- Search Bar:** search ENPDB, Info about field ID dropdown.
- Query Options:**
  - append wildcards to words
  - combine searches with
  - Number of entries to display per page: 30
  - Extended query form** (highlighted with a red box and arrow)
- Search Fields:** four dropdown menus for ID, each with a separate input field below it.
- Retrieval Options:** retrieve entries of type
- View Selection:**
  - Use predefined view: \* Names only \*
  - Create your own view: Select fields to display (ID, Header, Date, Title, Compound, Molecule, Synonym).

3. In the field 'Header' choose 'TNF' (1). In the field 'Molecule' input 'RECEPTOR'(2). Click the button 'Submit Query'(3).

The screenshot shows the ENPDB query interface with the following details:

- Top Navigation:** TOP PAGE, QUERY (highlighted), RESULTS, SESSIONS, VIEWS, DATABASES, HELP.
- Search Bar:** search ENPDB, Info about field ID dropdown.
- Query Options:**
  - append wildcards to words
  - combine searches with
  - Number of entries to display per page: 30
  - Standard query form**
  - Make default query page
- Search Fields (highlighted with red boxes):**

Field Name	Query	Include in View
ID		<input type="checkbox"/>
<u>Header</u>	TNF	<input type="checkbox"/> 1
Date	>= <input type="text"/> <= <input type="text"/>	<input type="checkbox"/>
Title		<input type="checkbox"/>
<u>Molecule</u>	RECEPTOR	<input type="checkbox"/> 2
Synonym		<input type="checkbox"/>
EC		<input type="checkbox"/>
BioUnit		<input type="checkbox"/>
- Buttons:** Submit Query (highlighted with a red box and arrow 3).

### III. Protein Integration Level. Chapter 1. Databases.

4. This will bring up the ‘Query Results’ page with the entries found in ENPDB. This page is illustrated in the Figure. To load the complete text of the entry, click the hyperlink of the entry.

The screenshot shows a web-based interface for querying a database. At the top, there's a navigation bar with links for TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, DATABANKS, and HELP. Below the navigation bar, a message says "Query "[enpdb-Header: TNF\*] & [enpdb-Molecule: RECEPTOR\*]" found 2 entries". On the left, a yellow sidebar contains buttons for "Perform operation" (radio buttons for "on all but selected" and "on selected", with "on all but selected" selected), "Link", "Save", and "View". It also has a dropdown menu set to "\* Names only \*". Below these are settings for "Number of entries to display per page" (set to 30) and a "Printer Friendly" link. The main content area lists two entries: "ENPDB:1CA4" and "ENPDB:1CA9", each preceded by a small checkbox. A red arrow points to the "ENPDB:1CA4" link. At the bottom of the page, it says "SRS 6.0.7.3 | feedback".

5. Click the name of the entry. This will bring up the page illustrated in the figure.

This screenshot shows the detailed information for the entry ENPDB:1CA4. The title "ENPDB:1CA4" is at the top in purple. Below it is the ID "1CA4 (RasMol, 3D Atlas)". The entry details include: HEADER TNF SIGNALING, DATE 23-FEB-1999, TITLE STRUCTURE OF TNF RECEPTOR ASSOCIATED FACTOR 2 (TRAF2), COMPOUND MOL\_ID: 1; MOLECULE: TNF RECEPTOR ASSOCIATED FACTOR 2; CHAIN: A, B, C, D, E, F; FRAGMENT: TRAF DOMAIN; SYNONYM: TRAF2; ENGINEERED: YES; BIOLOGICAL\_UNIT: TRIMER; MOL\_SOURCE BIOLOGY; SOURCE MOL\_ID: 1; ORGANISM\_SCIENTIFIC: HOMO SAPIENS; ORGANISM\_COMMON: HUMAN; SYNTHESIS\_MOL\_ID: 1; EXPRESSION\_SYSTEM: ESCHERICHIA COLI; EXPRESSION\_SYSTEM\_STRAIN: BL21; EXPRESSION\_SYSTEM\_VECTOR\_TYPE: PLASMID; EXPRESSION\_SYSTEM\_PLASMID: PET24D; KEYWORD TNF SIGNALING, TRAF, ADAPTER PROTEIN, CELL SURVIVAL; TECHNIQUE X-RAY DIFFRACTION; AUTHOR Y.C.PARK,V.BURKITT,A.R.VILLA,L.TONG,H.WU; JRNL AUTH Y.C.PARK,V.BURKITT,A.R.VILLA,L.TONG,H.WU; JRNL TITL STRUCTURAL BASIS FOR SELF ASSOCIATION AND RECEPTOR RECOGNITION OF HUMAN TRAF2; JRNL REF NATURE; JRNL VOLUME 398; JRNL YEAR 1999; RESOLUTION 2.2; CHAIN\_AMOUNT 6.

### Example 2.

The query is 'To find all protein complexes with the ordered quantitative parameters: the number of DNA/RNA molecules in the complex is at least 1; the number of proteins in the complex is at least 2; number of protein chains is at least 2; the protein chain length is at least 50; the number of alpha-helices in the protein is at least 2; the number of beta-sheets is at least 3; the number of ligands in the complex is at least 1; the ligand name is one of the following MAGNESIUM, ZINC, or INOSINE'

To make such a query, you should perform the following:

1 To input these parameters, in the field 'ChainAmount' select the sign “>=” and enter the value 2. In the field 'ChainSizes' choose the sign “>=” and input the value 50. In the field 'HelixAmount' choose the sign “>=” and insert the value 2. In the field SheetAmount choose the sign “>=” and enter the value 3. In the field 'DnaRnaAmount' choose the sign “>=” and input the value 1. In the filed 'ProteinAmount' choose the sign “>=” and input the value 2. In the field 'HetAmount' choose the sign “>=” and input the value 1. In the field Heterogen input MAGNESIUM|ZINC|INOSINE. Click the button “Submit Query”. This query is illustrated in the Figure.

Resolution	>=		<=		
<u>ChainAmount</u>	>=	2	<=		<input type="checkbox"/>
<u>ChainSizes</u>	>=	50	<=		<input type="checkbox"/>
<u>HelixAmount</u>	>=	2	<=		<input type="checkbox"/>
<u>SheetAmount</u>	>=	3	<=		<input type="checkbox"/>
<u>DnaRnaAmount</u>	>=	1	<=		<input type="checkbox"/>
<u>ProteinAmount</u>	>=	2	<=		<input type="checkbox"/>
<u>HetAmount</u>	>=	1	<=		<input type="checkbox"/>
<u>Heterogen</u>	MAGNESIUM ZINC INOSINE				<input type="checkbox"/>
<u>LinkEnbl</u>					<input type="checkbox"/>
<u>LinkPir</u>					<input type="checkbox"/>
<u>LinkSwissProt</u>					<input type="checkbox"/>
<u>LinkTransfac</u>					<input type="checkbox"/>
<u>LinkTrd4</u>					<input type="checkbox"/>

SRS 6.0.7.3 | [feedback](#)

### III. Protein Integration Level. Chapter 1. Databases.

2. After submitting the query, you will see the ‘Query Results’ page with the list of entries found in ENPDB database. To load the complete text of the entry, click the hyperlink with the entry name. This will bring up the page illustrated in the next Figure.

The screenshot shows a web-based search interface for the ENPDB database. At the top, there's a navigation bar with links for TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, DATABANKS, and HELP. Below the navigation bar, a search query is displayed: "Query "[enpdB-ChainAmount# 2] & [enpdB-ChainSizes# 50] & [enpdB-HelixAmount# 2] & [enpdB-SheetAmount# 3] & [enpdB-DnaRnaAmount# 1] & [enpdB-ProteinAmount# 2] & [enpdB-HetAmount# 1] & [enpdB-Heterogen: MAGNESIUM\*[ZINC\*|INOSINE\*]" found 37 entries". A red arrow points to the entry "ENPDB:1BSU" in the list of results.

3. Click the name of the entry. This will bring up the entry page illustrated in the Figure.

The screenshot shows the detailed view of the ENPDB entry for ID 1BSU. The page title is "ENPDB:1BSU". The entry details include:

- ID:** 1BSU (RasMol, 3D Atlas)
- HEADER:** COMPLEX (ENDONUCLEASE/DNA)
- DATE:** 31-AUG-1998
- TITLE:** ECORV/GAMITC/CA2+
- COMPOUND:** MOL\_ID: 1; MOLECULE: ENDONUCLEASE; CHAIN: A, B; EC: 3.1.21.4; ENGINEERED: YES; BIOLOGICAL\_UNIT: HOMODIMER; MOL\_ID: 2; MOLECULE: DNA; CHAIN: C, D; ENGINEERED: YES; BIOLOGICAL\_UNIT: DUPLEX; MOL\_SOURCE MIXED
- SOURCE:** MOL\_ID: 1; ORGANISM\_SCIENTIFIC: ESCHERICHIA COLI; MOL\_ID: 2; SYNTHETIC: YES;
- SYNTHESIS:** MOL\_ID: 1; EXPRESSION\_SYSTEM: ESCHERICHIA COLI; MOL\_ID: 2;
- KEYWORD:** COMPLEX (ENDONUCLEASE/DNA)
- TECHNIQUE:** X-RAY DIFFRACTION
- AUTHOR:** J.PERONA, A.MARTIN
- JRNL:** AUTH A.M.MARTIN, M.D.SAM, N.O.REICH, J.J.PERONA  
TITLE STRUCTURAL AND ENERGETIC ORIGINS OF INDIRECT READOUT IN SITE-SPECIFIC DNA CLEAVAGE BY A RESTRICTION ENDONUCLEASE  
REF TO BE PUBLISHED
- REMARK\_1:** REFERENCE 1  
AUTH N.C.HORTON, J.J.PERONA  
TITLE METAL ION MEDIATED SUBSTRATE-ASSISTED CATALYSIS IN TYPE II RESTRICTION ENDONUCLEASES  
REF TO BE PUBLISHED

### III. Protein Integration Level. Chapter 1. Databases.

4 Click the 3D Atlas link to visualize the entry by 3D Browser.

[ENPDB:1BSU](#)

ID [1BSU](#) ([RasMol](#), [3D Atlas](#))  
HEADER COMPLEX (ENDONUCLEASE/DNA)  
DATE 31-AUG-1998  
TITLE EC0RV/GAMITC/CA2+  
COMPOUND MOL\_ID: 1; MOLECULE: ENDONUCLEASE; CHAIN: A, B; EC: 3.1.21.4;  
ENGINEERED: YES; BIOLOGICAL\_UNIT: HOMODIMER; MOL\_ID: 2; MOLECULE: DNA; CHAIN:  
C, D; ENGINEERED: YES; BIOLOGICAL\_UNIT: DUPLEX;  
MOL\_SOURCE MIXED  
SOURCE MOL\_ID: 1; ORGANISM\_SCIENTIFIC: ESCHERICHIA COLI; MOL\_ID: 2; SYNTHETIC:  
YES;  
SYNTHESIS MOL\_ID: 1; EXPRESSION\_SYSTEM: ESCHERICHIA COLI; MOL\_ID: 2;  
KEYWORD COMPLEX (ENDONUCLEASE/DNA)  
TECHNIQUE X-RAY DIFFRACTION  
AUTHOR J.PERONA,A.MARTIN  
JRNL AUTH A.M.MARTIN,M.D.SAM,N.O.REICH,J.J.PERONA  
JRNL TITL STRUCTURAL AND ENERGETIC ORIGINS OF INDIRECT READOUT IN SITE-SPECIFIC  
DNA CLEAVAGE BY A RESTRICTION ENDONUCLEASE  
JRNL REF TO BE PUBLISHED  
REMARK\_1 REFERENCE 1  
REMARK\_1 AUTH N.C.HORTON,J.J.PERONA  
REMARK\_1 TITL METAL ION MEDIATED SUBSTRATE-ASSISTED CATALYSIS IN TYPE II  
RESTRICTION ENDONUCLEASES  
REMARK\_1 REF TO BE PUBLISHED

4.

[ENPDB:1CA4](#)

ID [1CA4](#) ([RasMol](#), [3D Atlas](#))  
HEADER TNF SIGNALING  
DATE 23-FEB-1999  
TITLE STRUCTURE OF TNF RECEPTOR ASSOCIATED FACTOR 2 (TRAF2)  
COMPOUND MOL\_ID: 1; MOLECULE: TNF RECEPTOR ASSOCIATED FACTOR 2; CHAIN: A, B, C,  
D, E, F; FRAGMENT: TRAF DOMAIN; SYNONYM: TRAF2; ENGINEERED: YES;  
BIOLOGICAL\_UNIT: TRIMER;  
MOL\_SOURCE BIOLOGY  
SOURCE MOL\_ID: 1; ORGANISM\_SCIENTIFIC: HOMO SAPIENS; ORGANISM\_COMMON: HUMAN;  
SYNTHESIS MOL\_ID: 1; EXPRESSION\_SYSTEM: ESCHERICHIA COLI;  
EXPRESSION\_SYSTEM\_STRAIN: BL21; EXPRESSION\_SYSTEM\_VECTOR\_TYPE: PLASMID;  
EXPRESSION\_SYSTEM\_PLASMID: PET24D;  
KEYWORD TNF SIGNALING, TRAF, ADAPTER PROTEIN, CELL SURVIVAL  
TECHNIQUE X-RAY DIFFRACTION  
AUTHOR Y.C.PARK,V.BURKITT,A.R.VILLA,L.TONG,H.WU  
JRNL AUTH Y.C.PARK,V.BURKITT,A.R.VILLA,L.TONG,H.WU  
JRNL TITL STRUCTURAL BASIS FOR SELF ASSOCIATION AND RECEPTORrecognition OF  
HUMAN TRAF2  
JRNL REF NATURE  
JRNL VOLUME 398  
JRNL YEAR 1999  
RESOLUTION 2.2  
CHAIN\_AMOUNT 6

Comments and questions are welcome to Vladimir Ivanisenko ([salix@bionet.nsc.ru](mailto:salix@bionet.nsc.ru))

## 2. PDBSite Database

Release 2003

**1. Database description:** PDBSite is a database on protein active sites and their spatial environment.

**2. Access to PDBSite database:**

<http://www.domain.com/mgs/gnw/pdbsite/>

**3. Database content**

SRS table	Description	Number of entries
PDBSITE	database on protein active sites and their spatial environment	2842

**4. List of biological tasks that could be solved by using the PDBSite database:**

- search for proteins according the data on protein function and structural characteristics of protein active sites;
- search for proteins based on the data about sites subjected to biochemical modifications;
- search for functional sites on the basis of data on biological classification of proteins.

**5. SRS table format**

### PDBSITE

PDBSite database contains the fields of two types: those that could be queried and searched or not. The fields that could not be queried contain additional information on the site structure and their surroundings.

Below is given the description of fields that could be queried.

Line code	Field name	Field description
ID	ID	Entry identifier. This identifier is unique within PDBSite.
PDBID	PDBID	PDB ID code. This identifier is unique within PDB.
HEADER	Header	It contains PDB classification for the entry. The field content corresponds to that in PDB.
TITLE	Title	It contains the title for experiment or analysis described in the entry. The field content corresponds to that in PDB.
KEYWORD	Keyword	It contains keywords describing the macromolecule. The content corresponds to that of KEYWDS of PDB.
MOLECULE	Molecule	It contains names of macromolecules from the COMPND of PDB and is designed to search for entries by the names of macromolecules.
NUM_SITE_CHAINS	NumSiteChains	It contains number of different chains, which the residues of the site belong to.
SITE_DESCR	SiteDescr	It contains description of the site. The content corresponds to that of SITE_DESCRIPTION sub-field of REMARK 800 field of PDB.
RESIDUE_NOTAA	ResidueNotAA	It contains the names of residues that are not amino acids but contained in the site.
SITE_CHAINS	LenSite	It contains number of residues in the site.
SURR_CHAINS	LenSurround	It contains number of residues in the site environment.

EXPOSURE	ExposureSite	It contains average exposure of residues of the site.
EXPOSURE	ExposureSurround	It contains average exposure of residues of the site environment.
DISCONTINUITY	Discontinuity	It characterises discontinuity of the site by its primary structure.

Below is the list of fields that could not be queried, but contain additional information on the site structure and its surroundings.

MOLECULE - contains names of macromolecules from the COMPND of PDB and is designed to search for entries by the names of macromolecules.

MOLCHAINS – contains identifier of chains of a macromolecule.

CHAIN\_ID - contains chain identifier for the site and its environment.

POS – indicates the positions of site residues and its surroundings in a protein sequence.

RESNAME - contains names of site residues and its environment.

EXPOSE - contains exposure of each of the residues of the site and its environment.

SITE\_SEQ\_PROFILE – Contains relative frequency of an amino acid type at each position of the site in homologous proteins.

Physico-chemical parameters of the site and its surroundings are listed in a special table. Types of physico-chemical characteristics are given in the columns of this table. The order of physico-chemical characteristics in the columns of the table is indicated in the line ORDER. In the lines, the type of a site and its surroundings is indicated, as well as the way of calculation of the physico-chemical parameter for the site and its surroundings. Three types of physico-chemical parameters are listed in the table: average, sum and module of spatial moment.

In the lines SITE, there are physico-chemical parameters calculated for the site. The lines FULL\_SURROUND contain physico-chemical parameters calculated for all residues of site surroundings. The lines EXPOSED\_SURROUND contain physico-chemical parameters calculated for exposed residues of the site surroundings. The lines BURIED\_SURROUND contain physico-chemical parameters calculated for buried residues of site surroundings.

The table is organized in the following way.

The first is the line ORDER.

Then follows the line AVERAGE that indicates that in lines SITE, FULL\_SURROUND, EXPOSED\_SURROUND and BURIED\_SURROUND the average values of physico-chemical parameters will be indicated. These lines are placed below the line AVERAGE.

Next is the line SUM indicating that physico-chemical parameters listed in the lines below were calculated as the sum values. Then follows the line SPATIAL MOMENT indicating that physico-chemical parameters listed in subsequent lines SITE, FULL\_SURROUND, EXPOSED\_SURROUND and BURIED\_SURROUND were calculated as module of spatial moment. This is the end of the table.

PAIRWISE - contains pairwise distances between residues of the site.

COORDINATES CA\_ATOMS - contains C-alpha atom co-ordinates of site residues.

COORDINATES CENTRE\_MASS - contains centre mass co-ordinates of site residues.

### Examples of SRS queries to the PDBSITE database

#### Example 1

The query is

'Find all proteins containing Zinc-binding site'.

To make such a query, you should execute the following:

1. Choose PDBSITE SRS table on the page 'SRS access'.

The Protein Data Bank (PDB) contains data on the spatial protein structures and their biologically active sites (i.e., ligand binding regions, enzyme catalytic centers, regions subjected to biochemical modifications, etc.). However, neither of the well known systems searching PDB does not provide the user with possibility to make the queries related with the active sites. A database PDBSITE storing the data on biologically active sites contained in the PDB database has been developed. PDBSITE accumulates amino acid content, structure features calculated by spatial protein structures, and physicochemical properties of sites and their spatial surroundings.

**ACCESS to PDBSITE**      SRS ACCESS: [PDBSITE](#)

<b>General information</b>	<b>User's guide</b>
<a href="#">How to cite PDBSITE</a>	<a href="#">Brief manual on the database PDBSITE</a>
<a href="#">PDBSITE publications</a>	<a href="#">PDBSITE publications</a>
<a href="#">Contact us</a>	<a href="#">Contact us</a>
<b>Current release</b>	<b>Additional information</b>
PDSITE is updated after each new PDB release. The current release has 4723 entries and was indexed 19-Mar-2003.	<a href="#">Links to other databases and programs</a>

2. Click the button 'Search'.

**TOP PAGE** | **QUERY** | **RESULTS** | **SESSIONS** | **VIEWS** | **DATABANKS** | **HELP**

**PDBSITE**

The current release has 4723 entries and was indexed 19-Mar-2003.

**Description**: PDSITE databank

**WWW**: [WWW site](#)

**Data-fields in SRS**

Name	Short Name	Type	No of Keys	No of References	Indexing Date	Status
<a href="#">ID</a>	<a href="#">id</a>	<a href="#">id</a>	4723	4723	19-Mar-2003	ok
<a href="#">PDBID</a>	<a href="#">pid</a>	<a href="#">index</a>	1956	4723	19-Mar-2003	ok
<a href="#">Header</a>	<a href="#">hdr</a>	<a href="#">index</a>	403	7527	19-Mar-2003	ok
<a href="#">Title</a>	<a href="#">ttl</a>	<a href="#">index</a>	3090	38713	19-Mar-2003	ok
<a href="#">Keyword</a>	<a href="#">kw</a>	<a href="#">index</a>	1890	31201	19-Mar-2003	ok
<a href="#">Molecule</a>	<a href="#">mol</a>	<a href="#">index</a>	992	11072	19-Mar-2003	ok
<a href="#">MolChains</a>	<a href="#">mch</a>	<a href="#">show</a>	0	0		not indexed
<a href="#">NumSiteChains</a>	<a href="#">sc</a>	<a href="#">num</a>	6	4723	19-Mar-2003	ok
<a href="#">SiteDescr</a>	<a href="#">sd</a>	<a href="#">index</a>	2206	31608	19-Mar-2003	ok
<a href="#">ResidueNotAA</a>	<a href="#">rma</a>	<a href="#">index</a>	153	5341	19-Mar-2003	ok
<a href="#">LenSite</a>	<a href="#">lsi</a>	<a href="#">num</a>	35	4723	19-Mar-2003	ok
<a href="#">LenSurround</a>	<a href="#">lsu</a>	<a href="#">num</a>	244	4723	19-Mar-2003	ok
<a href="#">ExposureSite</a>	<a href="#">esi</a>	<a href="#">real</a>	1155	4723	19-Mar-2003	ok

### III. Protein Integration Level. Chapter 1. Databases.

3. Click the field ‘SiteDescr’ from the drop-down menu. For querying Zinc-binding sites insert ‘Zinc’ in the text-box located to the right (1). Click the button ‘Submit Query’(2).

The screenshot shows the PDBSITE query interface. At the top, there are tabs for TOP PAGE, QUERY (which is selected), RESULTS, SESSIONS, VIEWS, DATABASES, and HELP. Below the tabs, there is a search bar with 'search PDBSITE' and an 'Info' button. A dropdown menu for 'about field ID' is open. On the left, there is a sidebar with buttons for 'Submit Query' (highlighted with a red arrow labeled 2), 'append wildcards to words' (with a checked checkbox), 'combine searches with AND', 'Number of entries to display per page' set to 30, and 'Extended query form'. The main search area has a dropdown for 'SiteDescr' set to 'Zinc' (highlighted with a red box labeled 1) and three empty 'ID' dropdowns. Below these is a button 'retrieve entries of type Entry'. There are also buttons for 'Use predefined view' and 'Create your own view'. A list of fields to display is shown on the right, including ID, PDBID, Header, Title, Keyword, Molecule, and MnfChains. A scroll bar is visible on the right side of the list.

4. The results of the querying operation will be a list of PDBSITE entries. To display the complete text of an entry, click the hyperlink with the entry name.

The screenshot shows the results of a query for 'PDBSITE-SiteDescr: Zinc\*'. The top bar indicates 'Query "[pdbname-SiteDescr: Zinc\*]" found 227 entries' and a 'next' button. The results list contains 227 entries, each preceded by a checkbox. The first entry, 'PDBSITE:1A1RZN1', is highlighted with a red arrow. The results list includes: PDBSITE:1A1RZN1, PDBSITE:1A1RZN2, PDBSITE:1AF0ACT, PDBSITE:1AGNZN1, PDBSITE:1AGNZN2, PDBSITE:1AGNZN3, PDBSITE:1AGNZN4, PDBSITE:1AGNZN5, PDBSITE:1AGNZN6, PDBSITE:1AGNZN7, PDBSITE:1AGNZN8, PDBSITE:1AYMZN, PDBSITE:1AYNZN, PDBSITE:1B4LZN, PDBSITE:1B4TZN, PDBSITE:1BA9ZN, PDBSITE:1BC2ZNA, PDBSITE:1BC2ZNB, PDBSITE:1BH5ZN1, PDBSITE:1BH5ZN2, and PDBSITE:1BH5ZN3.

5. Complete text of the entry found

```
PDBSITE:1A1RZN1

ID 1A1RZN1
PDBID 1A1R
HEADER VIRAL PROTEIN
TITLE HCV NS3 PROTEASE DOMAIN:NS4A PEPTIDE COMPLEX
KEYWORD VIRAL PROTEIN, SERINE PROTEASE, NONSTRUCTURAL PROTEINS, COFACTOR PEPTIDE, HEL
MOLECULE NS3 PROTEIN
MOL_CHAINS AB
NUM_SITE_CHAINS 1
SITE_DESCR ZINC BINDING SITE, MOLECULE A.
RESIDUE_NOTAA ZN
NUMBER_OF_AA 3 22
EXPOSURE 13.667 14.864
DISCONTINUITY 15.333
SITE_CHAINS AAA
POS 123 125 171
RESNAME CCC
EXPOSE 9 32 0
SITE_SEQ_PROFILE
  PDBNo V L I M F W Y G A P S T C H R K Q E N
  123 A 0 0 0 0 0 0 0 0 0 0 0 99 0 0 0 0 0 0 0
  125 A 0 0 0 0 0 0 0 0 0 0 0 0 99 0 0 0 0 0 0 0
  171 A 0 0 0 0 0 0 0 0 0 0 0 0 100 0 0 0 0 0 0 0
```

### Example 2

The query is

'Find all proteins containing the sites subjected to phosphorylation'.

To make such a query, you should perform the following:

1. Select from the drop-down menu the field SiteDescr. For searching for phosphorylation sites enter in the text-box to the right 'phosphorylation'(1). Click 'Submit Query' button (2).

The screenshot shows the PDB search interface with the following details:

- Top Bar:** Includes links for TOP PAGE, QUERY (highlighted in yellow), RESULTS, SESSIONS, VIEWS, DATABANKS, and HELP.
- Search Bar:** Contains "search PDBSITE" and "Info about field ID".
- Query Form:**
  - Submit Query:** Button labeled "Submit Query" with a red arrow pointing to it.
  - Advanced Options:** Includes "append wildcards to words" (checkbox checked) and "combine searches with AND".
  - Results Display:** "Number of entries to display per page" set to 30.
  - Extended Query Form:** "Extended query form" button.
  - Search Criteria:** A dropdown menu for "SiteDescr" is open, with "phosphorylation" entered in the text input field. This input field is highlighted with a red box and has a red arrow pointing to it.
  - Search Options:** "retrieve entries of type Entry" dropdown.
  - View Options:** "Use predefined view" dropdown set to "Names only".
  - Custom View:** "Create your own view" section with a scrollable list of fields to display: ID, PDBID, Header, Title, Keyword, Molecule, MolChains.

### III. Protein Integration Level. Chapter 1. Databases.

2. This brings up a list of PDBSite entries displayed in the ‘Query Results’ page. To display the complete text of an entry, click the hyperlink with the entry name.

The screenshot shows a web-based interface for querying protein phosphorylation sites. At the top, there's a navigation bar with links for TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, DATABANKS, and HELP. Below the navigation bar, there's a search bar with the query text: "Query '[pdbsite-SiteDescr: phosphorylation\*]' found 27 entries". On the left, there's a sidebar with options for 'Perform operation' (radio buttons for 'on all but selected' or 'on selected'), buttons for 'Link', 'Save', and 'View', and a dropdown menu set to '\* Names only \*'. It also includes a dropdown for 'Number of entries to display per page' set to 30, and a link to 'Printer Friendly'. The main content area lists 27 entries, each preceded by a checkbox. The entry 'PDBSITE:1A0BPHB' is highlighted with a red arrow pointing to it. Other entries listed include PDBSITE:1A041, PDBSITE:1A2OPON, PDBSITE:1A6JPH1, PDBSITE:1A6JPH2, PDBSITE:1AUZP, PDBSITE:1AX3PON, PDBSITE:1B89TK, PDBSITE:1BUZP, PDBSITE:1H4YSEA, PDBSITE:1H4YSEB, PDBSITE:1IBAS1, PDBSITE:1OPDAS6, PDBSITE:1OPDAS6, PDBSITE:1QKKPAD, PDBSITE:1RN1L1, PDBSITE:1SJSPOS, PDBSITE:2A0BPHB, PDBSITE:4PGMCIC, PDBSITE:5PGMCIA, and PDBSITE:5PGMCIB.

3. The complete text of the ‘phosphorilation site’ entry is shown in Figure.

```
PDBSITE:1A0BPHB

ID 1A0BPHB
PDBID 1AOE
HEADER HISTIDINE KINASE
TITLE HISTIDINE-CONTAINING PHOSPHOTRANSFER DOMAIN OF ARCB FROM ESCHERICHIA COLI
KEYWORD HISTIDINE KINASE, PHOSPHOTRANSFER, TWO-COMPONENT SYSTEM, FOUR-HELIX BUNDLE
MOLECULE AEROBIC RESPIRATION CONTROL SENSOR PROTEIN ARCB
MOL_CHAINS
NUM_SITE_CHAINS 1
SITE_DESCR PHOSPHORYLATION SITE.
NUMBER_OF_AA 1 8
EXPOSURE 35.000 17.250
DISCONTINUITY 0.000
SITE_CHAINS
POS 715
RESNAME H
EXPPOSE 35
SITE_SEQ_PROFILE
  PDBNO  V   L   I   M   F   W   Y   G   A   P   S   T   C   H   R   K   Q   E   N
    715   0   0   0   0   0   0   0   0   0   0   0   0   0 100   0   0   0   0
SURR_CHAINS
POS 711 713 714 716 717 718 719 737
RESNAME VEGKIKGQ
```

### Example 3

The query is

'Find all proteins containing discontinuous sites, with the length of 3 amino acid residues'.

To make such a query, you should perform the following:

1. Select from the drop-down menu the field 'Discontinuity'. In the text-box located to the right enter '0'. In menu below select the field 'LenSite'. In the field to the right enter the value '3'. Click the button 'Submit Query'.

TOP PAGE | QUERY | RESULTS | SESSIONS | VIEWS | DATABASES | HELP

Reset | search PDBSITE | Info about field ID

**Submit Query** 2

append wildcards to words

combine searches with **AND**

Number of entries to display per page 30

**Extended** query form

separate multiple values by & (and), | (or), ! (and not)

**1**

Discontinuity	0
LenSite	3
ID	
ID	

retrieve entries of type Entry

Use predefined view \* Names only \* Create your own view

Select fields to display:

- ID
- PDBID
- Header
- Title
- Keyword
- Molecule
- Macromolecules

2. The results of the querying operation will be a list of PDBSITE entries. To display the complete text of an entry, click the hyperlink with the entry name.

TOP PAGE | QUERY | RESULTS | SESSIONS | VIEWS | DATABASES | HELP

Reset | Query "[pdbname-Discontinuity# 0] & [pdbname-LenSite# 3]" found 12 entries

Perform operation

on all but selected

on selected

**Link** **Save** **View** **\* Names only \***

Number of entries to display per page 30

**Printer Friendly**

PDBSITE:1AWZNLS  
 PDBSITE:1B8FMIO  
**PDBSITE:1H9JMO1** ←  
 PDBSITE:1H9KWO1  
 PDBSITE:1H9MAC1  
 PDBSITE:1H9MAC2  
 PDBSITE:1H9MAC5  
 PDBSITE:1H9MAC6  
 PDBSITE:1MFNRGD  
 PDBSITE:2MFNRGD  
 PDBSITE:3LRIAR  
 PDBSITE:3LRINTR

SRS 6.0.7.3 | [feedback](#)

### III. Protein Integration Level. Chapter 1. Databases.

3. Complete text of the entry found is seen in Figure.

PDBSITE:1MFNRGD																																																																																																							
ID 1MFNRGD																																																																																																							
PDBID <u><a href="#">1MFN</a></u>																																																																																																							
HEADER CELL ADHESION PROTEIN																																																																																																							
TITLE SOLUTION NMR STRUCTURE OF LINKED CELL ATTACHMENT MODULES OF MOUSE FIBRONECTIN C																																																																																																							
KEYWORD CELL ADHESION PROTEIN, RGD, EXTRACELLULAR MATRIX, HEPARIN-BINDING, GLYCOPROTE																																																																																																							
MOLECULE FIBRONECTIN																																																																																																							
MOL_CHAINS																																																																																																							
NUM_SITE_CHAINS 1																																																																																																							
SITE_DESCR CELL ADHESION SITE.																																																																																																							
NUMBER_OF_AA 3 10																																																																																																							
EXPOSURE 46.000 40.900																																																																																																							
DISCONTINUITY 0.000																																																																																																							
SITE_CHAINS																																																																																																							
POS 168 169 170																																																																																																							
RESNAME RGD																																																																																																							
EXPPOSE 28 45 65																																																																																																							
SITE_SEQ_PROFILE																																																																																																							
<table border="1"> <thead> <tr> <th>PDBNo</th> <th>V</th> <th>L</th> <th>I</th> <th>M</th> <th>F</th> <th>W</th> <th>Y</th> <th>G</th> <th>A</th> <th>P</th> <th>S</th> <th>T</th> <th>C</th> <th>H</th> <th>R</th> <th>K</th> <th>Q</th> <th>E</th> <th>N</th> </tr> </thead> <tbody> <tr> <td>168</td> <td>3</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>3</td> <td>8</td> <td>5</td> <td>11</td> <td>3</td> <td>0</td> <td>3</td> <td>50</td> <td>8</td> <td>5</td> <td>0</td> <td>3</td> </tr> <tr> <td>169</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>53</td> <td>3</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>5</td> <td>5</td> <td>0</td> <td>34</td> <td>0</td> <td>0</td> </tr> <tr> <td>170</td> <td>3</td> <td>3</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>5</td> <td>5</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>5</td> <td>32</td> <td>11</td> <td>5</td> <td>0</td> </tr> </tbody> </table>																								PDBNo	V	L	I	M	F	W	Y	G	A	P	S	T	C	H	R	K	Q	E	N	168	3	0	0	0	0	0	0	3	8	5	11	3	0	3	50	8	5	0	3	169	0	0	0	0	0	0	0	53	3	0	0	0	0	5	5	0	34	0	0	170	3	3	0	0	0	0	0	0	5	5	0	0	0	0	5	32	11	5	0
PDBNo	V	L	I	M	F	W	Y	G	A	P	S	T	C	H	R	K	Q	E	N																																																																																				
168	3	0	0	0	0	0	0	3	8	5	11	3	0	3	50	8	5	0	3																																																																																				
169	0	0	0	0	0	0	0	53	3	0	0	0	0	5	5	0	34	0	0																																																																																				
170	3	3	0	0	0	0	0	0	5	5	0	0	0	0	5	32	11	5	0																																																																																				
SURR_CHAINS																																																																																																							

#### Example 4

The query is

'Find all sites for the protein family of histidine kinases'.

To make such a query, you should perform the following:

1. Select the option from the drop-down menu 'Header'. In the text-box enter 'HISTIDINE&KINASE'(1). Click the button 'Submit Query'(2).

The screenshot shows the PDB search interface with the following details:

- Header:** The text box contains "Header HISTIDINE&KINASE".
- Submit Query:** The button labeled "Submit Query" is highlighted with a yellow box and has a red arrow pointing to it.
- Other Fields:** There are three additional text boxes below the main one, each with "ID" selected in the dropdown menu.
- Buttons:** "Reset", "Info", and "Help" buttons are visible at the top.
- Views:** Options for "Use predefined view" and "Create your own view" are shown.
- Select fields to display:** A scrollable list includes "ID", "PDBID", "Header", "Title", "Keyword", and "Molecule".

### III. Protein Integration Level. Chapter 1. Databases.

2. Submitting the query will bring up the list of PDBSITES entries found. To display the complete text of an entry, click the hyperlink with the entry name.

The screenshot shows a software interface for querying databases. At the top, there's a navigation bar with links for TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, DATABASES, and HELP. Below the navigation bar, a yellow header bar displays the query "Query "[pdbsite-Header:HISTIDINE\*]" found 6 entries". On the left, there's a sidebar with options for "Perform operation" (radio buttons for "on all but selected" and "on selected", with "on selected" checked), buttons for Link, Save, View, and a dropdown menu for "Names only". Below this is a section for "Number of entries to display per page" with a dropdown set to 30. At the bottom of the sidebar is a "Printer Friendly" button. The main content area lists six entries: PDBSITE:1A0BPHB, PDBSITE:1A0BZNB (with a red arrow pointing to it), PDBSITE:1H5YAC1, PDBSITE:1H5YAC2, PDBSITE:1H5YAC5, and PDBSITE:1H5YAC6. At the very bottom of the interface, it says "SRS 6.0.7.3 | feedback".

3. Complete text of the entry found is shown in Figure.

PDBSITE:1A0BZNB	
ID	1A0BZNB
PDBID	<a href="#">1A0B</a>
HEADER	HISTIDINE KINASE
TITLE	HISTIDINE-CONTAINING PHOSPHOTRANSFER DOMAIN OF ARCB FROM ESCHERICHIA COLI
KEYWORD	HISTIDINE KINASE, PHOSPHOTRANSFER, TWO-COMPONENT SYSTEM, FOUR-HELIX BUNDLE MOLECULE AEROBIC RESPIRATION CONTROL SENSOR PROTEIN ARCB
MOL_CHAINS	
NUM_SITE_CHAINS	1
SITE_DESCR	ZN BINDING SITE.
NUMBER_OF_AA	4 39
EXPOSURE	37.500 19.897
DISCONTINUITY	6.750
SITE_CHAINS	
POS	728 746 754 758
RESNAME	HDEE
EXPPOSE	11 75 33 31
SITE_SEQ_PROFILE	
PDBNo	V L I M F W Y G A P S T C H R K Q E N
728	0 0 0 0 0 0 0 0 0 0 0 0 0 0 75 25 0 0 0 0
746	0 0 0 0 13 0 0 0 0 0 0 0 0 0 0 0 0 0 0 25 0
754	0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 13 13 75 0
758	0 13 0 0 0 0 13 0 0 0 0 0 0 0 0 0 0 0 0 0 75 0

Comments and questions are welcome to Vladimir Ivanisenko ([salix@bionet.nsc.ru](mailto:salix@bionet.nsc.ru))

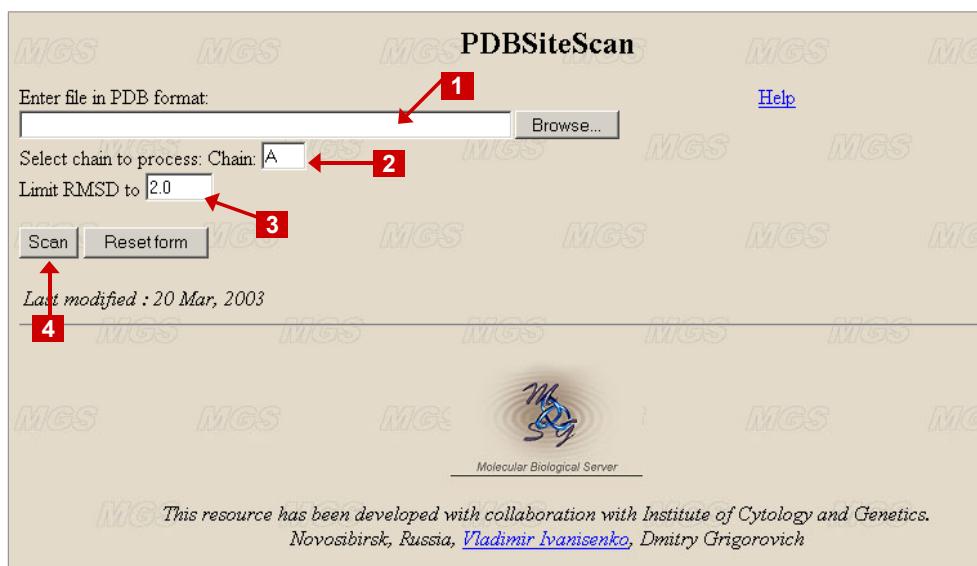
## 3. PDBSiteScan

Release 2003

### Program description:

The program PDBSiteScan automatically performs the best superposition of sites from PDBSite with the 3D structure of a protein under study.

The algorithm developed is based on exhaustion of all the possible combinations of protein positions to be compared with the site. The program realizes the following steps. At the first step, the amino acids of a fragment and the site are compared. If they are identical, their 3D structures are compared at the second step. If the RSMD obtained is lower than the user-specified value, the fragment examined, this fragment is added to the list of results.



### Access to PDBSiteScan:

<http://www.domain.com/mgs/gnw/pdbsitescan/> link 'Search for Functional Sites in Protein Tertiary Structures'

### List of biological tasks that could be solved by using the PDBSiteScan

- detection of potential active sites in the 3D structure;
- structural alignment with known functional sites stored in PDBSite.

### Data input

Input a filename of tertiary structure of a protein under study into the text window (1). The tertiary structure should be in PDB format.

### Program options

Input a chain ID of the protein to be analyzed (2).

The program selects all sites for which the root mean square deviation (RMSD) for atoms  $N$ ,  $C\alpha$ ,  $C$  does not exceed a level specified by the user in the text window (3).

### Program execution

Start the tools processing by clicking the button 'Scan' (4).

## Program output

The program output includes links to the PDBSite database, brief description of the site, RMSD value, positions of the protein fragment, and the list of residues.

### Example

Let us consider the application of the program by the example of recognition of catalytic site in the protein 1ELV (the family HYDROLASE).

This figure demonstrates the input data for PDBSiteScan.

The screenshot shows the PDBSiteScan web interface. At the top, there is a text input field labeled "Enter file in PDB format:" containing the path "C:\INSTALL\pdb1elv.ent". Next to it is a "Browse..." button. To the right is a "Help" link. Below the input field are two dropdown menus: "Select chain to process: Chain: A" and "Limit RMSD to 3.0". Underneath these are two buttons: "Scan" and "Reset form". A timestamp "Last modified : 20 Mar, 2003" is displayed. In the center, there is a logo for "Molecular Biological Server" featuring a stylized "M" and "B" with a blue ribbon-like effect. Below the logo, a note reads: "This resource has been developed with collaboration with Institute of Cytology and Genetics. Novosibirsk, Russia, Vladimir Ivanisenko, Dmitry Grigorovich".

This figure demonstrates the result of PDBSiteScan operation.

```
PDBSiteID 1AUOACA
PDBID 1AUO
SITE_DESCR CATALYTIC TRIAD.
RMSD 0.000
Chain for each residue: AAA
Positions: 114 168 199
Residues: S D H
END

PDBSiteID 1AUOACB
PDBID 1AUO
SITE_DESCR CATALYTIC TRIAD.
RMSD 0.265
Chain for each residue: AAA
Positions: 114 168 199
Residues: S D H
END

PDBSiteID 1AURACB
PDBID 1AUR
SITE_DESCR CATALYTIC TRIAD.
RMSD 0.465
Chain for each residue: AAA
Positions: 114 168 199
Residues: S D H
END
```

Comments and questions are welcome to Vladimir Ivanisenko ([salix@bionet.nsc.ru](mailto:salix@bionet.nsc.ru))

## 4. Artificial Selected Peptides/Proteins Database (ASPD)

Release 2003

### Database description:

ASPD is a curated database on selected from randomized pools proteins and peptides designed for accumulation of experimental data on protein functionality obtained by in vitro directed evolution methods (phage display, ribosome display, SIP etc.)

### Access to ASPD database:

<http://www.domain.com/mgs/gnw/aspd/>

### Database content

SRS table	Description	Number of entries
ASPD_ALIGN	Experiments	195
ASPD_REF	Literature	112

### Biological tasks that could be solved by using the ASPD database

- Find sequences of proteins and peptides selected by phage display by querying the database with biological keywords
- Annotate proteins by finding its homologs in ASPD

### SRS table format

#### ASPD\_ALIGN

Field name	Field description
Identifier	Unique identifier of an entry. It has the format PH1XXNNN, where X is a letter of latin alphabet, N is a number (0-9)
Lit_reference	Link to the literary reference in the ASPD_REF database (the identifier of the corresponding entry in ASPD_REF having the form PH4XXNNN)
Target	the substrate for binding to which the peptides or proteins were selected
Template	the native protein or peptide performing the function for which peptides from random pool were selected
Keywords	A list of keywords (separated with commas) describing the biological relevance of the entry.
Link	Links to external databases. Has the form Link DbName EntryId, where the DbName is one of the following: SwissProt, Prosite, PDB, Enzyme, and EntryId is the identifier of the corresponding entry in these databases.
Comment	The referent's comments concerning conditions of the experiment (free text)
Consensus	The general form of peptides retrieved by in vitro evolution (in PROSITE signature), usually provided by the authors of the original paper.
Number of sequences	Number of different amino acid sequences in the entry
Alignment	Author's alignment of amino acid sequences (those retrieved from random pools are denoted with numbers, others (native or constructed) with letters).

**ASPD\_REF**

Field name	Field description
Identifier	Unique identifier of an entry. It has the format PH4XXNNN, where X is a letter of latin alphabet, N is a number (0-9)
Authors	Authors of the paper
Title	Title of the paper
Journal	Journal name
Volume	Volume
Year	Year
Pages	Page number(s)
Medline	PubMed ID
Corresponding_Author	Full name of the corresponding author with e-mail address

**Comments and questions are welcome to Vadim Valuev ([valuev@bionet.nsc.ru](mailto:valuev@bionet.nsc.ru)).**

## 5. DCS - Database on residue coordination spheres

Release 2003

### 1. Database Description

DCS is database on structural and physico-chemical information of amino acid residue coordination spheres in protein 3D structures. Current release contain data from DNA-binding proteins structures (46 non-homologous protein chains). We concentrate on this specific subset of proteins, because they are involved in gene regulation function.

### 2. Access to DCS database:

<http://www.domain.com/mgs/gnw/dcs/>

### 3. Database content:

DCS database contains of three tables, related to environment class data (DCS\_CLUST), protein data (DCS\_PROT) and residues data (DCS\_RES). Brief description of tables is presented in Table 1.

**Table 1.** Description of data tables in DCS database.

SRS table	Description	Number of entries
DCS_CLUST	List of coordination sphere clusters with similar physico-chemical properties	9
DCS_PROT	List of proteins used for database derivation	46
DCS_RES	List of coordination spheres and their characteristics for residues in protein dataset	4240

### 4. Biological tasks that could be solved using DCS database

Database aimed to find residues that have similar physico-chemical environment in protein structures and annotate protein residues in terms of their environment classes.

### 5. How DCS database was constructed?

The construction of DCS database involved the following stages: (1) determining local spatial environment of residues; (2) calculating physicochemical characteristics of the local environment of individual residues; (3) clustering the local environment with respect to their physicochemical properties; (4) analyzing the resulting classes of physicochemical characteristics of the environment and their interrelations with the structural properties of residues; and (5) analyzing the similarity of spatial environments for various amino acid types.

*Data.* Spatial structures of the transcription factors displaying a degree of sequence similarity not exceeding 40% were used for the analysis. These structures included (chain names are indicated in parenthesis): 1AIS(B), 1BH9(A,B), 1BM8, 1BOR, 1BVO(A), 1C7U(A), 1CF7(A,B), 1CI6(A), 1CQT(I), 1D8J(A), 1DH3(A), 1DL6(A), 1DP7(P), 1ENW(A), 1EO0(A), 1EQF(A), 1EXE(A), 1F3U(A,B), 1F4S(P), 1F62(A), 1G2Y(A), 1GD2(E), 1HKS, 1HLO(A), 1I27(A), 1I4W(A), 1JFI(A,B), 1K99(A), 1MNM(A,C), 1NCS, 1PUE(E), 1QQH(A), 1SKN(P), 1SP1, 1TBA(A,B), 1TF3(A), 1TFI, 1YTF(B,C), and 3HSF. The total number  $N$  of the residues analyzed amounted to 4240.

*Structural properties of residues.* We used two parameters, calculated by the program DSSP (Kabsch and Sander, 1983), as structural properties, namely, secondary structure of residues and the surface area accessible to water. These characteristics were calculated for each protein chain

separately from others. Additionally, we determined the residues being in contact with DNA that have at least one atom at a distance less than 4 Å to any atoms of DNA.

*Determination of local spatial environment of residues.* It was assumed that the local spatial environment of the  $i$ th (central) residue is formed by the residues whose C $\alpha$  atoms were located at a distance not exceeding 7 Å from the C $\alpha$  atom of the  $i$ th residue. No exclusions were made for the immediate neighbor residues. We considered that the neighbors immediate in the primary structure also contribute to the local environment of a residue in question.

*Physicochemical characteristics of the local environment of individual residues.* Values of five amino acid properties, taken from (Bogardt *et al.*, 1980)—volume, polarity, isoelectric point, hydrophobicity, and the surface area accessible to water—were used while describing the physicochemical characteristics. The spatial environment of the  $i$ th residue was characterized with the vector  $f_i = \{f_k, k = 1, \dots, 5\}$ ; each component of the vector corresponded to the value of  $k$ th amino acid property averaged over the residues forming the local spatial environment.

*Classification of physicochemical properties of local environment.* To classify the local spatial environment of individual residues according to their physicochemical characteristics, we used the hierarchical cluster analysis (Sneath and Sokal, 1973). The squared Euclidian distance

$d_{ij} = [\sum_{k=1}^5 (f_{ik} - f_{jk})^2]$  was used as a measure of the distance between characteristics of the

environment for a pair of residues  $i, j$ . Unweighted paired grouping method with arithmetic mean (UPGMA) was used for constructing the similarity tree.

## 6. SRS table format

### DCS

Field name	Field description
<b>DCS CLUST</b>	
ID	Environment class identifier
SZ	Size of class (number of environments of this class)
AM	Average solvent accessibility for cluster members
FM	Average physico-chemical properties for environment residues of cluster members
SC	Secondary structure class counts for cluster residues
SF	Secondary structure class fraction for cluster
CM	Residues - members of this class
<b>DCS PROT</b>	
ID	Protein identifier
DB	Reference to protein structure in PDB (format is ID:Chain)
SZ	Size of protein (number of residues of this protein)
PM	Residues - members of this protein
<b>DCS RES</b>	
ID	Residue identifier
PR	Reference to protein
CR	Reference to environment class
NM	Residue name
AA	Amino Acid type of residue
PI	Residue index in protein structure
SI	Residue index in protein sequence
SS	Secondary structure according to DSSP program
AC	Solvent Accessibility according to DSSP
DC	Residue is in contact with DNA if DC=1, not in contact if DC=0, if DC=-1 information not available N/A.
CN	Coordination number (number of contact residues)
AF	Frequencies of different amino acid types in coordination sphere (list in order "ARNDCQEGHILKMFPTWYV")
FX	Average physico-chemical properties for environment residues

**Comments and questions to** Dmitry Afonnikov (ada@bionet.nsc.ru).

## 7. References

1. Bogardt, R.A., Jr, Jones, B.N., Dwulet, F.E., Garner, W.H., Lehman, L.D., and Gurd, F.R. (1980). Evolution of the amino acid substitution in the mammalian myoglobin gene. *J. Mol. Evol.*, **15**:197-218.
2. Kabsch, W. and Sander, C. (1983). Dictionary of protein secondary structure: pattern recognition of hydrogen bonded and geometrical features. *Biopolymers*, **22**:2577-2637.
3. Sneath, P.H.A. and Sokal, R.R. (1973). *Numerical Taxonomy. The Principles and Practice of Numerical Classification*. San Francisco: W.H. Freeman and Co.

## Examples of SRS queries to the DCS database

### Example 1

The query is 'To search for TBP protein (PDB identifier 1AIS), represented in DSC'.

To make such a query, you should perform the following:

1. From DCS start page follow the link 'DCS\_PROT' database (1).

DCS (Database on Coordination Spheres) contain structural and physico-chemical information on amino acid residue coordination spheres in protein 3D structures.

ACCESS to DCS → SRS ACCESS DCS RES **DCS PROT** DCS CLUST

**General information**  
How to cite DCS?  
Contact us

**User's guide**  
Brief manual on the database DCS

**Current release**  
The current release has 4240 entries and was indexed 01-Mar-2003.

**Additional information**  
Links to other databases and programs

**III. Protein Integration Level.** Chapter 1. Databases.

2. This will bring up the ‘DCS\_PROT’ page. This page is illustrated in the Figure. To perform database query, click 'PDBID' hyperlink (1).

Name	Short Name	Type	No of Keys	No of Entry References	Indexing Date	Status
<a href="#">ProteinIdentifier</a>	<a href="#">id</a>	id	46	46	25-Mar-2003	ok
<a href="#">PDBID</a>	<a href="#">pi</a>	index	39	46	25-Mar-2003	ok
<a href="#">ProteinSize</a>	<a href="#">sz</a>	num	38	46	25-Mar-2003	ok
<a href="#">ResidueList</a>	<a href="#">pm</a>	index	4240	4240	25-Mar-2003	ok

3. This will bring up the query page for ‘PDBID’ field of the DCS\_PROT table. Input '1ais' string into query field (1). Click the 'List Values' button (2).

4. This will bring up the ‘Query Results’ page with the entries found in DCS\_PROT. This page is illustrated in the Figure. Click the hyperlink of the entry.

**III. Protein Integration Level.** Chapter 1. Databases.

5. This will bring up the page illustrated in the figure. Click to the hyperlink to PROT1 entry.

The screenshot shows the SRS 6.0.7.3 software interface. At the top, there is a navigation bar with tabs: TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, DATABANKS, and HELP. Below the navigation bar, a yellow header bar displays the query "Query '[DCS\_PROT-PDBID:'1ais']' found 1 entries". On the left side, there is a sidebar with various buttons and dropdown menus. One of the buttons, "DCS PROT-PROT1", is highlighted with a red arrow pointing to it. The main content area shows the results of the query, which consists of a single line of text: "ID PROT1 DB 1AIS:B SZ 193 PM RES1 RES2 RES3 RES4 RES5 RES6 RES7 RES8 RES9 RES10 RES11 RES12 RES13 RES14 RES15 RES16 RES17 RES18 RES19 RES20 RES21 RES22 RES23 RES24 RES25 RES26 RES27 RES28 RES29 RES30 RES31 RES32 RES33 RES34 RES35 RES36 RES37 RES38 RES39 RES40 RES41 RES42 RES43 RES44 RES45 RES46 RES47 RES48 RES49 RES50 RES51 RES52 RES53 RES54 RES55 RES56 RES57 RES58 RES59 RES60 RES61 RES62 RES63 RES64 RES65 RES66 RES67 RES68 RES69 RES70 RES71 RES72 RES73 RES74 RES75 RES76 RES77 RES78 RES79 RES80 RES81 RES82 RES83 RES84 RES85 RES86 RES87 RES88 RES89 RES90 RES91 RES92 RES93 RES94 RES95 RES96 RES97 RES98 RES99 RES100 RES101 RES102 RES103 RES104 RES105 RES106 RES107 RES108 RES109 RES110 RES111 RES112 RES113 RES114 RES115 RES116 RES117 RES118 RES119 RES120 RES121 RES122 RES123 RES124 RES125 RES126 RES127 RES128 RES129 RES130 RES131 RES132 RES133 RES134 RES135 RES136 RES137 RES138 RES139 RES140 RES141 RES142 RES143 RES144 RES145 RES146 RES147 RES148 RES149 RES150 RES151 RES152 RES153 RES154 RES155 RES156 RES157 RES158 RES159 RES160 RES161 RES162 RES163 RES164 RES165 RES166 RES167 RES168 RES169 RES170 RES171 RES172 RES173 RES174 RES175 RES176 RES177 RES178 RES179 RES180 RES181 RES182 RES183 RES184 RES185 RES186 RES187 RES188 RES189 RES190 RES191 RES192 RES193 //".

6. This will bring up the resulting page of the PROT1 entry.

The screenshot shows the resulting page for the PROT1 entry. The title of the page is "DCS PROT-PROT1". The content of the page is a list of residues, starting with "ID PROT1" and "DB 1AIS:B", followed by "SZ 193" and "PM". The list of residues is as follows:

```
ID PROT1
DB 1AIS:B
SZ 193
PM
RES1 RES2 RES3 RES4 RES5 RES6 RES7 RES8 RES9 RES10
RES11 RES12 RES13 RES14 RES15 RES16 RES17 RES18 RES19 RES20
RES21 RES22 RES23 RES24 RES25 RES26 RES27 RES28 RES29 RES30
RES31 RES32 RES33 RES34 RES35 RES36 RES37 RES38 RES39 RES40
RES41 RES42 RES43 RES44 RES45 RES46 RES47 RES48 RES49 RES50
RES51 RES52 RES53 RES54 RES55 RES56 RES57 RES58 RES59 RES60
RES61 RES62 RES63 RES64 RES65 RES66 RES67 RES68 RES69 RES70
RES71 RES72 RES73 RES74 RES75 RES76 RES77 RES78 RES79 RES80
RES81 RES82 RES83 RES84 RES85 RES86 RES87 RES88 RES89 RES90
RES91 RES92 RES93 RES94 RES95 RES96 RES97 RES98 RES99 RES100
RES101 RES102 RES103 RES104 RES105 RES106 RES107 RES108 RES109 RES110
RES111 RES112 RES113 RES114 RES115 RES116 RES117 RES118 RES119 RES120
RES121 RES122 RES123 RES124 RES125 RES126 RES127 RES128 RES129 RES130
RES131 RES132 RES133 RES134 RES135 RES136 RES137 RES138 RES139 RES140
RES141 RES142 RES143 RES144 RES145 RES146 RES147 RES148 RES149 RES150
RES151 RES152 RES153 RES154 RES155 RES156 RES157 RES158 RES159 RES160
RES161 RES162 RES163 RES164 RES165 RES166 RES167 RES168 RES169 RES170
RES171 RES172 RES173 RES174 RES175 RES176 RES177 RES178 RES179 RES180
RES181 RES182 RES183 RES184 RES185 RES186 RES187 RES188 RES189 RES190
RES191 RES192 RES193
//
```

### Example 2

The query is 'To search for environment class with average solvent accessibility of central residue between 70 and 90'.

To make such a query, you should perform the following:

- From DCS start page follow the hyperlink 'SRS\_ACCESS' (1).

- This will bring up the 'SRS top page'. This page is illustrated in the Figure. To perform database query, select DCS\_CLUST database (1) and click extended query button (2).

### III. Protein Integration Level. Chapter 1. Databases.

3. This will bring up the query page. This page is illustrated in the Figure. To perform database query, input lower limit 70 for 'Mean Accessibility' field (1), upper limit 90 for the same field (2), click 'Include in View' button for this field (3). Click 'Submit Query button' (4).

Field Name	Query	Include in View
<u>ClassIdentifier</u>		<input type="checkbox"/>
<u>ClassSize</u>	>= <input type="text"/> 1 <= <input type="text"/>	<input type="checkbox"/>
<u>MeanAccessibility</u>	>= <input type="text"/> 70 <= <input type="text"/> 90	<input checked="" type="checkbox"/>
<u>Volume</u>	>= <input type="text"/> <= <input type="text"/>	<input type="checkbox"/>
<u>Polarity</u>	>= <input type="text"/> <= <input type="text"/>	<input type="checkbox"/>
<u>IsoelectricPoint</u>	>= <input type="text"/> <= <input type="text"/>	<input type="checkbox"/>

4. This will bring up the result page with the list of environment classes satisfying the above conditions. This page is illustrated in the Figure. To view environment class entry click to the corresponding hyperlink (1).

DCS_CLUST	MeanAccessibility
<input type="checkbox"/> <a href="#">DCS_CLUST:ENVCL7</a>	71.277
<input type="checkbox"/> <a href="#">DCS_CLUST:ENVCL8</a>	75.445

**III. Protein Integration Level.** Chapter 1. Databases.

5. This will bring up the resulting page of the ENVCL8 entry. This page is illustrated in the Figure.

The screenshot shows a web-based protein database interface. At the top, there is a navigation bar with icons for a paw print, and links for TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, and DATABANKS. A HELP button is also present. Below the navigation bar, there are buttons for Reset, View (with a dropdown menu showing 'Complete entries'), and a search input field. A message says 'This entry is from: DCS CLUSTENVCL8'. Below this, there is a section for 'DCS CLUST' with buttons for Save and Link, and a link to 'Printer Friendly'. The main content area displays detailed protein information for ENVCL8, including:

ID ENVCL8  
SZ 887  
AM 75.445  
FM  
Volume=47.008  
Polarity=44.392  
Isoelectric\_point=32.666  
Hydrophobicity=33.215  
Mean\_solvent\_acc=35.350  
SC  
X=177  
H=424  
T=117  
S=78  
G=24  
E=65  
B=2  
SF  
X=0.199549  
H=0.478016  
T=0.131905  
S=0.087937  
G=0.027057  
E=0.073281  
B=0.002255  
CM

## 6. Database on local conformations of protein chains -"conformons" (ConfDB)

Release 2003

### 1. Database description:

ConfDB contain structural information on conformational properties of short fragments (9 aa in length) in 3D structures of *C.elegans* proteins.

### 2. Access to ConfDB database:

<http://www.domain.com/mgs/gnw/confdb/>

### 3. Database content:

SRS table	Description	Number of entries
CONFDB_RES	List of peptides	790
CONFDB_PROT	List of proteins used for database derivation	5
CONFDB_CLUST	List of clusters of peptides with similar conformations	320

### 4. Biological tasks that could be solved using ConfDB database:

Find peptides that have similar local conformations in protein structures.

Annotate protein residues in terms of local conformations of polypeptide chains.

### 5. How ConfDB was constructed?

Database constructed from short polypeptide segments. They are collected from set of protein structures by sliding window of 9 residues. These segments are overlapped. All of peptide segments were classified in conformational space in conformational clusters. Nearest neighbour method was used with r.m.s.d. cutoff 1.2 Å for CA-atoms.

### 6. SRS table format:

#### CONFDB\_RES

Field name	Field description
ID	Peptide identifier
PR	Reference to protein
CR	Reference to CLUSTER table
SP	Index of start position in protein structure
SS	Index of start position in protein sequence
SQ	peptide sequence
PC	peptide coordinates for CA atoms (columns: SP;SS;Amino Acid 1-letter code; X coordinate; Y coordinate; Z coordinate.)

#### CONFDB\_PROT

Field name	Field description
ID	Protein identifier
DB	Reference to protein structure in PDB (format is ID:Chain)
SZ	Size of protein (number of peptides in this protein)

#### CONFDB\_CLUST

Field name	Field description
ID	Conformation cluster identifier
SZ	Size of cluster (number of peptides)
CM	Peptides - members of this class

Comments and questions to Dmitry Afonnikov (ada@bionet.nsc.ru).

### Example of SRS query to the ConfDB database

The query is 'To search for signal transduction protein (PDB identifier 1SEM), represented in ConfDB'.

To make such a query, you should perform the following:

- From ConfDB start page follow the link 'SRS\_ACCESS' (1).

- This will bring up the 'SRS top page'. This page is illustrated in the Figure. To perform database query, select CONFDB\_PROT database (1), input query string '1sem' into query field (2) and click 'Quick Search' button (3).

### III. Protein Integration Level. Chapter 1. Databases.

3. This will bring up the page illustrated in the figure. Click to the hyperlink to PROT3 entry.

The screenshot shows a web-based search interface with a navigation bar at the top: TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, DATABANKS, and HELP. Below the navigation bar, there is a search bar containing the query "Query "[confdb\_prot-ALLTEXT: 1sem\*]" found 1 entries". A red arrow points to the link labeled "CONFDB PROT.PROT3". On the left side, there is a sidebar with options for "Perform operation" (radio buttons for "on all but selected" and "on selected", with "on all but selected" selected), "Link", "Save", "View", and a dropdown menu set to "\* Names only \*". Below this is a section for "Number of entries to display per page" with a dropdown set to "30". At the bottom of the sidebar are "Printer Friendly" and "SRS 6.0.7.3 | feedback" links.

4. This will bring up the resulting page of the PROT3 entry.

The screenshot shows the detailed view of the PROT3 entry. The top navigation bar is identical to the previous screenshot. The main content area starts with "This entry is from: CONFDB PROT". It includes "Save" and "Link" buttons and a "Printer Friendly" link. To the right, there is a "view" button and a dropdown menu set to "\* Complete entries \*". The main body of the page displays the following text:  
CONFDB PROT.PROT3  
ID PROT3  
DB 1SEM:A  
SZ 50  
PM  
PEPT525 PEPT526 PEPT527 PEPT528 PEPT529  
PEPT530 PEPT531 PEPT532 PEPT533 PEPT534  
PEPT535 PEPT536 PEPT537 PEPT538 PEPT539  
PEPT540 PEPT541 PEPT542 PEPT543 PEPT544  
PEPT545 PEPT546 PEPT547 PEPT548 PEPT549  
PEPT550 PEPT551 PEPT552 PEPT553 PEPT554  
PEPT555 PEPT556 PEPT557 PEPT558 PEPT559  
PEPT560 PEPT561 PEPT562 PEPT563 PEPT564  
PEPT565 PEPT566 PEPT567 PEPT568 PEPT569  
PEPT570 PEPT571 PEPT572 PEPT573 PEPT574  
//

# CHAPTER 2. SOFTWARE

## 1. CRASP

### 1.1. Analysis of pairwise positional correlations

The program package CRASP, module P\_CRASP.

Release 2003

#### Program description:

The program package CRASP has been developed for analysis of co-ordinated substitutions of amino acid residues in aligned sequences of protein families. The module P\_CRASP is designed for evaluation of pairwise dependency of physico-chemical properties of amino acid residues.

#### Access to CRASP package:

<http://www.domain.com/mgs/gnw/crasp/> link 'Analysis of pairwise positional correlations'

#### List of biological tasks that could be solved by using the CRASP package

- evaluation of the mutual interdependency of physico-chemical property values in the pairs of positions of multiple alignment;
- searching for conserved (variable) physico-chemical properties of a protein

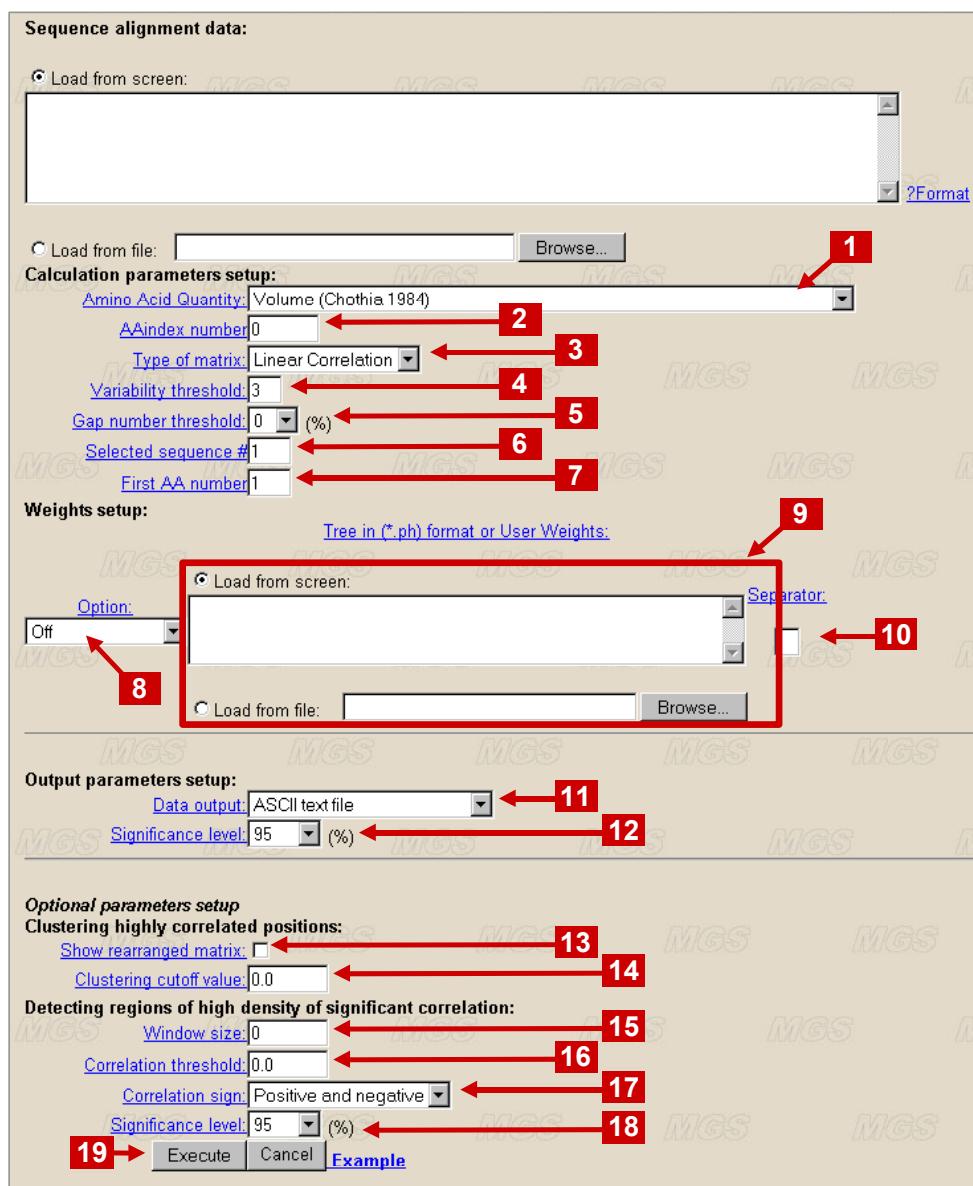
#### Data input

The data could be entered either from the browser screen in the appropriate text-box for data input (1) (option "Load from screen") or from the file (2) from the user's computer (option "Load from file"). The sequences should be aligned and entered in the FASTA format. There are no restrictions for the sequence length in alignment and the number of sequences.

The screenshot shows the CRASP software interface with three main sections:

- Sequence alignment data:** Contains two options:
  - Load from screen: A large text area for pasting sequence data, with a red box labeled 1 pointing to its top-left corner.
  - Load from file: A text input field and a "Browse..." button, with a red box labeled 2 pointing to the input field.
- Calculation parameters setup:** Includes:
  - Amino Acid Quantity: Volume (Chothia 1984)
  - AAindex number: 0
  - Type of matrix: Linear Correlation
  - Variability threshold: 3
  - Gap number threshold: 0 (%)
  - Selected sequence #: 1
  - First AA number: 1
- Output parameters setup:** Includes:
  - Data output: ASCII text file
  - Significance level: 95 (%)
  - Optional parameters setup:
    - Clustering highly correlated positions:
      - Show rearranged matrix:
      - Clustering cutoff value: 0.0
    - Detecting regions of high density of significant correlation:
      - Window size: 0
      - Correlation threshold: 0.0
      - Correlation sign: Positive and negative
      - Significance level: 95 (%)
  - Execute, Cancel, Example buttons.

### Program options.



(1) Choose one of 36 physico-chemical amino acid properties. See also (2).

(2) Choose amino acid properties by the number in the database AAIndex. Analysis of amino acid indices and mutation matrices for sequence comparison and structure prediction of proteins. The limits of this parameter range within 0-434. In case the parameter value is chosen to be equal to "0", then the choice of the property is determined by the menu (1). By default, the parameter equals to "0".

(3) Choose the type of the matrix to be calculated:

- 'Covariation' - covariation matrix calculation
- 'Linear correlation' - linear correlation coefficients calculation
- 'Partial correlation' - partial correlation coefficients calculation

By default, the option 'Linear correlation' is installed.

(4) Installation of the threshold of variability of positions. Here the minimal admissible number of different types of residues in a column of alignment is ordered. If the number of different types of residues is less than the ordered one, then such alignment column is excluded from calculation of a matrix. This parameter varies within the range 1-20. By default, the value “3” is installed.

(5) This option sets the threshold for the number of gaps in alignment column (expressed in percentage to the number of sequences). If the number of gaps exceeds the ordered value, then this column will be excluded from calculations of a matrix. This parameter is chosen from the drop-down menu within the range “0”-“15” %.

(6) This option sets the number of a sequence that will be the reference sequence in alignment. The reference sequence will be used for denoting positions of alignment. The admissible values vary from “1” to the number of sequences in alignment.

(7) This option orders the number of the first amino acid in alignment reference sequence. This value will be used for numeration when denoting alignment positions. Any whole positive numbers are accessible.

(8) This option is designed for choosing the regime of data weighting. See also description of the parameter (9).

- “Off” - all the sequence weights are equal to 1;
  - “Vingron & Argos” - the method suggested by Vingron and Argos (1985);
  - “User defined” - the weight coefficients are introduced by the user;
  - “Felsenstein” - the method is suggested by Felsenstein (1985) and its calculation is based on phylogenetic tree data. If you these data are available, this weighting scheme is recommended.
- By default, the option “Off” is installed.

(9) This option is designed for ordering weighting parameters if the option (8) is installed as “User defined” or for phylogenetic tree, if the option (8) is set as “Felsenstein”. The weights determined by user are entered in a text format with separator. The symbol of separator is ordered by the parameter (10). Phylogenetic tree is entered in a text form in the format PHYLIP (\*.ph). It is possible to enter data either from browser screen or from file.

(10) This option is for denoting the separator symbol. This parameter is installed in case the option (8) is set as “User defined”. It is possible to use any symbol, for example, “;”, “,”, “:”. By default, separator is set as the symbol “semicolon”.

(11) this option sets the output format of the matrix calculated:

- “ASCII text file” - ASCII text file (with HTML header);
- “HTML table” - HTML-table;
- “Matrix color diagram” - color diagram for matrix elements in GIF-format;
- “Significant pairs” - color diagram for statistically significant correlation coefficients in GIF-format (not defined for covariation matrix)

By default, the option “ASCII text file” is on.

(12) This option denotes the significance level of correlation coefficient that is chosen from menu.

(13) This is an additional option – graphical output of the matrix so that positions forming the cluster are closely located in the matrix. By default, this option is off. See also (14).

(14) This parameter is absolute value of correlation coefficient that determines the level of clusterisation. This parameter ranges within the limits “0” – “1”. If this parameter equals to “0”, then as the threshold is taken the critical value of correlation coefficient under given level of significance (parameter (12)).

(15) Parameters (15)-(18) determine visualisation of the matrix regions with predominance of correlation coefficients that exceed the threshold  $r_t$ . The matrix region is ordered as the rectangle window. The parameter (15) determines its size. In case this parameter is set as “0”, then this option is not fulfilled. The window size could be from “1” to the size of the matrix itself. By default, the value “0” is installed.

### III. Protein Integration Level. Chapter 2. Software

(16) This parameter sets the threshold  $r_t$  of correlation coefficient ranging within the limits “0” – “1”. If this parameter is set as “0”, then the threshold value is the critical value of correlation coefficient under significance level given (parameter (12)).

(17) This parameter sets the sign of correlation coefficients.

- “Only positive” – positive relationship only;
- “Only negative” – negative ones;
- “Positive and negative” – both sign;

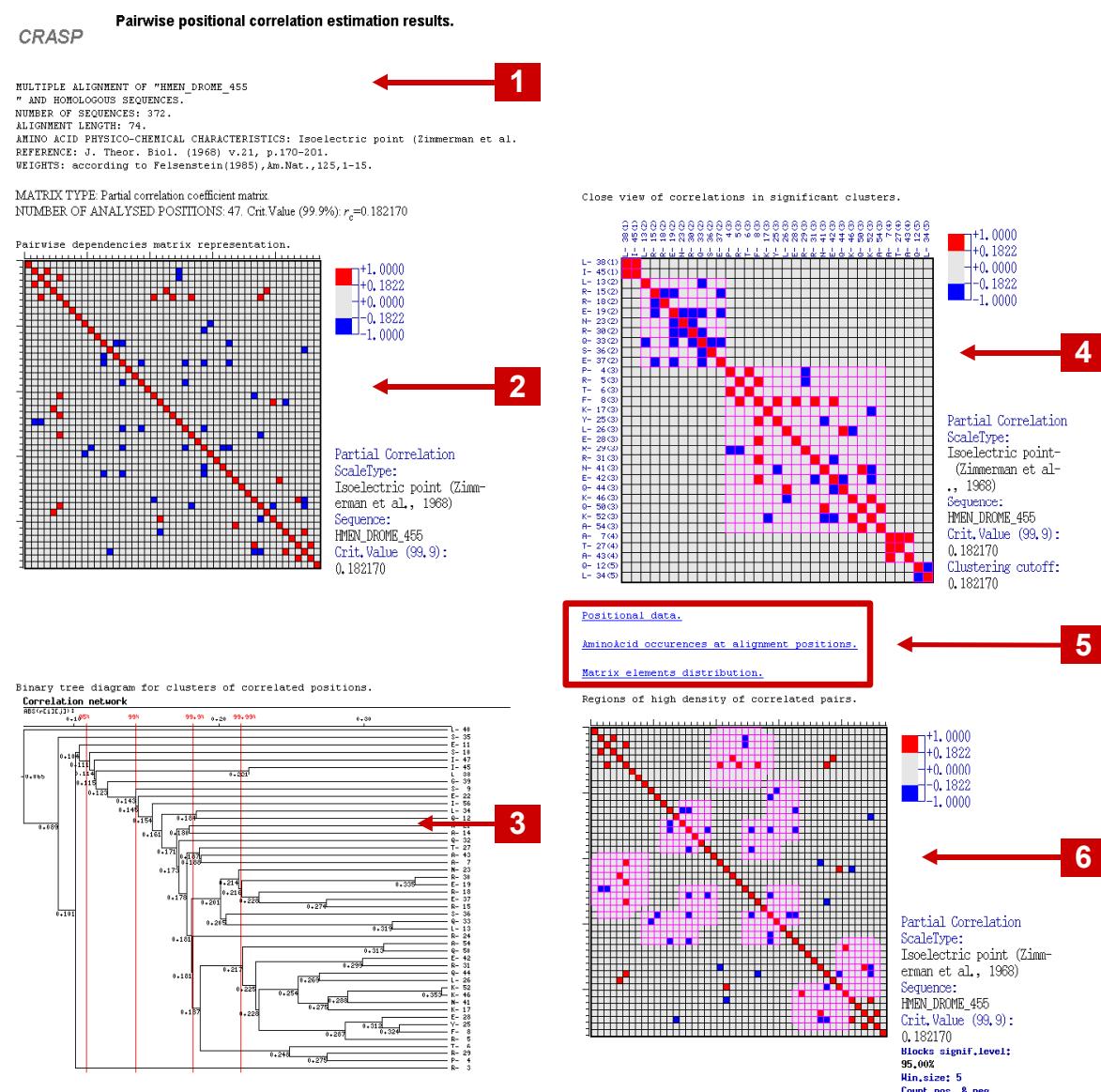
default option is “Positive and negative”.

(18) Significance level. This value should be chosen from the menu.

#### Program execution.

(19) For the program execution, click the button “Execute”.

#### Program output.



(1) This option is for displaying calculation parameters.

- (2) This option is for displaying correlation matrix.
- (3) This option is for displaying the diagram of hierarchical clusterisation of positions.
- (4) This option sets displaying of the ordered correlation matrix such that cluster positions in this matrix are closely located.
- (5) This option enables to display additional information: parameters of distribution of physico-chemical property in columns (Positional data), amino acid occurrences at alignment columns (AminoAcid occurrences at alignment positions), matrix elements distribution.
- (6) This option enables to display the regions of a matrix with predominance of significant correlation coefficients that are marked by violet colour.

**Example:**

The examples of data input for each program module could be displayed by the hyperlink (1):



**Comments and questions are welcome to Dmitry Afonnikov (ada@bionet.nsc.ru).**

## **1.2. Analysis of protein integral physico-chemical characteristics**

The program package CRASP, module F\_CRASP.

**Release 2003**

**Program description:**

The program package CRASP has been developed for analysis of co-ordinated substitutions of amino acid residues in aligned sequences of protein families. The module F\_CRASP evaluates the impact of co-ordinated substitutions into the constancy (or variability) of chemical and physical protein characteristics.

**Access to the CRASP package:**

<http://www.domain.com/mgs/gnw/crasp/> link 'Analysis of protein integral physico-chemical characteristics'

**List of biological tasks that could be solved by using the CRASP package:**

evaluation of the impact of co-ordinated substitutions into the constancy (or variability) of chemical and physical protein characteristics.

The program package CRASP consists of two modules. The first module (P\_CRASP) is designed for evaluation of the pairwise correlation of physico-chemical properties of the residues.

## Data input

The data could be entered either from the browser screen in the appropriate text-box for data input (1) (option “Load from screen”) or from the file (2) from the user’s computer (option “Load from file”). The sequences should be aligned and entered in the FASTA format. There are no restrictions for the sequence length in alignment and the number of sequences.

**Sequence alignment data**

Load from screen:  
[Empty text area] 1

Load from file: [File input field]  2

[?Format](#)

**Calculation parameters setup:**

- Amino Acid Quantity: Volume (Chothia 1984)
- AAindex number: 0
- Random samples number: 1000

**Weights setup:**

Tree in (\*.ph) format or User Weights:  
 Load from screen:  
Option: Off

Load from file: [File input field]

**Physico-chemical linear characteristics setup:**

Characteristic	Name	Description
F1	[Text input field]	
F2	[Text input field]	
F3	[Text input field]	
F4	[Text input field]	

**Output parameters setup:**

Data output: Text

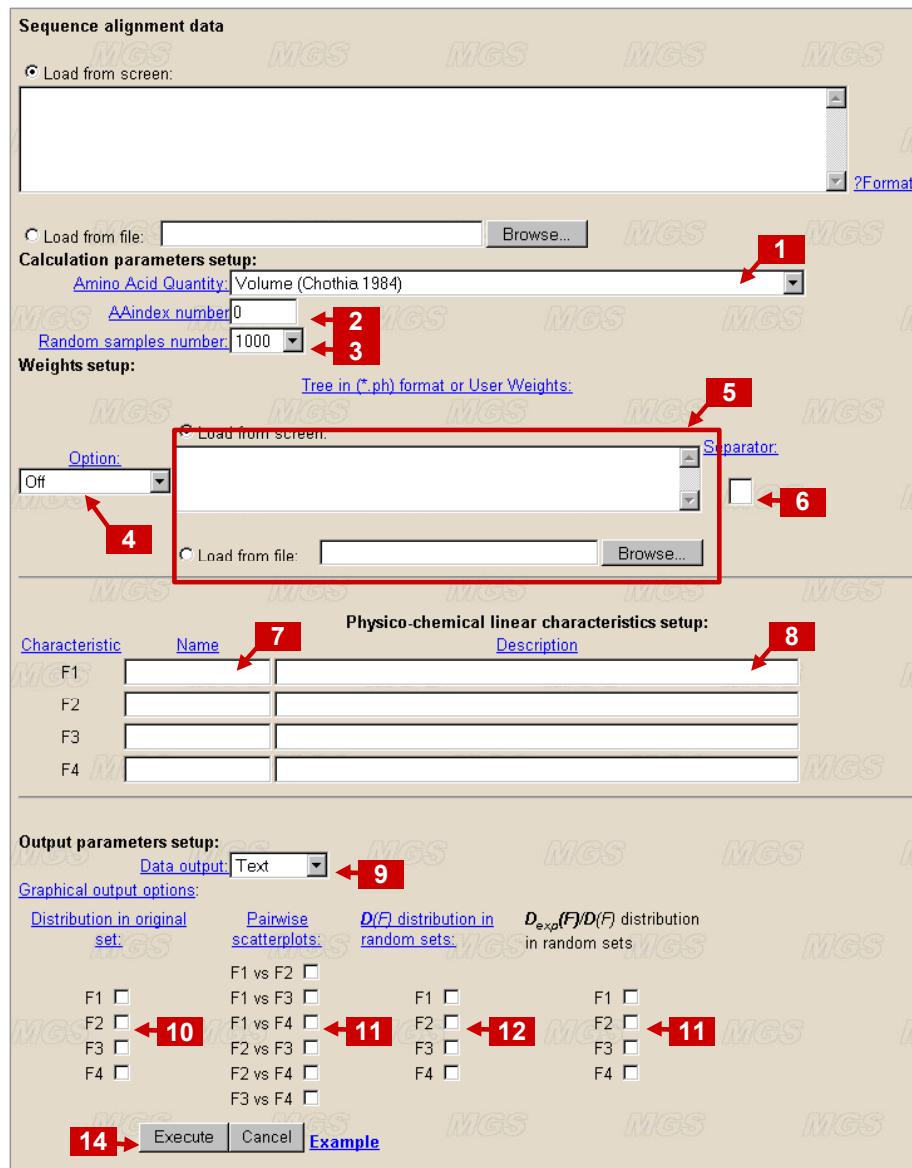
**Graphical output options:**

Distribution in original set:	Pairwise scatterplots:	D(F) distribution in random sets:	D <sub>exp</sub> (F)/D(F) distribution in random sets:
F1 <input type="checkbox"/>	F1 vs F2 <input type="checkbox"/>	F1 <input type="checkbox"/>	F1 <input type="checkbox"/>
F2 <input type="checkbox"/>	F1 vs F3 <input type="checkbox"/>	F2 <input type="checkbox"/>	F2 <input type="checkbox"/>
F3 <input type="checkbox"/>	F1 vs F4 <input type="checkbox"/>	F3 <input type="checkbox"/>	F3 <input type="checkbox"/>
F4 <input type="checkbox"/>	F2 vs F3 <input type="checkbox"/>	F4 <input type="checkbox"/>	F4 <input type="checkbox"/>
	F2 vs F4 <input type="checkbox"/>		
	F3 vs F4 <input type="checkbox"/>		

[Example](#)

## Program options.

### III. Protein Integration Level. Chapter 2. Software



(1) Choose one of 36 physico-chemical amino acid properties. See also (2).

(2) Choose amino acid properties by the number in the database AAindex (Tomii, K. and Kanehisa, M. (1996) Analysis of amino acid indices and mutation matrices for sequence comparison and structure prediction of proteins. *Protein Eng.* 9, 27-36.). The limits of this parameter range within 0-434. In case the parameter value is chosen to be equal to "0", then the choice of the property is determined by the menu (1). By default, the parameter equals to "0".

(3) This option determines the number of random sets generated for testing statistical significance of the impact of co-ordinated substitutions into the constancy (variability) of the conserved parameter. The value of this parameter is chosen from the pull-down menu.

(4) This option determines the regime of data weighting. See also description of the parameter (9).

- "Off" - all the sequence weights are equal to 1;
- "Vingron & Argos" - the method suggested by Vingron and Argos (1985);
- "User defined" - the weight coefficients are introduced by the user;
- "Felsenstein" - the method is suggested by Felsenstein (1985) and its calculation is based on phylogenetic tree data. If you these data are available, this weighting scheme is recommended.

By default, the option “Off” is installed.

(5) The parameters of data weighting should be entered if the option (8) is set as “User defined” or, for phylogenetic tree, if the option (8) is ordered as “Felsenstein”. The weights determined by a user are entered in a textual format with separators. The symbol of a separator is ordered by the parameter (10). Phylogenetic tree is entered in a textual form, in the PHYLIP (\*.ph) format. It is possible to enter the data either from screen, or from file.

(6) This option determines the symbol-separator, which is ordered if the option (8) is set as “User defined”. It is possible to use any symbol, i.e., “;”, “,”, “.”. By default, the symbol-separator is the “ semicolon”.

(7) This option determines the name of integral characteristics. The symbol line should be at most of 50 symbols.

(9) This option is designed for setting parameters of integral characteristics. To setup integral physico-chemical characteristics user should use the following format: “x1(npos1); x2(npos2); ...; xn(nposn);” xi - arbitrary numbers in a floating point format. nposi - corresponding positions of alignment (not a protein positions). They can be enumerated in an arbitrary form (using ',' and '-' symbols), i.e., “(1-3,4,5,30-44)” denotes positions from 1 to 5 and from 30 to 44. For example, net value of certain amino acid characteristic at the protein positions corresponding to alignment positions 6-8 is described as “1.(6-8);” Projection of alpha helical momentum (for helix positions 1 to 5 ) could be expressed as “1.(1); -0.17(2); -0.94(3); 0.5(4); 0.77(5);” where cos( $0^\circ$ )=1; cos( $100^\circ$ )=-0.17; cos( $200^\circ$ )=-.94; .... In total, up to 4 different integral protein characteristics related to one and the same physico-chemical property of residues (parameters (1) or (2)) could be ordered.

(9) This option orders the data output format:

- “Text” displays the textual data output;
- “Graphics” displays the resulted data as a set of images.

By default, the option “Text” is entered.

(10) This option enables to output information about distribution of integral characteristics value in the set analysed.

(11) This option enables to display information about dependency between two integral characteristics.

(12) This option enables to display information about distribution of dispersions of integral characteristics in random sets.

(13) This option enables to display information about the ratio between dispersion of characteristics in random sets and dispersion, which is expected in case the substitutions are independent (non-correlating).

### **Program execution.**

(14) For the program execution, click the button “Execute”.

### **Program output.**

### III. Protein Integration Level. Chapter 2. Software

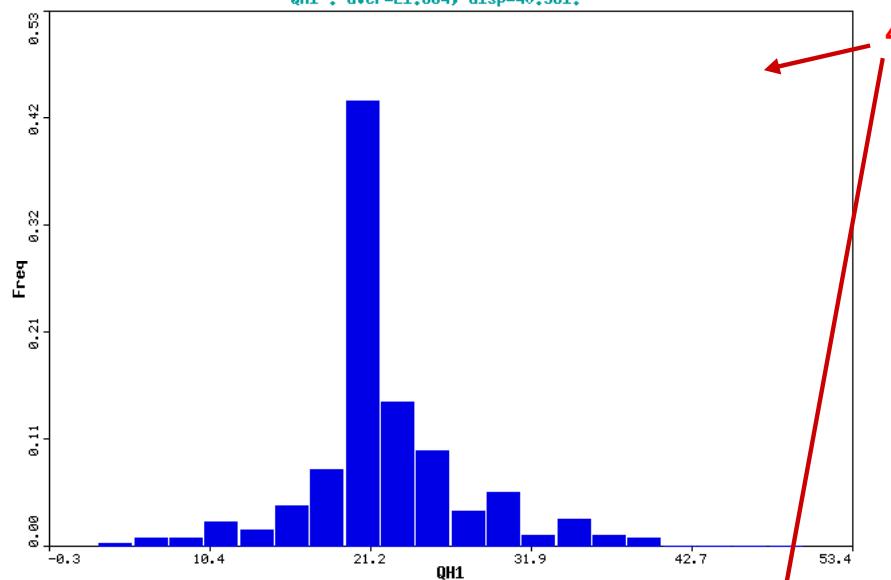
#### CRASP      Integral characteristic calculation results

MULTIPLE ALIGNMENT OF "HMEN\_DROME\_455" AND HOMOLOGOUS SEQUENCES.  
 NUMBER OF SEQUENCES: 372.  
 ALIGNMENT LENGTH: 74.  
 AMINO ACID PHYSICO-CHEMICAL CHARACTERISTICS: Isoelectric point (Zimmerman et al., 1968).  
 REFERENCE: J. Theor. Biol. (1968) v.21, p.170-201.  
 WEIGHTS: according to Felsenstein(1985), Am.Nat., 125, 1-15.  
 Mean(M) and variance(V):  
 1: QH1            21.68394    40.56061  
 2: QH2            26.49128    57.51518 ← 2  
 3: QH1+QH2      48.17522    80.70230  
 4: QH1-QH2      -4.80733    115.44930

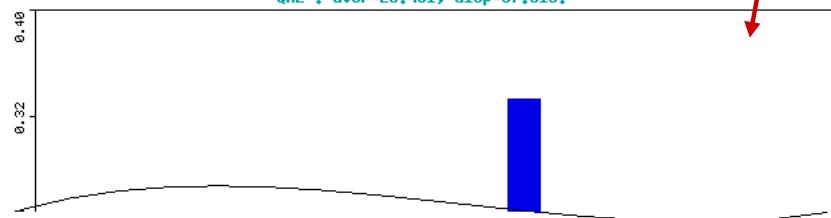
! WARNING:  
 ! All graphs corresponding to original sample data represent ← 3  
 ! contrasts calculated according to Felsenstein(1985), Am.Nat., 125, 1-15.

===== F(i) values distribution in original sample =====

Multiple alignment of "HMEN\_DROME\_455" and homologous sequences.  
 of sequences 372; Align. length 74. AA property: Isoelectric point (Zimmerman et al., 1968)  
 QH1 : aver=21.684; disp=40.561.



Multiple alignment of "HMEN\_DROME\_455" and homologous sequences.  
 of sequences 372; Align. length 74. AA property: Isoelectric point (Zimmerman et al., 1968)  
 QH2 : aver=26.491; disp=57.515.

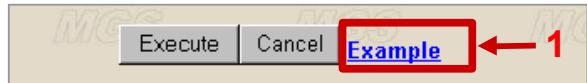


- (1) This option displays the output of main parameters of calculations.
- (2) This option displays the mean values and dispersions of characteristics ordered.
- (3) This warning becomes visible if the method of data weighting by Felsenstein(1985) was used.
- (4) Below the histograms and the scattering plots are displayed that were ordered by the parameters (10)-(13). In the heading of each plot, the alignment parameters are indicated (number of sequences and the length of alignment); type of the physico-chemical characteristics of amino acid residues;

mean and dispersion of characteristics considered; regression and correlation coefficients and the intervals of their significance.

**Example:**

The examples of data input for each program module could be displayed by the hyperlink (1) shown in the picture fragment below:



**Comments and questions are welcome to Dmitry Afonnikov (ada@bionet.nsc.ru).**

# PART IV. GENE NETWORKS INTEGRATION LEVEL

## CHAPTER 1. GENENET DATABASE

Release 2003

### 1. Database description:

The GeneNet database is designed to accumulate the information about structure and functional organisation of gene networks

### 2. Access to GeneNet database:

<http://www.domain.com/mgs/gnw/genenet/>

### 3. Database content

SRS table	Description	Number of entries
GN GENE	Genes	1006
GN RNA	RNAs	342
GN PROTEIN	proteins and their complexes	1766
GN SUBSTANCE	nonproteinaceous substances	241
GN RELATION	relationships between entities	3634
GN SCHEME	description of gene networks diagrams	37
GN SCHEME ENTITY	entities (elementary structures)	3530
GN SCHEME RELATION	relationships in gene networks	3634
GN COMPARTMENT	cell compartments	126
GN ORGANISM	species	93
GN PROCESS	input and output processes	128
GN CELL	cells, cell lines, tissues and organs	393
GN BIBLIOGRAPHY	references to the papers	1980
GN EXPERT	GeneNet annotators	32

### 4. List of biological tasks that could be solved by using the GeneNet database

- to extract the list of entities that are involved in functioning of a particular gene network and select the items by species, compartment, type of an entity;
- to extract the list of all reactions and regulatory relations for a particular gene network;
- to browse information about all relationships that involve the protein of interest;
- to extract the list of genes, transcription of which is induced by a particular transcription factor;
- to view reactions that involve a protein, as well as its role in them.

### 5. SRS table format

#### GN GENE

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
IC	IdCode	<Species abbreviation>:<gene abbreviation>
DT	EntryInfo	<date>; <annotator name>; <created/updated>
OS	Species	<Latin name> (<English name>)

SN	ShortName	Short gene name
NM	FullName	Full gene name
SY	Synonym	Synonym(s) of a gene name
SO	CellTableLink	Link to the GN CELL table <Cell IC>
CH	Chromosome	Chromosomal localisation
RE	Inducer Repressor	Inducer/repressor name
PN	Protein	Link to the GN PROTEIN table
DR	DatabaseLink	<database name>; <database first AC>; <database ID>
RF	LiteratureReference	Link to the GN BIBLIOGRAPHY table
CC	Comments	Comments

#### GN RNA

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
IC	IdCode	<Species abbreviation>;<RNA abbreviation>
DT	EntryInfo	<date>; <annotator name>; <created/updated>
OS	Species	<Latin name> (<English name>)
SN	ShortName	Short RNA name
NM	FullName	Full RNA name
SY	Synonym	Synonym(s) of a RNA name
RF	LiteratureReference	Link to the GN BIBLIOGRAPHY table
CC	Comments	Comments

#### GN PROTEIN

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
IC	IdCode	<Species abbreviation>;<protein abbreviation>
DT	EntryInfo	<date>; <annotator name>; <created/updated>
OS	Species	<Latin name> (<English name>)
SN	ShortName	Short protein name
NM	FullName	Full protein name
SO	CellTableLink	Link to the GeneNet CELL table <Cell IC>
FN	FunctionalState	Functional state <active/inactive/no data>
MM	Multimerization Level	Multimerization state <monomer/homodimer/heterodimer/multimer/no data>
MD	Modifications	Protein modifications <phosphorylated/non-phosphorylated/no data>
GN	Gene	Link to the GeneNet GENE table <Gene IC>
RF	Reference	Link to the GN BIBLIOGRAPHY table
CC	Comments	Comments

#### GN SUBSTANCE

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
IC	IdCode	Substance abbreviation
DT	EntryInfo	<date>; <annotator name>; <created/updated>
SN	ShortName	Short substance name
NM	FullName	Full substance name
CC	Comments	Comments

#### GN RELATION

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
DT	EntryInfo	Entry info: <date>; <annotator name>; <created/updated>
RE	Relation	Relation code <<type of the component entering into reaction>><component identifier>>^<component localisation>-> <<type of the component that is product of the reaction>><component identifier>>^<component localisation>

IN	Input	Input entity
TY	RelationType	Relation class <control/reaction>
OU	Output	Output entity
CO	RegTarget	Code of the controlled relation <<type of the component entering into reaction>><component identifier>^<component localisation> -> <<type of the component that is product of the reaction>><component identifier>^<component localisation>
EF	ReactionType	Reaction type <direct/indirect>
AT	InfluenceType	Type of regulatory event <increase/decrease/switch on/switch off>
RF	Reference	Link to the GN BIBLIOGRAPHY table
CC	Comments	Comments

#### **GN SCHEME**

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
NM	SchemeName	Gene network name
DT	EntryInfo	Entry info: <date>; <annotator name>; <created/updated>
MD	GeneNetDynamicModel	Link to the GeneNet dynamic model
TP	RepresentationLevel	Level of the gene network representation <cell/organism>
EN	EntityList	Link to the list of entities included in the gene network
RA	ReactionList	Link to the list of reactions in the gene network
RE	RegulatoryEventList	Link to the list of regulatory events in the gene network
DE	Description	Gene network description
OS	Species	<Latin name> (<English name>)
RR	Reference	Link to the GN BIBLIOGRAPHY table
DR	DatabaseLink	Links to other databases

#### **GN SCHEME ENTITY**

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
SC	SchemeId	Link to the GN SCHEME table
NM	SchemeName	Gene network name
ET	EntityType	Entity type <gene/process/protein/RNA/Substance>
OS	Species	<Latin name> (<English name>)
EN	EntityName	Entity name
SU	Compartment	Compartment localisation

#### **GN SCHEME RELATION**

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
SC	SchemeId	Link to the GN SCHEME table
NM	SchemeName	Gene network name
RT	RelationType	Relation type <reaction/regulatory event>
IN	Input	Input entity
OU	Output	Output entity
CO	RegTarget	Code of the controlled relation <<type of the component entering into reaction>><component identifier>^<component localisation> -> <<type of the component that is product of the reaction>><component identifier>^<component localisation>

#### **GN COMPARTMENT**

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
IC	IdCode	Compartment code
AC	IdCode2	Compartment name
DT	EntryInfo	Entry info: <date>; <annotator name>; <created/updated>
CC	Comments	Comments

### GN ORGANISM

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
IC	Species Latin	Species Latin name
AC	Species English	Species English name
SN	Abbreviation	Species abbreviation
OC	Classification	Classification

### GN PROCESS

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
IC	IdCode	Process code
DT	EntryInfo	Entry info: <date>; <annotator name>; <created/updated>
SN	ShortName	Process name
NM	FullName	Complete process name
CC	Comments	Comments

### GN CELL

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
IC	IdCode	<Species abbreviation>;<item abbreviation>
DT	EntryInfo	Entry info: <date>; <annotator name>; <created/updated>
OS	Species	<Latin name> (<English name>)
SN	ShortName	Abbreviated name
NM	FullName	Complete name
SY	Synonym	Synonymous name
RF	Reference	Link to the GN BIBLIOGRAPHY table
CC	Comments	Comments

### GN BIBLIOGRAPHY

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
IC	IdCode	<Species abbreviation>;<item abbreviation>
AU	Author	Authors of the paper
TI	Title	Title of the paper
SO	Journal	Journal name
VL	Volume	Volume
IS	Issue	Issue
YR	Year	Year
PG	Pages	Pages
ML	Medline UI	Link to the MEDLINE

### GN EXPERT

Line code	Field name	Field description
ID	Identifier	GeneNet identifier
IC	IdCode	Abbreviation of the annotator
NM	Name	Name of the annotator
LB	Laboratory	Laboratory
OR	Organization	Organization
CO	Country	Country
EM	EMail	E-mail address

### Example 1

To search for all relationships described in the GeneNet database, in which IRF gene or protein participate.

To make such a query, you should perform the following:

- 1 Choose GeneNet SRS table (RELATION) on the page 'SRS access'.

The screenshot shows the GeneNet website interface. At the top, there is a navigation bar with links for HOME, DNA, RNA, PROTEIN, GENENETWORKS, and MAP. Below the navigation bar, there is a logo for 'GENENET' with a yellow spider icon. To the right of the logo, a section titled 'GeneNet DATABASE (SRS)' displays a hierarchical tree diagram of biological concepts. A red arrow points to the 'RELATION' node in this tree. Other nodes include ORGANISM, COMPARTMENT, NUCLEUS, CELL, GENE, RNA, PROTEIN, SUBSTANCE, PROCESS, hypoxia, SCHEME, SCHEME\_ENTITY, SCHEME\_RELATION, BIBLIOGRAPHY, and EXPERT. To the right of the tree diagram, there is a magnifying glass icon over a circular network visualization labeled 'GeneNet VIEWER'. Below the viewer, there is a blue oval containing the 'GeneNet' logo and the text 'System for formalized description, visualization, and modelling of gene networks'. At the bottom right, there is a section titled 'GeneNet MODELLING' featuring a circular network visualization and two line graphs labeled 'Oncoprotein interaction model' and 'Growth hormone release by insulin'.

2. This will bring up the home page of the chosen GeneNet SRS table (RELATION). Click the button 'Search'.

The screenshot shows the 'GN\_RELATION' table homepage. At the top, there is a navigation bar with links: TOP PAGE, QUERY, RESULTS, SESSIONS, VIEWS, DATABANKS, and HELP. Below the navigation bar, there is a sidebar with sections: Home, Status, Description, Literature, and WWW. The 'WWW' section contains a link to the 'WWW site'. The main content area displays the table structure and data. The table has columns: Name, Short Name, Type, No of Keys, No of References, Indexing Date, and Status. The data rows are:

Name	Short Name	Type	No of Keys	No of References	Indexing Date	Status
<a href="#">Identifier</a>	<a href="#">id</a>	id	3254	3254	05-Jan-2003	ok
<a href="#">EntryInfo</a>	<a href="#">dt</a>	index	441	11139	05-Jan-2003	ok
<a href="#">Relation</a>	<a href="#">re</a>	index	3242	3254	05-Jan-2003	ok
<a href="#">Input</a>	<a href="#">in</a>	index	1810	6508	05-Jan-2003	ok
<a href="#">RelationType</a>	<a href="#">ty</a>	index	2	3254	05-Jan-2003	ok
<a href="#">Output</a>	<a href="#">ou</a>	index	1740	2422	05-Jan-2003	ok
<a href="#">RegTarget</a>	<a href="#">co</a>	index	523	832	05-Jan-2003	ok
<a href="#">ReactionType</a>	<a href="#">ef</a>	index	2	3254	05-Jan-2003	ok
<a href="#">InfluenceType</a>	<a href="#">at</a>	index	4	1760	05-Jan-2003	ok
<a href="#">Reference</a>	<a href="#">rf</a>	index	633	1276	05-Jan-2003	ok
<a href="#">Comments</a>	<a href="#">CC</a>	show	0	0		not indexed

3. Select the mode 'OR' from the list of the combine searches (1).

Select the fields to be searched for from the lists (2). You will need the following fields: 'Input', 'Output', 'RegTarget'.

Type term \*IRF\* to be searched in the text windows (3). (An asterisk marks "any symbol")

Then choose the parameter 'Complete entries' for displaying the query results (4) and click the button 'Submit Query' (5).

The screenshot shows the GeneNet QUERY interface. The search bar contains "search GN\_RELATION". Below it, the query string is entered as "[gn\_relation-Input: \*IRF\*] || [gn\_relation-Output: \*IRF\*] || [gn\_relation-RegTarget: \*IRF\*]". The interface includes several dropdown menus and input fields:

- Submit Query** button (1)
- Input dropdown set to "Input" with value "\*IRF\*" (2)
- Output dropdown set to "Output" with value "\*IRF\*" (3)
- RegTarget dropdown set to "RegTarget" with value "\*IRF\*" (3)
- Identifier dropdown (4)
- Search parameters: "append wildcards to words" checked, "combine searches with OR" selected (1), "Number of entries to display per page" set to 30.
- View options: "Use predefined view" set to "\* Complete entries \*" (4), "Create your own view", and a list of fields to display: Identifier, EntryInfo, Relation, Input, RelationType.

4. The query result is shown below.

The screenshot shows the GeneNet RESULTS interface displaying the query results for the search performed in the previous step. The results are listed under three sections:

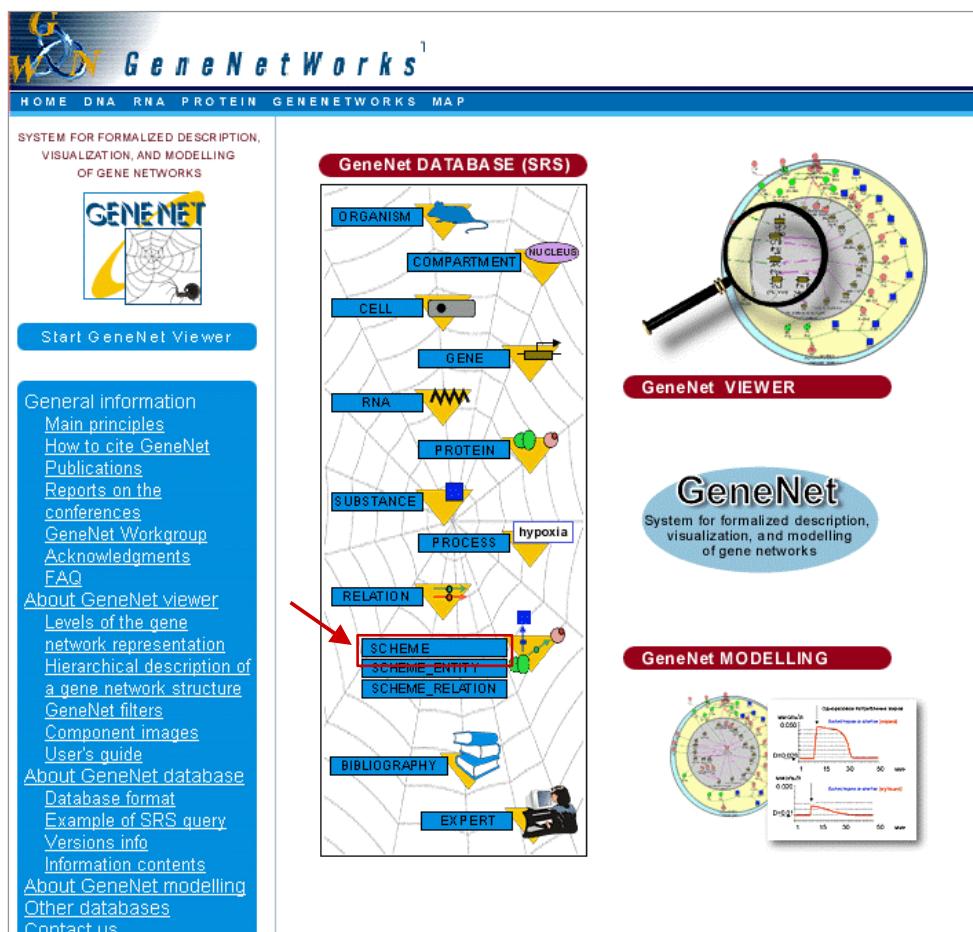
- GN\_RELATION:R127**: ID R127, IN <RNA>[Rn:IRF-1^nucleus](#), TY Reaction, OU <protein>[Rn:IRF-1^nucleus](#), DT 14-Dec-01.; [Suslov V.V.](#); created, EF direct, //
- GN\_RELATION:R345**: ID R345, IN <gene>[Hs:IRF-1^nucleus](#), TY Reaction, OU <protein>[Hs:IRF-1^nucleus](#), DT 06.4.1999.; [Ananko E.](#); created, EF indirect, //
- GN\_RELATION:R346**: ID R346, IN <gene>[Hs:IRF-2^nucleus](#), TY Reaction, OU <protein>[Hs:IRF-2^nucleus](#), DT 30.3.1999.; [Ananko E.](#); created, EF indirect, //

## Example 2

To visualise description of the diagram 'Macrophage activation' and the list of all entities and relationships entering this diagram.

To make such a query, you should perform the following:

- 1 Choose GeneNet SRS table (SCHEME) on the page 'SRS access'.



2. In the Query Form for GN\_SCHEME select the field 'SchemeName' to be searched for from the list (1).

Type term \*macrophage\* to be searched in the text windows (2). (An asterisk marks "any symbol").

Then choose the parameter 'Complete entries' for displaying the query results (3) and click the button 'Submit Query' (4).

The screenshot shows the 'QUERY' tab selected in the top navigation bar. The main search bar contains 'GN SCHEME' and 'search GN SCHEME'. A red box labeled '1' points to the search bar. A red box labeled '2' points to the input field containing '\*macrophage\*'. A red box labeled '3' points to the dropdown menu 'Complete entries'. A red box labeled '4' points to the 'Submit Query' button.

3. This will bring up the description of the gene network of macrophage activation in the GeneNet database.

By clicking the links 'List of Entities', 'List of Reactions', 'List of Regulatory Events' (1), you may extract the complete lists of Entities, Reactions, and Regulatory Events for this gene network. Also, here are the links to GeneNet viewer (2).

The screenshot shows the results of the query. The top message says 'Query "[gn\_scheme-SchemeName: \*macrophage\*]" found 1 entries'. A red box labeled '2' points to the link 'GN\_SCHEME:SCH17'. Three red boxes labeled '1' point to the links 'List of Entities', 'List of Reactions', and 'List of Regulatory Events' listed under the entry details. The entry details themselves are as follows:

ID SCH17  
**GeneNet viewer:** [Macrophage activation \(model\)](#)  
 DT 29.10.2000; [Nedosekina E.A.](#); created.  
 DT 19.6.2002.; [Ananko E.](#); updated.  
 DT 11.12.2002; [Nedosekina E.A.](#); updated.  
 TP cell  
 EN [List of Entities](#) 1  
 RA [List of Reactions](#) 2  
 RE [List of Regulatory Events](#) 3  
 DE Macrophage activation by the IFN-gamma and lipopolysaccharides ()  
 OS [Gallus gallus](#) (chicken)  
 OS [Homo sapiens](#) (human)  
 OS [Mus musculus](#) (mouse)  
 OS [Rattus norvegicus](#) (rat)  
 RR [Aijan R.A. et al., 1998a](#)  
 RR [Amura C.R. et al., 1998](#)  
 RR [Aoudjit F. et al., 1994](#)  
 RR [Arend W.P., 1997](#)  
 RR [Arias-Negrete S. et al., 1995](#)  
 RR [Barrios-Rodiles M. and Chadee K., 1998](#)  
 RR [Bayon Y. et al., 1998](#)  
 RR [Briken V. et al., 1995](#)  
 RR [Byrnes A.A. et al., 2001](#)  
 RR [Celada A. et al., 1996](#)  
 RR [Chan E.D. et al., 1999](#)  
 RR [Chen F. et al., 1995](#)  
 RR [Chen F. et al., 1995a](#)  
 RR [Chen Y.Q. et al., 1998](#)  
 RR [Cockerill P.N. et al., 1995](#)  
 RR [Cockerill P.N. et al., 1996](#)  
 RR [Cottin V. et al., 1999](#)

Comments and questions are welcome to Elena Ananko (eananko@bionet.nsc.ru).

## CHAPTER 2. GENENET SOFTWARE

### 1. GeneNet Viewer.

Release 2003

#### 1. Program description:

Software for graphical representation of the information stored in the GeneNet database

#### 2. Access to GeneNet Viewer

[http://www.domain.com/mgs/gnw/systems/genenet/applet\\_genenet\\_viewer.shtml](http://www.domain.com/mgs/gnw/systems/genenet/applet_genenet_viewer.shtml)

#### 3. List of biological tasks that could be solved by using the GeneNet Viewer:

- ◆ To display the general visualisation of the structure and functional organisation of the gene network described in the GeneNet database;
- ◆ To view, which data were obtained for particular species or organisms.

#### 4. Data input

Data input is available only via the special program Data Input GUI (<http://www.domain.com/mgs/gnw/systems/genenet/>). To use this program via the Internet, it is necessary to register by the address <http://wwwmgs.bionet.nsc.ru/mgs/systems/genenet/registration.html>.

#### 5. Program options

##### System of filters

The data obtained from different species are summarised in the scheme of the gene network. As a result, the diagram may contain several equivalent objects (for example, homologous genes or proteins from different species). The equivalent components of the gene network are displayed as a single image.

By default, the gene network diagram is drawn from all the information given in the corresponding entry from the scheme table.

A system of filters enable a user to select for visualisation only the entities and relations that have been experimentally identified for a specified organism or specified cell types as well as those specific to the cell response to specified external stimuli. For this purpose, the GeneNet is provided with filters of three types according to:

1. Species
2. Cell type
3. Inducer

All the three filters can be used at the same time. As a result, only the objects that meet the requirements of all the filters, will be displayed in the diagram.

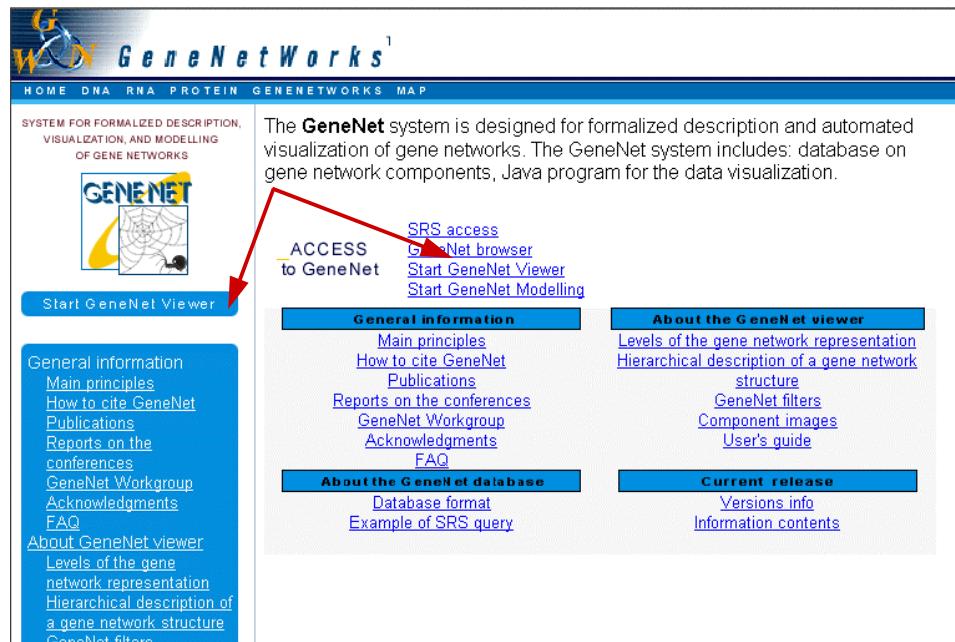
The system of filters is available via the menu 'Filters' (see Example).

##### Zooming

Zooming of a diagram is ordered through the menu 'Diagram', 'Zoom" (see Example).

## 6. Program execution

The program is executed by clicking the respective links at the GeneNet Home Page (they are marked by red arrows in the figure), as well as from SRS entries in GN\_SCHEME table via the special link in the field 'NM GeneNet viewer.'



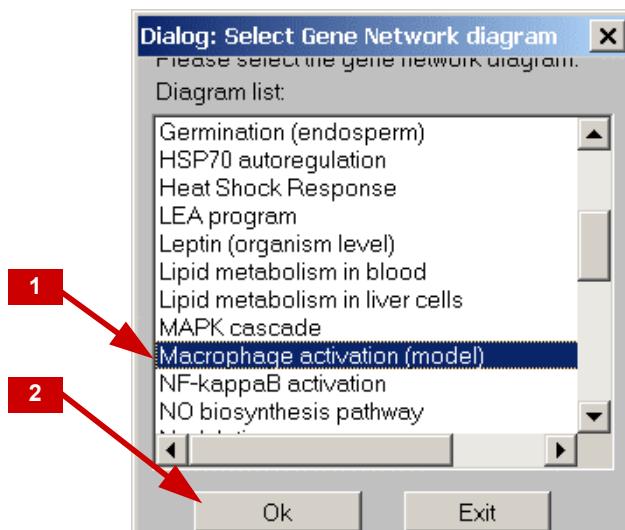
## 7. Data output

Data are displayed as a graphical diagram of structure and functional organisation of a gene network.

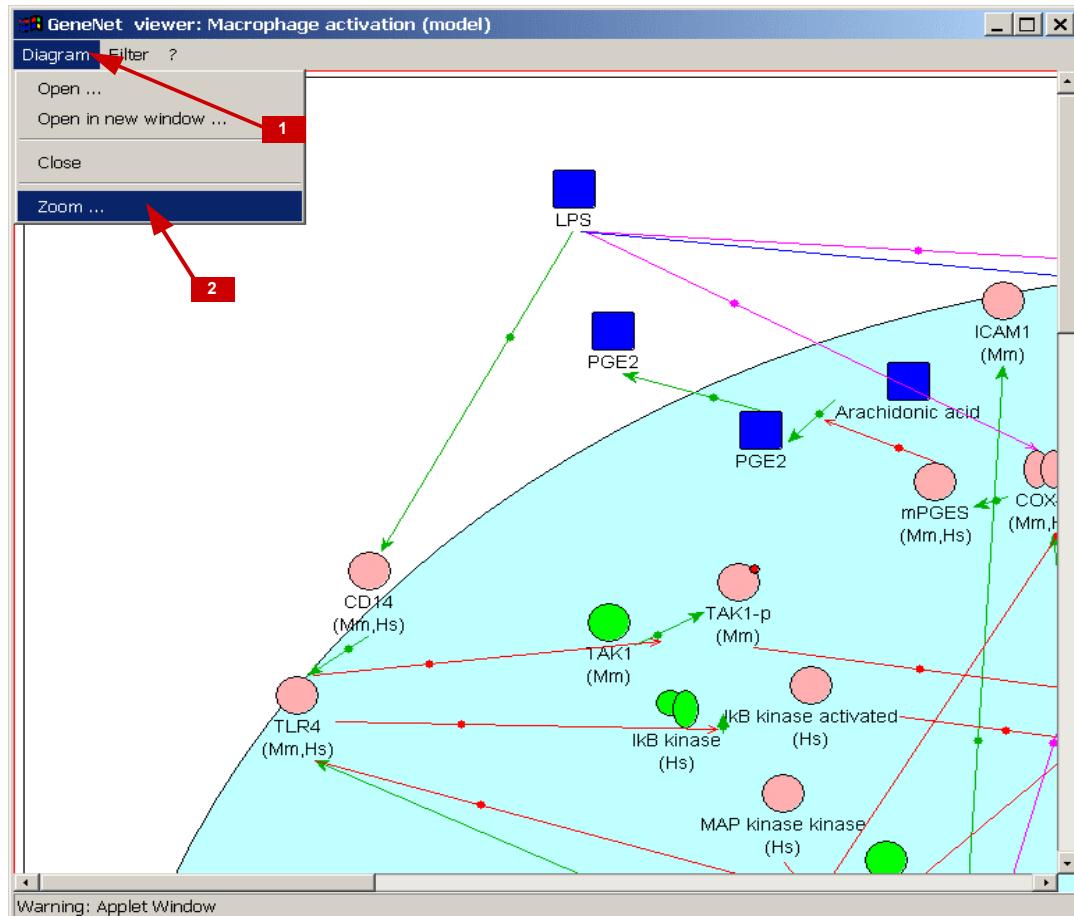
### Example

Visualisation of the diagram 'Macrophage activation (model)' by means of GeneNet Viewer

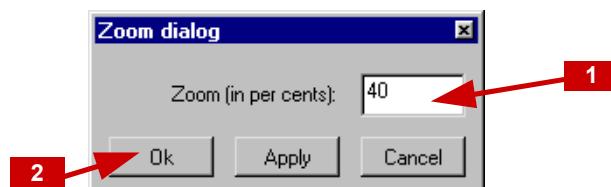
1. Start GeneNet Viewer. After data loading into the dialog window, select the option 'Macrophage activation (model)' (1) and click the button 'OK' (2)



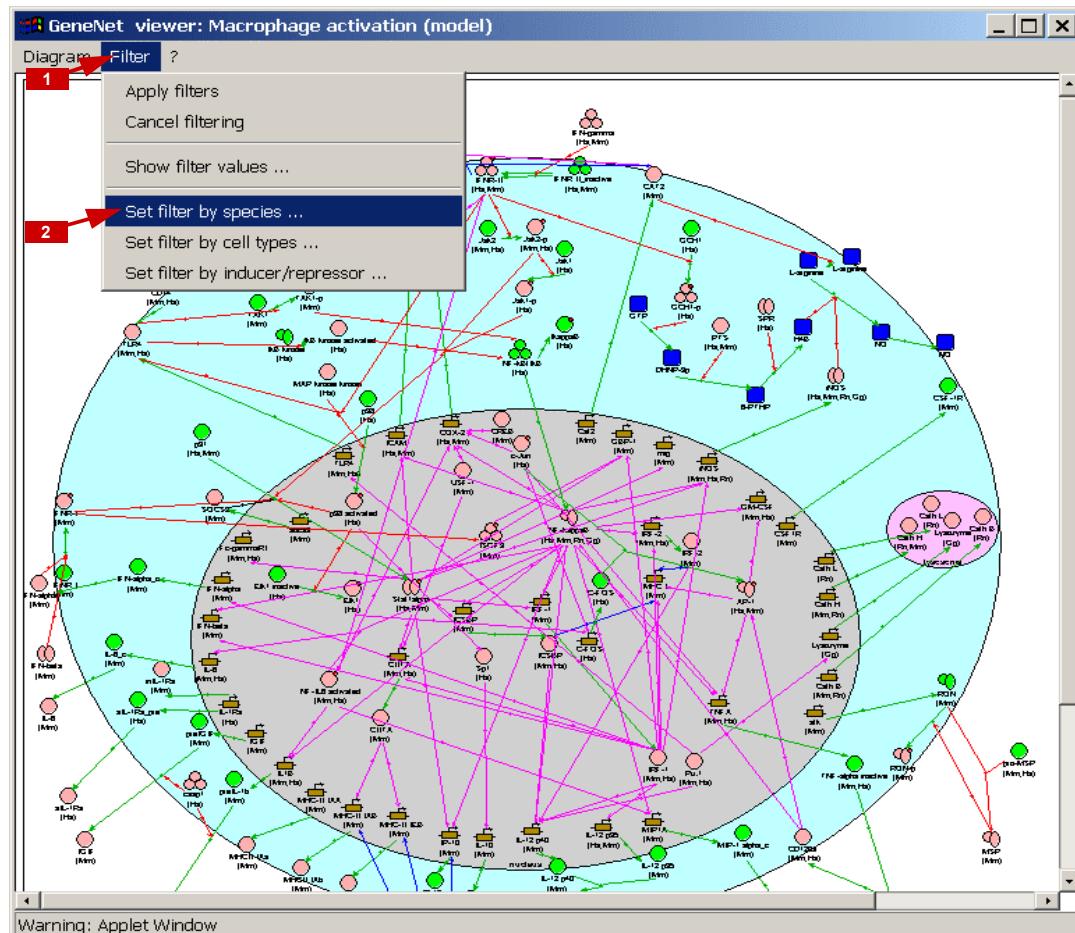
2. This will bring up the graphical diagram of the gene network Macrophage activation (model). In order to modify the zoom of the diagram, select in the menu 'Diagram' (1) the option 'Zoom' (2).



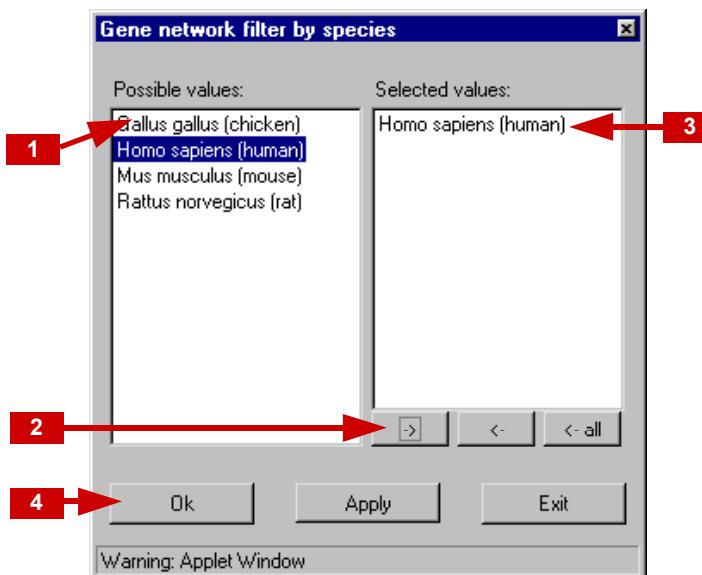
3. Enter the percentage value '40' into the textual window (1) and click the button 'OK' (2).



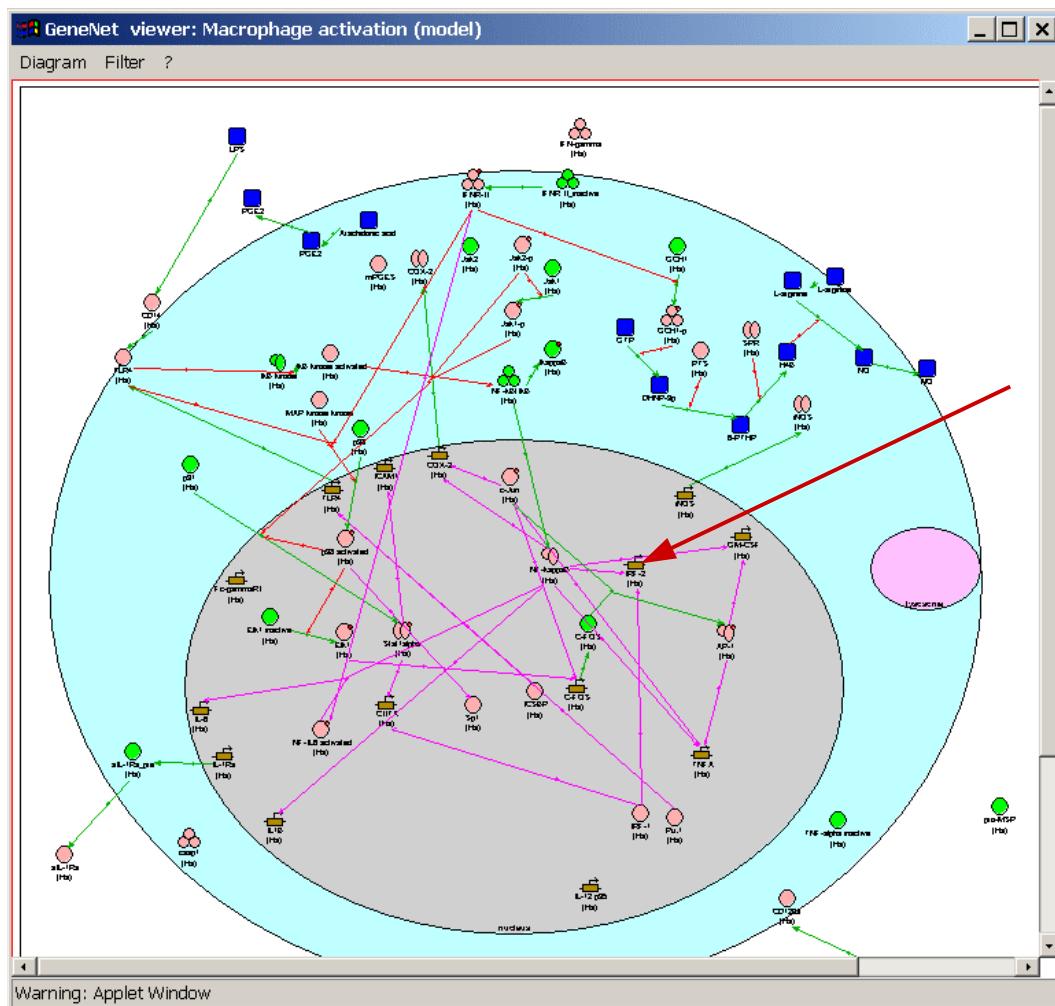
4. This will bring up the window with the altered zooming of the diagram. In order to display the visualisation of data obtained only for human cells, select in the menu 'Filter' (1) the value 'Set filter by species' (2).



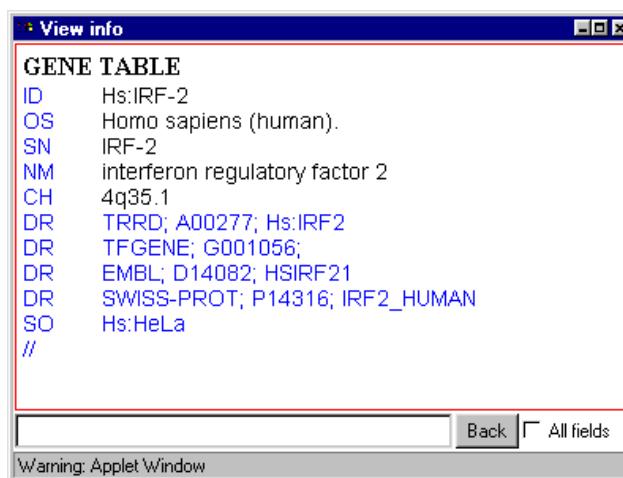
5. In the dialog window for setting the filters, choose the value 'Homo sapiens' (1), click the button ' $\rightarrow$ ' (2). This will bring up the option 'Homo sapiens' in the window from the right (3). Click the button 'OK' (4).



6. This will bring up the diagram displaying the data obtained by experiments only with human cells. The data obtained for the cells of the other organisms will disappear from the diagram. To display textual description of an object in the GeneNet database, click this object by mouse twice.



7. Textual description of the chosen object will be displayed in the special text window, in which the hyperlinks to the other databases are also available.



**Comments and questions are welcome to Elena Ananko (eananko@bionet.nsc.ru).**

## **2. GeneNet Modelling**

**Release 2003**

### **1. Program description:**

GeneNet Modeling is a software for operating with mathematical models of gene networks dynamics. The current version of the GeneNet Modelling operates with three models:

- ◆ Model of cholesterol biosynthesis regulation
- ◆ Model of erythrocyte maturation regulation
- ◆ Model of macrophage activation under the action of interferon gamma and/or lipopolysaccharides (LPS).

### **2. Access to GeneNet Modeling**

[http://www.domain.com/mgs/gnw/gn\\_model/](http://www.domain.com/mgs/gnw/gn_model/)

### **3. List of biological tasks that could be solved by using the GeneNet Modeling:**

- ◆ viewing the reaction of a gene network in response to alteration of initial concentration values of the system's components;
- ◆ viewing the reaction of a gene network in response to alteration of the values of concentrations of the system's components during fixed periods of time;
- ◆ viewing the gene network behavior in response to alteration of parameters determining the rates of elementary processes;
- ◆ computer simulation of gene mutation (disruption of a gene or enzyme function, alteration of rates of translation, transcription, transportation, utilisation processes, etc.) by varying parameters of a gene network mathematical model.

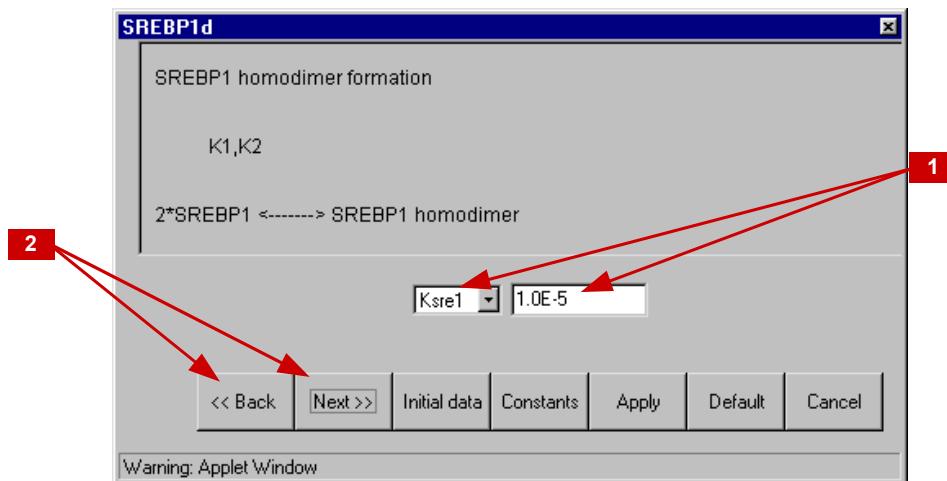
### **4. Data input**

Novel structure components and interactions might be added into existing model only by means of the special program Data Input GUI (<http://wwwmgs.bionet.nsc.ru/mgs/systems/genenet/>), which could be used via Internet after registration by the address <http://wwwmgs.bionet.nsc.ru/mgs/systems/genenet/registration.html>.

### **5. Program options**

In the regime of the graphical interface it is possible to modify reaction rates and initial concentrations of some components. By clicking by mouse the component chosen, user displays the dialog window, in which the parameters available for a given gene network component could be modified (1). It is possible to view all the parameters of this component available for alteration by clicking the buttons 'Next' and 'Back' (2).

Note: if the dialog window is not open by clicking the object, this means that alteration of parameters of this object is prohibited for user.

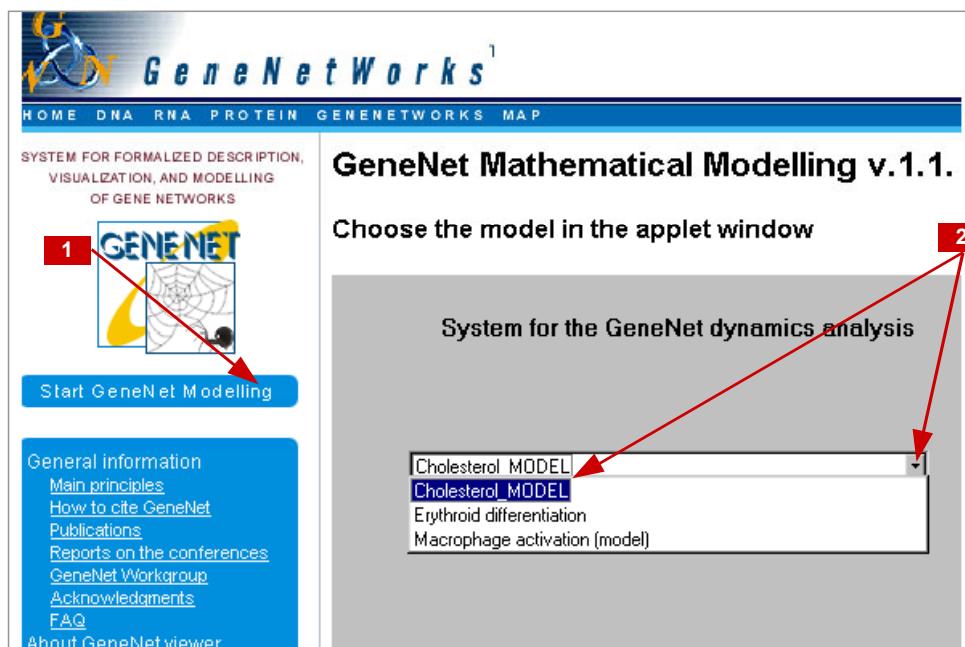


## 6. Program execution

Click 'Start GeneNet Modelling' (1)

This will be your working window of the GeneNet Modelling package.

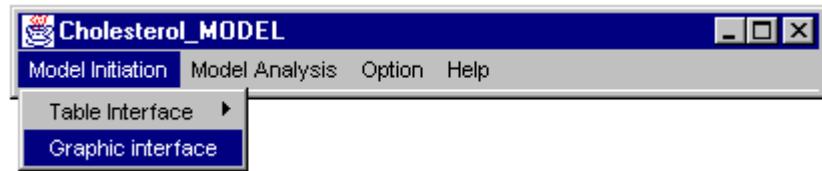
Choose the model to operate with (2).



If you click the option 'Cholesterol\_MODEL' in the drop-down menu 'Select a model', this will bring up the starting page of cholesterol biosynthesis gene network modeling.



Choose on the menu bar the options 'Model Initiation' and 'Graphic interface' to initiate the graphical interface of the model of cholesterol biosynthesis.



## 7. Data output

The data could be output both in graphical representation and in a form of a table (see Example).

### Example 1

Operating with dynamic model of inner cellular cholesterol concentration.

Use it in parallel with the tutorial window. Tutorial pages will guide you how to work with mathematical models of gene networks.

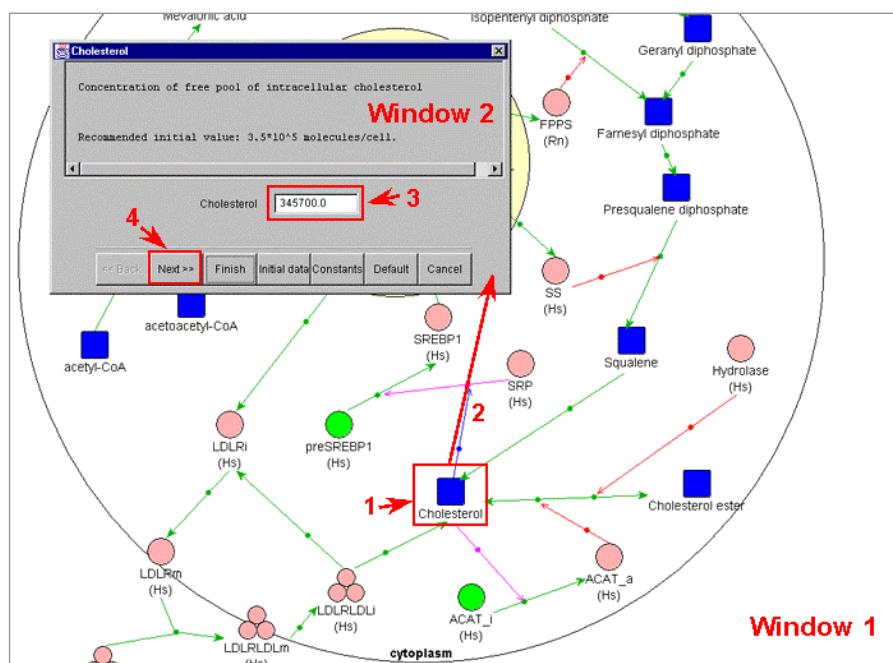
#### 1. Modification of parameters of a component.

Select the component 'Cholesterol' at the diagram (1).

By clicking the component you will get the dialog window (2). You may change the value of initial concentration (3) of the component selected.

If you want to restore initial values that are input by default, click the button 'Default'.

For modifying the next parameter, click the button 'Next' (4).



2. This brings up the dialog window with description of reaction(s) in which the selected component is involved. You may change the reaction constant(s) value(s).

If you would like to return to initial values of the constants in this window installed by default, click the button 'Default' (1).

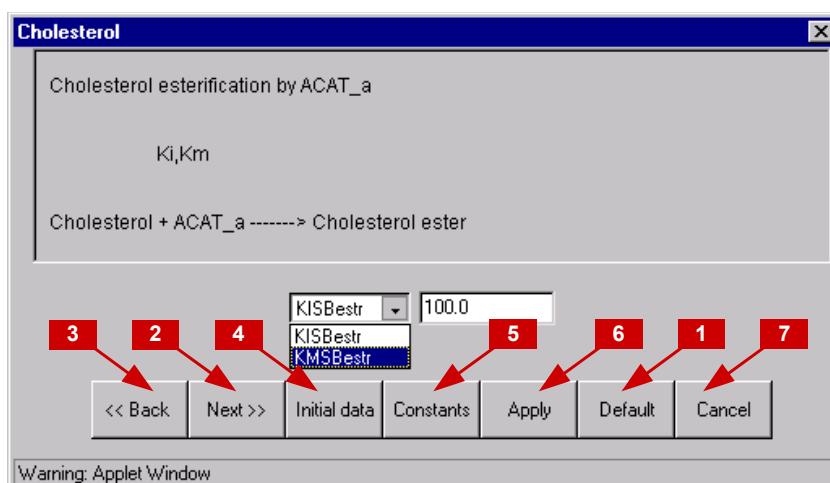
After finishing the editing of reaction, click the button 'Next' (2) if you want to edit the following reaction in which the selected component is involved. For returning to preceding window, click the button 'Back' (3).

If you want to return the window for editing initial values, click the button 'Initial data' (4).

If you want to edit constant(s), click the button 'Constants' (5).

After finishing the editing of initial value(s) and constants of reactions in which the selected component is involved, click the button 'Apply' (6) to save the changes.

Click the button 'Cancel' (7) to return back to initial data.

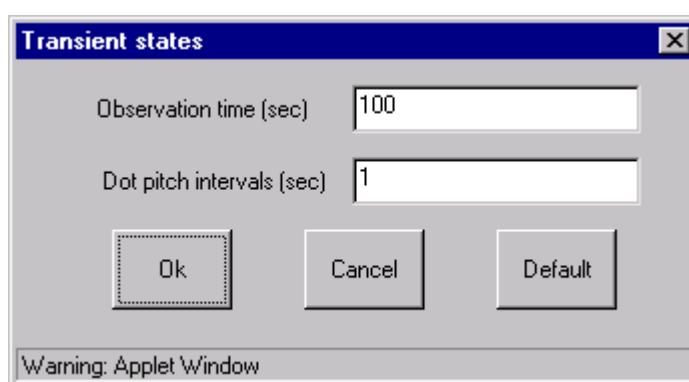


Analogously, user may change the parameters (initial values and reaction constants) of the other components involved in modeling of the gene network.

3. Go to the model window and click from the drop-down menu 'Model Analysis' the view 'Transient states'.



This will bring up the dialog window for editing the observation time and dot pitch of the output data.



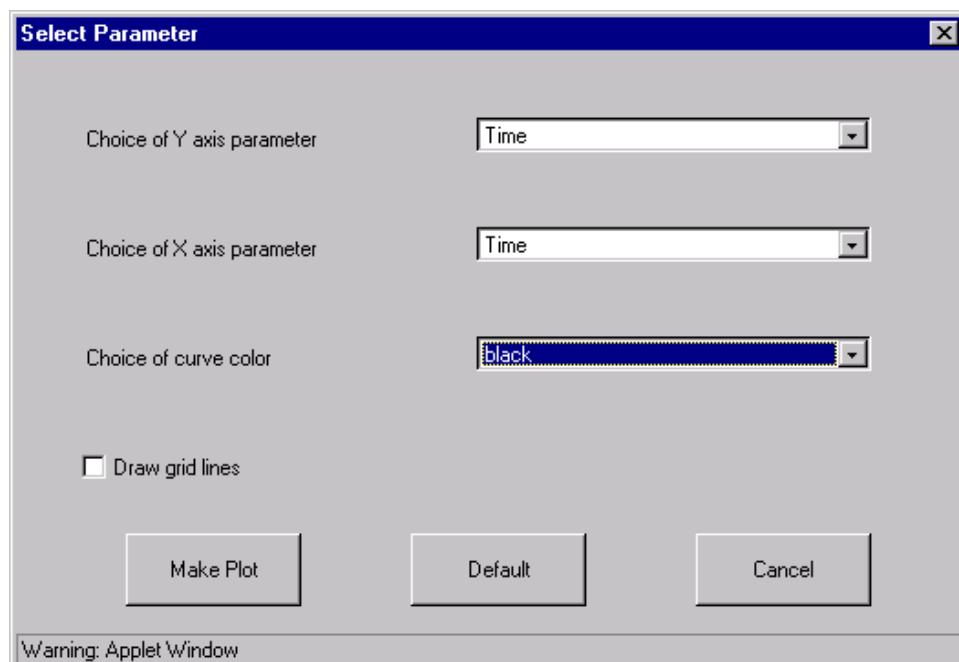
For default settings, click the button 'Default'.

To go back to editing of reaction constants and (or) initial values of components of mathematical model of a gene network, click the button 'Cancel'.

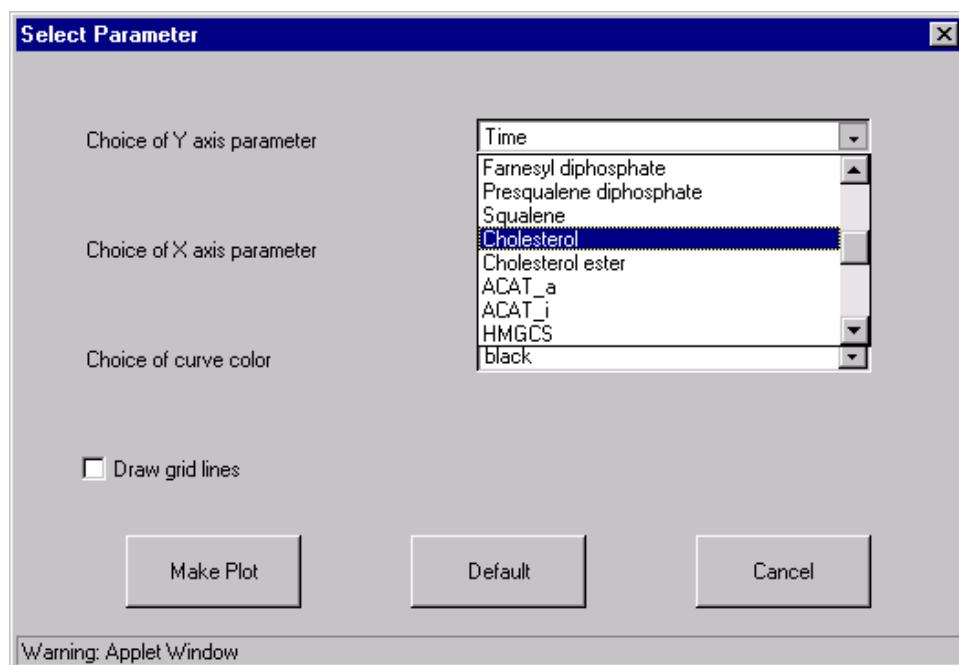
Select the observation parameters and click 'OK'.

Wait for calculation of the model.

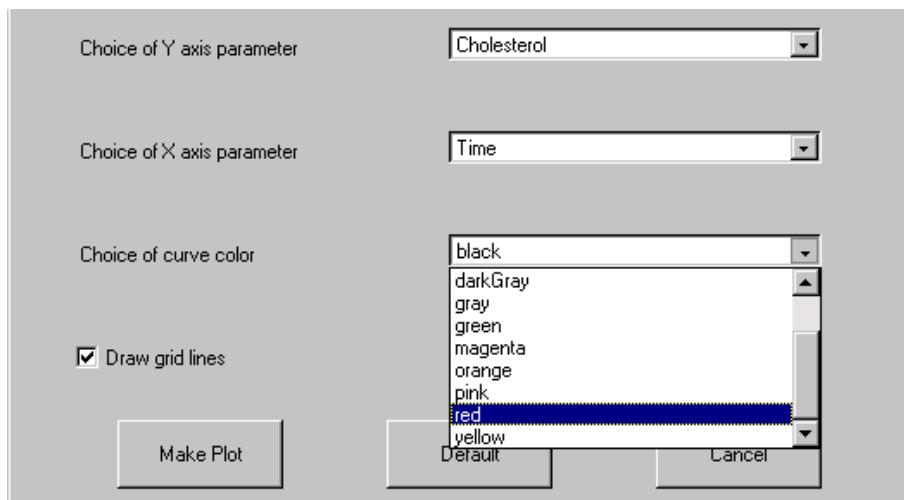
Next, you will get the dialog window for ordering the view of graphical representation data output.



4. Choose the parameters for X and/or Y axis from the drop down menu 'Choice of Y axis parameter' and 'Choice of X axis parameter'.



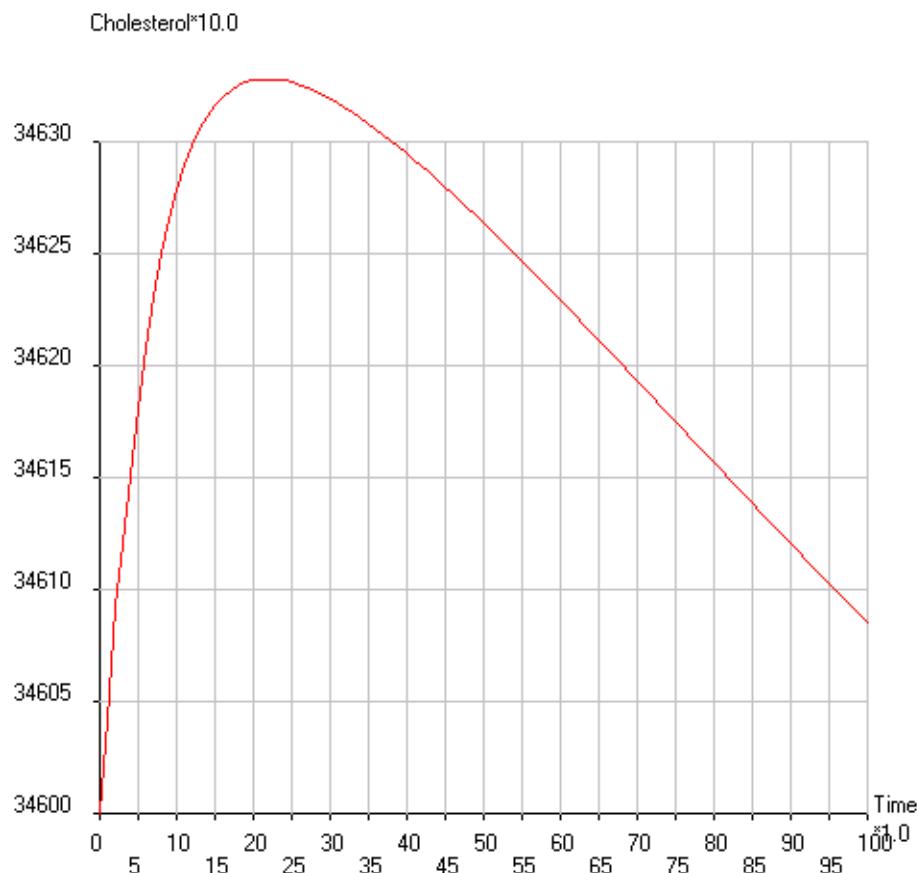
5. Select the colour of the curve at the plot from the drop down menu ‘Curve Color’.



For visualization of grid lines at the plot, tick off the check box ‘Draw grid lines’.

When the button 'Default' is clicked, the first component from the drop down menu ‘Choice of Y axis parameter’ and the option ‘Time’ from the drop down menu ‘Choice of X axis parameter’ will be selected.

6. Click the button 'Make Plot'.  
This brings up the plot resulted.



If necessary, you may view the lists of reaction constants and final component concentrations. For this purpose, go to the model window, click the drop-down menu below the 'Model Analysis' button, click the view 'Stationary states' and find the view you want, 'Data' or 'Constant'.



If you click the view 'Constant', this will display the following table:

Constants List	
Kcoaacet	1000.0 1/sec
KISBacet	1000.0 1/sec
KMSBacet	1.0E7 molecules/cell
KISBCoA	1000.0 1/sec
KutiACoA	0.1 1/sec
KMSBCoA	1.0E7 molecules/cell
Kuti2CoA	0.1 1/sec
KISBeduc	980.0 1/sec
KMSBeduc	1.2E7 molecules/cell
KISBme5P	1.0 1/sec
KISBm5PP	1.0 1/sec
KISBpePP	1.0 1/sec
KdimetPP	10.0 1/sec
KgeraPP1	1.0E-6 cell/(molecules*sec)
KgeraPP2	0.01 1/sec
KISBynth	40.0 1/sec
KMSBynth	200000.0 molecules/cell
Kutisop	1.0 1/sec

Cancel

Warning: Applet Window

This table lists the names of constants that are involved in mathematical model, their values and units.

If you selected the view 'Data', this will bring up the following table:

Data	
acetyl-CoA	7995308.158 molecules/cell
acetoacetyl-CoA	2.487072197E7 molecules/cell
HMG-CoA	2.536485297E7 molecules/cell
Mevalonic acid	1408774.93 molecules/cell
Mevalonic acid 5-phosphate	1408774.93 molecules/cell
5-Diphosphomevalonate	1408774.93 molecules/cell
Isopentenyl diphosphate	45444.35258 molecules/cell
Dimethylallyl diphosphate	1.001800329E7 molecules/cell
Geranyl diphosphate	81814.78419 molecules/cell
Farnesyl diphosphate	288258.082 molecules/cell
Presqualene diphosphate	408476.8396 molecules/cell
Squalene	83092.72186 molecules/cell
Cholesterol	345732.4696 molecules/cell
Cholesterol ester	175310.8088 molecules/cell
ACAT_a	9855.395495 molecules/cell
ACAT_i	10144.6045 molecules/cell
HMGCS	2117.732403 molecules/cell
HMGR	2117.732403 molecules/cell

Warning: Applet Window

This table lists the names of dynamic variables of the gene network, values of concentrations of dynamic variables at the end of calculations and measurement units.

### Example 2.

Modelling of a mutation by the example of mutation decreasing cell concentration of LDL receptor mRNA.

For modelling of a mutation, it is necessary to do the following:

1. Select the model of interest. In our case, click the option ‘Cholesterol\_MODEL’ in the drop-down menu ‘Select a model’, this will bring up the starting page of cholesterol biosynthesis gene network modeling (see Example 1);
2. Load the graphical interface of the model selected;
3. Click the component named LDLR<sub>i</sub> that denotes the concentration of LDL receptors, which are located inside the cell and have not reached the cell membrane. A part of this LDL pool consists of the native LDL synthesized from mRNA;
4. By clicking the button ‘Next’, select the following reaction: *LDLR gene + SREBP1 homodimer*  $\xrightarrow{K_i, K_m}$  *LDLR*;
5. For decreasing the LDL synthesis rate, for example, for 30%, input the value ‘0.7’ in the text-box ‘KISBldlr’;
6. Click the button ‘Enter’;
7. Calculate the model.

In order to compare the results obtained to the results of the model with initial parameters or some other set of parameters, it is necessary to close the windows ‘Select Parameter’ and ‘Transient states’ and input the necessary set of parameters (reaction constant values and initial concentrations of components of a gene network). Then calculate the model again. Make necessary plots and compare them with the plots for the previous set of parameters.

By analogy, mutations described in the Section 'Application of the models: Computer simulation of mutations of the Low Density Lipoprotein (LDL) receptor gene' could be modeled, as well as many other mutations influencing the system of cell cholesterol biosynthesis.

**Example 3.**

Modeling of the EKLF gene mutation that modifies the functioning of erythrocyte maturation gene network.

This mutation completely arrests the synthesis of the EKLF transcription factor. For simulation of this mutation, do the following:

1. Select the model ‘Erythroid\_differentiation’.
2. Upload graphic interface of the model chosen.
3. Click the component named GATA1-p, which determines the inner-cellular concentration of the key transcription factor GATA-1;
4. Choose the reaction  $EKLF\ gene \xrightarrow{Ki, Km, GATA1-p} EKLF\ mRNA$  by clicking the button ‘Next’;
5. Enter the value of the constant KISBgmre equaling to 0;
6. Click the button ‘Enter’;
7. Calculate the model.

In order to compare the results obtained to the results of the model with initial parameters or some other set of parameters, it is necessary to close the windows ‘Select Parameter’ and ‘Transient states’ and input the necessary set of parameters (reaction constant values and initial concentrations of components of a gene network). Then calculate the model again. Make necessary plots and compare them with the plots for the previous set of parameters.

### 3. Mathematical models of gene networks in SBML format

#### 3.1 About mathematical models of gene networks in SBML format

[This web page](#) contains 38 mathematical models in SBML format (Level 1 and Level 2) of gene networks represented in [GeneNet](#) database.

What is the SBML?

SBML - [Systems Biology Markup Language](#) is a common representation language for storing biochemical models. SBML is based on [XML](#), and contains structures for representing compartments, species and reactions, as well as optional unit definitions, parameters and rules (constraints). All information about SBML you can find on [this web page](#).

All compartments, species and reactions in the models of gene networks are the same as in [GeneNet](#) database.

There are some software packages supporting SBML Level 1:

[The Systems Biology Workbench \(SBW\)](#) and [SBW-enabled modules \(Jarnac\)](#), a biochemical simulation package for Windows; [JDesigner](#), a visual biochemical network layout tool and [others](#).

[Gepasi](#) is a software package for modeling biochemical systems. It simulates the kinetics of systems of biochemical reactions and provides a number of tools to fit models to data, optimize any function of the model, perform metabolic control analysis and linear stability analysis.

It is necessary to download the software packages from web sites given above for working with [the mathematical models of gene networks](#).

#### 3.2 List of mathematical models in SBML format of gene networks represented in [GeneNet](#) database

1.	Adipocyte	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
2.	Antiviral response	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
3.	Apoptosis (general scheme)	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
4.	Cell Cycle (G0/G1-S transition)	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
5.	Cholesterol metabolism (intracellular) 1	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
6.	Cholesterol metabolism (intracellular) 2	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
7.	Cholesterol metabolism (intracellular) 3	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
8.	Erythroid differentiation	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
9.	Germination (endosperm)	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
10.	Heat Shock Response	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
11.	HSP70 autoregulation	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
12.	LEA program	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
13.	Leptin (organism level)	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
14.	Lipid metabolism in blood	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
15.	Lipid metabolism in liver cells	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
16.	Lipid metabolism (integral diagram)	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
17.	Macrophage activation (model)	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
18.	MAPK cascade	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
19.	NF-kappaB activation	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>

20.	NO biosynthesis pathway	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
21.	Nodulation	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
22.	Oxidative stress response (glutathione)	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
23.	Photomorphogenesis	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
24.	Plant-pathogen	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
25.	Principal cell of CCD	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
26.	REDOX-REGULATION	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
27.	Seed reserve mobilisation (1): carbohydrates	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
28.	Seed reserve mobilisation (2): lipids and phosphates	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
29.	Seed reserve mobilisation (3): proteins	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
30.	Seed reserve mobilisation (4): regulatory relationships	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
31.	Seed reserve mobilisation (5): the general diagram	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
32.	Seed reserve mobilisation (organism level)	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
33.	Steroidogenesis (adrenal cortex)	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
34.	Steroidogenesis (sex steroids)	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
35.	Storage protein biosynthesis (dicot)	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
36.	Storage protein biosynthesis (monocot)	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
37.	Thermotolerance	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>
38.	Thyroid system	<a href="#">SBML Level 1</a>	<a href="#">SBML Level 2</a>

**Comments and questions are welcome to Alexander V. Ratushny ([ratushny@bionet.nsc.ru](mailto:ratushny@bionet.nsc.ru)).**