





Genomes to Hits: The Emerging Assembly Line in Silico

Prof. B . J ayaram

Department of Chemistry &

Supercomputing Facility for Bioinformatics & Computational Biology &

School of Biological Sciences

Indian Institute of Technology Delhi



The Dream @ SCFBio

From Genome to Drug : Establishing the Central Dogma of Modern Drug Discovery



Develop In Silico Suggestions of Personalized Medicine



A Case Study



http://www.scfbio-iitd.res.in/publication/CHAPTER-3-B%20Jayaram-LATEST.pdf

Hepatitis B virus (HBV) is a major blood-borne pathogen worldwide. Despite the availability of an efficacious vaccine, chronic HBV infection remains a major challenge with over 350 million carriers.

No.	HBV ORF	Protein	Function
1	ORF P	Viral polymerase	DNA polymerase, Reverse transcriptase and RNase H activity ^[36,48] .
2	ORF S	HBV surface proteins (HBsAg, pre-S1 and pre-S2)	Envelope proteins: three in-frame start codons code for the small, middle and the large surface proteins ^[36,49,50] . The pre-S proteins are associated with virus attachment to the hepatocyte ^[51]
3	ORF C	Core protein and HBeAg	HBcAg: forms the capsid ^[36] . HBeAg: soluble protein and its biological function are still not understood. However, strong epidemiological associations with HBV replication ^[52] and risk for hepatocellular carcinoma are known ^[42] .
4	ORF X	HBx protein	Transactivator; required to establish infection <i>in vivo</i> ^[53,54] . Associated with multiple steps leading to hepatocarcinogenesis ^[45] .





United States FDA approved agents for anti-HBV therapy

Agent	Mechanism of action / class of drugs	
Interferon alpha	Immune-mediated clearance	
Peginterferon alpha2a	Immune-mediated clearance	
Lamivudine	Nucleoside analogue	
Adefovir dipivoxil	Nucleoside analogue	
Tenofovir	Nucleoside analogue	
Entecavir	Nucleoside analogue	
Telbivudine	Nucleoside analogue	

Resistance to nucleoside analogues have been reported in over 65% of patients on long-term treatment. It would be particularly interesting to target proteins other than the viral polymerase.

Wanted: New targets and new drugs





Input the HBV Genome sequence to *ChemGenome 3.0*:

Hepatitis B virus, complete genome NCBI Reference Sequence: NC_003977.1 >gi|21326584|ref|NC_003977.1| Hepatitis B virus, complete genome

ChemGenome 3.0 output Five protein coding regions identified

Gene 2 (BP: 1814 to 2452) predicted by the *ChemGenome 3.0* software encodes for the HBV precore/ core protein (Gene Id: 944568) >gi|77680741|ref|YP_355335.1| precore/core protein [Hepatitis B virus]



MQLFPLCLIISCSCPTVQASKLCLGWLWGMDIDPYKE FGASVELLSFLPSDFFPSIRDLLDTASALYREALESPEH CSPHHTALRQAILCWGELMNLATWVGSNLEDPASREL VVSYVNVNMGLKIRQLLWFHISCLTFGRETVLEYLVS FGVWIRTPPAYRPPNAPILSTLPETTVVRRRGRSPRRR TPSPRRRRSQSPRRRRSQSRESQC

Input Amino acid sequence to Bhageerath







Input Protein Structure to Active site identifier of Sanjeevini (AADS) 10 potential binding sites identified A quick scan against a million compound library Sanjeevini (RASPD) calculation with an average cut off binding affinity to limit the number of candidates.

RASPD output

2057 molecules were selected with good binding energy from one million molecule database corresponding to the top 5 predicted binding sites.





Out of the 2057 molecules, top 40 molecules are given as input to *Sanjeevini* (ParDOCK) for atomic level binding energy calculations. Out of this 40, (with a cut off of -7.5 kcal/mol), 24 molecules are seen to bind well to precore/core protein target. These molecules could be tested in the Laboratory.

Molecule ID	Binding Energy (kcal/mol)	
0001398	-10.14	
0004693	-8.78	
0007684	-10.05	
0007795	-9.06	
0008386	-8.38	
0520933	-8.21	
0587461	-10.22	
0027252	-8.39	
0036686	-8.33	
0051126	-8.73	
0104311	-9.3	
0258280	-7.8	
0000645	-7.89	
0001322	-8.23	
0001895	-9.49	
0002386	-8.53	
0003092	-8.35	
0001084	-8.68	
0002131	-8.07	
0540853	-11.08	
1043386	-10.14	
0088278	-9.16	
0043629	-7.5	
0097895	-8.04	



....

....





From Genome to Hits



Genome





X Teraflops Chemgenome Bhageerath Sanjeevini

Hits





www.scfbio-iitd.res.in

•Genome Analysis - ChemGenome

A novel *ab initio* Physico-chemical model for whole genome analysis

•Protein Structure Prediction – *Bhageerath*

A *de novo* energy based protein structure prediction software

•Drug Design – Sanjeevini

A comprehensive target directed lead molecule design protocol





Arabidopsis Thaliana (Thale Cress)



Gene Prediction Accuracies

Software	Method	Sensitivity*	Specificity*
GeneMark.hmm http://www.ebi.ac.uk/genemark/	5th-order Markov model	0.82	0.77
GenScan http://genes.mit.edu/GENSCAN.html	Semi Markov Model	0.63	0.70
MZEF http://rulai.cshl.org/tools/genefinder/	Quadratic Discriminant Analysis	0.48	0.49
FGENF http://www.softberry.com/berry.phtml	Pattern recognition	0.55	0.54
Grail http://grail.lsd.ornl.gov/grailexp/	Neural network	0.44	0.38
FEX http://www.softberry.com/berry.phtml	Linear Discriminant analysis	0.55	0.32
FGENESP http://www.softberry.com/berry.phtml	Hidden Markov Model	0.42	0.59

*Desirable: A sensitivity & specificity of unity => While it is remarkable that these methods perform so well with very limited experimental data to train on, more research, new methods and new ways of looking at DNA are required.





ChemGenome:

Build a three dimensional physico-chemical vector which, as it walks along the genome, distinguishes genes from non-genes





"A Physico-Chemical model for analyzing DNA sequences", Dutta S, Singhal P, Agrawal P, Tomer R, Kritee, Khurana E and Jayaram B,J.Chem. Inf. Mod., 46(1), 78-85, 2006.







 $\mathbf{E}_{\mathbf{HB}} = \mathbf{E}_{\mathbf{i}-\mathbf{l}} + \mathbf{E}_{\mathbf{j}-\mathbf{m}} + \mathbf{E}_{\mathbf{k}-\mathbf{n}}$

$$E_{\text{Stack}} = (E_{i-m} + E_{i-n}) + (E_{j-l} + E_{j-n}) + (E_{k-l} + E_{k-m}) + (E_{i-j} + E_{i-k} + E_{j-k}) + (E_{l-m} + E_{l-n} + E_{m-n})$$

Hydrogen bond & Stacking energies for all 32 unique trinucleotides were calculated from 50 ns long **Molecular Dynamics Simulation Trajectories on 39 sequences encompassing all possible tetranucleotides in the #ABC database* and the data was averaged out from the multiple copies of the same trinucleotide. The resultant energies were then linearly mapped onto the [-1, 1] interval giving the x & y coordinates for each of the 64 codons.

*Beveridge et al. (2004). *Biophys J* 87, 3799-813. #Dixit et al. (2005). *Biophys J* 89, 3721-40. Lavery et al. (2009) *Nucl. Acid Res.*, *38*(1), 299-313.



Dinucleotide	Hydrogen bond	Stacking Energy	Strength parameter
AA	-6.92	-26.92	-33.84
AC	-9.64	-27.87	-37.51
AG	-8.78	-26.91	-35.69
AT	-7.05	-27.34	-34.38
СА	-9.34	-27.23	-36.57
CC	-11.84	-26.33	-38.17
CG	-11.37	-27.83	-39.20
СТ	-8.78	-26.91	-35.69
GA	-10.12	-26.98	-37.10
GC	-12.03	-28.27	-40.30
GG	-11.84	-26.33	-38.17
GT	-9.64	-27.87	-37.51
ТА	-7.16	-27.15	-34.31
ТС	-10.12	-26.98	-37.10
TG	-9.34	-27.23	-36.57
TT	-6.92	-26.92	-33.84

Energy parameters (in kcal) for dinucleotides derived from molecular dynamics simulations

Tm (C) = {(-8.69 x E) + $[6.07 x \ln(\text{Len})] + [4.97 x \ln(\text{Conc})] + [1.11 x \ln(\text{dna})]$ } -233.45

G. Khandelwal, J. Gupta, B. Jayaram, J. Bio Sc., 2012, 37, xxx-xxx.





Melting temperatures of ~ 200 oligonucleotides: Prediction versus Experiment



The computed (MD derived hydrogen bond + stacking) energy (E) correlates very well with experimental melting temperatures of DNA oligonucleotides





Solute-Solvent Interaction Energy for Genes/Non-genes



Coding and non-coding frames have different solvation characteristics which could be used to build the third parameter (z) besides hydrogen bonding (x) & stacking (y)





TTAConjugate
rule acts as a
good
constraint on
the 'z'ATA
ATC
ATA
ATC
TGO
TGO
TGO
TGO
TGO
TGO
TGO
TGO
AGTConjugate
rule acts as a
good
constraint on
the 'z'
parameter of
Chemgenome
or one could
simply use
+1/-1 as in the
Table for 'z'!!

TTT Phe -1	GGT Gly +1	TAT Tyr -1	GCT Ala +1
TTC Phe -1	GGC Gly+1	TAC Tyr -1	GCC Ala +1
TTA Leu -1	GGA Gly+1	TAA Stop -1	GCA Ala +1
TTG Leu -1	GGG Gly +1	TAG Stop -1	GCG Ala +1
ATT Ile -1	CGT Arg+1	CAT His +1	ACT Thr -1
ATC Ile +1	CGC Arg -1	CAC His -1	ACC Thr +1
ATA Ile +1	CGA Arg -1	CAA Gln -1	ACA Thr +1
ATG Met -1	CGG Arg+1	CAG Gln +1	ACG Thr -1
TGT Cys -1	GTT Val +1	AAT Asn -1	CCT Pro +1
TGC Cys -1	GTC Val +1	AAC Asn +1	CCC Pro -1
TGA Stop -1	GTA Val +1	AAA Lys +1	CCA Pro -1
TGG Trp -1	GTG Val +1	AAG Lys -1	CCG Pro +1
AGT Ser -1	CTT Leu +1	GAT Asp +1	TCT Ser -1
AGC Ser +1	CTC Leu -1	GAC Asp +1	TCC Ser -1
AGA Arg+1	CTA Leu -1	GAA Glu +1	TCA Ser -1
AGG Arg -1	CTG Leu +1	GAG Glu +1	TCG Ser -1

Extent of Degeneracy in Genetic Code is captured by *Rule of Conjugates*: A_{1,2} is the conjugate of $C_{1,2}$ & $U_{1,2}$ is the conjugate of $G_{1,2}$:(A₂ x C₂ & G₂ x U₂) With 6 h-bonds at positions 1 and 2 between codon and anticodon, third base is inconsequential

With 4 h-bonds at positions 1 and 2 third base is essential

With 5 h-bonds middle pyrimidine renders third base inconsequential;

middle purine requires third base.

B. Jayaram, "Beyond Wobble: The Rule of Conjugates", J. Molecular Evolution, 1997, 45, 704-705.

Codons with $G_1 \rightarrow +1$; C_1G_3 or $C_1T_3 \rightarrow +1$; C_1A_3 or $C_1C_3 \rightarrow -1$



ChemGenome

A Physico-Chemical Model for identifying signatures of functional units on Genomes



(1) "A Physico-Chemical model for analyzing DNA sequences", Dutta S, Singhal P, Agrawal P, Tomer R, Kritee, Khurana E and Jayaram B, J.Chem. Inf. Mod., 46(1), 78-85, 2006; (2) "Molecular Dynamics Based Physicochemical Model for Gene Prediction in Prokaryotic Genomes ", P. Singhal, B. Jayaram, S. B. Dixit and D. L. Beveridge, *Biophys. J.*, 2008, 94, 4173-4183; (3) "A phenomenological model for predicting melting temperatures of DNA sequences", G. Khandelwal and B. Jayaram, PLoS ONE, 2010, 5(8): e12433. doi:10.1371/journal.pone.0012433





Distinguishing Genes (blue) from Non-Genes (red) in ~ 900 Prokaryotic Genomes



Three dimensional plots of the distributions of gene and non-gene direction vectors for six best cases (A to F) calculated from the genomes of

(A) Agrobacterium tumefaciens (NC_003304), (B) Wolinella succinogenes (NC_005090),

(C) Rhodopseudomonas palustris (NC_005296), (D) Bordetella bronchiseptica (NC_002927),

(E) *Clostridium acetobutylicium* (NC_003030), (F) *Bordetella pertusis* (NC_002929)



SCFBio

http://www.scfbio-iitd.res.in/chemgenome/index.jsp

ChemGenome 1.1 GENE EVALUATOR ChemGenome is a physico-chemical method [1] which accepts DNA sequence in FASTA format and characterizes it as gene or nongene based on hydrogen bonding energy, stacking energy and groove potentials for each trinucleotide (codon). Agrobacterium Wolinella Bordetella Clostridium Rhodopseudomonas Bordetella pertusis succinogenes acetobutylicium tumefaciens bronchiseptica palustris (NC_005296) (NC_002929) (NC_003304) (NC_005090) (NC_002927) (NC_003030) Above is a pictorial representation of the separation of genes(blue) from non-genes(red). ChemGenome is ab initio in nature and has been tested on 294786 experimentally verified genes in 331 prokaryotic genomes. The observed average sensitivity, specificity & correlation-coefficient are found to be 96.9% (min: 90%, max: 100%), 86.0% & 85.0% respectively. Preliminary studies on eukaryotic genomes show that the model successfully separates the exonic regions from the non-coding regions.A software for whole genome analysis is available at www.scfbio-iitd.res.in/chemgenome2 ChemGenome Please specify the E-mail id : ailesh@scfbio-iitd.res.in Insert the Nucleotide sequence (in FASTA format)* : Help >Gene Name (This comment line is necessary) ATGTTGGTGTCCGCAAGGGTAGAGAAACAÁAGCGTGTTGCTTATCAGGGGAAGGCGACAGTGCTTGCTCTCGG TAAGG CCTTGCCGAGCAATGTTGTTTCCCAGGAGAATCTCGTGGAGGAGTATCTCCGTGAAATCAAATGCGATAACCTTTC TAT CAAAGACAAGCTGCAACACTTGTGCAAAAGCACAACTGTCAAGACACGCTACACAGTCATGTCACGGGAGACG CTGCAC AAATACCCTGAACTAGCAACCGAGGGTTCCCCCAACCATCAAACAGAGGCTTGAGATTGCAAACGATGCGGTTGT GCAGA SUBMIT RESET Upload Browse.

Instructions for using the Tool

- The tool takes DNA sequence in FASTA format as input file.
- Browse to select the input file and upload.
- The input file can contain multiple sequences, each sequence being in FASTA format.
- For multiple sequences, please specify the E-mail address or wait for a few minutes to get the on-line result.
- Click on Submit to get the result
- For further information, please see the Help file.

Suggestions and Comments

We will be glad to receive your suggestions and comments/feedback at scfbio@scfbio-iitd.res.in. References

[1] "A Physico-Chemical model for analyzing DNA sequences", Dutta S, Singhal P, Agrawal P, Tomer R, Kritee, Khurana E and Jayaram B, *J. Chem. Inf. Mod.*, 46 (1), 78-85, 2006. [ABSTRACT].

[2] "Beyond the Wobble : The rule of conjugates", Jayaram B, Journal of Mol. Evol., 1997.45.704.

Copyright 2004-2006, Prof B. Jayaram & Co-workers

The ChemGenome2.0 WebServer

http://www.scfbio-iitd.res.in/chemgenome/chemgenomenew.jsp

	CHENIGENOME 2.0
	ab-initio dens Prediction Software
Chemgenome is an ab-inti reading frames. The meth prokaryotic genomes. Rea	io gene prediction software, which find genes in prokaryotic genomes in all six adology follows a physico-chemical approach and has been validated on 372 ad more about ChemGenome
	A.
Download CHEMGENOME	2.0 for Linux environment from here 斗
	[General Info] [Data Set] [Validated Result Set] [Help] [Hume]
Input File	Browse
OR paste Genome Sequer	nce in FASTA format
Bun Chemogenores CD	
Additional Parameters	
Threshold Values : 100 🖻	Start Codon: ATG 🖻 ETG 🗔 6TG 🖻 TTG 🖻
Method: 💿 DNA 🛇 Prot	cein V Swissprot
Method : ⓒDNA ○Prot E-mail ID :	cein 🔍 Swissprot (Optional)
Method : ③DNA 〇Prot E-mail ID :	cein 🔍 Swissprot (Optional)
Method : ⑦DNA ○Prot E-mail ID : Thresthold Value: If you You have large genomes	cein CSwissprot (Optional) have small genome you can specify lower threshold value to find smaller genes. [f you can specify higher threshold value to weed out false positives
Method : DNA Prot E-mail ID : Threshold Value: If you You have large genomes Start Codon: You can spe	cin Cytional] (Optional] have small genome you can specify lower threshold value to find smaller genes. [F you can specify higher threshold value to weed out false positives edify what should be the start codon with which you want to find genes.
Method : DNA Prot E-mail ID : Threshold Value: If you you have large genomes Start Codon: You can spo Method : DNA Space: The method b input file. It searches for (DNA).	Coptional] (Optional] have small genome you can specify lower threshold value to find smaller genes. [F you can specify higher threshold value to weed out false positives edify what should be the start codon with which you want to find genes. takes complete or part of genome sequence of prokaryotic species in FASTA format as genes based on physico-chemical properties of double-helical deoxyribonucleic acid
Method : DNA Prot E-mail ID : Threshold Value: If you you have large genomes Start Codon: You can spo Method : DNA Space: The method b input file. It searches for (DNA). Arotein Space: The metho on stereochemical proper	Coptional] (Optional] have small genome you can specify lower threshold value to find smaller genes. [F you can specify higher threshold value to weed out false positives edify what should be the start codon with which you want to find genes. cakes complete or part of genome sequence of prokaryotic species in FASTA format as genes based on physico-chemical properties of double-helical deoxyribonucleic acid d takes the result generated from DNA space as input file and works as a filter based rises of protein sequences to reduce false positives.
Method : DNA Prot E-mail ID : Threshold Value: If you you have large genomes Start Codon: You can spo Method : DNA Space: The method to input file. It searches for (DNA). Anotein Space: The method on stereochemical proper Swissprot Space : The method standard deviation of a of based on the frequency of based on the freq based on based on the frequency of based on the frequency of	Coptional] (Optional] have small genome you can specify lower threshold value to find smaller genes. If you can specify higher threshold value to weed out false positives edify what should be the start codon with which you want to find genes. cakes complete or part of genome sequence of prokaryotic species in FASTA format as genes based on physico-chemical properties of double-helical deoxyribonucleic add d takes the result generated from DNA space as input file and works as a filter based rise of protein sequences to reduce false positives. hod takes the result generated from protein space as input file and calculates the juery nucleotide sequence (predicted gene sequence) with the swissprot proteins of protein at maximum.
Method : DNA Prot E-mail ID : Threshold Value: If you you have large genomes Start Codon: You can spe Method : DNA Space: The method b input file. It searches for (DNA). Arotein Space: The method on spereochemical proper Swissprot Space : The method on stareochemical proper Standard daviation of a g based on the frequency of false positives at minimum There is no file size limita with us. If the program cr	Coptional] (Optional] have small genome you can specify lower threshold value to find smaller genes. If you can specify higher threshold value to weed out false positives ecify what should be the start codon with which you want to find genes. cakes complete or part of genome sequence of prokaryotic species in FASTA format as genes based on physico-chemical properties of double-helical deoxyribonucleic acid d takes the result generated from DNA space as input file and works as a filter based rise of protein sequences to reduce false positives. had takes the result generated from protein space as input file and calculates the guery nucleotide sequence (predicted gene sequence) with the swissprot proteins of occurrence of aminoacids. A threshold standard deviation is chosen to keep the m and precision at maximum. tion for the genomes. We have tested on more than 5 MB genome file size available rashes on large genome size, more than 5 MB, please intimate us.





Arabidopsis Thaliana

(Thale Cress)

Software	Method	Sensitivity	Specificity
<i>ChemGenome</i> www.scfbio-iitd.res.in/chemgenome	Physico-chemical model	0.87	0.89
GeneMark.hmm http://www.ebi.ac.uk/genemark/	5th-order Markov model	0.82	0.77
GenScan http://genes.mit.edu/GENSCAN.html	Semi Markov Model	0.63	0.70
MZEF http://rulai.cshl.org/tools/genefinder/	Quadratic Discriminant Analysis	0.48	0.49
FGENF http://www.softberry.com/berry.phtml	Pattern recognition	0.55	0.54
Grail http://grail.lsd.ornl.gov/grailexp/	Neural network	0.44	0.38
FEX http://www.softberry.com/berry.phtml	Linear Discriminant analysis	0.55	0.32
FGENESP http://www.softberry.com/berry.phtml	Hidden Markov Model	0.42	0.59

A simple physico-chemical model works just as well as any of the sophisticated knowledge base driven methods and has scope for further systematic improvements









Chemgenome methodology enables detection of not only protein coding regions on a genome but also promoters (top panel) and introns (bottom panel) etc..

Solvation energies of DNA distinguish mRNA genes from tRNA genes



Relative solvation energy per base pair of DNA sequences coding for 2063537 mRNAs (Blue) and 56251 tRNAs (Pink) from 1531 genomic sequences; calculated from MD data. The X-axis denotes the index of the genome, the Y-axis depicts the solvation energy of the sequence relative to the average for that genome.

G. Khandelwal, B. Jayaram, J. Am. Chem. Soc., 2012, 134 (21), 8814-8816, DOI:10.1021/ja3020956.



Supercomputing Facility for Bioinformatics & Computational Biology IITD Some day, it should be possible to read the book of Human Genome like a novel !!! 3000 Mb **Gene & Gene related Sequences Extra-genic DNA** 2100 Mb 900 Mb **Repetitive DNA Unique & low copy number Coding DNA Non-coding DNA** 420 Mb 1680 Mb 90 Mb (3%) !!! 810 Mb

Tandemly repeated DNA

Satellite, micro-satellite, mini-satellite DNA

LTR elements, Lines, Sines, DNA Transposons

Interspersed genome wide repeats





www.scfbio-iitd.res.in

•Genome Analysis - *ChemGenome* A novel *ab initio* Physico-chemical model for whole genome analysis

•**Protein Structure Prediction** – *Bhageerath* A *de novo* energy based protein structure prediction software

•Drug Design – Sanjeevini

A comprehensive target directed lead molecule design protocol





Bhageerath Protein Tertiary Structure Prediction

.....GLU ALA GLU MET LYS ALA SER GLU ASP LEU LYS LYS HIS GLY VAL THR VAL LEU THR ALA LEU GLY ALA ILE LEU LYS LYS GLY HIS HIS GLU ALA GLU LEU LYS PRO LEU ALA GLN SER HIS ALA THR LYS HIS LYS ILE PRO ILE LYS TYR LEU GLU PHE ILE SER GLU ALA ILE ILE HIS LEU HIS.....







Protein Folding Problem







WHY FOLD PROTEINS ?

Pharmaceutical/Medical Sector



- Active site directed drug-design
- Mapping the functions of proteins in metabolic pathways.





PROTEIN FOLDING LANDSCAPE



Native structure at the bottom of the rugged free energy well is the folded protein.





Protein Structure Prediction Approaches

Comparative Modeling

Homology

Similar sequences adopt similar fold is the basis.

Alignment is performed with related sequences. (SWISS-MODEL-www.expasy.org, 3D JIGSAW-www.bmm.icnet.uk etc).

Threading

Sequence is aligned with all the available folds and scores are assigned for each alignment according to a scoring function. (Threader - bioinf.cs.ucl.ac.uk)





Computational Requirements for *ab initio* **Protein Folding**

Strategy A

• Generate all possible conformations and find the most stable one.

- For a protein comprising 200 AA assuming 2 degrees of freedom per AA
- 2^{200} Structures => 2^{200} Minutes to optimize and find free energy.

```
2^{200} Minutes = 3 \times 10^{54} Years!
```

Strategy B

- Start with a straight chain and solve F = ma to capture the most stable state
- A 200 AA protein evolves
- $\sim 10^{\text{--}10} \sec$ / day / processor
- $10^{-2} \sec => 10^{8} \text{ days}$
 - $\sim 10^6$ years

With 10^6 processors ~ 1 Year





From Sequence to Structure: The Bhageerath Pathway

AMINO ACID SEQUENCE

Bioinformatics Tools

EXTENDED STRUCTURE WITH PREFORMED SECONDARY STRUCTURAL ELEMENTS

TRIAL STRUCTURES (~10⁶ to 10⁹)

SCREENING THROUGH BIOPHYSICAL FILTERS

- 1. Persistence Length
- 2. Radius of Gyration
- 3. Hydrophobicity
- 4. Packing Fraction

MONTE CARLO OPTIMIZATIONS AND MINIMIZATIONS OF RESULTANT STRUCTURES (~10³ to 10⁵)

ENERGY RANKING AND SELECTION OF 100 LOWEST ENERGY STRUCTURES

STRUCTURE EVALUATION (Topology & ProRegIn) & SELECTION OF 5 LOWEST ENERGY STRUCTURES

NATIVE-LIKE STRUCTURES

Narang P, Bhushan K, Bose S and Jayaram B 'A computational pathway for bracketing native-like structures for small alpha helical globular proteins.' *Phys. Chem. Chem. Phys.* 2005, 7, 2364-2375.




Sampling 3D Space







Filter-Based Structure Selection

Persistence Length Analysis of 1,000 Globular Proteins



Frequency vs Hydrophobicity Ratio of 1,000 Globular Proteins



Radius of Gyration vs N^{3/5} of 1,000 Globular Proteins



N^{3/5} (N= number of amino acids)

N^{3/5} plot incorporates excluded volume effects (Flory P. J., *Principles of Polymer Chemistry*, Cornell University, New York, 1953).

Frequency vs Packing Fraction of 1,000 Globular Proteins



Globular proteins are known to exhibit packing fractions around 0.7





Removal of Steric Clashes in Selected Structures (Distance Based Monte Carlo)



A CONTRACTOR OF A CONTRACTOR O



Supercomputing Facility for Bioinformatics & Computational Biology IITD

Validation of Empirical Energy Based Scoring Function



Narang, P., Bhushan, K., Bose, S., and Jayaram, B. *J. Biomol.Str.Dyn*, **2006**,*23*,385-406; Arora N.; Jayaram B.; *J. Phys. Chem. B.* **1998**, *102*, 6139-6144; Arora N, Jayaram B, *J. Comput. Chem.*, **.1997**, *18*, 1245-1252.







Bhageerath is currently implemented on a 280 processor (~3 teraflops) cluster

Jayaram et al., Bhageerath, 2006, Nucleic Acid Res., 34, 6195-6204













Performance of Bhageerath on 70 Small Globular Proteins

			No. of	Lowest	Energy rank of
S No	DDDID	No of Amino	Secondary	RMSD Å	lowest RMSD
S.INU.	FDDID	Acids	Structure	(from	structure in top
			elements	native)	5 structures
1	1E0Q	17	2E	2.5	2
2	1B03	18	2E	4.4	2
3	1WQC	26	2H	2.5	3
4	1RJU	36	2H	5.9	4
5	1EDM	39	2E	3.5	2
6	1AB1	46	2H	4.2	5
7	1BX7	51	2E	3.2	4
8	1B6Q	56	2H	3.8	5
9	1ROP	56	2H	4.3	2
10	1NKD	59	2H	3.9	1
11	1RPO	61	2H	3.8	2
12	1QR8	68	2H	3.9	4
13	1FME	28	1H,2E	3.7	5
14	1ACW	29	1H,2E	5.3	3
15	1DFN	30	3E	5	1
16	1Q2K	31	1H,2E	4.8	4
17	1SCY	31	1H,2E	3.1	5
18	1XRX	34	1E,2H	5.6	1
19	1ROO	35	3H	2.8	5
20	1YRF	35	3H	4.8	4
21	1YRI	35	3H	4.6	3
22	1VII	36	3H	3.7	2
23	1BGK	37	3Н	4.1	3
24	1BHI	38	1H,2E	5.3	2





			No. of		Energy rank of
C N-	DDDID	No of Amino	Secondary	Lowest	lowest RMSD
S.NO .	PDBID	Acids	Structure	RMSD Å	structure in top 5
			elements		structures
25	10VX	38	1H,2E	4	1
26	1I6C	39	3E	5.1	2
27	2ERL	40	3Н	4	3
28	1RES	43	3Н	4.2	2
29	2CPG	43	1E,2H	5.3	2
30	1DV0	45	3Н	5.1	4
31	1IRQ	48	1E,2H	5.5	3
32	1GUU	50	3Н	4.6	4
33	1GV5	52	3Н	4.1	2
34	1GVD	52	3Н	5.1	4
35	1MBH	52	3Н	4	4
36	1GAB	53	3Н	4.9	1
37	1MOF	53	3Н	2.9	5
38	1ENH	54	3Н	4.6	3
39	1IDY	54	3Н	3.6	5
40	1PRV	56	3Н	5	5
41	1HDD	57	3Н	5.5	4
42	1BDC	60	3Н	4.8	5
43	1I5X	61	3H	3.6	3
44	1I5Y	61	3H	3.4	5
45	1KU3	61	3H	5.5	4
46	1YIB	61	3H	3.5	5
47	1AHO	64	1H,2E	4.5	4
48	1DF5	68	3H	3.4	1
49	1QR9	68	3H	3.8	2
50	1AIL	70	3Н	4.4	3





S.No.	PDBID	No of Amino Acids	No. of Secondary Structure elements	Lowest RMSD Å	Energy rank of lowest RMSD structure in top 5 structures
51	2G7O	68	4H	5.8	2
52	2OCH	66	4H	6.6	3
53	1WR7	41	3E,1H	5.2	2
54	2B7E	59	4H	6.8	4
55	1FAF	79	4H	6.4	4
56	1PRB	53	4H	6.9	4
57	1DOQ	69	5H	6.8	3
58	1I2T	61	4H	5.4	4
59	2CMP	56	4H	5.6	1
60	1BW6	56	4H	4.2	1
61	1X4P	66	4H	5.2	3
62	2K2A	70	4H	6.1	1
63	1TGR	52	4H	6.8	2
64	2V75	90	5H	7.0	3
65	1HNR	47	2E,2H	5.2	2
66	2KJF	60	4H	5.0	4
67	1RIK	29	2E,2H	4.4	4
68	1JEI	53	4H	5.8	5
69	2HOA	68	4H	6.3	4
70	2DT6	62	4H	5.9	3





Predicted Structures with Bhageerath

for 70 Globular Proteins superposed on their corresponding experimental structures

A		25	27		-		-	2 👣	
1e0q	1b03	1wqc	1rju	1edm	1ab1	1bx7	1fme	1acw	1ail
and the second s		States -	C.T. T. S.	and a	82	33	5.35	2	V
160	lrop	1.pkd		1	1yrf	1yri	2erl	1res	1gvd
podr	тор		Irpo	Idr8	2 C	230	1	and and a second	n
S		2 P	300		Samo		23	1df5	[]
1dfn	1q2k	1scy	1xrx	1roo	Imbn			IUIS	Le S
25	E ??	G		R			to	AL SO	·
	*	De	2	W	2g7o	2och	1wr7	2b7e	1faf
1vii	1bgk	1bhi	lovx	160	SE	- 1	1	10	
-	A S	1 and			185	(Com	Ea.	AND	25
2cpg	1dv0	1irq	1guu	1gv5	lprb	1doq	1i2t	2cmp	1x4p
E	100 A	Sec.	Ste		87 C	375		100	C.
53F	1		Sec.	100 m	1bw6	2k2a	1tgr	2v75	1hnr
1gab	1mof	1enh	1idy	1prv		de		~~~~	122
and the second second	AND	See.	r	K	36	5	1 Alexandree	A start	3.00
1i5x	1i5y	1ku3	1yib	1aho	2kjf	1rik	1jei	2hoa	2dt6
			Native stru	ucture	Predic	ted structu	ıre		





Bhageerath versus Homology modeling

No	Protein PDB ID	CPHmodels RMSD(Å)	ESyPred3D RMSD(Å)	Swiss-model RMSD(Å)	3D-PSSM RMSD(Å)	Bhageerath# RMSD(Å)
1.	1IDY (1-54)*	3.96 (2-54)*	3.79 (2-51)*	5.73 (1-51)*	3.66 (1-51)*	3.36
2.	1PRV (1-56)*	5.66 (2-56)*	5.56 (3-56)*	6.67 (3-56)*	5.94 (1-56)*	3.87

*Numbers in parenthesis represent the length (number of amino acids) of the protein model. #Structure with lowest RMSD bracketed in the 5 lowest energy structures.

The above two proteins have maximum sequence similarity of 38% and 48% respectively.

In cases where related proteins are not present in structural databases Bhageerath achieves comparable accuracies.

Homology models are simply superb where the similarities between query sequence and template in the protein data bank are high. Where there is no match/similarity ab initio methods such as Bhageerath are the only option.



Residues



The Protein Structure Prediction Olympics: CASP9 (May 3rd to July 17th, 2010: 129 Targets)



Bhageerath vs other servers for Template free prediction in CASP9

				TASSER	ROBETTA	SAM-T08
Target	No.of		Bhageerath	RMSD Å	RMSD Å	RMSD Å
No.	residues	PDBID	RMSD Å			
T0531	65	2KJX	7.1	11.0	11.9	12.6
T0553	141	2KY4	9.6	6.0	11.5	8.6
T0581	136	3NPD	15.8	11.6	5.3	15.1
T0578	164	3NAT	19.2	11.6	15.5	19.1

While *Bhageerath* – an *ab initio* method - works well for small proteins (<100 residues), improvements are necessary to tackle larger proteins





Development of a homology - ab initio hybrid server Bhageerath-H Protocol





Homology *ab initio* hybrid methods are getting better in tertiary structure prediction *Bhageerath-H* is a participant in CASP10 (May-July, 2012) Stay Tuned to (<u>http://predictioncenter.org/casp10/</u>) for further progresses

BHAGEERATH : An Energy Based Protein Structure Prediction Server

The present version of "Bhageerath" accepts amino acid sequence and secondary structure information to predict 10 candidate structures for the native. It is anticipated that at least one native like structure (RMSD < 6Å without end loops) is present in the final structures. The server has been validated on 50 small globular proteins. Know about Protein Folding

Download BHAGEERATH 1.0 for Solaris 10.0 environment from here.

		[Repository]	[General Info]	[Links]	[Help]	[Home]
Process ID	56703599					
E-mail Address:		(Optional)				

Input Amino acid sequence in FASTA format OR Click on the Amino acid to add to the sequence

ALA VAL LEU ILE PRO
MET PHE TRP GLY SER
 THR CYS ASN GLN TYR
ASP GLU LYS ARG HIS

Secondary Structure Information

ullet Auto Secondary Structure Prediction $igodot$ Enter Secondary Structure Information					
Helix 😪 Residue Range 🛛 - 🔹 Add Clear					
SUBMIT RESET					

Retrieve previous results		
Job I);	Get Status

In case of any Suggestions/Exceptions, Please contact us at scfbio@scfbio-iitd.res.in

Bhageerath-H WebServer http://www.scfbio-iitd.res.in/bhageerath/bhageerath_h.jsp

BHAGEERATH-H: A Homology ab-intio Hybrid Web server for Protein Tertiary Structure Prediction						
"Bhageerath-H" accepts amino acid sequence to predict 5 candidate structures for the native. Here user has the flexibility to mention reference PDB(s) for modeling. Method has been fielded in CASP9 Experiment and has been improved since.						
[Repository] [Tutorial] [Sample File] [Links] [Help] [Home]						
Process ID 1764624 E-mail Address:						
Unload sequence in EASTA format Choose File No file choose						
OR Input Amine poid coquence in EASTA format						
ALA VAL LEU ILE PRO MET PHE TRP GLY SER THR CYS ASN GLN TYR ASP GLU LYS ARG HIS						
Template Information						
Auto Template Searching User Defined Template PDB ID - Chain ID Add Clear SUBMIT RESET						

In search of rules of protein folding:

Margin of Life: Amino acid compositions in proteins have a tight distribution

The average percentage occurrence of each amino-acid

The average percentage occurrence of each aminoacid for folded proteins gives the "Chargaff's rules" for protein folding and the standard deviations give the "margin of life".

	Folded Proteins	from the ExPASy Server.			
Amino Acid	Margin of Life (mean ± std, n = 3718)	Amino Acid	Protein sequences confirmed by annotation and experiments (mean \pm std, n = 131855)		
A	7.8 ± 3.4	A	7.2 ± 3.0		
v	7.1 ± 2.4	V	6.3 ± 2.1		
I	5.8 ± 2.4	I	5.1 ± 2.2		
L	9.0 ± 2.9	L	9.6 ± 2.9		
Y	3.4 ± 1.7	Y	3.0 ± 1.5		
F	3.9 ± 1.8	F	3.9 ± 1.8		
w	1.3 ± 1.0	W	1.2 ± 0.9		
P	4.4 ± 2.0	Р	5.4 ± 2.6		
M	2.2 ± 1.3	M	2.2 ± 1.3		
C	1.8 ± 1.5	С	1.9 ± 2.3		
Т	5.5 ± 2.4	Т	5.5 ± 1.8		
S	6.0 ± 2.5	S	7.9 ± 2.8		
Q	3.8 ± 2.0	Q	4.3 ± 2.0		
N	4.3 ± 2.2	N	4.2 ± 1.9		
D	5.8 ± 2.0	D	5.2 ± 1.9		
E	7.0 ± 2.7	E	6.8 ± 2.8		
н	2.3 ± 1.4	H	2.4 ± 1.3		
R	5.0 ± 2.3	R	5.3 ± 2.9		
K	6.3 ± 2.8	К	6.0 ± 2.9		
G	7.2 ± 2.8	G	6.6 ± 2.8		

The average percentage occurrence of each amino acid, their STD as observed and as calculated from the binomial distribution.

	P(%)	STD (observed)	STD (random)
A	7.8	3.4	7:2
V	7.1	2.4	6.6
I	5.8	2.4	5.5
L	9.0	2.9	8.2
Y	3.4	1.7	3.3
F	3.9	1.8	3.7
W	1.3	1.0	1.3
Р	4.4	2.0	4.2
м	2.2	1.3	2.2
С	1.8	1.5	1.8
Т	5.5	2.4	5.2
S	6.0	2.5	5.6
Q	3.8	2.0	3.7
N	4.3	2.2	4.1
D	5.8	2.0	5.5
Е	7.0	2.7	6.5
H	2.3	1.4	2.2
R	5.0	2.3	4.8
к	6.3	2.8	5.9
G	7.2	2.8	6.7

Mezei (2011), JBSD

In search of rules of protein folding:

Cα atoms of proteins of varying sequences and sizes follow a single (universal) spatial distribution



All 400 Cα spatial distributions (above) collapse into one narrow band (below) irrespective of the chemical nature of the amino acids when their percentage occurrences are considered => A Stoichiometric Hypothesis for Protein Folding.



Mittal & Jayaram et al., (2010) *JBSD*, 28, 133-142; (2011), *JBSD*, 28, 443-454; (2011), *JBSD*, 28, 669-674. While structure prediction attempts are progressing well, rules of folding are still elusive.

Functional Implications of SNPs



Cartoon representation of the structure of Human Angiogenin (PDB entry 1B1I) showing its functional sites; catalytic triad residues are represented as stick models, nuclear localization signal is represented in magenta color and receptor binding site is represented in orange color.

Mutations in the coding region of the ANG (angiogenin) gene have been found in patients suffering from Amyotrophic Lateral Sclerosis (ALS). Neurodegeneration results from the loss of angiogenic ability of ANG (protein coded by ANG gene). This is one of the several examples where SNPs could lead to disease/disorder or a predisposition. We performed extensive molecular dynamics (MD) simulations of wild-type ANG protein and disease associated ANG variants to elucidate the mechanism behind the loss of ribonucleolytic activity and nuclear translocation activity, functions needed for angiogenesis. MD simulations can yield information on structural and dynamic differences in the catalytic site and nuclear localization signal residues between WT-ANG (Wild-type ANG) and six mutants. Variants K17I, S28N, P112L and V113I have confirmed association with ALS, while T195C and A238G single nucleotide polymorphisms (SNPs) at the gene level encoding L35P and K60E mutants respectively, have not been associated with the disease. Our results show that the loss of ribonucleolytic activity in K17I is caused by conformational switching of the catalytic residue His114 by 99°. The loss of nuclear translocation activity of S28N and P112L is caused by changes in the folding of the residues ³¹RRR³³ that result in the reduction in solvent accessible surface area. Based on the results obtained, we predict that V133I mutant will exhibit loss of angiogenic properties by loss of nuclear translocation activity and L35P mutant by loss of both ribonucleolytic activity and nuclear translocation activity. No functional loss was inferred for K60E. This is just an illustration of how molecular simulations on protein tertiary structures can be used to infer functional implications of mutations. MD simulations on a series of mutants are time consuming. Faster methods are required for genomic scans.

(A. K. Padhi, H. Kumar, S. V. Vasaikar, B. Jayaram and James Gomes, "Mechanisms of Loss of Functions of Human Angiogenin Variants Implicated in Amyotrophic Lateral Sclerosis", *PLoS One*, 2012, 7(2): e32479. doi:10.1371/journal.pone.0032479)





www.scfbio-iitd.res.in

•Genome Analysis - *ChemGenome* A novel *ab initio* Physico-chemical model for whole genome analysis

•Protein Structure Prediction – *Bhageerath* A *de novo* energy based protein structure prediction software

•Drug Design – Sanjeevini

A comprehensive target directed lead molecule design protocol





Target Directed Lead Design



Given the structure of the drug target, design a molecule that will bind to the target with high affinity and specificity





COST & TIME INVOLVED IN DRUG DISCOVERY



Source: PAREXEL's Pharmaceutical R&D Statistical Sourcebook, 2001, p96.; Hileman, Chemical Engg. News, 2006, 84, 50-1.





Pharmaceutical R&D is Expensive

New Chemical Entities (NCEs) need to be continuously developed since income from older drugs gets gradually reduced on account of increasing competition from other products, generics as well *drug resistance*.

Drug Development is an Uphill Task

1035 new drugs approved by FDA between 1989 to 2000361 (35%) were New Molecular Entities (NME).Only 15% were deemed to provide significant improvement over existing medicines.

http://www.seniors.gov/articles/0502/medicine-study.htm

SCFBio

Structure Based Lead Molecule Design







Present Scenario of Drug Targets



BLUE: Number of targets in each class. (Imming P, Sinning C, Meyer A. *Nature Rev Drug Discov* **2006**;5: 821) (Total 218 targets & 8 classes) GREEN: Number of 3D structures available in each class (Total: 130) (Protein Data Bank)

Shaikh SA, Jain T, Sandhu G, Latha N, Jayaram, B. Current Pharmaceutical Design, 2007, 13, 3454-3470.





Some Concerns in Lead Design In Silico

- Novelty and Geometry of the Ligands
- Accurate charges and other Force field parameters
- Ligand Binding Sites
- Flexibility of the Ligand and the Target
- Solvent and salt effects in Binding
- Internal energy versus Free energy of Binding
- Druggability
- Computational Tractability

De novo LEAD-LIKE MOLECULE DESIGN: THE SANJEEVINI PATHWAY

SCFBio



Jayaram, B., Latha, N., Jain, T., Sharma, P., Gandhimathi, A., Pandey, V.S., Indian Journal of Chemistry-A. 2006, 45A, 1834-1837. Tanya Singh, Goutam Mukherjee, Abhinav Mathur, B. Jayaram, *Sanjeevini* – A User Manual, 2012, manuscript in preparation







Supercomputing facility for bioinformatics and computational biology IIT Delhi

Molecular Descriptors / Drug-like Filters

Lipinski's rule of five

Molecular weight	≤ 500
Number of Hydrogen bond acceptors	s <u><</u> 10
Number of Hydrogen bond donors	<u><</u> 5
logP	≤ 5

Additional filters

Molar Refractivity	≤ 140
Number of Rotatable bonds	<u><</u> 10

http://www.scfbio-iitd.res.in/utility/LipinskiFilters.jsp



http://www.scfbio-iitd.res.in/dock/ActiveSite_new.jsp





Rank of the cavity points vs. cumulative percentage prediction Top ten cavity points capture the active site 100 % of time in 640 protein targets



Prediction accuracies of the active site by different softwares

Sl. No	Softwares	Top1	Тор3	Top5	Top10
1	SCFBIO (Active	73	92	95	100
	Site Finder)				
2	Fpocket	83	92	-	
3	PocketPicker	72	85	-	
4	LIGSITE ^{cs}	69	87	-	
5	LIGSITE	69	87	_	
6	CAST	67	83	-	
7	PASS	63	81	-	
8	SURFNET	54	78	-	
9	LIGSITEcsc	79	-	-	

http://www.scfbio-iitd.res.in/software/drugdesign/raspd.jsp



Supercomputing Facility for Bioinformatics & Computational Biology, IIT Delhi



Home | Drug Design Software

RASPD for Preliminary Screening of Drugs

The challenge for computer aided drug discovery is to achieve this specificity - with small molecule inhibitors - in binding to target proteins, at reduced cost and time while ensuring synthesizability, novelty of the scaffolds and proper ADMET profiles. RASPD is a computationally fast protocol for identifying good candidates for any target protein. The binding pocket of the input target protein is scanned for the number of hydrogen bond donors, acceptors, number of hydrophobic groups and number of rings. A QSAR type equation combines the aforementioned properties of the target protein and the candidate molecule and an estimate of the binding free energy is generated if the target protein were to complex with the candidate. The most interesting feature of this methodology is that it takes fraction of a second for calculating the binding affinities of the protein-candidate molecule complexes as opposed to several minutes in known art today for regular docking and scoring method, whereas the accuracy of this method in sorting good candidates is comparable with the conventional techniques. We have also created million molecules database. This database is prepared to include chemical formula, structure, topological index, number of hydrogen bond donors and acceptors, number of hydrophobic groups, number of rings, logP values for each of the million molecules. Scoring of 1 million small molecule database by RASPD method to identify hits for a particular protein target is also web enabled for free access at the same site.

Know	more	about	RASPD	Screeing.	Click	here	to	see	'How	to	Use	Tool'.	Click	here	to	see
'Comj	putatio	onal Flo	w Chart'.													

Browse
Enter Drug Id: DRG
your job





Quantum Chemistry on Candidate drugs for Assignment of Force Field Parameters



G. Mukherjee, N. Patra, P. Barua and B. Jayaram, (2011), JCC, 32,893-907.
http://www.scfbio-iitd.res.in/software/drugdesign/charge.jsp







MONTE CARLO DOCKING OF THE CANDIDATE DRUG IN THE ACTIVE - SITE OF THE TARGET www.scfbio-iitd.res.in/dock/pardock.jsp









RMSD between the crystal structure and one of the top five docked structuresT. Singh, D. Biswas and B. Jayaram, AADS - An automated active site identification, docking and scoring protocol for
protein targets based on physico-chemical descriptors, (2011), JCIM, 51 (10), 2515-2527

ENERGY BASED SCORING FUNCTION

 $\Delta G^{\circ}_{bind} = \Delta H^{\circ}_{el} + \Delta H^{\circ}_{vdw} - T\Delta S^{\circ}_{rtvc} + \Delta G^{\circ}_{hpb}$



Correlation between experimental & calculated binding free energy for 161 protein-ligand complexes (comprising 55 unique proteins)

SCFBio

Jain, T & Jayaram, B, *FEBS Letters*, **2005**, 579, 6659-6666 www.scfbio-iitd.res.in/software/drugdesign/bappl.jsp

Correlation between experimental ΔT_m and calculated free energy of interaction for DNA-Drug Complexes

S.A Shaikh and B.Jayaram, J. Med.Chem. , 2007, 50, 2240-2244

www.scfbio-iitd.res.in/software/drugdesign/preddicta.jsp





Comparative Evaluation of Scoring Functions

S	Scoring Function	Method	Dataset		Correlation	Reference			
No.			Training Test		Coefficient				
1.00	1 unction				(r)				
1.	Present	Force field /	61	100	r = 0.92	FEBS Letters, 2005, 579, 6659			
	Work(BAPPL*)	Empirical							
2.	DOCK	Force field	-	-	-	J. ComputAided Mol. Des. 2001, 15, 411			
3.	EUDOC	Force field	-	-	-	J. Comp. Chem. 2001, 22, 1750			
4.	CHARMm	Force field	-	-	-	J. Comp. Chem. 1992, 13, 888			
5.	AutoDock	Force field	-	-	-	J. Comp. Chem. 1998, 19, 1639			
6.	DrugScore	Knowledge	-	-	-	J. Mol. Biol. 2000, 295, 337			
7.	SMoG	Knowledge	-	36	r = 0.79	J. Am. Chem. Soc. 1996, 118, 11733			
8.	BLEEP	Knowledge	-	90	r = 0.74	J. Comp. Chem. 1999, 202, 1177			
9.	PMF	Knowledge	-	77	r = 0.78	J. Med. Chem. 1999, 42, 791			
10.	DFIRE	Knowledge	-	100	r = 0.63	J. Med. Chem. 2005, 48, 2325			
11.	SCORE	Empirical	170	11	r = 0.81	J. Mol. Model. 1998, 4, 379			
12.	GOLD	Empirical	-	-	-	J. Mol. Biol. 1997, 267, 727			
12	LUDI	Empirical	82	12	r = 0.83	J. ComputAided Mol. Des. 1994, 8, 243 &			
15.		Empiricai				1998, 12, 309			
14.	FlexX	Empirical	-	-	-	J. Mol. Biol. 1996, 261, 470			
15.	ChemScore	Empirical	82	20	r = 0.84	J. ComputAided Mol. Des. 1997, 11, 425			
16.	VALIDATE	Empirical	51	14	r = 0.90	J. Am. Chem. Soc. 1996, 118, 3959			
17.	Ligscore	Empirical	50	32	r = 0.87	J. Mol. Graph. Model. 2005, 23, 395			
18.	V CSCODE	Empirical	200	30	r = 0.77	J. ComputAided Mol. Des. 2002, 16, 11			
	A-CSCUKE	(consensus)							
19.	CLIDE	Force field /	-	-	-	J. Med. Chem. 2004, 47, 1739			
	GLIDE	Empirical							





Binding Affinity Analysis on Zinc Containing Metalloprotein-Ligand Complexes



Correlation between the predicted and experimental binding free energies for 90 zinc containing metalloprotein-ligand complexes comprising 5 unique targets

T. Jain & B. Jayaram, *Proteins: Struct. Funct. Bioinfo.* 2007, 67, 1167-1178.

www.scfbio-iitd.res.in/software/drugdesign/bapplz.jsp

Comparative evaluation of some methodologies reported for estimating binding affinities of zinc containing metalloproteinligand complexes

S. No.	Contributing Group	Method	Protein Studied	Training Set	Test Set	R ²
1.	Donini et al	MM-PBSA	MMP	-	6	
2.	Raha <i>et al</i>	QM	CA & CPA	-	23	0.69
3.	Toba <i>et al</i>	FEP	MMP	-	2	-
4.	Hou, et al	LIE	MMP	-	15	0.85
5.	Hu et al	Force Field	MMP	-	14	0.50
6.	Rizzo et al	MM-GBSA	MMP	-	6	0.74
7.	Khandelwal et al	QM/MM	MMP	-	28	0.76
8.	Present Work	Force Field / Empirical	CA, CPA, MMP, AD & TL	40	50	0.77



BAPPL server



HIV-I Protease complexed with U75875 (1hiv.pdb)

Welcome to the BAPPL server

Binding Affinity Prediction of Protein-Ligand (BAPPL) server computes the binding free energy of a nonmetallo protein-ligand complex using an all atom energy based empirical scoring function [1] & [2].







ParDOCK

Automated Server for Protein Ligand Docking





Logarithm of the frequencies of the occurrence of base sequences of lengths 4 to 18 base pairs in *Plasmodium falciparum* and in humans embedding a regulatory sequence TGCATGCA (shown in green), GTGTGCACAC (blue) and GCACGCGTGC (orange) or parts thereof, of the plasmodium. The solid lines and the dashed lines correspond to humans and plasmodium, respectively. Curves lying between 0 and 1 on the log scale indicate occurrences in single digits.

One needs to cover at least 18 bp for uniqueness of the drug target



PreDDICTA

Predict DNA-Drug Interaction strength by Computing Δ Tm and Affinity of binding.



S.A Shaikh and B.Jayaram, J. Med. Chem., 2007, 50, 2240-2244





Supercomputing facility for bioinformatics and computational biology IIT Delhi

Binding Affinity Analysis



P. Kalra, T. V. Reddy, and B. Jayaram, "Free energy component analysis for drug design: A case study of HIV-1 protease-inhibitor binding", *J. Med. Chem.*, 2001, *44*, 4325-4338.



Shaikh, S., Jain. T., Sandhu, G., Latha, N., <u>Jayaram., B</u>., *A physico-chemical pathway from targets to leads*, 2007, *Current Pharmaceutical Design*, 13, 3454-3470.

	Drug1	Drug2	Drug3	Drug4	Drug5	Drug6	Drug7	Drug8	Drug9	Drug10	Drug11	Drug12	Drug13	Drug14
Target1														
Target2														
Target3														
Target4														
Target5														
Target6														
Target7														
Target8														
Target9														
Target10														
Target11														
Target12														
Target13														
Target14														

BLUE: HIGH BINDING AFFINITY

GREEN: MODERATE AFFINITY

ORANGE: POOR AFFINITY

Diagonal elements represent drug-target binding affinity and off-diagonal elements show drug-non target binding affinity. Drug 1 is specific to Target 1, Drug 2 to Target 2 and so on. Target 1 is lymphocyte function-associated antigen LFA-1 (CD11A) (1CQP; Immune system adhesion receptor) and Drug 1 is lovastatin. Target 2 is Human Coagulation Factor (1CVW; Hormones & Factors) and Drug 2 is 5-dimethyl amino 1-naphthalene sulfonic acid (dansyl acid). Target 3 is retinol-binding protein (1FEL; Transport protein) and Drug 3 is n-(4-hydroxyphenyl)all-trans retinamide (fenretinide). Target 4 is human cardiac troponin C (1LXF; metal binding protein) and Drug 4 is 1-isobutoxy-2-pyrrolidino-3[n-benzylanilino] propane (Bepridil). Target 5 is DNA {1PRP; d(CGCGAATTCGCG)} and Drug 5 is propamidine. Target 6 is progesterone receptor (1SR7; Nuclear receptor) and Drug 6 is mometasone furoate. Target 7 is platelet receptor for fibrinogen (Integrin Alpha-11B) (1TY5; Receptor) and Drug 7 is n-(butylsulfonyl)-o-[4-(4-piperidinyl)butyl]-l-tyrosine (Tirofiban). Target 8 is human phosphodiesterase 4B (1XMU; Enzyme) and Drug 8 is 3-(cyclopropylmethoxy)-n-(3,5-dichloropyridin-4-yl)-4-(difluoromethoxy)benzamide (Roflumilast). Target 9 is Potassium Channel (2BOB; Ion Channel) and Drug 9 is tetrabutylammonium. Target 10 is {2DBE; d(CGCGAATTCGCG)} and Drug 10 is Diminazene aceturate (Berenil). Target 11 is Cyclooxygenase-2 enzyme (4COX; Enzymes) and Drug 13 is carboxyatractyloside. Target 14 is Glutamate Receptor-2 (2CMO; Ion channel) and Drug 14 is 2-({[(imethylamino)(dihydroxy)-lambda~4~-sulfanyl]phenyl}-8-methyl-2-oxo-6,7,8,9-tetrahydro-1H-pyrrolo[3,2-H]isoquinolin-3(2H)-ylidene]amino}oxy)-4-hydroxybutanoic acid. The binding affinities are calculated using the software made available at http://www.scfbio-iitd.res.in/software/drugdesign/bappl.jsp and <u>http://www.scfbio-iitd.res.in/software/drugdesign/bappl.jsp</u> and <u>http://www.scfbio-iitd.res.in/preddicta</u>.





Supercomputing facility for bioinformatics and computational biology IIT Delhi

Supercomputing Facility For Bioinformatics and Computational Biology, IIT Delhi





Future of Drug Discovery: Towards a Molecular View of

ADMET (Absorption, Distribution, Metabolism, Excretion & Toxicity)



The distribution path of an orally administered drug molecule inside the body is depicted. Black solid arrows: Complete path of drug starting from absorption at site of administration to distribution to the various compartments in the body, like sites of metabolism, drug action and excretion. Dashed arrows: Path of the drug after metabolism. Dash-dot arrows: Path of drug after eliciting its required action on the target. Dot arrows: Path of the drug after being reabsorbed into circulation from the site of excretion.





From Genome to Hits



Genome





X Teraflops Chemgenome Bhageerath Sanjeevini

Hits





SCFBio Team



~ 6 teraflops of computing; 20 terabytes of storage + huge brain power





BioComputing Group, IIT Delhi (PI : Prof. B. Jayaram)

Present

Shashank Shekhar Tanya Singh Avinash Mishra Abhilash Jayaraj Sahil Kapoor Sanjeev Kumar Garima Khandelwal Priyanka Dhingra Ashutosh Shandilya Anjali Soni Nagarajan Goutam Mukherjee Vandana Satyanarayan Rao Navneet Tomar Preeti Bisht

Dr. Achintya Das Dr. Tarun Jain Dr. Kumkum Bhushan Dr. Nidhi Arora Pankaj Sharma A.Gandhimathi Neelam Singh Dr. Sandhya Shenoy

Former

- Dr. N. Latha Dr. Saher Shaikh Dr. Poonam Singhal Dr. E. Rajasekaran Praveen Agrawal Gurvisha Sandhu Shailesh Tripathi Rebecca Lee
- Dr. Pooja Narang Dr. Parul Kalra Dr. Surjit Dixit Surojit Bose Vidhu Pandey Anuj Gupta Dhrubajyoti Biswas Bharat Lakhani

Collaborators: Dr. Aditya Mittal & Prof. D.L. Beveridge

Lead Invent

Drug Design Solutions



Biospectrum Award 2011 Asia Pacific Emerging Company of the Year

Technologies

Novel Drug Discovery

Mr. Pankaj Sharma Mr. Surojit Bose Mr. Praveen Aggarwal Ms. Gurvisha Sandhu

Incubated at IIT Delhi (2007-2009) www.leadinvent.com



Under Incubation at IITD (since April, 2011) Received TATA NEN Award 2012 for being one of the best Upcoming Start-Up companies







Acknowledgements

Department of Biotechnology

Department of Science & Technology

Ministry of Information Technology

Council of Scientific & Industrial Research

Indo-French Centre for the Promotion of Advanced Research (CEFIPRA)

HCL Life Science Technologies

Dabur Research Foundation

Indian Institute of Technology, Delhi





A Few Key References Genome Annotation

(a) S. Dutta, P. Singhal, P. Agrawal, R. Tomer, Kritee, E. Khurana and B. Jayaram. "A Physico-Chemical Model for Analyzing DNA sequences", Journal of Chemical Information & Modelling, 2006, 46(1), 78-85. (b) P. Singhal, B. Jayaram, S. B. Dixit and D. L. Beveridge. Molecular Dynamics Based Physicochemical Model for Gene Prediction in Prokaryotic Genomes, 2008, Biophysical Journal, 94, 4173-4183; (c) G. Khandelwal and B. Jayaram. "A Phenomenological Model for Predicting Melting Temperatures of DNA Sequences", PLoS One, 2010, 5(8): e12433. doi:10.1371/journal.pone.0012433; (d) G. Khandelwal, B. Jayaram, "DNA-water interactions distinguish messenger RNA genes from transfer RNA genes", J. Am. Chem. Soc., 2012,134 (21), 8814-8816, DOI:10.1021/ja3020956; (e) G. Khandelwal, J. Gupta, B. Jayaram, "DNA energetics based analyses suggest additional genes in prokaryotes"; J. Bio Sc., 2012, 37, xxx-xxx.

Proten Structure prediction

2. Bhageerath: (a) P. Narang, K. Bhushan, S. Bose and B. Jayaram. "A computational pathway for bracketing native-like structures for small alpha helical globular proteins", *Phys. Chem. Chem. Phys.*, 2005, 7, 2364-2375; (b) B. Jayaram et al., "*Bhageerath..*", *Nucleic Acid Res.*, 2006, 34, 6195-6204; (c) S. R. Shenoy and B. Jayaram, "Proteins: Sequence to Structure and Function – Current Status", *Curr. Prot. Pep. Sci.*, 2010, 11, 498-514; (d) A. Mittal, B. Jayaram, S. R. Shenoy and T. S. Bawa, "A Stoichiometry driven universal spatial organization of backbones of folded proteins: Are there Chargaff's rules for protein folding ?", J. Biomol. Struc. Dyn., 2010, 28, 133-142; (e) Aditya K. Padhi, Hirdesh Kumar, Suhas V. Vasaikar, B. Jayaram and James Gomes, "Mechanisms of Loss of Functions of Human Angiogenin Variants Implicated in Amyotrophic Lateral Sclerosis", *PLoS One*, 2012, 7(2): e32479. doi:10.1371/journal.pone.0032479; (f) B. Jayaram, P. Dhingra, B. Lakhani and S. Shekhar, "*Bhageerath* - Targeting the Near Impossible: Pushing the Frontiers of Atomic Models for Protein Tertiary Structure Prediction", *Journal of Chemical Sciences*, 2012, 124, 83-91.

Drug Design

3. Sanjeevini: (a) T. Jain and B. Jayaram. "An all atom energy based computational protocol for predicting binding affinities of proteinligand complexes", *FEBS Letters*, 2005, 579, 6659-6666; (b) T. Jain and B. Jayaram. "A computational protocol for predicting the binding affinities of zinc containing metalloprotein-ligand complexes", *Proteins: Structure, function & Bioinformatics*, 2007, 67, 1167-1178; (c) S. Shaikh and B. Jayaram. "A swift all atom energy based computational protocol to predict DNA-Drug binding affinity andDT_m", *J. Med. Chem.*, 2007, 50, 2240-2244; (d) S. Shaikh et al.. "A physico-chemical pathway from targets to leads", *Current Pharmaceutical Design*, 2007, *13*, 3454-3470; (e) G. Mukherjee, N. Patra, P. Barua and B. Jayaram, "A fast empirical GAFF compatible partial atomic charge assignment scheme for modeling interactions of small molecules with biomolecular targets", *J. Computational Chemistry*, 2011, *32*,893-907; (f) T. Singh, D. Biswas and B. Jayaram, "AADS - An automated active site identification, docking and scoring protocol for protein targets based on physico-chemical descriptors", *Journal of Chemical Information & Modeling*, 2011, *51* (*10*), 2515-2527, DOI: 10.1021/ci200193z



Let us go climb Everest, cure cancer



Visit Us at www.scfbio-iitd.res.in

Thank You