

Helix-Turn-Helix Motif Detection in Protein Sequences

Giri Narasimhan¹, Changsong Bu², Yuan Gao³, Tom Milledge¹,
Xuning Wang⁴, Ning Xu⁵, Gaolin Zheng¹, And Kalai Mathee⁵



¹School of Computer Science, Florida International University, ²Idax Inc., ³IBM T.J. Watson Research,
⁴Parke Davis, ⁵University of Memphis, ⁶Department of Biology, Florida International University.

ABSTRACT

We use methods from Data Mining and Knowledge Discovery to design an algorithm for detecting motifs in protein sequences. The algorithm assumes that a motif is constituted by the presence of a "good" combination of residues in appropriate locations of the motif. The algorithm attempts to compile such good combinations into a "pattern dictionary" by processing an aligned training set of protein sequences. The dictionary is subsequently used to detect motifs in new protein sequences. Statistical significance of the detection results are ensured by statistically determining the various parameters of the algorithm.

Based on this approach, we have implemented a program called GYM. The Helix-Turn-Helix (HTH) Motif was used as a model system on which to test our program. The program was also extended to detect Homeodomain motifs. The detection results for the two motifs compare favorably with existing programs.

Previous Methods

Profile Method: [Gribskov '90]

- Build a **Profile** matrix based on frequencies of occurrence of amino acids in specified locations within the motif. $Weight(i, AA) = 100 \log(p(i, AA) / (m(AA) N))$
- Use the profile to perform **detection** of the motif in new sequences.

Hidden Markov Models Neural Networks

New Algorithm: Basic Assumptions

- Combinations of residues in specific locations (may not be contiguous) contribute towards stabilizing a structure.
- Some **reinforcing** combinations are relatively rare.

Patterns: Examples

Loc	Protein Name	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20		
14	Cro	F	G	L	K	T	A	K	D	L	L	L	Y	Q	S	A	W	K	A	I	H		
16	434 Cro	M	T	D	T	E	L	A	T	K	A	G	V	K	Q	O	S	L	Q	L	I	E	A
11	P22	G	T	D	R	A	V	A	K	A	L	G	I	S	D	A	A	V	S	O	W	K	E
31	Rep	L	S	E	S	V	A	D	K	M	M	G	O	S	G	V	G	A	L	F	N		
16	434 Rep	L	N	A	E	L	A	Q	K	V	G	T	Q	S	I	E	D	L	E	K			
19	P22 Rep	I	R	D	A	A	L	G	K	M	V	G	V	S	N	V	A	I	S	O	W	E	R
24	CI1	L	G	T	E	K	T	E	A	V	G	V	D	K	S	O	L	S	R	W	K	R	
4	LacR	V	T	L	Y	D	V	R	E	F	A	G	V	S	Y	Q	T	V	S	R	V	N	
167	CAP	I	T	R	O	E	T	G	O	I	V	G	C	S	R	E	T	V	G	R	I	L	K
66	TrpR	M	S	O	R	E	L	K	N	E	L	A	G	I	A	T	I	T	R	G	S	N	
22	BlaA P9	L	N	F	T	K	A	A	L	E	L	V	V	T	O	G	A	V	S	O	O	V	
23	TrpR P9	N	S	V	S	D	A	A	E	L	H	V	T	H	G	A	V	S	R	O	L	K	

- Q1 G9 N20
- A5 G9 V10 I15

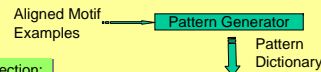
Motifs in Protein Sequences

Motifs are combinations of secondary structures in proteins with a specific **structure** and a specific **function**. Protein families are often characterized by one or more such motifs. Motif detection in proteins is thus an important problem since motifs carry out and regulate various functions, and the presence of specific motifs may help classify a protein.

Examples: Helix-Turn-Helix, Zinc-finger, Homeobox domain, Hairpin-beta motif, Calcium-binding motif, Beta-alpha-beta motif, Coiled-coil motifs.

New Algorithm: Outline

Pattern Mining:



Motif Detection:



Experimental Results

Motif	Protein Family	Number Tested	GYM = DE Agree	Number Annotated	GYM = Annot.
HTH Motif (22)	Master	88	88 (100%)	13	13
	Sigma	314	284 + 23 (98%)	96	82
	Negates	93	86 (92%)	0	0
	LysR	130	127 (98%)	95	93
	AraC	68	57 (84%)	41	34
	Rreg	116	99 (85%)	57	46
Total		675	653 + 23 (94%)	289	255 (88%)

Helix-Turn-Helix Motifs

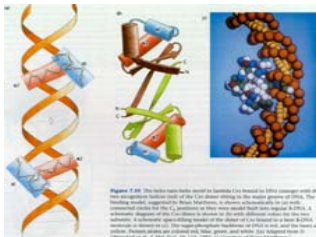
Structure

- 3-helix complex
- Length: 22 amino acids
- Turn angle



Function

Gene regulation by binding to DNA



GYM Algorithm: Pattern Mining

Algorithm Pattern-Mining

Input: Motif length m , support threshold T , list of aligned motifs M .
Output: Dictionary L of frequent patterns.

- $L_1 :=$ All frequent patterns of length 1
- for $i = 2$ to m do
- $C_i :=$ Candidates(L_{i-1})
- $L_i :=$ Frequent candidates from C_i
- if $(|L_i| \leq 1)$ then
- return L as the union of all $L_j, j < i$.

Conclusions

- Pattern Discovery is a powerful way to detect motifs in protein sequences.
- A Pattern Dictionary is a composite descriptor of a motif or a domain, and is better than a consensus sequence or a motif signature.
- The GYM program is accurate, sensitive, and efficient. It also provides lot of useful information along with the motif detection.
- Identical patterns occurring in different proteins have been observed to often share near-identical structures.
- Pattern discovery can also be used as a basis for local and global structure prediction

GYM Algorithm: Motif Detection

Algorithm Motif-Detection

Input: Motif length m , threshold score T , pattern dictionary L , and input protein sequence $P[1..n]$.
Output: Information about motif(s) detected.

- for each location i do
- $S :=$ MatchScore($P[i..i+m-1], L$).
- if $(S > T)$ then
- Report it as a possible motif

Future Work

- Automate Motif Detection Process for other targeted motifs.
- Automating the Choice of a Training Set.
- Negative Training Set.
- Mutational Studies.
- Protein Structure Prediction.
- Functional Annotations.

Website for GYM Online: www.cs.fiu.edu/~giri/bioinf/GYM2/welcome.html