# **Motif Extraction and Protein Classification**

Vered Kunik School of Computer Science Tel Aviv University Tel Aviv 69978, Israel kunikver@tau.ac.il Zach Solan School of Physics and Astronomy Tel Aviv University Tel Aviv 69978, Israel zsolan@tau.ac.il Shimon Edelman Department of Psychology Cornell University Ithaca, NY 14853, USA se37@cornell.edu

Eytan Ruppin School of Computer Science Tel Aviv University Tel Aviv 69978, Israel ruppin@tau.ac.il David Horn School of Physics and Astronomy Tel Aviv University Tel Aviv 69978, Israel horn@tau.ac.il

### Abstract

We present a novel unsupervised method for extracting meaningful motifs from biological sequence data. This de novo motif extraction (MEX) algorithm is data driven, finding motifs that are not necessarily over-represented in the data. Applying MEX to the oxidoreductases class of enzymes, containing approximately 7000 enzyme sequences, a relatively small set of motifs is obtained. This set spans a motif-space that is used for functional classification of the enzymes by an SVM classifier. The classification based on MEX motifs surpasses that of two other SVM based methods: SVMProt, a method based on the analysis of physical-chemical properties of a protein generated from its sequence of amino acids, and SVM applied to a Smith-Waterman distances matrix. Our findings demonstrate that the MEX algorithm extracts relevant motifs, supporting a successful sequence-to-function classification.

keywords motif extraction, enzyme classification

### Introduction

It is commonly accepted that high sequence similarity guarantees functional similarity of proteins. A contemporary analysis of enzyme function conservation by Tian and Skolnick [14] suggests that 40% pairwise sequence identity can be used as a threshold to certify functional similarity, i.e. the first three digits of the Enzyme Commission (EC) number are identical <sup>1</sup>. Using pairwise sequence similarity, and combining it with the Support Vector Machine (SVM) classification method [15, 10], Liao and Noble [7] have argued that they obtain a significantly improved remote homology detection relative to existing state-of-the-art algorithms.

There are alternative sequence-based approaches to the task of protein classification. One is based on general characteristics of the sequence, such as the number of specific amino-acids within it, as suggested in [6]. A recent variation of this approach represents the amino-acid sequence as a sequence of physical-chemical features [3, 4], such as hydrophobicity, normalized Van der Waals volume, polarity, polarizability, charge, surface tension, secondary structure and solvent accessibility. Cai *et al.* [3, 4] have applied SVM to these feature vectors and reported that the SVM-Prot technique reaches a high degree of accuracy, at a level of two digits of the EC number hierarchy, on various enzyme subclasses.

An alternative to the straightforward sequence similarity approach is the usage of motifs. Appropriately chosen sequence motifs may be expected to reduce noise in the data and indicate active regions of the protein, hence improving predictability of its function. A protein can then be represented as a 'bag of motifs' [1] (i.e. neglecting their particular order on the linear sequence), or a vector in a space spanned by these motifs. A recent work by Ben-Hur and Brutlag [2], based on the eMOTIF approach [9, 8], led to very good results. Starting out with 5911 enzyme se-

<sup>&</sup>lt;sup>1</sup>The function of an enzyme is specified by a name and a number given to it by the Enzyme Commission (EC). The EC number consists of four

numbers, n1:n2:n3:n4, corresponding to four levels of classification. The oxidoreductases class discussed in this paper corresponds to n1=1, one of the six main divisions. For this class, n2 (subclass) specifies electron donors, n3 (sub-subclass) specifies electron acceptors and n4 indicates the exact enzymatic activity.

quences of the oxidoreductases class, which consisted 129 EC subclasses, they based their analysis on 59783 regularexpression eMOTIFs. By using an appropriate feature selection method they obtained success rates well over 90% for a variety of classifiers.

The approach presented in this work is motif based. Its novelty is the employed motif extraction algorithm (MEX). Conventional approaches [5] construct motifs in terms of position specific weight matrices, or else use hidden Markov models and Bayesian networks, hence are supervised to some extent. MEX extracts motifs from proteins sequential data in an **unsupervised** manner, without requiring over-representation of its amino-acid motifs in the data set. MEX motifs are explicit strings in contradistinction to position-specific weight matrices or regular expressions. In the application described below, 3165 MEX motifs are extracted. This is a low number of motifs in comparison with the 59783 regular-expression eMOTIFs used by Ben-Hur and Brutlag [2].

In what follows, we demonstrate that an SVM analysis of oxidoreductases enzymes based on MEX motifs leads to results that are better than those obtained by an SVM based on pairwise sequence similarity. Furthermore, it outperforms SVMProt on the class of oxidoreductases enzymes, even though the latter is based on physical and chemical properties of the amino-acid sequence. Moreover, our algorithm is highly predictive of function, down to the third level (subsubclass) of the EC hierarchy.

### The Motif Extraction Algorithm (MEX)

MEX is a motif extraction algorithm that serves as the basic unit of ADIOS [12, 13], an unsupervised method for extraction of syntax from linguistic corpora. We apply it to the task of finding sequence-motifs within biological data. Consider a data set of sequences of variable length, each such sequence expressed in terms of an alphabet of finite size N (e.g. N=20 amino-acids in proteins). The N letters form vertices of a graph on which the sequences are placed as ordered paths. Each sequence defines such a path over the graph. In terms of all  $p(e_i|e_i)$  the graph defines a Markov model. Moreover, using any path on the graph, to be called henceforth a search-path, we find a particular instantiation of a variable order Markov model up to order k, where k is the length of the search-path. For each such search-path  $(e_1; e_k) = e_1 e_2 \cdots e_k$  we define a right-moving probability function, whose value at  $i, j \leq k$  is

$$P_R(e_i; e_j) = p(e_j | e_i e_{i+1} e_{i+2} \dots e_{j-1}) = \frac{l(e_i; e_j)}{l(e_i; e_{j-1})} \quad (1)$$

where  $l(e_i; e_j)$  is the number of occurrences of sub-paths  $(e_i; e_j)$  in the graph. Starting from the other end of the path

we define a left-moving probability function

$$P_L(e_j; e_i) = p(e_i | e_{i+1}e_{i+2}...e_{j-1}e_j) = \frac{l(e_j; e_i)}{l(e_j; e_{i+1})}.$$
 (2)

Fig. 1 demonstrates the type of structures that we expect to find in our graph - an assimilation of paths over a subsequence of the search-path. Such a subsequence is a candidate motif. The criteria for motif selection are defined by local maxima of  $P_L$  and  $P_R$  signifying, respectively, the beginning and ending of a motif.



Figure 1. The definition of a motif within the MEX algorithm. Note that the maxima of  $P_L$  and  $P_R$  defines the beginning and ending of the motif, respectively. Descents in  $P_L$  and  $P_R$  following the maxima signify divergence of paths.

Let us define the drop in probability functions as:

$$D_R(e_i; e_j) = P_R(e_i; e_j) / P_R(e_i; e_{j-1})$$
(3)

$$D_L(e_j; e_i) = P_L(e_j; e_i) / P_L(e_j; e_{i+1})$$
(4)

A threshold parameter  $\eta$  is introduced. The location  $e_{j-1}$  is declared as the ending of the motif if  $D_R(e_i; e_j) < \eta$ . Analogously,  $e_{i+1}$  is declared as the beginning of the motif if  $D_L(e_j; e_i) < \eta$ . Since the experimental probabilities,  $P_R(e_i; e_j)$  and  $P_L(e_j; e_i)$ , are determined by finite numbers of paths, a statistical measure is introduced in order to avoid erroneous results. Hence, we calculate the significance values of both  $D_R(e_i; e_j) < \eta$  and  $D_L(e_j; e_i) < \eta$  and require that their maximum be smaller than a parameter  $\alpha < 1$ . In the following application we have set  $\eta = 0.9$  and  $\alpha = 0.01$ . Once the algorithm reaches the stop criteria (e.g. ceases to locate new patterns) they are sorted in a length-significance descending order, by which their loci are identified on the original data.

# SVM functional classification based on MEX motifs

We have concentrated our analysis on the oxidoreductases class of enzymes. 7095 protein sequences and their EC number annotations were extracted from the SwissProt database Release 40.0. These proteins served as the data-set to which MEX was applied. The algorithm identified 3165 motifs of various lengths.

Classification was tested on levels 2 (subclass) and level 3 (sub-subclass) of the EC number. Subclasses were required to have a sufficient number of elements to ensure reasonable statistics. Protein sequences were represented as 'bags of MEX-motifs'. A linear SVM classifier (SVM-Light package, available online at http://svmlight.joachims.org/) was trained on each subclass separately, taking the protein sequences of the subclass as positive examples and the protein sequences of other subclasses as negative examples. 75% of the examples were used for training and the remaining examples for testing. The train-test procedure was repeated on six different random choices of train-test sets in order to accumulate statistics. We have tested various subsets of MEX motifs and discovered that the subset of motifs longer than five amino-acids leads to optimal results in the classification task. There are 1222 such motifs, spanning the space in which we represent all enzymes. The enzymes are classified into 16 subclasses of level 2 and 39 sub-subclasses of level 3.

Our results are compared to those of two other approaches. The first, SVMProt [3, 4], uses a performance measurement parameter defined as

$$Q = \frac{TP + TN}{TP + TN + FP + FN},$$
(5)

where TP, TN, FP and FN denote the number of true positive, true negative, false positive, and false negative outcomes respectively. The SVMProt results presented below are obtained from their published results. However, since the large negative set used in each classification task quickly yields a high Q value, we have chosen to use the Jaccard score

$$J = \frac{TP}{TP + FP + FN} \tag{6}$$

instead. Not taking into account TN, this performance measurement is more discriminative than Q.

The second approach, the Smith-Waterman algorithm [11], is based on a one-versus-all sequence similarity ap-

proach. This algorithm has been applied to the same set of 7095 oxidoreductases sequences analyzed by MEX. The ariadne tool has been used (written by R. Mott, available online at http://www.well.ox.ac.uk/ariadne) in order to obtain the p-values distances matrix,  $M_{SW}$ , defining the feature space of the SVM classifier. A minimal p-value threshold of 10<sup>-6</sup> was imposed in order to allow usage of p-values logarithm, defining a normalized distances matrix  $D_{SW}$ . This procedure is similar to the approach described in [7], however, the entire vector of  $D_{SW}$  has been used in our analysis for specifying an enzyme. The classification task has been performed with the same SVM classifier (linear kernel) employed to the data driven by MEX. The dataset has been preprocessed in order to produce an appropriate input file for the learning task. A random 75%:25% partition of the data into a training set and a testing set, respectively, has been used for each learning task. The train-test procedure was repeated on three different random choices of data sets in order to accumulate statistics.

Fig. 2 shows a comparison of the Jaccard score obtained by MEX, Smith-Waterman analysis and SVMProt (error deviations are not presented for the latter as they were not included in their published results). The scores obtained by MEX are clearly higher than those obtained by the other methods. The average J-scores are  $0.89 \pm 0.06$  for MEX,  $0.74 \pm 0.13$  for SVMProt and  $0.79 \pm 0.12$  for the Smith-Waterman method. Noticeably, there is no correlation between the size of the subclass and the J-scores obtained by the various methods. Clearly, if the size of the subclass is too small, i.e. the number of the positive examples is small, a large variance in the train/test different divisions may exist, resulting in large error deviations. However, in most cases, the average J-scores are high, independent of the tested subclass.

Third level classification results were not compared to SVMProt as none were included in their published results. Table 1 presents a comparison of the Jaccard scores obtained by MEX and Smith-Waterman analysis. The scores obtained by MEX are clearly higher. The average J-scores are  $0.89 \pm 0.08$  for MEX and  $0.78 \pm 0.15$  for Smith-Waterman. These findings attest to the meaningful information embodied in MEX selected motifs, facilitating a fine tuned classification of these proteins.

### **Motif selection**

Motifs of various lengths were extracted by applying the MEX algorithm. The enzyme function classification capabilities of the motifs were tested using various length-dependent subsets of motifs. The results attest that the classification task performed by the subset of 601 motifs of length 6 obtain remarkable J-scores, comparable to those achieved by using the entire set of motifs longer than 5. In



Figure 2. Jaccard scores for second-level EC subclasses obtained by MEX (upper panel), Smith-Waterman (second panel) and SVMProt (third panel). The bottom panel depicts the size of each subclass. The subclasses are labeled according to their EC number and are ordered according to size.

order to gain additional insights regarding the predictive capabilities of motifs of length 6, we analyzed which of these motifs are unique of a single subclass. Statistics are presented in Fig. 3.

Evidently, motifs of length 6 are both abundant and, concomitantly, comprise a large fraction of motifs unique to a single subclass. Out of the 601 motifs of length 6, 493 are unique (e.g., belong to a single EC subclass at the second level). This group of 493 motifs are not sufficient for the classification task as their coverage of all proteins within their EC subclass is limited.

The 125 unique motifs of length 6 of the relatively large sub-subclass 1.1.1 (comprised of 1699 proteins) span only 63% of the protein sequences. Nonetheless, a classification task based on the entire set of motifs of length 6 obtains a

Jaccard score of  $0.91 \pm 0.03$ , hence the additional 108 nonunique MEX motifs span the rest of the protein sequences.

A level 3 classification task performed solely with motifs of length 6 yielded a J-score of  $0.89 \pm 0.08$ , which is lower than the J-score quoted in the previous section. Apparently, the space spanned by the 601 motifs of length 6 is not as comprehensive for this refined classification task, as the space spanned by the 1222 motifs of length 5 and longer.

Fig. 3 introduces an additional interesting insight, clarifying the reason an SVM analysis based on the set of MEX motifs longer than 4 results in relatively lower J-scores (Average J-scores are  $0.83\pm0.09$  for a level 2 classification task and  $0.83\pm0.14$  for a level 3 classification task). Apparently, the large fraction of non-unique motifs of length 5 impairs the predictive power of the unique motifs.

class	# of elements	MEX J	SW J
1.1.1	1699	$0.91\pm0.03$	$0.85\pm0.04$
1.1.99	59	$0.92\pm0.2$	$0.80\pm0.11$
1.10.2	69	$0.94\pm0.14$	$0.52\pm0.00$
1.10.3	38	$0.78\pm0.17$	$0.77\pm0.11$
1.11.1	310	$0.98\pm0.02$	$0.89\pm0.01$
1.12.99	26	$0.92\pm0.09$	$0.83\pm0.00$
1.13.11	112	$0.90\pm0.06$	$0.62\pm0.08$
1.14.11	47	$0.87\pm0.14$	$0.69\pm0.10$
1.14.13	101	$0.82\pm0.12$	$0.71\pm0.12$
1.14.14	233	$0.93\pm0.02$	$0.91\pm0.07$
1.14.15	38	$0.91\pm0.12$	$0.85\pm0.13$
1.14.16	28	$0.93 \pm 0.1$	$0.80\pm0.08$
1.14.19	26	$0.89\pm0.14$	$0.94\pm0.10$
1.14.99	72	$0.89\pm0.07$	$0.85\pm0.09$
1.15.1	233	$0.92\pm0.06$	$0.96\pm0.00$
1.16.1	21	1	$0.60\pm0.20$
1.17.4	113	$0.86\pm0.04$	$0.90\pm0.02$
1.18.1	47	$0.77\pm0.31$	$0.69 \pm 0.14$
1.18.16	123	$0.88\pm0.08$	$0.93\pm0.03$
1.2.1	512	$0.88\pm0.03$	$0.89\pm0.03$
1.2.4	66	$0.83\pm0.06$	$0.91\pm0.03$
1.3.1	156	$0.84 \pm 0.1$	$0.68\pm0.03$
1.3.3	139	$0.96\pm0.04$	$0.88\pm0.05$
1.3.5	18	1	1
1.3.99	73	$0.76\pm0.09$	$0.61\pm0.09$
1.4.1	83	$0.86\pm0.07$	$0.82\pm0.03$
1.4.3	89	$0.93\pm0.11$	$0.68\pm0.07$
1.4.99	31	$0.92\pm0.13$	$0.80\pm0.08$
1.5.1	167	$0.67\pm0.19$	$0.68\pm0.10$
1.6.1	21	1	$0.87 \pm 0.12$
1.6.2	20	$0.81\pm0.16$	$0.67\pm0.12$
1.6.5	814	$0.87\pm0.02$	$0.84 \pm 0.01$
1.6.99	177	$0.70\pm0.09$	$0.63 \pm 0.04$
1.7.1	58	$0.91\pm0.08$	$0.76\pm0.15$
1.7.2	26	1	$0.72 \pm 0.10$
1.7.99	43	1	$0.\overline{40 \pm 0.20}$
1.8.1	138	$0.91 \pm 0.03$	$0.\overline{86\pm0.04}$
1.8.4	137	$0.93\pm0.05$	$0.88\pm0.13$
1.9.3	552	$0.94\pm0.02$	$0.90\pm0.03$

Table 1. J-values derived from MEX and Smith-Waterman analysis, corresponding to level 3 classification tasks.

# Discussion

Applying the MEX algorithm on a group of 7095 enzymes, it has been shown that the extracted motifs form an excellent basis for classifying these enzymes into small classes known to have different functional roles. In particular, the classification from sequence to function based on these motifs of this enzymes class was demonstrated to outperform any of the alternative methods.



Figure 3. distribution of MEX motifs of lengths 5-10 according to their length. The three sets correspond to (left) entire set of MEX motifs, (middle) set of MEX motifs unique to a single level 2 subclasses and (right) set of MEX motifs unique to a single level 3 sub-subclass.

Our results are compared with two approaches: (i) Classification based on pairwise sequence similarity, analogous to the one employed by [7], using the same SVM procedure that was employed for MEX. As demonstrated, MEX derived motifs form a better basis for classification, indicating that MEX selected motifs improve the signal to noise ratio inherent in the original sequences. (ii) The SVMProt method introduced by [3, 4] on level 2 data (using their published results). Although their method is based on semantic information, i.e. physical and chemical properties of the sequence of amino-acids, the results obtained by MEX are better, again indicating that the MEX selected motifs carry relevant information.

It should be noted that the MEX based classification is accomplished by using only 1222 motifs of length 6 or longer. Moreover, similar results were obtained by using only the 601 motifs of length 6. Considering the 55 classification tasks for about 7000 proteins, the number of features allowing a successful classification by the MEX algorithm is surprisingly small. Furthermore, as opposed to the regular-expression motifs used by other methods, MEX motifs are all deterministic consecutive amino-acid sequences.

Such regular-expression motifs approach was presented

by [2]. They have used regular-expression motifs of average length of 21 amino-acids (termed eMOTIFs) derived in a supervised manner. Applying a feature-selection procedure to select approximately 1000 eMOTIFs out of their original very large set of eMOTIFs, they have achieved impressive classification results. However, while the small number of selected eMOTIFs is comparable to the 1222 motifs used by our approach, it should be noted that the deterministic, consecutive motif sequences extracted by MEX spans a much smaller sequence space than the one spanned by the eMOTIFs, yet, achieving successful classification. Unfortunately, a direct comparison with this work could not be made due to insufficient data.

The application of the MEX algorithm studied here applies only a single level of feature extraction. Higher level patterns may be extracted by iteratively applying MEX, where each MEX iteration uses the observed sequence-motifs as vertices in the MEX graph. Moreover, utilizing the full extent of the ADIOS approach [13] may further reveal higher syntactic structures in biological sequence data.

## Acknowledgment

This research was partially supported by the US Israel Binational Science Foundation. ZS has a student fellowship of the Horowitz Foundation for Complexity Sciences. We thank Asa Ben-Hur and Doug Brutlag for helpful conversations. Additionally, we thank Elhanan Borenstein for his helpful remarks.

## References

- Ben-Hur,A., Brutlag, D. (2003) Remote homology detection: a motif based approach. *Bioinformatics*, 19, Suppl. 1, i26-i33.
- [2] Ben-Hur,A., Brutlag, D. (2004) Sequence motifs: highly predictive features of protein function. *Neural Information Processing Systems 2004.*
- [3] Cai,C. Z., Han,L. Y., Ji,Z. L., Chen, Y. Z. (2003) SVM-Prot: web-based support vector machine software for functional classification of a protein from its primary sequence. *Nuclear Acids Research*, **31**, 3692-3697.
- [4] Cai,C. Z., Han,L. Y., Ji,Z. L., Chen, Y. Z. (2003) Enzyme family classification by support vector machines. *PROTEINS: Structure, Function and Bioinformatics*, 55, 66-76.
- [5] Durbin, R., Eddy, S., Krogh, A., Mitchison, G. (1998). Biological sequence analysis Probabilistic models of proteins and nucleic acids, Cambridge University Press.
- [6] des Jardin, M., Karp, P. D., Krummenacker, M., Lee, T. J. and Ouzounis, C. A. (1997). Prediction of enzyme classification from protein sequence without the use of sequence similarity. Proceedings of ISMB.

- [7] Liao, L., Noble, W. S., (2003) Combining pairwise sequence analysis and support vector machines for detecting remote protein evolutionary and structural relationships. *J. of Comp. Biology*, **10**, 857-868.
- [8] Huang, J. Y., Brutlag, D. L., (2001) The eMOTIF database. Nuclear Acids research, 29, 202-204.
- [9] Neville-Manning, C. G., Wu, T. D., Brutlag, D. L., (1998) Highly specific protein sequence motifs for genome analysis. *Proc. Natl. Acad. Sci. USA* 95, 5865-5871.
- [10] Schölkopf, B., (1997) Support Vector Learning. R. Oldenburg Verlag, Munich.
- [11] Smith, T., Waterman, M., (1981) Identification of common molecular subsequences. J. of Mol. Biology 147, 195-197.
- [12] Solan, Z., Ruppin, E., Horn, D., Edelman, S., (2003) Automatic acquisition and efficient representation of syntactic structures. In S. Becker, S. Thrun and K. Obermayer, editors, *Advance in Neural Information Processing Systems* 15, 91-98, MIT Press, Cambridge, MA.
- [13] Solan, Z., Horn, D., Ruppin, E., Edelman, S., (2004) Unsupervised context sensitive language acquisition from a large corpus. In Sebastian Thrun and Lawrence Saul and Bernhard Schölkopf, editors, *Advances in Neural Information Processing Systems* bf 16 MIT Press, Cambridge, MA.
- [14] Tian, W., Skolnick, J., (2003) How well is enzyme function conserved as a function of pairwise sequence identity? *J. Mol. Biol.* 333, 863 - 882.
- [15] Vapnik, V. (1995) The Nature of Statistical Learning Theory. Springer Verlag, NY.