

Journal of Bioinformatics and Computational Biology
© Imperial College Press

Profile-based string kernels for remote homology detection and motif extraction

Rui Kuang¹, Eugene Ie^{1,3}, Ke Wang¹, Kai Wang², Mahira Siddiqi²,
Yoav Freund^{1,3,4}, Christina Leslie^{1,3,4*}

¹*Department of Computer Science*, ²*Department of Biomedical Informatics*,

³*Center for Computational Learning Systems*,

⁴*Center for Computational Biology and Bioinformatics*
Columbia University

Received (September 01, 2004)

Revised (Day Month Year)

Accepted (Day Month Year)

We introduce novel profile-based string kernels for use with support vector machines (SVMs) for the problems of protein classification and remote homology detection. These kernels use probabilistic profiles, such as those produced by the PSI-BLAST algorithm, to define position-dependent mutation neighborhoods along protein sequences for inexact matching of k -length subsequences (“ k -mers”) in the data. By use of an efficient data structure, the kernels are fast to compute once the profiles have been obtained. For example, the time needed to run PSI-BLAST in order to build the profiles is significantly longer than both the kernel computation time and the SVM training time. We present remote homology detection experiments based on the SCOP database where we show that profile-based string kernels used with SVM classifiers strongly outperform all recently presented supervised SVM methods. We further examine how to incorporate predicted secondary structure information into the profile kernel to obtain a small but significant performance improvement. We also show how we can use the learned SVM classifier to extract “discriminative sequence motifs”—short regions of the original profile that contribute almost all the weight of the SVM classification score—and show that these discriminative motifs correspond to meaningful structural features in the protein data. The use of PSI-BLAST profiles can be seen as a semi-supervised learning technique, since PSI-BLAST leverages unlabeled data from a large sequence database to build more informative profiles. Recently presented “cluster kernels” give general semi-supervised methods for improving SVM protein classification performance. We show that our profile kernel results also outperform cluster kernels while providing much better scalability to large datasets.

Supplementary website: <http://www.cs.columbia.edu/compbio/profile-kernel>.

Keywords: protein classification; support vector machine; kernels; protein motifs.

1. Introduction

There has been much recent work on support vector machine (SVM)⁴ approaches for the classification of protein sequences into functional and structural families and for remote

*Corresponding author. Mailing address: 1214 Amsterdam Ave, MC 0401, New York, NY 10027. Email: cleslie@cs.columbia.edu. Telephone: 1-212-939-7043. Fax: 1-212-666-0140

2 Rui Kuang, Eugene Ie, Ke Wang, Kai Wang, Mahira Siddiqi, Yoav Freund and Christina Leslie

homology detection. Most of this research effort focuses on finding useful representations of protein sequence data for SVM training, either using explicit feature vector representations or *kernel* functions—specialized sequence similarity functions that define an inner product in an implicit feature space for the SVM optimization problem. Among the approaches that have been presented are the Fisher-SVM method¹³, which represents each protein sequence as a vector of Fisher scores extracted from a profile hidden Markov model (HMM) for a protein family, and kernels that extend the Fisher kernel method³⁰; families of efficient string kernels^{22,21,23}, such as the mismatch kernel, which are based on inexact-matching occurrences of k -length subsequences (“ k -mers”); the SVM-pairwise approach²⁴, which uses a feature vector of pairwise alignment scores between the input sequence and a set of training sequences; the eMOTIF kernel³, where the feature vector represents counts of occurrences of eMOTIF patterns in the sequence; feature vectors defined by structure-based I-sites motifs¹⁰; and string alignment kernels²⁹, which detect sequence similarity by approximating the behavior of the Smith-Waterman score. These studies show that most of the methods achieve comparable classification performance on benchmark datasets, though there are significant differences in computational efficiency²¹. Interestingly, except for the Fisher kernel method and its extensions, these representations do not make intrinsic use of standard tools for protein sequence analysis such as profiles⁹ and profile HMMs^{19,6,2}—more commonly, they use scores based on alignment or probabilistic models to construct a large set of features. It is perhaps surprising that very general k -mer based string kernels perform as well as the Fisher kernel approach, which makes well-motivated use of profile HMMs²¹.

In this paper, we define a natural extension of the k -mer based string kernel framework to define kernels on protein sequence profiles, such as those produced by PSI-BLAST¹. We choose to use profiles (rather than more complex models) because they can be calculated by PSI-BLAST in a tractable amount of time and because, once the profiles are obtained, we can efficiently compute string kernel values using an appropriate data structure; in fact, the time needed to compute the profile kernel matrix and the SVM training time are significantly shorter than the time needed by PSI-BLAST to compute profiles. From a machine learning point of view, use of PSI-BLAST profiles can be viewed as a *semi-supervised* approach—that is, a method that learns both from labeled training examples (sequences whose structural classification is known) and unlabeled examples—an important consideration given the relatively small amount of labeled data in this problem. Through iterative heuristic alignment, PSI-BLAST leverages unlabeled data from a large sequence database to obtain a much richer profile representation of each sequence. Intuitively, this richer data representation, made available to an SVM through a profile-based kernel, should greatly improve classification performance. Also, profile-based kernels are a significantly different semi-supervised approach than the Fisher-SVM method: with the Fisher kernel, unlabeled data in the form of domain homologs are used to train a model for a protein family of sequences in the training set, and then each sequence is represented by sufficient statistics with respect to the learned model; in our approach, unlabeled data is used to produce a profile model for each training sequence independently, and then the kernel is defined on the profiles. Our experimental results for the remote homology detection task, using

a benchmark based on the SCOP database, show that our profile-based kernel used with SVM classifiers strongly outperform all the recent purely supervised SVM methods that we compared against.

There have been several attempts to use predicted local structure, including secondary structure, for improving remote homology detection^{18,7}, generally resulting in modest performance improvement. In this paper, we discuss one variation of the profile kernel that incorporates additional predicted secondary structure profiles to help remote homology detection. Our experimental results show that secondary structure profiles can help the profile kernel achieve better performance, but given the current prediction accuracy for secondary structures, we obtain only limited improvement.

Usually, SVM methods are treated as a “black box” method, since in general it is difficult to interpret the SVM classification rule. For the case of profile string kernels, we show how we can use the trained SVM classifiers to define positional scores along the protein profiles that define a smoothed contribution to the positive classification decision. In general, we find that a low percentage of positions in the profile is responsible for a high percentage of the positive contribution to the total discrimination score for positive training sequences, and thus we can extract distinguished regions that we call “discriminative sequence motifs”. We give examples from our SCOP dataset to show that these discriminative motifs correspond to meaningful structural features, giving a proof of principle that the SVM-profile kernel approach allows us to extract useful sequence information.

Recently presented “cluster kernels” approaches³¹ give general semi-supervised methods for improving SVM protein classification performance of a base kernel using unlabeled data together with a similarity measure on input examples. These cluster kernels were successfully applied to the protein classification problem using the mismatch kernel as a base kernel for sequence data and BLAST or PSI-BLAST to define similarity scores. However, for large amounts of unlabeled data, these more general methods do not scale as well as our profile kernel approach. We show that our profile kernel results also outperform cluster kernels while providing much better scalability to large datasets.

The current paper is an expanded version of work that originally appeared in a conference proceedings²⁰. We have added new results on incorporating predicted secondary structure information into the profile kernel as well as additional examples of superfamilies where our extracted discriminative motif regions correspond to conserved structural features. We also found and corrected a small bug in our previously reported results²⁰: there was an error in mapping one of the amino acids in our profile kernel code, resulting in a kernel that was mathematically valid but deviated slightly from the one we defined in the text. The corrected classification results that we report here are stronger than those we obtained previously. Statistics and results relating to our discriminative motif regions, in particular the overlap with eMOTIF and I-sites motifs and with structural features in our case studies, have also slightly changed but display the same trends as before. Finally, we add experiment-by-experiment results on the percentage of positions in the positive training sequences needed to account for 90% of the total positive contribution to the discriminant score. Again, the trend that was previously reported—that relatively few positions account for most of the positive discrimination—was correct, but the mean percentage reported was

4 Rui Kuang, Eugene Ie, Ke Wang, Kai Wang, Mahira Siddiqi, Yoav Freund and Christina Leslie

inaccurate, and there is considerable variation across superfamilies in coverage by discriminative motif regions.

2. The Profile Kernel

A key feature of the SVM optimization problem is that it depends only on the inner products of the feature vectors representing the input data, allowing us to use *kernel techniques*. If we define a feature map Φ from the input space of protein sequences into a (possibly high-dimensional) vector space called the *feature space*, we obtain a *string kernel*—that is, a kernel on sequence data—defined by $K(x, y) = \langle \Phi(x), \Phi(y) \rangle$.

We first show how to define a feature mapping for protein sequence profiles—more precisely, we consider input examples to be profiles $P(x)$, where x is a sequence $x = x_1x_2 \dots x_N$ from the alphabet Σ of amino acids ($|\Sigma| = 20$, and the length $N = |x|$ depends on the sequence), and $P(x) = \{p_i(a), a \in \Sigma\}_{i=1}^N$ is a profile for sequence x , with $p_i(a)$ denoting the emission probability of amino acid a in position i and $\sum_{a \in \Sigma} p_i(a) = 1$ for each position i . We then show how to efficiently and directly compute the profile-based string kernel values $K(P(x), P(y))$ without storing the feature vector representation.

2.1. Profile-defined Mapping to k -mer Feature Space

Following the framework of k -mer based string kernels^{22,21,23}, our profile-based kernels will depend on a feature mapping to the $|\Sigma|^k$ -dimensional feature space indexed by the set of all possible k -length subsequences (“ k -mers”) of amino acids, where k is a small positive integer.

Previous string kernels relied on defining an inexact-matching neighborhood of k -mers around each k -length contiguous subsequence in the input sequence. For example, for the (k, m) -mismatch kernel, one defines the “mismatch neighborhood” around k -mer $\alpha = a_1a_2 \dots a_k$ to be the set of all k -length sequences β from Σ that differ from α by at most m mismatches. For a k -mer α , the mismatch feature map is defined as

$$\Phi_{(k,m)}^{\text{Mismatch}}(\alpha) = (\phi_\beta(\alpha))_{\beta \in \Sigma^k}, \quad (1)$$

where $\phi_\beta(\alpha) = 1$ if β belongs to $N_{(k,m)}(\alpha)$, and $\phi_\beta(\alpha) = 0$ otherwise, and one extends additively to full-length sequences x by summing the feature vectors for all the k -mers in x :

$$\Phi_{(k,m)}^{\text{Mismatch}}(x) = \sum_{k\text{-mers } \alpha \text{ in } x} \Phi_{(k,m)}^{\text{Mismatch}}(\alpha). \quad (2)$$

Thus each coordinate of the feature map is a count of the inexact-matching occurrences of a particular k -mer, where mismatching is used to define inexact matching.

For the profile kernel, we use the probabilistic profile $P(x)$ to define a mutation neighborhood for each k -length segment in the input sequence x . Therefore, unlike previous string kernels, the inexact-matching neighborhood k -mers are not the same for all the data but instead vary from sequence to sequence and within different regions of the same sequence. For each k -length contiguous subsequence $x[j+1 : j+k] = x_{j+1}x_{j+2} \dots x_{j+k}$

in x ($0 \leq j \leq |x| - k$), the *positional mutation neighborhood* is defined by the corresponding segment of the profile $P(x)$:

$$M_{(k,\sigma)}(P(x[j+1:j+k])) = \{\beta = b_1 b_2 \dots b_k : -\sum_{i=1}^k \log p_{j+i}(b_i) < \sigma\}. \quad (3)$$

Note that the emission probabilities $p_{j+i}(b)$, $i = 1 \dots k$, come from the profile $P(x)$ —for notational simplicity, we do not explicitly indicate the dependence on x . Typically, the profiles are estimated from close homologs found in a large sequence database and may be too restrictive for our purposes. Therefore, we smooth the estimates using background frequencies $q(b)$, $b \in \Sigma$, of amino acids in the training dataset via

$$\tilde{p}_i(b) = \frac{p_i(b) + Cq(b)}{1 + C}, i = 1 \dots |x|, \quad (4)$$

where C is a smoothing parameter, and we use the smoothed emission probabilities $\tilde{p}_i(b)$ in place of $p_i(b)$ in defining the mutation neighborhoods.

We now define the profile feature mapping as

$$\Phi_{(k,\sigma)}^{\text{Profile}}(P(x)) = \sum_{j=0 \dots |x|-k} (\phi_{\beta}(P(x[j+1:j+k])))_{\beta \in \Sigma^k}, \quad (5)$$

where the coordinate $\phi_{\beta}(P(x[j+1:j+k])) = 1$ if β belongs to the mutation neighborhood $M_{(k,\sigma)}(P(x[j+1:j+k]))$, and otherwise the coordinate is 0.

The profile kernel is simply defined by the inner product of feature vectors:

$$K_{(k,\sigma)}^{\text{Profile}}(P(x), P(y)) = \langle \Phi_{(k,\sigma)}^{\text{Profile}}(P(x)), \Phi_{(k,\sigma)}^{\text{Profile}}(P(y)) \rangle. \quad (6)$$

2.2. Efficient Computation of the Kernel Matrix

Rather than storing sparse feature vectors in high-dimensional k -mer space, we directly and efficiently compute the kernel matrix using a *trie* data structure, similar to the mismatch tree approach presented in our previous work^{22,21,23}. The difference for the profile kernels is that instead of matching k -mers along the path to a leaf, we pass k -length profiles down the tree branches.

Our new (k, σ) -profile trie is a rooted tree of depth k where each internal node has $|\Sigma| = 20$ branches, each labeled with an amino acid (symbol from Σ). A leaf node still represents a fixed k -mer in our feature space, obtained by concatenating the branch symbols along the path from root to leaf. We perform a depth-first traversal of the data structure and store, at a node of depth d , a set of pointers to all k -length profiles $P(x[j+1:j+k])$ from the sample data set, whose current cumulative substitution scores, up to depth d , are less than the σ threshold, that is, $-\sum_{i=1}^d \log p_{j+i}(b_i) < \sigma$, where $b_1 \dots b_d$ is the prefix of the current node. As we pass from a parent node at depth d to a child node at depth $d+1$ along a branch with symbol label b , we add for each k -length profile $P(x[j+1:j+k])$ a score $-\log p_{j+d+1}(b)$. Only those profile segments whose cumulative substitution scores are still less than σ will be passed to the child node. At the leaf node, we update the kernel by computing the contribution of active profile segments to the corresponding k -mer feature.

6 Rui Kuang, Eugene Ie, Ke Wang, Kai Wang, Mahira Siddiqi, Yoav Freund and Christina Leslie

The complexity of computing each value $K(x, y)$ depends on the size of the positional mutation neighborhood of k -length profiles. If M_σ represents the maximum number of k -mers found in any mutation neighborhood defined by a k -length profile from the input data for threshold σ , then we can bound the kernel computation complexity by $O(kM_\sigma(|x| + |y|))$. In our experiments, we choose σ so that the typical mutation neighborhood defined by a k -length profile allows about $m = 1$ or 2 mismatches relative to the original k -mer. Thus we can estimate that the running time is bounded by that of the (k, m) -mismatch kernel, which is $O(k^{m+1}|\Sigma|^m(|x| + |y|))$, with $m \leq 2$. More details on the complexity analysis for k -mer based string kernels can be found in Leslie and Kuang²³. See Section 3 for actual running times in benchmark experiments.

2.3. Extending the Profile Kernel with Secondary Structure Information

A natural approach for including secondary structural information in the profile kernel is to associate with each sequence an additional secondary structure profile, using the true secondary structures for the training set and predicted secondary structure profiles—produced with existing methods such as PSI-PRED¹⁶ and PHD²⁸—for the test set. Thus each k -length segment of the sequence profile has an associated k -length secondary structure profile over symbols from Σ_{str} , the alphabet of secondary structure elements (typically the three symbol alphabet representing alpha helix, beta sheet, and coil). We expand the original feature space to a new one, with dimension $k \times |\Sigma_{str}|$ times larger, by associating to every original k -mer feature β a vector of $k \times |\Sigma_{str}|$ features in the new space. When a k -mer β falls within the mutation neighborhood of a particular k -length segment of the sequence profile, we add the vectorized secondary structure profile for this segment to the $k \times |\Sigma_{str}|$ vector of features associated with k -mer β . In this way, the feature with index (j, s) associated to the k -mer β ($j = 1 \dots k$, $s = 1 \dots |\Sigma_{str}|$) is an expected count of occurrences of secondary structure symbol s in position j of a k -length segment, over all segments such that the k -mer β is in mutation neighborhood defined by the sequence profile. Let $P_{str}(x)$ represent the probabilistic profile of secondary structures. A formal definition of the extended profile kernel is given by

$$\Phi_{(k,\sigma)}^{\text{Profile-Str}}(P(x)) = \sum_{j=0 \dots |x|-k} (\phi_\beta^{str}(P(x[j+1:j+k])))_{\beta \in \Sigma^k}, \quad (7)$$

where the vector-valued coordinate function $\phi_\beta^{str}(P(x[j+1:j+k])) = P_{str}(x[j+1:j+k])$ if β belongs to the mutation neighborhood $M_{(k,\sigma)}(P(x[j+1:j+k]))$, and otherwise is a 0 vector of length $k \times |\Sigma_{str}|$.

The computation of the extended profile kernel is the same as before except that at leaf nodes in the trie traversal, we increment the kernel value with the dot product between the vectorized structural profiles of each pair of instance k -mers. This modification introduces an additional multiplicative constant of $k \times |\Sigma_{str}|$ to the computational complexity of original profile kernel.

One can also consider other strategies for adding structural information. For example, one can put a probabilistic threshold on the k -length structure profiles and redefine the

mutation neighborhood to be the set of pairs of k -mers and k -length structure sequences that satisfy a substitution probability both for the sequence profile and the secondary structure profile. However, two problems may arise: the double threshold approach can cause a diagonally dominant gram matrix, and one introduces an additional parameter-dependant constant to the computational complexity, which may make the kernel too expensive for practical use.

2.4. Extraction of Discriminative Motifs

Using the PSI-BLAST sequence profiles and the learned SVM weights, we can do a positional analysis to determine which regions of the positive training sequence contribute most to the classification score and thus extract “discriminative” protein motif regions. For a training set of protein sequence $\{x_i\}_{i=1}^n$, the normal vector to the SVM decision hyperplane is given by

$$\mathbf{w} = \sum_{i=1}^n y_i c_i \Phi_{(k,\sigma)}^{\text{Profile}}(P(x_i)), \quad (8)$$

where the c_i are learned weights and $y_i \in \{\pm 1\}$ are training labels. For each k -length profile segment of sequence x , its contribution to the classification score is (up to a constant):

$$S(x[j+1:j+k]) = \langle \phi_{(k,\sigma)}^{\text{Profile}}(P(x[j+1:j+k])), \mathbf{w} \rangle. \quad (9)$$

We are mainly interested in discriminative motifs that contribute to the positive decision of the classifier, so we define a positional score for each position j in a (positive) training sequence by summing up positive contributions of k -length segments containing the position:

$$\sigma(x[j]) = \sum_{q=1}^k \max(S(x[j-k+q:j-1+q]), 0). \quad (10)$$

We now sort these positional scores (for all positions in all positive training sequences) in decreasing order $\sigma(x[j_1]) \geq \sigma(x[j_2]) \geq \dots \geq \sigma(x[j_N])$, and we find the first index M such that cumulative sum $\sum_{i=1}^M \sigma(x[j_i])$ is greater than 0.9 times the total sum $\sum_{i=1}^M \sigma(x[j_i])$. Thus positions j_1, \dots, j_M constitute 90% of the positionally averaged positive classification scores. We will see in the Section 3.3 that these positions tend to fall in short segments of the protein sequence; we call these segments “discriminative motif regions”.

3. Experiments

We test SVM classification performance of profile-based string kernels against other recently presented SVM methods on a SCOP benchmark dataset. Methods are evaluated on the ability to detect members of a target SCOP family (positive test set) belonging to the

8 *Rui Kuang, Eugene Ie, Ke Wang, Kai Wang, Mahira Siddiqi, Yoav Freund and Christina Leslie*

same SCOP superfamily as the positive training sequences; no members of the target family are available during training. We use the same experimental set-up that has been used in several previous studies of remote homology detection algorithms^{13,24}.

We use the same 54 target families and the same test and training set splits as in the remote homology experiments in Liao and Noble²⁴. The sequences are 7329 SCOP domains obtained from version 1.59 of the database after purging with `astral.stanford.edu` so that no pair of sequences share more than 95% identity. Compared to Liao and Noble²⁴, we reduce the number of available labeled training patterns by roughly a third. Data set sequences that were neither in the training nor test sets for experiments from Liao and Noble are considered to be additional unlabeled data, used for cluster kernel method we compare against. All methods are evaluated using the receiver operating characteristic (ROC) score and the ROC-50, which is the ROC score computed only up to the first 50 false positives⁸.

We computed the profiles needed for our kernels by running PSI-BLAST¹ from the nonredundant protein database with default search parameters and with background frequencies, used for smoothing, estimated from the full dataset of 7329 SCOP domains. We tried two options for the maximum number of iterative database searches, 2 iterations and 5 iterations, to show the tradeoff between computational efficiency and classification performance. We used smoothing parameter corresponding to $\frac{1}{1+C} = .8$ in the profile kernel computation. The time needed to compute PSI-BLAST profiles for all sequences was approximately 36 hours on a 2.2 GHz Linux server using at most 2 iterative database searches or about 3 days using 5 iterative searches; on the same CPU, the time required to compute the 7329 x 7329 kernel matrix was 10 hours, and all 54 SVM experiments were completed in 30 minutes.

3.1. SCOP Experiments: Comparison with Supervised and Semi-Supervised Methods

We compared the results of profile kernels with five recently presented SVM methods, using different representations of protein sequence data—the eMOTIF kernel³, the SVM-pairwise method²⁴, the mismatch kernel²¹, the string alignment kernel²⁹, and the Fisher kernel¹³—as well as recent semi-supervised cluster kernel methods³¹. We also compared the SVM methods to PSI-BLAST, used directly as a method for ranking test sequences relative to positive training sequence queries (see below).

We evaluated the first three SVM methods—the eMOTIF kernel, SVM-pairwise, and the mismatch kernel—on the SCOP 1.59 benchmark dataset described above. We used the eMOTIF database extracted from eBlocks and packaged with eBAS version 3.7^{11,26}, and we obtained code for computing eMOTIF feature vectors from the authors³. For the SVM-pairwise method, we used PSI-BLAST E-values as pairwise similarity scores (see Weston et al.³¹ for details on this representation). We note that this use of PSI-BLAST with the SVM-pairwise method is not fully-supervised, because the PSI-BLAST scores themselves make use of unlabeled data. For the mismatch kernel, we use $(k, m) = (5, 1)$ as presented in the original paper²².

We include results for PSI-BLAST, used directly as a ranking method, in order to pro-

vide a baseline comparison with a widely used remote homology detection method and also to demonstrate the added benefit of combining PSI-BLAST with our SVM string kernel approach. The PSI-BLAST algorithm, which iteratively builds up a probabilistic profile for a query sequence by searching a large database, also gives a faster approximation to the iterative training method of profile HMMs. (We do not test profile HMMs here due to computational expense, but in previous benchmark results for the remote homology problem, SVM string kernel and Fisher kernel methods were both found to outperform profile HMMs^{22,13}.) Since PSI-BLAST is not a family-based method, we report results by averaging over queries: for each experiment, we use PSI-BLAST with each of the positive training sequences as the query and search against the nonredundant protein database in order to produce a set of profiles, and then we use these profiles to rank the test set sequences by their PSI-BLAST E-values. The ROC (ROC-50) score that we report for the experiment is the average of all ROC (ROC-50) scores from these rankings. (We note that a more sophisticated PSI-BLAST training procedure that uses all positive training sequences at once might be possible, but it is not clear how best to do this given the diverse positive training set.) For the PSI-BLAST ranking method, we use PSI-BLAST with the default parameters, allowing a maximum of 10 iterative searches against the nonredundant protein database in order to build the profiles.

In our main experiments, we computed the profile kernel with two sets of PSI-BLAST profiles, one set built using at most 5 iterative searches for better accuracy and the other set using 2 iterative searches for reduced PSI-BLAST computational cost. We tested profile kernels with $(k, \sigma) = (4, 6.0), (5, 7.5)$ and $(6, 9.0)$. These parameters were chosen using the following heuristics: we took the same range of $k = 4 \dots 6$ that we found useful in our previous string kernel work, and we chose the parameter σ to allow one or two mismatches from the input k -mer in a typical k -length profile. All three parameter choices yield similar results. Note that we did not try to exhaustively optimize parameter choices, since the benchmark dataset does not include a cross-validation set. Figure 1 shows the comparison of SVM performance of the $(5, 7.5)$ -profile kernel against the PSI-BLAST ranking method, the eMOTIF kernel, the mismatch kernel, and the SVM-pairwise method using PSI-BLAST across the 54 experiments in the benchmark. A signed rank test with Bonferroni correction for multiple comparisons concludes that the profile kernel significantly outperforms the mismatch kernel (p-value $1.3e^{-09}$), SVM-pairwise kernel ($7.2e^{-09}$), eMOTIF kernel ($1.3e^{-09}$), and mean PSI-BLAST ranking ($1.1e^{-09}$). Average ROC and ROC-50 scores across the experiments for all methods are reported in Table 1. Although running PSI-BLAST for 5 iterations instead of 2 iterations increased the PSI-BLAST computation time by a multiplicative factor, the results demonstrate that we can have significant improvement in ROC and ROC-50 scores for the profile kernel method by improving the profiles. In our subsequent motif analysis, we refer to the second set of SVM classifiers, which use profiles based on up to 5 PSI-BLAST iterations.

We note that the original authors of the SVM-pairwise used Smith-Waterman scores (SW) for pairwise comparison scores; however, on a similar benchmark with more training data than the current dataset, results for SVM-pairwise with SW scores were weaker than the PSI-BLAST results reported here, and ROC performance was only slightly better (ROC

10 Rui Kuang, Eugene Ie, Ke Wang, Kai Wang, Mahira Siddiqi, Yoav Freund and Christina Leslie

= 0.893, ROC-50 = 0.434). Thus the semi-supervised PSI-BLAST scores do indeed give a richer and more effective representation for SVM-pairwise; however, using PSI-BLAST profiles to define a profile-based string kernel is clearly more effective than SVM-pairwise with PSI-BLAST.

Our SCOP dataset is different from and larger than the benchmark on which the eMOTIF kernel was originally tested³. In cases where a superfamily has a common eMOTIF pattern or set of patterns, the eMOTIF kernel should achieve good specificity. We speculate that in our 54 experiments, fewer superfamilies are characterized by common eMOTIF patterns and that accordingly the eMOTIF kernel achieves weaker performance.

To compare with the two remaining kernel representations—the Fisher kernel¹³ and the string alignment kernel²⁹—we also tested the (5,7.5)-profile kernel on a second benchmark dataset, which is derived from an earlier version of the SCOP database (SCOP 1.53). In Table 2, we report the average ROC and ROC-50 scores for the profile kernel with published results for the other two methods²⁹. The (5,7.5)-profile kernel shows significantly stronger performance on this dataset over the Fisher kernel (p-value $1.8e^{-10}$) and the string alignment kernel (p-value $2.3e^{-07}$).

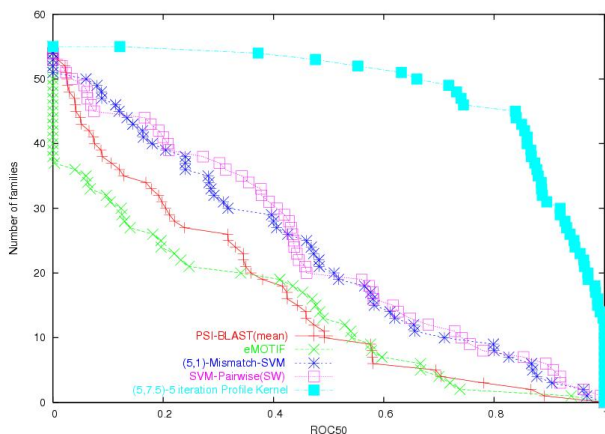


Fig. 1. Comparison of recent SVM-based homology detection methods for the SCOP 1.59 benchmark dataset. The graph plots the total number of families for which a given method exceeds an ROC-50 score threshold. Each series corresponds to one of the homology detection methods described in the text.

Finally, we also compared our profile kernel against recently presented cluster kernel methods³¹. These methods use “clustering” of additional unlabeled sequence data to improve the base representation. Here, sequences from the original SCOP dataset of 7329 domains that are not used in the training or test sets of any experiment provide the unlabeled data. For simplicity, we give results for only one of the two novel cluster kernel methods from Weston et al.³¹, the neighborhood kernel. (Results for the bagged kernel are very similar but more time-consuming to compute.) The neighborhood kernel uses the (5,1)-mismatch kernel as the base kernel and uses PSI-BLAST to define “neighborhood

sets' $\text{Nbd}(x)$ around each input sequence x , consisting of labeled or unlabeled sequences x' with similarity score to x below E-value threshold of 0.05, together with x itself. Then the implicit feature vector is $\Phi_{\text{nb}}(x) = \frac{1}{|\text{Nbd}(x)|} \sum_{x' \in \text{Nbd}(x)} \Phi_{\text{Mismatch}}(x')$.

We see from figure 2 that the profile kernel outperforms the neighborhood kernel (the preference to profile kernel is significant by a signed rank test with p-value threshold of $1.73e^{-05}$). We also note that our profile kernel is making use of more unlabeled data than the neighborhood kernel, since the neighborhoods are based on a smaller unlabeled database. However, as we scale up, computing the neighborhood kernel for extremely large neighborhood sets of sequences becomes expensive (computation time scales linearly with the size of the neighborhood). One can randomly select sequences from the neighborhood, but then one still has to devise an appropriate way of computing a sample without storing many thousands of sequences. (The bagged kernel from Weston et al. ³¹ has similar scalability issues as the database gets large.) By comparison, the profile-based string kernel approach achieves good SVM performance and computational efficiency while only representing the sequence profiles.

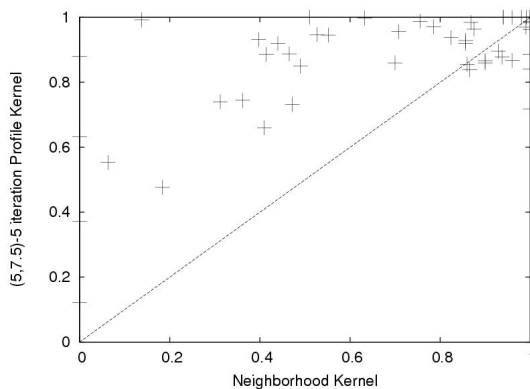


Fig. 2. Comparison of profile kernel (using 2 PSI-BLAST iterations) with recent cluster kernel approaches on the SCOP 1.59 benchmark dataset. The graph plots ROC-50 scores of the profile kernel (y-axis) versus the neighborhood kernel (x-axis), a recent cluster kernel method, for the 54 experiments in the SCOP benchmark.

3.2. Incorporating Secondary Structure into the Profile Kernel

To test the extended version of the profile kernel that incorporates secondary structure profiles, we performed one more experiment on the same SCOP dataset with 5-iteration PSI-BLAST sequence profiles. The true secondary structures are parsed from PDB formatted files with the DSSP program ¹⁷. The predicted secondary structure profiles are produced with PSI-PRED ¹⁶. For the training set, true secondary structures were encoded with a bit representation, where the corresponding match of a secondary structure element has a probability 1 and other probabilities are 0. Both true and predicted secondary structure pro-

12 *Rui Kuang, Eugene Ie, Ke Wang, Kai Wang, Mahira Siddiqi, Yoav Freund and Christina Leslie*

Table 1. Mean ROC and ROC-50 scores over 54 target families (SCOP version 1.59 benchmark).

Kernel	ROC	ROC-50
eMOTIF	0.711	0.247
PSI-BLAST(mean)	0.743	0.293
Mismatch(5,1)	0.870	0.416
SVM-pairwise(PSI-BLAST)	0.866	0.533
Neighborhood	0.923	0.699
Profile(4,6.0)-2 iterations	0.962	0.767
Profile(5,7.5)-2 iterations	0.973	0.821
Profile(6,9.0)-2 iterations	0.974	0.814
Profile(4,6.0)-5 iterations	0.974	0.837
Profile(5,7.5)-5 iterations	0.984	0.874
Profile(6,9.0)-5 iterations	0.987	0.866

Table 2. Mean ROC and ROC-50 scores over 54 target families (SCOP version 1.53 benchmark).

Kernel	ROC	ROC-50
Profile(5,7.5)-5 iterations	0.971	0.796
String alignment	0.923	0.661
Fisher	0.773	0.25

Table 3. Mean ROC and ROC-50 scores over 54 target families (SCOP version 1.59 benchmark).

Kernel	Original Profile		With Structure Profile	
	ROC	ROC-50	ROC	ROC-50
Profile(4,6.0)-5 iterations	0.975	0.834	0.983	0.861
Profile(5,7.5)-5 iterations	0.985	0.873	0.989	0.883
Profile(6,9.0)-5 iterations	0.987	0.863	0.989	0.869

files were also smoothed by a background frequency over secondary structure elements, as in equation 4. Results of adding secondary structure profile to original profile kernels of different parameters^a are shown in Table 3. With p-value threshold 0.05, a signed rank test suggests the improvement of the extended profile kernel over the original kernel is small but significant for the (4,6.0)-profile kernel (p-value 0.0023) but insignificant for (5,7.5)-profile kernel and (6,9.0)-profile kernel. Given the current accuracy of secondary structure prediction, we find that combining sequence information and local structural information leads only to modest improvements in remote homology detection.

3.3. Motif Extraction from SVM Predictions

We next calculated positional contribution scores $\sigma(x[j])$ for our trained SVM classifiers, as outlined in Section 2.4, to analyze which parts of the positive training sequences were most important for positive classification. Typically, we found peaky distributional plots of $\sigma(x[j])$ along positive training sequences, as shown for one experiment in Figure 3: the peaks in these plots correspond to “discriminative motif regions”. In Table 4 we show by

^aDue to post-processing of secondary structure, there are slight discrepancies between the k -length segments of SCOP sequences used in the two versions of the experiments. In the secondary structure experiments, we used only those sequence regions for which secondary structure can be determined, and we reproduced all results with the original profile kernels for consistency.

Table 4. Coverage of positive training sequences by discriminative motif regions in 54 experiments. Each experiment is identified by its SCOP target (test) family (SCOP version 1.59 identifiers). The coverage statistic indicates the fraction of positions in the positive training sequences that contribute 90% of the positive SVM discrimination score. Experiments for which the coverage statistic is small correspond to positive training superfamilies where the discriminative sequence information is concentrated in a small percentage of sequence positions.

SCOP Target	Coverage	SCOP Target	Coverage	SCOP Target	Coverage	SCOP Target	Coverage
a.26.1.1	0.701	b.1.1.5	0.583	c.1.8.1	0.317	c.47.1.1	0.421
a.26.1.2	0.659	b.10.1.2	0.583	c.1.8.3	0.336	c.47.1.10	0.444
a.35.1.2	0.665	b.10.1.3	0.632	c.2.1.2	0.325	c.47.1.5	0.379
a.35.1.5	0.601	b.10.1.4	0.651	c.2.1.3	0.391	g.3.11.1	0.685
a.39.1.2	0.46	b.29.1.1	0.54	c.2.1.4	0.397	g.3.6.2	0.703
a.39.1.5	0.412	b.29.1.3	0.526	c.2.1.5	0.412	g.3.7.1	0.661
a.4.1.1	0.588	b.40.4.1	0.586	c.2.1.6	0.407	g.3.7.2	0.69
a.4.1.2	0.609	b.40.4.3	0.577	c.2.1.7	0.4	g.3.7.5	0.662
a.4.1.3	0.6	b.40.4.5	0.557	c.3.1.2	0.295	g.39.1.2	0.687
a.43.1.2	0.654	b.47.1.2	0.44	c.3.1.5	0.209	g.39.1.3	0.683
b.1.1.1	0.558	b.55.1.2	0.417	c.37.1.1	0.411	g.41.5.1	0.764
b.1.1.2	0.568	b.6.1.1	0.403	c.37.1.11	0.38	g.41.5.2	0.696
b.1.1.3	0.583	b.6.1.3	0.521	c.37.1.13	0.4		
b.1.1.4	0.626	b.60.1.2	0.224	c.37.1.8	0.347		

superfamily the percentage of the positions in the positive training sequences give a cumulative total of 90% of the SVM classification scores for these sequences. We found that for some superfamilies low percentage (20%-30%) of positions contributed 90% of the SVM classification scores and on average below 50% of positions contribute 90% of classification scores. We manually examined the motif candidates for positive training sequence sets in 13 experiments (2 sets from all- α class, 5 from all- β class, 5 from $\alpha+\beta$ class, and 1 from small proteins class) with high ROC scores. By comparing them with PDB annotations, we tried to identify common functional and structural characteristics captured by motif candidates for these superfamilies. We found results of four experiments to be of particular interest. For all examples we take the expected number of top-scored positions as motif regions for each sequence given the percentage of positions contributing 90% of SVM classification scores for the superfamily. We describe these four experiments below.

The first interesting example came from the homology detection experiment for PH domain-like protein superfamily (SCOP 1.59 superfamily b.55.1). Proteins in this superfamily share a conserved fold made up of a beta-barrel composed of two roughly perpendicular, anti-parallel beta-sheets and a C-terminal alpha helix. Previous studies have shown that PH domains bind to their inositol phosphate ligands via a binding surface composed primarily of residues from the $\beta 1/\beta 2$, $\beta 3/\beta 4$, and $\beta 6/\beta 7$ loops¹². The motif candidates we extracted correspond well with the C-terminal alpha helix and the ligand-binding regions on the beta-sheets but not in the loop regions at the $\beta 1/\beta 2$, $\beta 3/\beta 4$, and $\beta 6/\beta 7$. This may suggest those binding sites on the main structural components are more conserved than those in loop regions. In Figure 4, we show the motif regions for one member of this superfamily, mouse beta-spectrin protein, together with structural and functional annotations.

The second example was the EF-hand calcium-binding protein superfamily (SCOP version 1.59 superfamily a.39.1). The motif candidates that we extracted correspond well with

the two calcium-binding loops with adjacent helices, forming a local helix-loop-helix structure. In PDB annotations, these regions are labeled as EF-hand PROSITE patterns, which are important for calcium coordination. As an example, we show the SVM-extracted motif regions for one member of this superfamily, shark parvalbumin protein, in the Figure 5, together with structural and functional annotations.

In the third example, the homology detection experiment for the scorpion toxin-like superfamily (SCOP 1.59 superfamily g.3.7), we found a common motif region that forms a beta-hairpin with two adjacent disulphides. Previous studies have found that this hairpin structure might be structurally important in interacting with membrane receptors and ionic channels for proteins in this superfamily, and the disulphide bridges can help to stabilize the toxin protein. Figure 6 gives an example from this superfamily, the scorpion OSK1 toxin protein, to demonstrate the structure of the motif candidate ¹⁴.

The last example was the superfamily of homeodomain-like proteins (SCOP version 1.59 superfamily a.4.1). A common structural feature of the superfamily consists of 3 helices containing a helix-turn-helix (HTH) DNA-binding motif (also referred to as a homeodomain in eukaryotes). This structural motif is believed to interact with the major groove of DNA double strands, which facilitates the binding of many transcription regulators ²⁷. Motif candidates extracted from all sequences in our training dataset from this superfamily align well with the HTH region as indicated by the PDB annotations. In Figure 7, we show an example from the homeodomain-like protein superfamily, the MarA protein from *E. Coli*, with functional and structural annotations of the HTH motif regions.

3.4. Discriminative regions versus protein motif databases

To analyze our discriminative motif candidates further, we consider whether the discriminative regions that we found coincide with known protein motifs from the eMOTIF database (Version 3.6) ³ or structural motifs from the I-sites library (Version 16.2) ⁵. For a simple comparison, we compute the extent to which eMOTIF and I-sites motifs contribute to the overall positive discriminative scores for positive training sequences. We calculate accumulated discriminative scores falling into the sequence regions matched by any motif from the eMOTIF database or the I-sites database, and then we compare it with the expected contribution based on motif coverage, which is estimated by the ratio between total length of motif regions and the sequence length. We also compute the ratio of the eMOTIF/I-sites contribution to the expected contribution.

Interestingly, we found that on average, the eMOTIF/I-sites contribution to the discriminative score is slightly, but not dramatically, higher than expected. We show this comparison in Table 5 for eMOTIF and Table 6 for I-sites, giving results for different confidence thresholds using the eBAS and I-sites software, respectively. We conclude that our discriminative motif regions provide information that is complementary or additional to eMOTIF or I-sites motifs in many experiments.

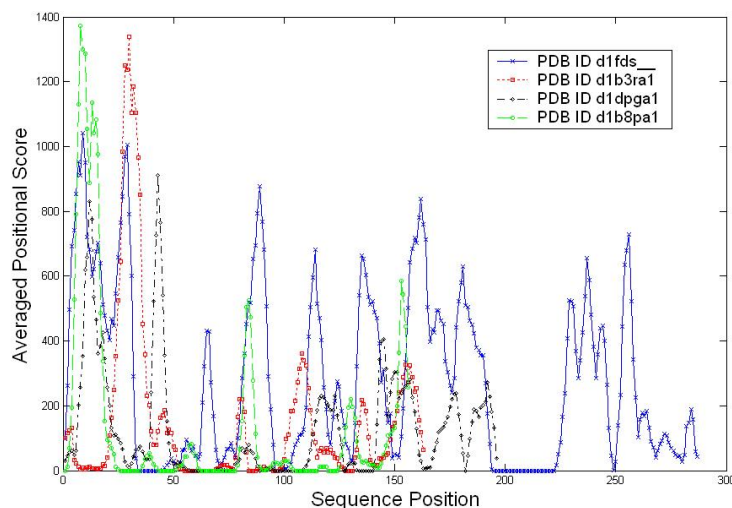


Fig. 3. **Positional contribution analysis of SVM classification score for SCOP superfamily 3.2.1 (target family 3.2.1.7).** The plot shows the contribution of each position along the sequence, obtained by averaging k -mer profile SVM scores for all k -mers containing the position, for positive training sequences in the experiment.

Table 5. **Comparison of eMOTIF motifs versus SVM discriminative scores.**

eBAS cutoff threshold	Average eMOTIF contribution	Average expected contribution	Average ratio of eMOTIF over expected
-1	0.476	0.362	1.316
-4	0.334	0.255	1.311
-8	0.271	0.207	1.306
-15	0.190	0.146	1.302

Table 6. **Comparison of I-sites motifs versus SVM discriminative scores.**

I-sites confidence threshold	Average I-sites contribution	Average expected contribution	Average ratio of I-sites over expected
0.7	0.541	0.467	1.159
0.8	0.358	0.318	1.126
0.9	0.122	0.125	0.971

4. Discussion

We have presented a novel string kernel based on protein sequence profiles, such as those produced by PSI-BLAST. The profile kernel extends the framework of k -mer based string kernels but dramatically improves SVM classification and remote homology detection over these earlier kernels. In our SCOP benchmark experiments, the SVM-profile kernel also outperformed other recently presented SVM approaches such as the eMOTIF kernel and SVM-pairwise and gave far better performance than PSI-BLAST used directly as a ranking method. Furthermore, the profile kernel is competitive with recent semi-supervised cluster

16 Rui Kuang, Eugene Ie, Ke Wang, Kai Wang, Mahira Siddiqi, Yoav Freund and Christina Leslie

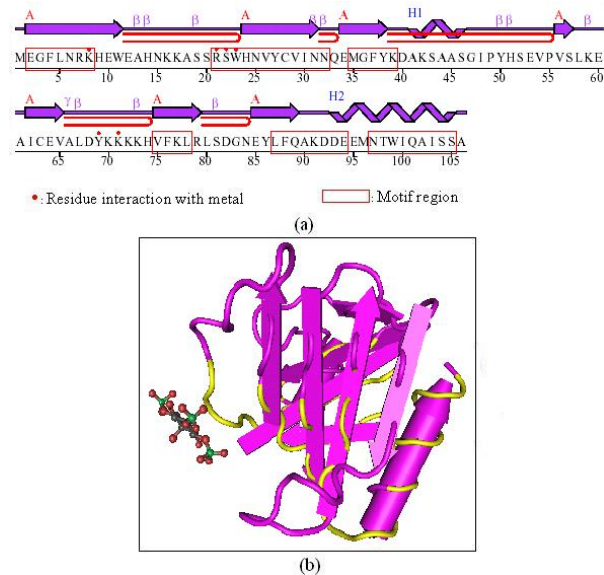


Fig. 4. Motif regions on the Mouse beta-spectrin protein that belongs to the PH domain-like protein superfamily. (a) PDB sequence annotation (PDB id 1btm) and SVM-extracted motif regions. (b) 3D structure of the mouse beta-spectrin showing the SVM-extracted motif regions on the protein structure. The yellow regions are the motif regions; the molecule is shown in pink and the ligand in green.

kernels, such as the neighborhood kernel, while achieving much better scalability to large datasets. We note that the cluster kernel approaches are general methods that can be used for a variety of applications, while the profile kernel is specialized for protein sequence data; profiles are often computed and stored for other kinds of protein sequence analysis, so profile-based kernels are particularly convenient. We then extend the profile kernel with predicted secondary structure information to obtain further small but significant improvement in remote homology detection performance.

We also show how to compute positional scores along profiles for the positive training sequences and thus extract discriminative sequence motifs. As a proof of principle, we give examples from preliminary analysis where these discriminative regions indeed map to important functional and structural features of the corresponding superfamilies. These discriminative motifs may be of use to structural biologists for improving comparative models. Moreover, we observed that motifs from known protein motif libraries like eMOTIF and I-sites were only slightly over-represented in our discriminative regions, suggesting that discriminative motifs for structural categories provide information that is complementary or supplementary to known motif databases. Moreover, in cases where the protein classification to be learned is a functional category, such as enzymatic activity, the method could be used to find discriminative sites associated with protein function.

Several authors have recently defined kernels on probabilistic models like HMMs^{25,15}. One may be able to extend the semi-supervised methodology we introduce here to use these

Remote protein homology detection and motif extraction using profile kernels 17

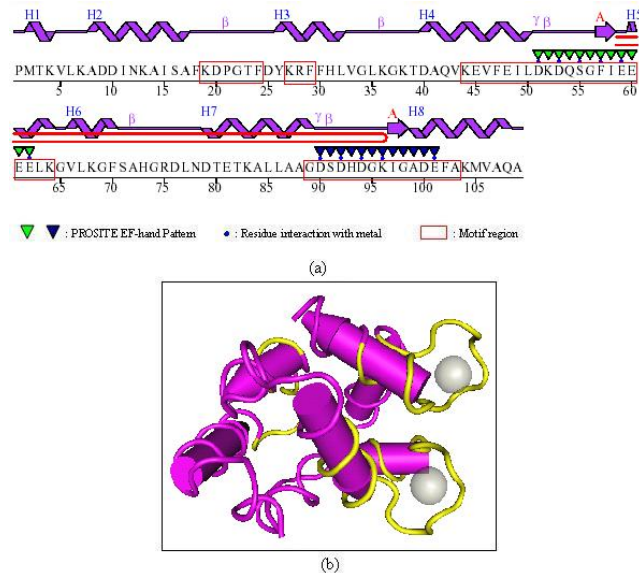


Fig. 5. Motif regions on the shark parvalbumin protein, from the EF-hand calcium-binding protein superfamily. (a) PDB sequence annotation (PDB id 5pal) and SVM-extracted motif sequences. (b) 3D structure of the parvalbumin protein showing the motif regions on the protein structure. The yellow regions are the motif regions; the two silver balls represent the calcium molecules.

kernels on suitable probabilistic models for protein sequences; for example, one could use PSI-BLAST profiles to estimate profile HMMs and then use kernels defined on HMMs for SVM training. However, these probability kernels are more computationally expensive than the profile kernel, which scales linearly with sequence length.

One significant finding from the analysis of our method is that for some superfamilies, only 20%-30% of the positions in the positive training sequences give a cumulative total of 90% of the SVM classification score for these sequences in remote homology detection experiments. This result may suggest that the multiple alignment of protein domain sequences from a superfamily—which would be used, for example, in a superfamily-based profile HMM approach—might be unnecessary for this problem, since the discriminative information is concentrated in relatively short subregions of the protein sequences. Our profile-based string kernel approach does implicitly use heuristic alignment via PSI-BLAST, but this is only to build a local profile model around each sequence, not to build a model for all the positive sequences at once. We find that local profile information, when combined with an effective profile-based string kernel representation and a powerful classification algorithm, allows us to implement a new and compelling alternative approach to remote homology detection.

18 Rui Kuang, Eugene Ie, Ke Wang, Kai Wang, Mahira Siddiqi, Yoav Freund and Christina Leslie

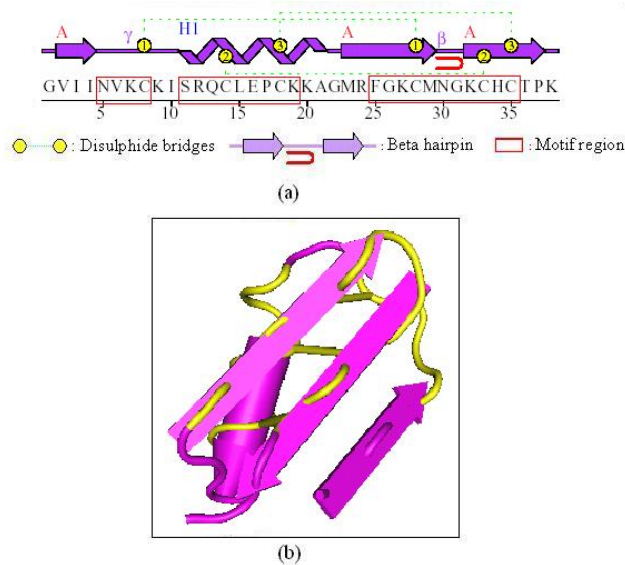


Fig. 6. Motif regions on the scorpion OSK1 Toxin from the Scorpion toxin-like superfamily. (a) PDB sequence annotation (PDB id 1sco) and SVM-extracted motif regions. (b) 3D structure of the OSK1 toxin showing the SVM-extracted motif regions on the protein structure. The yellow regions are the motif regions; the yellow bars represent the disulphide bridges.

Acknowledgments

This work is supported by an Award in Informatics from the PhRMA Foundation, NIH grant LM07276-02, NSF grant ITR-0312706, and NSF grant CCR-0325463. We thank Asa Ben-Hur for providing his eMOTIF kernel code and Chris Bystroff for providing and helping us with the I-sites package.

References

1. S. F. Altschul, T. L. Madden, A. A. Schaffer, J. Zhang, Z. Zhang, W. Miller, and D. J. Lipman. Gapped BLAST and PSI-BLAST: A new generation of protein database search programs. *Nucleic Acids Research*, 25:3389–3402, 1997.
2. P. Baldi, Y. Chauvin, T. Hunkapiller, and M. A. McClure. Hidden markov models of biological primary sequence information. *Proceedings of the National Academy of Sciences of the United States of America*, 91(3):1059–1063, 1994.
3. A. Ben-Hur and D. Brutlag. Remote homology detection: a motif based approach. *Proceedings of the Eleventh International Conference on Intelligent Systems for Molecular Biology*, 2003.
4. B. E. Boser, I. Guyon, and V. Vapnik. A training algorithm for optimal margin classifiers. In *Fifth Annual Workshop on Computational Learning Theory*, pages 144–152, Pittsburg, PA, 1992. ACM Press.
5. C. Bystroff and D. Baker. Prediction of local structure in proteins using a library of sequence-structure motifs. *Journal of Molecular Biology*, 281:565–577, 1998.
6. S. R. Eddy. Multiple alignment using hidden markov models. In *Proceedings of the Third Inter-*

Remote protein homology detection and motif extraction using profile kernels 19

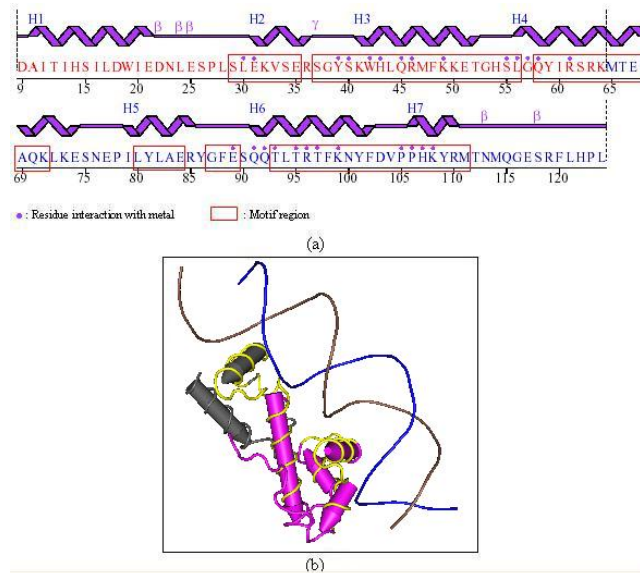


Fig. 7. Motif regions on the *E. Coli* MarA protein, from the Homeodomain-like protein superfamily. (a) PDB sequence annotation (PDB id 1bl0) and SVM-extracted motif sequences. The sequences labeled in red and blue are two domains on the same protein, both of which belong to this superfamily. (b) 3D structure of the MarA protein showing the motif regions on the protein structure. The yellow regions are the motif regions; the blue and brown strands are DNA double helices; the two protein domains are shown in black and pink respectively.

national Conference on Intelligent Systems for Molecular Biology, pages 114–120. AAAI Press, 1995.

7. V. Di Francesco, P.J. Munson, and J. Garnier. Foresst - fold recognition from secondary structures. *Bioinformatics*, 15(2):131–140, 1999.
8. M. Gribskov and N. L. Robinson. Use of receiver operating characteristic (roc) analysis to evaluate sequence matching. *Computers and Chemistry*, 20(1):25–33, 1996.
9. Michael Gribskov, Andrew D. McLachlan, and David Eisenberg. Profile analysis: Detection of distantly related proteins. *Proceedings of the National Academy of Sciences of the United States of America*, 84:4355–4358, 1987.
10. Y. Hou, W. Hsu, M. L. Lee, and C. Bystroff. Efficient remote homology detection using local structure. *Bioinformatics*, 19(17):2294–2301, 2003.
11. J. Y. Huang and D. L. Brutlag. The EMOTIF database. *Nucleic Acids Research*, 29:202–204, 2001.
12. M. Hyvonen, M. J. Macias, M. Nilges, H. Oschkinat, M. Saraste, and M. Wilmanns. Structure of the binding site for inositol phosphates in a ph domain. *EMBO Journal*, 14:4676, 1995.
13. T. Jaakkola, M. Diekhans, and D. Haussler. A discriminative framework for detecting remote protein homologies. *Journal of Computational Biology*, 7(1/2):95–114, 2000.
14. V. A. Jaravine, D. E. Nolde, M. J. Reibarkh, Y. V. Korolkova, S. A. Kozlov, K. A. Pluzhnikov, E. V. Grishin, and A. S. Arseniev. Three-dimensional structure of toxin osk1 from orthochirus scrobiculosus scorpion venom. *Biochemistry*, 36(6):1223–32, 1997.
15. T. Jebara, R. Kondor, and A. Howard. Probability product kernels. *Journal of Machine Learning*, 5:819–844, 2004.

20 Rui Kuang, Eugene Ie, Ke Wang, Kai Wang, Mahira Siddiqi, Yoav Freund and Christina Leslie

16. D. T. Jones. Protein secondary structure prediction based on position-specific scoring matrices. *Journal of Molecular Biology*, 292:195–202, 1999.
17. W. Kabsch and C. Sander. Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers*, 22:2577–2637, 1983.
18. R. Karchin, M. Cline, Y. Mandel-Gutfreund, and K. Karplus. Hidden markov models that use predicted local structure for fold recognition: alphabets of backbone geometry. *Proteins*, 51(4):504–514, 2003.
19. A. Krogh, M. Brown, I. Mian, K. Sjolander, and D. Haussler. Hidden markov models in computational biology: Applications to protein modeling. *Journal of Molecular Biology*, 235:1501–1531, 1994.
20. R. Kuang, E. Ie, K. Wang, K. Wang, M. Siddiqi, Y. Freund, and C. Leslie. Profile-based string kernels for remote homology detection and motif extraction. In *Computational Systems Bioinformatics Conference*, 2004.
21. C. Leslie, E. Eskin, A. Cohen, J. Weston, and W. S. Noble. Mismatch string kernels for discriminative protein classification. *Bioinformatics*, 2004. To appear.
22. C. Leslie, E. Eskin, J. Weston, and W. S. Noble. Mismatch string kernels for SVM protein classification. *Advances in Neural Information Processing Systems 15*, pages 1441–1448, 2003.
23. C. Leslie and R. Kuang. Fast kernels for inexact string matching. *Sixteenth Annual Conference on Learning Theory and Seventh Kernel Workshop*, pages 114–128, 2003.
24. L. Liao and W. S. Noble. Combining pairwise sequence similarity and support vector machines for remote protein homology detection. In *Proceedings of the Sixth Annual International Conference on Research in Computational Molecular Biology*, pages 225–232, Washington, DC, 2002. ACM Press.
25. R. Lyngsø, C. N. S. Pedersen, and H. Nielsen. Metrics and similarity measures for hidden Markov measures. *Proceedings of the Seventh International Conference on Intelligent Systems for Molecular Biology*, 1999.
26. C. G. Nevill-Manning, T. D. Wu, and D. L. Brutlag. Highly specific protein sequence motifs for sequence analysis. *Proceedings of the National Academy of Sciences 95*, pages 5865–5871, 1998.
27. S. Rhee, R. G. Martin, J. L. Rosner, and D. R. Davies. A novel dna-binding motif in mara; the first structure for an arac family transcriptional activator. *Proceedings of the National Academy of Sciences 95*, page 10413, 1998.
28. B. Rost and C. Sander. Prediction of protein secondary structure at better than 70% accuracy. *Journal of Molecular Biology*, 232:584–599, 1993.
29. H. Saigo, J. Vert, N. Ueda, and T. Akutsu. Protein homology detection using string alignment kernels. *Bioinformatics*, 2004.
30. K. Tsuda, M. Kawanabe, G. Rtsch, S. Sonnenburg, and K. Müller. A new discriminative kernel from probabilistic models. *Neural Computation*, 14:2397–2414, 2002.
31. J. Weston, C. Leslie, D. Zhou, A. Elisseeff, and W. S. Noble. Cluster kernels for semi-supervised protein classification. *Advances in Neural Information Processing Systems 17*, 2003.