

# Motif-Based Protein Sequence Classification Using Neural Networks

KONSTANTINOS BLEKAS, DIMITRIOS I. FOTIADIS, and ARISTIDIS LIKAS

## ABSTRACT

We present a system for multi-class protein classification based on neural networks. The basic issue concerning the construction of neural network systems for protein classification is the sequence encoding scheme that must be used in order to feed the neural network. To deal with this problem we propose a method that maps a protein sequence into a numerical feature space using the matching scores of the sequence to groups of conserved patterns (called motifs) into protein families. We consider two alternative ways for identifying the motifs to be used for feature generation and provide a comparative evaluation of the two schemes. We also evaluate the impact of the incorporation of background features (2-grams) on the performance of the neural system. Experimental results on real datasets indicate that the proposed method is highly efficient and is superior to other well-known methods for protein classification.

**Key words:** protein sequence classification, neural networks, probabilistic motifs, MEME algorithm, motif-based features.

## 1. INTRODUCTION

PROTEIN SEQUENCE CLASSIFICATION CONSTITUTES an important problem in biological sciences for annotating new protein sequences and detecting close evolutionary relationships among sequences. It deals with the assignment of sequences to known categories based on homology detection properties (sequence similarity). In several studies, protein classification has been examined at various levels, according to a top-down hierarchy in molecular taxonomy, consisting of superfamilies, families, and subfamilies (Dayhoff *et al.*, 1978). Throughout this paper, we will use the terms *family* (or *subfamily*) and *class* interchangeably to denote any collection of sequences that are presumed to share common characteristics and belong to the same category.

Various approaches have been developed for solving the protein classification problem. Most of them are based on appropriately modeling protein families, either directly or indirectly. Direct modeling techniques use a training set of sequences to build a model that characterizes the family of interest. Hidden Markov models (HMMs) are a widely used probabilistic modeling method for protein families (Durbin *et al.*, 1998) that provides a probabilistic measurement (score) of how well an unknown sequence fits to a family. Indirect techniques use direct models as a preprocessing tool in order to extract useful sequence features. In this way,

---

Department of Computer Science and Biomedical Research Institute - FORTH, University of Ioannina, GR-45110 Ioannina, Greece.

sequences of variable length are transformed into fixed-length input vectors that are subsequently used for training discriminative models, such as neural networks.

In protein sequences, *motifs* or *patterns* enclose significant homologous attributes, since they correspond to conserved regions in protein families holding useful structural and functional biological properties. They can be considered as islands of amino acids conserved in the same order of a given family. Therefore, they can be seen as local features characterizing the sequences. Motifs can be either deterministic or probabilistic (Brázma *et al.*, 1998; Rigoutsos *et al.*, 2000). Deterministic motifs follow grammatical inference properties in order to syntactically describe conserved regions of homologous sequences. The PROSITE database (Hofmann *et al.*, 1999) represents a large collection of such motifs used to identify protein families. On the other hand, probabilistic motifs are more flexible models, and they provide a probabilistic matching score of a sequence to a motif. The BLOCKS database (Henikoff and Henikoff, 1994) is an example of ungapped probabilistic motifs. In any case, motif models are suitable for constructing efficient similarity score functions that can be subsequently used as local features for protein classification. An example is presented by Ma and Wang (2000), and by Wang *et al.* (2001) where motif-based local features are produced based on the minimum description length (MDL) principle for the case of deterministic motif models.

The *background* information also constitutes another source for extracting features from sequence data. The determination of the background features, also defined as *global* features, is usually made by using the *2-gram* encoding scheme that counts the occurrences of two consecutive amino acids in protein sequences (Wang *et al.*, 2001). In the case of protein sequences (generated from the alphabet of the 20 aminoacids), there are 400 possible 2-grams that produce a large feature space. A recent approach (Almeida and Vinga, 2002) proposes a scheme for globally encoding sequences, where each amino acid character is initially represented as a unique binary number with  $n$  bits ( $n = 5$  for the 20 aminoacids) and then each sequence is mapped into a position inside the  $n$ -dimensional hypercube.

In this paper, we focus on building efficient neural classifiers for discriminating multiple protein families by using appropriate local features that have been extracted by efficient probabilistic motif models. As motifs constitute family diagnostic signatures, our aim is to exploit this information by constructing a neural network scheme that exploits motif-based (local) features.

The proposed method can be considered as combining an unsupervised with a supervised learning technique. Starting by applying a motif-discovery (unsupervised) algorithm, we identify probabilistic motifs in a training set of multiclass sequences. This can be achieved in two alternative ways: applying the algorithm for motif discovery either to each family training set separately (*class-dependent* motifs), or to the whole dataset of training sequences (*class-independent* motifs). The discovered motifs are then used to convert each sequence to a numerical input vector that subsequently can be applied to a typical feed-forward neural network. Using a Bayesian regularization training technique, the neural network parameters are adjusted, and therefore a classifier is obtained suitable for predicting the family of an unlabeled sequence.

The next section provides a brief overview of statistical and neural techniques proposed for classifying biological sequences, while Section 3 describes the proposed method. Experimental results obtained using several sets of protein families are presented in Section 4, along with a comparison with other known protein classification approaches. Finally, Section 5 summarizes the proposed classification scheme and specifies directions for future research.

## 2. PROTEIN CLASSIFICATION METHODS

One class of methods for protein sequence classification work directly with sequence information and establish classification criteria based on sequence homology properties. In the general scheme, a representative set of sequences is selected for each protein family and used to build an appropriate model for each family. The classification of an unknown sequence is made by selecting the family that best matches according to the model homology mechanism. This can be considered as a simple *nearest neighbor* scheme that ranks sequence similarities and selects the best ranking.

The popular BLAST tool (Altschul *et al.*, 1990) represents the simplest nearest neighbor approach and exploits pairwise local alignments to measure sequence similarity. The BLAST technique compares protein

queries with a protein database of labeled sequences and produces normalized alignment scores for each comparison by calculating the corresponding expectation values ( $E$ -values). The classification procedure is based on the selection of the label of the sequence that produces the best pairwise alignment score (i.e., minimum  $E$ -value).

Another type of direct modeling methods is based on hidden Markov models (HMMs) (Durbin *et al.*, 1998; Karplus *et al.*, 1998). After constructing an HMM for each family, protein queries can be easily scored against all established HMMs by calculating the log-likelihood of each model for the unknown sequence and then selecting the class label of the most likely model.

The Motif Alignment and Search Tool (MAST) (Bailey and Gribskov, 1998) is based on the combination of multiple motif-based statistical score values. According to this scheme, groups of probabilistic motifs discovered by the MEME algorithm (Bailey and Elkan, 1994), are used to construct protein profiles for the families of interest. The MAST algorithm successively estimates the significance of the match of a query sequence to a family model as the product of the  $p$ -values of each motif match score. This measure (called  $E$ -value) can then be used to select the family of the unknown sequence.

Neural network schemes for protein classification consist of two stages: the *encoding* stage, where discriminative numerical features are computed for each training sequence, and the *decision* stage, where the created feature vectors are used as input vectors to a neural network classifier. Various encoding schemes have been proposed in the literature that produce numerical features in the encoding stage based on the calculation of background features (global sequence homology) and local features (locally conserved family information) embedded in motifs. In the decision stage, feed-forward neural networks have been used trained either through back-propagation (Wu *et al.*, 1996) or using Bayesian regularization (Ma and Wang, 2000; Wang *et al.*, 2001). These approaches are characterized by the enormous size of the extracted input vectors, the imbalance between global and local features (more emphasis on global features), and the need for large training sets (since the number of network inputs is very large). For example, in Ma and Wang (2000) and in Wang *et al.* (2001) only one feature was responsible for carrying local information, while all the others were produced by the 2-grams encoding scheme (background features).

Support vector machines (SVMs) (Vapnik, 1979) have been also applied to protein homology detection problems. Such an approach, which has been introduced by Logan *et al.* (2001), feeds probabilistic score values from all motifs available (nearly 10,000) in the BLOCKS database (Henikoff and Henikoff, 1994) into an SVM classifier. Obviously, this scheme uses only local features, but the dimensionality of the input space is extremely high. Another method has been proposed by Jaakkola *et al.* (2000) and by Karchin *et al.* (2002) that combines hidden Markov models (HMMs) and SVMs for detecting remote protein homologies. In particular, an HMM is first trained to model a protein family, and then the observed probabilities (in the log space) of each sequence with respect to each parameter of the HMM are calculated. The obtained gradient-log-probability vectors are applied to an SVM to identify the decision boundary between the family and the rest of the protein universe.

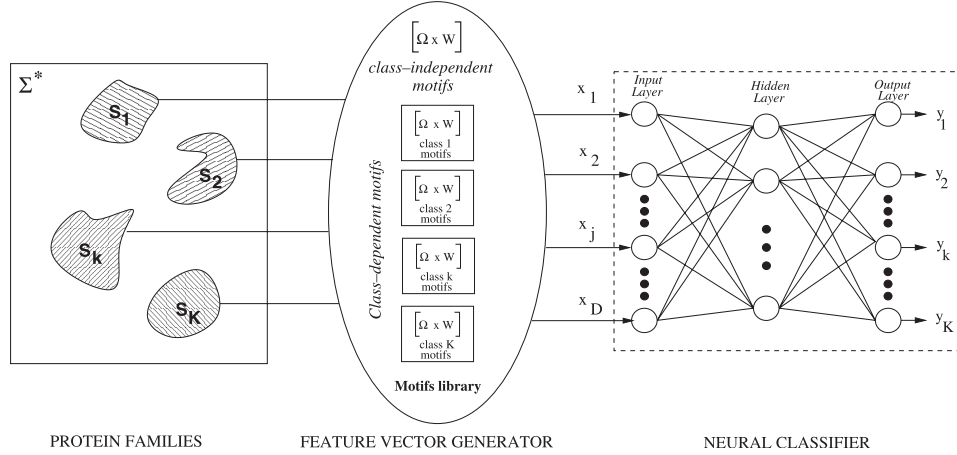
### 3. THE PROPOSED METHOD

This paper studies the problem of classifying a set of  $N$  protein sequences  $\mathbf{S} = \{S_i, i = 1, \dots, N\}$  into  $K$  classes. The set  $\mathcal{S}$  is a union of positive example datasets  $\mathcal{S}_k$  from  $K$  different classes, i.e.,  $\mathcal{S} = \{\mathcal{S}_1 \cup \dots \cup \mathcal{S}_K\}$ , and can be seen as a subset of the complete set of all possible sequences over the amino acid alphabet ( $\mathcal{S} \subseteq \Sigma^*$ ).

Figure 1 illustrates the architecture of the proposed protein classification scheme. It consists of a search tool (unsupervised learning) for discovering probabilistic motifs in a set of  $K$  protein families, a feature vector generator that converts protein sequences into feature vectors, and a decision module (neural network) for assigning a protein family to each input sequence. The following subsections describe in detail the major building blocks of the proposed architecture.

#### 3.1. Using motifs for feature generation

Consider a finite alphabet consisting of set of characters  $\Sigma = \{\alpha_1, \dots, \alpha_\Omega\}$  ( $\Omega = 20$  for protein sequences). We can probabilistically model a contiguous (ungapped) motif  $M_j$  of length  $W_j$  using a



**FIG. 1.** The architecture of the proposed classification scheme.

position weight matrix ( $PWM_j$ ) that follows a multinomial character distribution. Each column ( $l$ ) of the matrix corresponds to a position  $l$  in the motif sequence ( $l = 1, \dots, W_j$ ), where the column elements provide the probability of each character of the alphabet  $p_{\alpha_{\xi},l}$  ( $\xi = 1, \dots, \Omega$ ) to appear in that position.

Let  $s_p = a_{p,1} \dots a_{p,W_j}$  denote a segment of a sequence  $S$  beginning at position  $p$  and ending at position  $p + W_j - 1$ . This represents a subsequence of length  $W_j$ . Totally, there are  $L - W_j + 1$  such subsequences for a sequence  $S$  of length  $L$ . Then, we can define the probability that  $s_p$  matches the motif  $M_j$ , or alternatively, has been generated by the model  $PWM_j$  corresponding to that motif, using the following equation:

$$P(s_p|M_j) = \prod_{l=1}^{W_j} p_{a_{p,l},l}. \quad (1)$$

A major advantage of using the probabilistic matrix  $PWM_j$  is the ability to compute the corresponding position-specific score matrix ( $PSSM_j$ ) in order to score a sequence. The  $PSSM_j$  is a log-odds matrix calculating the logarithmic ratio  $r_{\alpha_{\xi},l}$  of the probabilities  $p_{\alpha_{\xi},l}$  suggested by the  $PWM_j$  and the corresponding general relative frequencies of aminoacids  $\rho_{\alpha_{\xi}}$  in the family<sup>1</sup>. According to the definition of  $PSSM_j$ , the score value  $f_j(s_p)$  of a subsequence  $s_p$  of a sequence  $S$  can be defined as

$$f_j(s_p) = \sum_{l=1}^{W_j} \log \left( \frac{p_{a_{p,l},l}}{\rho_{a_{p,l}}} \right) = \sum_{l=1}^{W_j} r_{a_{p,l},l}. \quad (2)$$

At the sequence level, the score value of a protein sequence  $S$  against a motif  $M_j$  can be determined as the maximum value among all scores of the possible subsequences of  $S$ , i.e.,

$$f_j(S) = \max_{1 \leq p \leq L - W_j + 1} f_j(s_p). \quad (3)$$

It must be noted that it is possible to adopt other definitions for scoring a sequence, such as setting scores below a certain threshold equal to zero (Bailey and Gribskov, 1998).

If we assume that we have discovered a group of  $D$  motifs in the set of sequences  $\mathbf{S}$ , we can generate a  $D$ -dimensional numerical feature space and map each sequence  $S_i$  into a vector  $\mathbf{x}_i$  in the  $D$ -dimensional feature space by calculating the score values  $x_{ij} = f_j(S_i)$  ( $j = 1, \dots, D$ ) for each of the  $D$  motif models.

<sup>1</sup>The general relative frequencies of amino acids indicate the background information in a protein family and can be presented as a probabilistic vector  $\rho$  of size  $\Omega = 20$ .

### 3.2. Finding probabilistic motifs in protein sequences

Several approaches have been proposed for discovering probabilistic motifs in a set of unaligned biological sequences. CONSENSUS (Hertz and Stormo, 1999), the Gibbs sampler (Lawrence *et al.*, 1993), and MEME (Bailey and Elkan, 1994) are examples of such methods that identify multiple shared motifs in protein families. We have selected the MEME approach for the motif identification component of our strategy, since it has been widely used in biological applications and directly extracts position-specific score matrices. Below, we briefly describe this algorithm and propose two ways to integrate it in our classification system.

The MEME algorithm follows an iterative procedure, which applies at each iteration a two-component mixture model to discover one motif of length  $W$ . In the two-component model, one component describes the motif (ungapped common subsequences of length  $W$ ) while the other component models the background information. Multiple motifs can be found by sequentially fitting the two-component model to the set of sequences that remain after removing the sequences containing occurrences of the already identified motifs.

In particular, MEME (Bailey and Elkan, 1994) uses the Expectation Maximization (EM) algorithm (Dempster *et al.*, 1977) to maximize the log-likelihood function of the two-component mixture model, i.e., to estimate the elements of the corresponding position weight matrix<sup>2</sup>. Furthermore, MEME provides a strategy for locating efficient initial parameter values in order to prevent the EM algorithm from getting stuck in local optima (Bailey and Elkan, 1994). The  $D$  motif models  $PWM_j$  ( $j = 1, \dots, D$ ) discovered by MEME can be of either fixed or variable length  $W_j$ . In our experimental studies, both types of motifs will be examined to evaluate the impact of this decision on the performance of the neural classifier.

In order to discover a group of motifs from a multiclass training set of sequences (containing sequences of  $K$  classes), two alternative approaches can be followed. The first approach is to apply the MEME algorithm  $K$  times, *separately* to the training sequences of each protein family. Then, putting all the discovered  $K$  family profiles together, we can form the final group of  $D$  motifs. An alternative approach is to apply the motif-discovery algorithm only once to the total training set  $\mathcal{S}$ , ignoring class labels. In this way, we do not allow the algorithm to directly create  $K$  protein family profiles, but rather to discover  $D$  *class-independent* motifs.

The advantage of the second approach is the ability of taking into account local similarity measurements in the whole training set, without restricting the search procedure to a single class. Therefore, possible partial homologies among sequences from different families can be defined that may prove helpful for the classification task. On the other hand, a disadvantage of the class-independent approach is that the  $D$  discovered motifs may not be equally distributed among the  $K$  families. This may result in insufficient modeling of some families, thus leading to performance deterioration. During experiments, both motif-discovery strategies will be considered and evaluated.

### 3.3. Construction of a neural classifier

After discovering  $D$  motifs and constructing the  $D$ -dimensional feature space, the last stage in our methodology is to implement and train a feed-forward neural network that will be able to map the input vectors into the protein classes of interest. A typical network architecture is illustrated in Fig. 1. To construct the neural classifier, we use the training set  $\mathbf{X} = \{\mathbf{x}_i, \mathbf{t}_i\}$ ,  $i = 1, \dots, N$  consisting of positive examples  $\mathbf{x}_i$  from the set of  $K$  protein families. The target vector  $\mathbf{t}_i$  is a binary vector of size  $K$  indicating the class label of input  $\mathbf{x}_i$ ; i.e.,  $t_{ik} = 1$  if  $\mathbf{x}_i$  corresponds to a sequence  $S_i$  belonging to class  $k$ , and 0 otherwise. The output of the classifier is represented by the  $K$ -dimensional vector  $\mathbf{y}_i$  where component  $y_{ik}$  corresponds to class  $k$ . Based on this scheme, the predicted class  $h(\mathbf{x}_i)$  of an unlabeled feature vector  $\mathbf{x}_i$  corresponding to a query sequence  $S_i$  is given by the index of the output node with the largest value  $y_{ic}$ ; i.e.,

$$h(\mathbf{x}_i) = c : y_{ic} = \max_{1 \leq k \leq K} y_{ik} . \quad (4)$$

<sup>2</sup>The model used in our experiments assumes that there are zero or more nonoverlapping occurrences of the motif in each sequence of the dataset. Alternative models that can be used are the exactly one-occurrence-per-sequence and the zero-or-one-occurrence-per-sequence models.

Setting a threshold value  $\theta$  ( $\in [0, 1]$ ), we can restrict the classifiers' decision to only those input vectors whose maximum output value surpasses this threshold. In this case, we can write

$$h(\mathbf{x}_i, \theta) = c : y_{ic} = \max_{1 \leq k \leq K} y_{ik} \wedge y_{ic} \geq \theta . \quad (5)$$

Parameter  $\theta$  can be used to specify the sensitivity of the classifier.

In order to train the neural network, we used the Gauss–Newton Bayesian Regularization (GNBR) learning algorithm (Foresse and Hagan, 1997). This algorithm applies Bayesian regularization and implements a Gauss–Newton approximation to the Hessian matrix of the objective function.

In the Bayesian regularization framework, the objective function is formulated as the weighted sum of two terms: the sum of the squared errors ( $E_X$ ) and the sum of squares of the network weights ( $E_W$ ). Using Bayes' rule, the posterior probability distribution for the weights  $\mathbf{w}$  of the network given a training set  $\mathbf{X}$  can be written as follows:

$$P(\mathbf{w}|\mathbf{X}) = \frac{P(\mathbf{X}|\mathbf{w})P(\mathbf{w})}{P(\mathbf{X})} . \quad (6)$$

By properly choosing the prior distribution  $P(\mathbf{w})$  and the likelihood function  $P(\mathbf{X}|\mathbf{w})$ , we can obtain the following expression (Bishop, 1995; Foresse and Hagan, 1997) for the posterior distribution:

$$P(\mathbf{w}|\mathbf{X}) = \frac{1}{Z_F} \exp(-\beta E_X - \alpha E_W) = \frac{1}{Z_F} \exp(-F(\mathbf{w})), \quad (7)$$

where the  $Z_F$  corresponds to the normalizing factor that is independent of the weights.

Maximizing the above posterior distribution is equivalent to minimizing the regularized objective function  $F(\mathbf{w})$ :

$$F(\mathbf{w}) = \frac{\beta}{2} \sum_{i=1}^{N_X} \{y_i - \mathbf{t}_i\}^2 + \frac{\alpha}{2} \sum_{j=1}^{N_W} w_j^2 , \quad (8)$$

where  $N_X$  and  $N_W$  represent the number of input vectors and network parameters, respectively. In order to estimate the normalizing factor  $Z_F$ , a Gaussian approximation can be used for the posterior distribution (MacKay, 1992) as obtained by the Taylor expansion of function  $F(\mathbf{w})$  around the minimum value of the posterior,  $\mathbf{w}_{MP}$ . This gives the following estimation (Bishop, 1995):

$$Z_F^*(\alpha, \beta) = \exp(-F(\mathbf{w}_{MP}))(2\pi)^{N_W/2} |\mathbf{H}|^{-1/2} , \quad (9)$$

where  $\mathbf{H}$  corresponds to the Hessian matrix of the regularized objective function and, therefore, optimal values for parameters  $\alpha$  and  $\beta$  at the minimum point  $\mathbf{w}_{MP}$  can be computed as follows:

$$\hat{\alpha} = \frac{\gamma}{2E_W(\mathbf{w}_{MP})} \quad \text{and} \quad \hat{\beta} = \frac{\gamma N_X}{2E_X(\mathbf{w}_{MP})} . \quad (10)$$

The quantity  $\gamma$  represents the effective number of network parameters  $\mathbf{w}$  and can be defined using the eigenvalues of  $H^{-1}$  as  $\gamma = N_W - 2\alpha \text{Tr} \mathbf{H}^{-1}$ . In cases where the number of effective parameters is equal to the actual ones ( $\gamma \approx N_W$ ), more hidden units must be added to the network. Furthermore, the GNBR algorithm follows a Gauss–Newton approximation method (Foresse and Hagan, 1997) for calculating the Hessian matrix of  $F(\mathbf{w})$  at the minimum point  $\mathbf{w}_{MP}$ , using the Levenberg–Marquardt optimization algorithm (Bishop, 1995). It must be noted that in our experiments, the best results for the GNBR algorithm were obtained by scaling the network inputs in the range  $[-1, 1]$ .

#### 4. EXPERIMENTAL RESULTS

Several experiments were conducted to evaluate the proposed method. The classification accuracy was measured by counting the sensitivity and specificity rates. In all  $K$ -class classification problems, each

TABLE 1. THE TWO PROSITE FAMILIES USED IN THE EXPERIMENTAL STUDY

<i>Problem: PROSITE 1 (K = 6)</i>			<i>Problem: PROSITE 2 (K = 7)</i>		
<i>PROSITE family</i>	<i>Positive data</i>	<i>Training set (avg length of seqs)</i>	<i>PROSITE family</i>	<i>Positive data</i>	<i>Training set (avg length of seqs)</i>
PS00030	302	20 (370)	PS00070	129	15 (558)
PS00038	289	20 (359)	PS00077	155	15 (502)
PS00061	317	20 (299)	PS00118	168	15 (127)
PS00198	300	20 (284)	PS00180	123	15 (408)
PS00211	574	30 (478)	PS00215	123	15 (321)
PS00301	386	20 (517)	PS00217	148	15 (490)
			PS00338	173	15 (212)

protein family  $\mathcal{S}_k$  ( $k = 1, \dots, K$ ) was randomly partitioned into training and test sequences, with the training set being only a small percentage (5–10%) of the family dataset. Using the training datasets, experiments have been carried out using the MEME algorithm to discover groups of motifs. Two cases were considered: in the first case, the MEME algorithm has been applied separately to each training set providing a group of  $D_k = 5$  class-dependent motifs for each family  $\mathcal{S}_k$ .<sup>3</sup> In the second case, the MEME algorithm was applied only once to the total training dataset (ignoring the class labels) to provide a group of  $D = 5 \times K$  class-independent motifs.

In any case, the obtained final group of  $D$  motifs were used to transform each sequence of the dataset into a dataset with numerical  $D$ -dimensional feature vectors, denoted  $\mathbf{X}_s$  for the class-dependent case and  $\mathbf{X}_g$  for the class-independent case. Furthermore, we also experimented with the effect of the length  $W$  of the discovered motifs to the performance of the proposed classifier, by applying the MEME algorithm with either fixed or variable motif length. We selected  $W = 20$  for the first case and the range  $[10, 30]$  for the second case. In summary, we have considered four distinct cases considering the application of MEME: discovering either class-dependent or class-independent motifs with either fixed or variable motif length. Therefore, for each classification problem, four distinct neural classifiers will be constructed and tested.

To evaluate classification performance, ROC (receiver operating characteristic) analysis was used. More specifically, we used the ROC<sub>50</sub> curve which is a plot of the sensitivity as a function of false positives for various decision threshold values until 50 false positives are found.

For our experimental study, three real datasets were selected. In particular we have used protein families from the PROSITE database (Hofmann *et al.*, 1999), which is a large collection of protein families together with their characteristic (deterministic) motifs. Two datasets with  $K = 6$  (PROSITE 1) and  $K = 7$  (PROSITE 2) classes from the PROSITE database (Hofmann *et al.*, 1999) were selected, summarized in Table 1. Moreover, experiments have also been conducted on a dataset of G-protein coupled receptors (GPCR) (Horn *et al.*, 1998), that is, a superfamily of cell membrane proteins. The GPCR database is hierarchically classified into five major classes and their subfamilies (Horn *et al.*, 1998). We studied the problem of classifying subfamilies within the class A, since it dominates the whole GPCR database. As indicated by Karchin *et al.* (2002), the difficulty of recognizing GPCR subfamilies arises from the fact that the classification of the subfamilies has been made based on chemical properties rather than sequence homology. Therefore, members from different subfamilies may share strong homology, thus making their discrimination hard. Among 15 subfamilies consisting of class A, seven of them have been selected in our experimental study described in Table 2. The remaining eight subfamilies are of very small size, and it is difficult to construct an effective system for their discrimination. Details of the three datasets (family/subfamily names and their protein ID's) used in our experiments are given in the appendix.

#### 4.1. Local versus global features

In this series of experiments, we assessed the impact of using 2-grams (background features) on the performance of the proposed classification scheme. For a sequence  $S_i$  with length  $L_i$ , we define the feature

<sup>3</sup>Experiments with a greater number of motifs did not yield better classification performance.

TABLE 2. SEVEN FAMILIES FROM THE GPCR CLASS A USED IN THE EXPERIMENTAL STUDY

<i>Problem: GPCR (K = 7)</i>		
<i>GPCR Class A subfamily</i>	<i>Positive data</i>	<i>Training set (avg length of seqs)</i>
Amine	306	20 (485)
Peptide	654	30 (383)
Hormone	43	10 (378)
Rhodopsin	270	20 (358)
Olfactory	325	20 (317)
Prostanoid	43	10 (721)
Nucleotide-like	58	10 (348)

value  $g_{iq}$  for each 2-gram  $q$  with respect to this sequence as

$$g_{iq} = \frac{\mathcal{N}(q|S_i)}{L_i - 1}, \quad (11)$$

where  $\mathcal{N}(q|S_i)$  denotes the number of occurrence of the 2-gram feature  $q$  in the sequence  $S_i$ . Obviously, the above equation gives the relative frequency of a 2-gram feature in a sequence. In a training set  $\mathbf{S} = \{S_1, S_2, \dots, S_N\}$  of  $N$  sequences, we can ignore *redundant* 2-grams and consider only the  $N_g$  features  $g_{iq}$  that correspond to the most frequently occurring 2-grams. We select the  $N_g$  2-grams occurring in at least half of the training sequences and by computing the corresponding  $g_{iq}$  ( $q = 1, \dots, N_g$ ) values for each sequence  $S_i$ , we construct the corresponding feature vectors to be fed in the neural classifier.

Table 3 presents the dimensionality of the feature spaces obtained using 2-grams and motifs for each dataset used in the experiments. It must be noted that we can further reduce the dimensionality of the 2-gram feature vectors using standard dimension reduction techniques, such as principal component analysis (PCA).

To assess the impact of the several feature types on the performance of the classification system, we have considered five different datasets:

- $\mathbf{X}_S$ :  $D$  motif-based features separately identified for each family (class-dependent),
- $\mathbf{X}_G$ :  $D$  motif-based class-independent features,
- $\mathbf{X}_S \cup \mathbf{G}$ :  $D$  motif-based class-dependent features along with  $N_g$  2-gram features,
- $\mathbf{X}_G \cup \mathbf{G}$ :  $D$  motif-based class-independent features, along with  $N_g$  2-gram features
- $\mathbf{G}$ :  $N_g$  2-gram features.

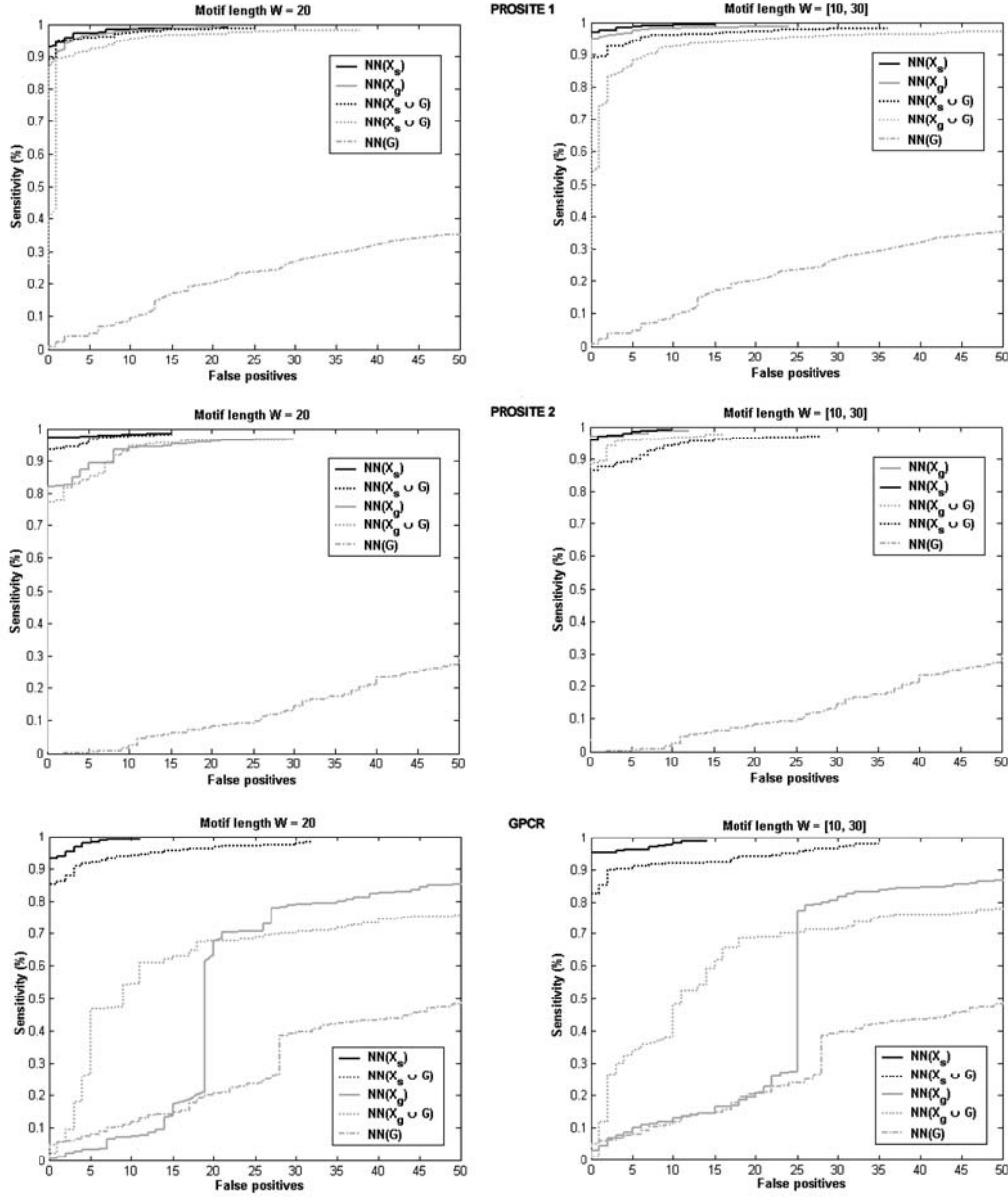
The neural network architecture had one hidden layer of either 10 (for the cases  $\mathbf{X}_S$  and  $\mathbf{X}_G$ ) or 20 nodes (for the other three cases).

Figure 2 displays the ROC<sub>50</sub> curves obtained after training the five neural classifiers in each of the three classification problems, respectively. For each problem, two different graphs are presented concerning

TABLE 3. THE NUMBER OF THE EXTRACTED MOTIF-BASED ( $D$ ) AND 2-GRAM ( $N_g$ ) FEATURES THAT CORRESPONDS TO EACH DATASET

<i>Problem</i>	$N_g$ <i>2-gram features</i>	$D$ <i>motif-based features</i>
PROSITE 1	174	$5 \times 6 = 30$
PROSITE 2	285	$5 \times 7 = 35$
GPCR	152	$5 \times 7 = 35$

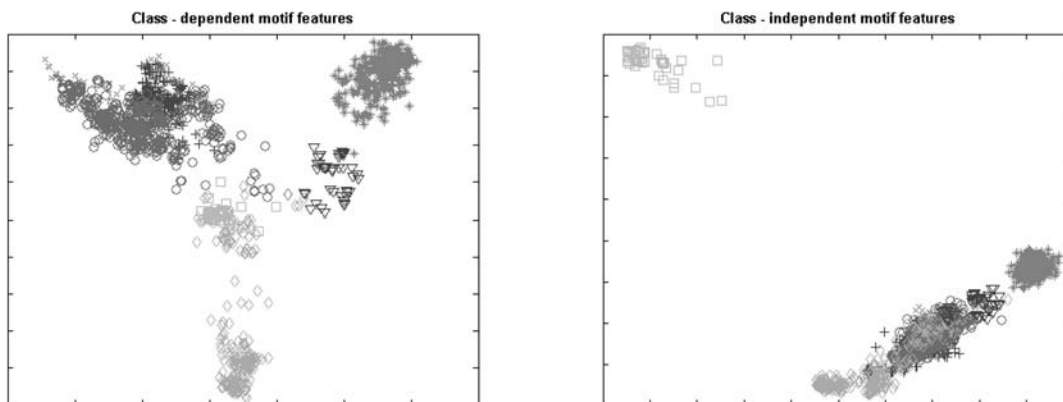




**FIG. 2.** ROC<sub>50</sub> curves illustrating the performance of the neural classifier on the three datasets using the five different feature vectors.

motifs of fixed length ( $W = 20$ ) and of variable length  $W \in [10, 30]$ . Obviously, motif-based features themselves constitute an excellent source of information able to generate significant features and lead to the construction of efficient classifiers. In all cases, the neural networks trained by mixed features (e.g.,  $NN(X_s \cup G)$ ) exhibit lower classification accuracy compared to the corresponding classifier trained with only motif-based features (e.g.,  $NN(X_s)$ ). Furthermore, the 2-grams features alone (case  $NN(G)$ ) do not seem to contain significant discriminant information.

Another observation that can be made from the ROC<sub>50</sub> curves in Fig. 2 is related to the performance of the neural classifier with class-dependent motifs (network  $NN(X_s)$ ) compared to that obtained with class-independent motifs (network  $NN(X_g)$ ). In almost all cases, we obtained better classification results with the network  $NN(X_s)$ . One explanation for this behavior is that, when searching for a specific number  $D$  of motifs in the whole training set (ignoring class labels), the algorithm may focus on some of the families



**FIG. 3.** The seven class regions in the GPCR dataset in the case of class-dependent and class-independent features. The data have been projected in two dimensions using PCA.

and leave the other families explored only partially. This possibly affects the satisfactory modeling of some families, since the discovered class-independent motifs may not be sufficient for describing them (only a few individual motifs are dedicated to this family). Experiments in the  $\mathbf{X}_g$  datasets with MEME have shown that the allocation of motifs in most cases was not equal for all the  $K$  families.

An example is shown in Fig. 3 that illustrates the constructed feature space of the  $\mathbf{X}_s$  and  $\mathbf{X}_g$  datasets in the case of the GPCR problem (seven classes), after projecting the 35-dimensional numerical to a two-dimensional space using PCA. It can be observed that in the case of class-dependent motifs the protein classes exhibit less overlap while in the reduced feature space of class-independent motifs there is a significant overlapping among class regions, thus making the discrimination harder. A selection of higher values of  $D$  probably would lead to better results for the class-independent case, but would simultaneously result in larger feature spaces or to the overestimation of some families.

#### 4.2. Comparison with other approaches

We have also compared the neural classifier (with class-dependent motif-based features) with two other protein classification methods, namely, the MAST homology detection algorithm (Bailey and Gribskov, 1998) and the profile HMMs built using SAM, (Hughey and Krogh, 1996). In both MAST and SAM, each protein family (or subfamily) is transformed (indirectly or directly) into a probabilistic model-profile, and the test sequences are classified using the class of the profile with the best score value.

More specifically, the MAST procedure (Bailey and Gribskov, 1998) initially uses the MEME algorithm to discover groups of motifs separately for each one of the  $K$  protein families. For each sequence in the testing set, the MAST algorithm combines the calculated  $p$ -values and estimates the significance of the observed match (called  $E$ -value) of the sequence to each of the  $K$  groups of motifs.<sup>4</sup> Then the query sequence is assigned to the class with the minimum  $E$ -value. The SAM method (Hughey and Krogh, 1996) works in a similar way by building an HMM for each one of the  $K$  protein families (or subfamilies) instead of discovering groups of motifs.<sup>5</sup>

Figure 4 provides comparative results from the application of the proposed neural classifier, MAST and SAM, to the three datasets. We have created five ROC curves for each method (number of false positives versus sensitivity for several threshold values) until 25 false positives were found ( $\text{ROC}_{25}$ ). The performance of the neural classifier and MAST was given by two curves respectively<sup>6</sup> concerning motifs of fixed ( $W = 20$ ) and variable length ( $W = [10, 30]$ ), while the last one corresponds to SAM performance.

<sup>4</sup>We use the *meme* and *mast* commands from the available MEME package v.3.0.4.

<sup>5</sup>We use the *buildmodel* and *hmmscore* commands from the available SAM package v.3.3.1.

<sup>6</sup>The curves for the neural classifier performance were the best plots from the corresponding  $\text{ROC}_{50}$  diagrams in Fig. 2.

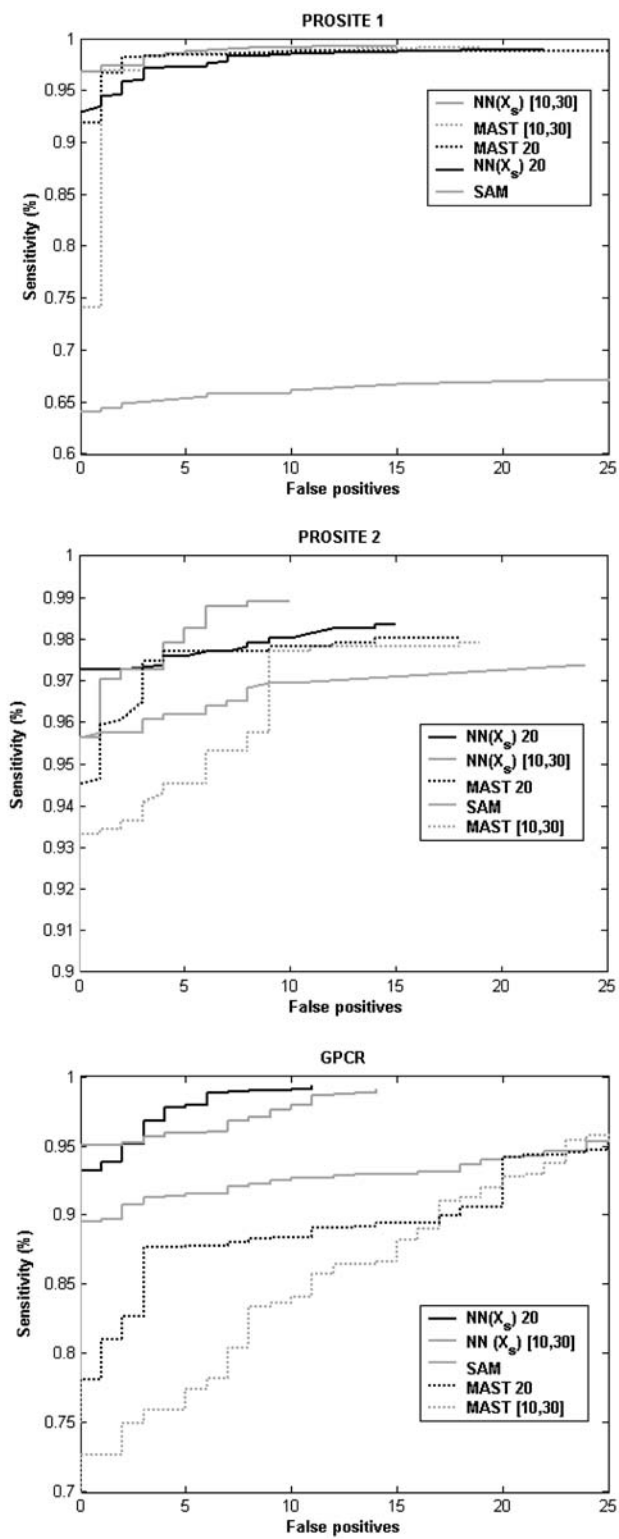


FIG. 4. ROC<sub>25</sub> curves for the three methods (neural (NN), MAST, and SAM) on the three datasets.

In the case of MAST and SAM methods, ROC curves were obtained by setting several  $E$ -value thresholds. When the lowest estimated  $E$ -value for a query sequence was greater than the threshold, then the test sequence was considered unclassified.

The superior classification of the proposed neural approach is obvious from the plotted curves in all problems, offering greater sensitivity rates with perfect specificity (zero false positives). For the GPCR dataset, which is more difficult to discriminate, the classification improvement is more clear: a sensitivity rate of 99.30% was measured with only 11 false positives, while the corresponding results for MAST and SAM are (95.76%, 25) and (95.38%, 25), respectively. It is also important to stress the higher accuracy that the neural scheme achieves compared with MAST (dot lines). Although these two methods use the same groups of motifs, our method seems to offer a more efficient scheme for combining the motif match scores compared to the combination of their  $p$ -values as suggested by MAST. In addition, the neural classifier achieves fewer false positives with higher sensitivity rates in all datasets concerning either fixed or variable motif length. Again, the improvement is more clear in the plots corresponding to the GPCR dataset.

Regarding more carefully the three selected datasets, they can be considered as three different types of protein sequence classification problems. In particular, the PROSITE 1 dataset consists of *diverse protein families* in the sense that their corresponding PROSITE motifs are not very specific (such as in the case of PS00030 and PS00198) and they can be found in sequences from a large number of protein families. Hence, this application can be seen as a diverse protein family recognition problem. On the other hand, the PROSITE 2 dataset consists of protein families with more specific PROSITE motifs that can be distinguished more easily. Finally, the third dataset, GPCR, is related to the recognition of protein subfamilies within a broader protein family domain sharing strong homology.

In all the above three types of protein sequences classification problems, our approach has shown a superior classification performance providing better results in comparison with the two other approaches. As illustrated in Fig. 4, the SAM method seems to be unsuccessful in recognizing diverse protein families (PROSITE 1 case), and the obtained classification rate was low (the individual classification error for each diverse family was about 50%). On the other hand, the performance of the MAST method was lower in the case of the GPCR subfamily recognition problem where sequences from different subfamilies share strong homology. Finally, in the case of recognizing simple protein families (PROSITE 2 dataset), all the three approaches provide similar classification rates, with the proposed neural scheme offering slightly better results.

## 5. CONCLUSIONS

In this paper, we have presented a neural network approach for the classification of protein sequences. The proposed methodology is motivated by the principle that in biological sequence analysis motifs can provide major diagnostic features for determining the class label of the unknown sequences. The method is implemented in two steps, where a preprocessing step (based on the MEME algorithm) is initially applied for discovering a group of probabilistic motifs appearing in the sequences. We have suggested and evaluated two alternative ways for motif discovery in a set of  $K$ -class sequences depending on whether the class labels are taken into account. Using the discovered motifs, a numerical feature vector is generated for each sequence by computing the matching score of the sequence to each motif. At the second stage of the proposed method, the extracted feature vectors are used as inputs to a feed-forward neural network trained using the Gauss–Newton Bayesian Regularization algorithm that provides the class label of a sequence.

Experiments were conducted on real datasets (using very small training sets), and comparisons were made with the MAST and SAM probabilistic methods. ROC curves were used as a performance indicator, and the experimental results clearly illustrate the superiority of the proposed neural system. In addition we have shown that background features do not constitute a useful source of information for the classification task since they do not lead to performance improvement.

In future work, more extensive experiments could be conducted to assess the performance of the method on specific protein superfamilies of important biological functions, as was the case with the GPCR dataset. Also, alternative methods could be implemented and tested, both in the classification stage (mixture models, SVMs, etc.) and in the motif discovery stage.

APPENDIX: DATASETS

In the next tables proteins with bold ID's correspond to the training examples and the rest of them to the test set.

TABLE 4. DESCRIPTION OF THE PROSITE 1 DATASET

Table with 2 columns: Family and Protein ID's. It lists various protein families such as PS00030, PS00038, and PS00061, along with their corresponding protein identifiers.

(continued)

TABLE 4. (Continued)

Family	Protein ID's
PS00198	<p> <b>DHSB-CYACA DHSB-PARDE DHSB-RICCN FER3-PLEBO FER-ALIAE FER-CLOST FIXG-RHIME FIXX-BRAJA HMC6-DESVH MAUM-METEX NIFJ-ECOLI NUIC-ARATH NUIM-NEUCR PORD-METJA PSAC-ORNSA RNF8-PASMU RNF6-ECOS7 Y208-METJA YD49-METJA YFHL-ECOLI AEGA-ECOLI ASRA-SALTY ASRC-SALTY COOF-RHORA DCA1-METMA DCA2-METMA DCA-METJA DCA-M-SETS DCA-M-METTE DCA-M-METTE DCMG-METTE DHSB-SCHPO DHSB-BACSU DHSB-CAEEL DHSB-CHOCR DHSB-COXBU DHSB-DROME DHSB-ECOLI DHSB-HUMAN DHSB-MYCGR DHSB-PORPU DHSB-RAT DHSB-RECAM DHSB-RICPR DHSB-SCHPO DHSB-USTMA DHSB-YEAST DMSB-ECOLI DMSB-HAEIN DPYD-CAEEL DPYD-HUMAN DPYD-PIG DSRB-ARCFU DSVB-DESGI DSVB-DESVH FDHB-METFO FDHB-METJA FDHB-METTF FDHB-WOLSU FDNH-ECOLI FDOH-ECOLI FDXX-HAEIN FDXN-ANASP FDXN-ANAVA FDXN-AZOCHE FDXN-BRAJA FDXN-RHILT FDXN-RHIME FDXN-RHISN FDXN-RHOCA FER1-AZOFI FER1-CAUCR FER1-CHILLI FER1-DESFAF FER1-DESIN FER1-DESVM FER1-METJA FER1-RHOPE FER1-RHORU FER1-SULTO FER2-CHILLI FER2-DESDN FER2-DESVM FER2-METJA FER2-RHOCA FER2-RHORU FER2-SULTO FER2-THEAC FER3-ANASP FER3-ANAVA FER3-DESFAF FER3-METJA FER3-RHISN FER3-RHOCA FER4-METJA FER5-METJA FER6-METJA FER7-METJA FER8-METJA FER9-AZOFI FERVA-AZOFI FER-ACIAM FER-BACSC FER-BACST FER-BACSU FER-BACTH FER-BUTME FER-CHILLI FER-CHIRVI FER-CLOAC FER-CLOBU FER-CLOPA FER-CLOPE FER-CLOSP FER-CLOTM FER-CLOTS FER-DESGI FER-ENTHI FER-MEGL FER-METBA FER-METTE FER-METT FER-MOOTH FER-MYCSM FER-MYCTU FER-PEPAS FER-PSEPK FER-PSEST FER-PYRAB FER-PYRFU FER-PYRIS FER-RICPR FER-SACER FER-STROR FER-SULAC FER-THEAC FER-THELI FER-THEMA FER-THEM FIXX-AZOCA FIXX-AZOFI FIBX-ECOLI FIXX-RHLEF FIXX-RHILP FIXX-RHILT FIXX-RHIME FIXX-RHISN FRPB-MYCLE FRPB-MYCTU FRD1-AQUAE FRD2-AQUAE FRDB-ECOLI FRDB-HAEIN FRDB-HELPI FRDB-HELPE FRDB-MYCTU FRDB-PROVU FRDB-WOLSU FRHG-METJA FRHG-METTH FRHG-METVO GLCP-ECOLI GLPC-ECOLI GLPC-HAEIN HMC2-DESVH HYBA-ECOLI HYCB-ECOLI HYCE-ECOLI HYDN-ECOLI HYFA-ECOLI HYFH-ECOLI IORA-ARCFU IORA-METTH IORA-PYRAB IORA-PYRHO IORA-PYRKO MAUM-METFL MAUM-METME MAUM-PARDE MAUN-METEX MAUN-METTL MAUN-PARDE NAF-ECOLI NAFP-HAEIN NAFG-ECOLI NAFG-HAEIN NAPH-ECOLI NAPH-HAEIN NIF-ANASP NIF-ENTAG NIF-KLEPN NIF-RHORU NIF-SYNY3 NOO9-PARDE NOO9-THEH NRFC-ECOLI NRFC-HAEIN NUG2-RHIME NUIC-MAIZE NUIC-MARPO NUIC-MESVI NUIC-ORYSA NUIC-PLEBO NUIC-SPIOL NUIC-SYNY3 NUIC-TOBAC NUIC-WHEAT NUIC-WHEAT NUIM-ARATH NUIM-BOVIN NUIM-CAEEL NUIM-HUMAN NUIM-RECAM NUIM-SOLTU NUIM-TOBAC NUIM-TRYBB NUOI-BCUAI NUOI-ECOLI NUOI-MYCTU NUOI-RHOCA NUOI-RICCN NUOI-RICPR PHFI-CLOPA PHFI-DESZH PHFI-DESVM PHFI-SALTY PORD-METTH PORD-PYRAB PORD-PYRFU PORD-PYRHO PORD-THEMA PSAC-ANASP PSAC-ANTSP PSAC-ARATH PSAC-CHLRE PSAC-CHLVU PSAC-CYAPA PSAC-PSAC PSAC-EUGGR PSAC-FREDI PSAC-GUITH PSAC-MAIZE PSAC-MARPO PSAC-MASLA PSAC-MESVI PSAC-ODOSI PSAC-PEA PSAC-PINTH PSAC-PORPU PSAC-SKECO PSAC-SPIOL PSAC-SYNEL PSAC-SYNP2 PSAC-SYNP6 PSAC-SYNY3 PSRB-WOLSU RDXA-RHOSH RDXB-RHOSH RNF8-BUCAI RNF8-ECOS7 RNF8-ECOLI RNF8-HAEIN RNF8-PSEAE RNF8-RHOCA RNF8-VIBCH RNF8-BUCAI RNF8-ECOLI RNF8-HAEIN RNF8-PASMU RNF8-PSEAE RNF8-RHOCA RNF8-VIBCH VORD-METTH VORD-PYRAB VORD-PYRFU VORD-PYRHO Y092-METJA Y264-METJA Y492-MYCTU Y578-METJA Y726-METJA Y870-METJA YA43-HAEIN YCCM-ECOLI YCXI-PORPU YDIJ-ECOLI YDIJ-HAEIN YDIT-ECOLI YEIA-ECOLI YFHL-HAEIN YFRA-PROVU YG84-METTH YGFS-ECOLI YGFT-ECOLI YGL5-BACST YJES-ECOLI YJWJ-ECOLI YKGF-ECOLI YNFG-ECOLI YSAA-ECOLI         </b> </p>
PS00211	<p> <b>ABC2-HUMAN APPD-BACSU FTSE-HAEIN HISP-SALTY KST1-ECOLI LCCL-LACLA LMRA-LACL LOLD-BUCAI MDLB-BUCAI MKL-MYCTU MODC-HAEIN MRP2-RABIT NKD-ECOLI NODI-AZOCA NODI-RHISN NOSF-PSEST NRKD-SYNY3 OPPF-LACLA OPPF-MYCPN POTA-MYCGE RFBM-MYXXA SUFC-ECOLI UVRA-BRUAB UVRA-STRMU VEXC-SALTI WHIT-ANOAL Y348-CHLPN YF08-METJA YJKH-HAEIN YXDL-BACSU AAPP-RHILV AB11-HUMAN AB11-HUMAN AB11-HUMAN AB11-RABIT AB11-RAT ABC1-HUMAN ABC1-MOUSE ABC1-SCHPO ABC2-MOUSE ABC3-HUMAN ABC6-HUMAN ABC7-HUMAN ABC7-MOUSE ABC8-HUMAN ABCA-ARSA ABCR-HUMAN ABCX-ANTSP ABCX-CYACA ABCX-CYAPA ABCX-GALSU ABCX-GUITH ABCX-ODOSI ABCX-PORPU ABCX-STRMU ABCX-METN-HAEIN ABD2-HUMAN ABD3-HUMAN ABD3-MOUSE ABD3-RAT ABD4-HUMAN ABD4-MOUSE ABF2-HUMAN ABG1-HUMAN ABG1-MOUSE ABG2-HUMAN ABG3-MOUSE ABG4-HUMAN ABG5-HUMAN ABG5-MOUSE ABG5-RAT ABG8-HUMAN ABG8-MOUSE ABG8-RAT ACC8-CRICR ACC8-HUMAN ACC8-RAT ACC9-HUMAN ACC9-MOUSE ACC9-RABIT ACC9-RAT ADCC-STRMP ADP1-YEAST FBPC-ACFPI FBPC-ECOLI FBC1-HAEIN AGLK-RHIME ALD-HUMAN ALD-MOUSE ALSA-ECOLI AMIE-STRPN AMIF-STRPN AOTF-PSEAE APPF-BACSU APRD-PSEAE ARAG-ECOLI ARTP-ECOLI ARTP-HAEIN ATMI-YEAST BCRA-BACLI BEXA-HAEIN BFERI-SCHPO BPTI-YEAST BRAP-PSEAE BRAG-PSEAE BROW-DROME BROW-DROVI BTUD-ECOLI BZTD-RHOCA CBIO-SALTY CBRD-ERWCH CCMA-BRAJA CCMA-ECOLI CCMA-HAEIN CCMA-PARDE CCMA-RHOCA CDRI-CANAL CDR2-CANAL CDR3-CANAL CDR4-CANAL CFTR-BOVIN CFTR-CAVPO CFTR-HUMAN CFTR-MACMU CFTR-MOUSE CFTR-RABIT CFTR-RAT CFTR-SHEEP CFTR-SOUAC CFTR-XENLA CHVA-AGRT5 CHVD-AGRTU COMA-STRPN CTDR-NEIMA CTDR-NEIMB CVAB-ECOLI CVAB-BORPE CYDC-BACSU CYDC-ECOLI CYDC-HAEIN CYDD-BACSU CYDD-ECOLI CYDD-HAEIN CYSA-CHLVU CYSA-ECOLI CYSA-MARPO CYSA-MESVI CYSA-SALTY CYSA-SYNP7 CYSA-SYNY3 DPPP-BACSU DPPP-ECOLI DPPP-HAEIN DPPP-ECOLI DPPP-HAEIN DRRA-STRPE ECSA-BACSU EF3A-YEAST EF3B-YEAST EF3-CAEEL EF3-PNECA EF3-SCHPO EGO-ECOLI EXP8-STRPN FECE-ECOLI FEPC-ECOLI FHUC-BACSU FHUC-ECOLI FTSE-ECOLI GC20-YEAST GLNQ-BACST GLNQ-ECOLI GLTL-ECOLI CLUA-CORGL HEP-ANASP HFAC-CAUCR HISP-ECOLI FBC2-HAEIN HLY2-ECOLI HLYB-ACTAC HLYB-ECOLI HLYB-PASHA HLYB-PASSP HLYB-PROVU HMT1-SCHPO HNUV-YERPE HST6-CANAL KST5-ECOLI LACK-AGRRD LCN3-LACLA LCN3-LACLA LIVF-ARCFU LIVF-ECOLI LIVF-METJA LIVF-SALTY LIVG-ARCFU LIVG-ECOLI LIVG-METIA LIVG-SALTY LIVG-SALTY LOLD-BUCAP LOLD-ECOLI LOLD-NEIMA LOLD-NEIMB LOLD-NEIMB LOLD-XYLEFA MACB-ECOLI MALK-ECOLI MALK-ENTAE MALK-PHOLI MALK-SALTY MAMI-SCHPO MCHF-ECOLI MDLI-CANAL MDLI-YEAST MDL2-YEAST MDLA-BUCAI MDLA-ECOLI MDLB-ECOLI MDRI-CAEEL MDRI-CRIGR MDRI-ENTHI MDRI-HUMAN MDRI-LEIEN MDRI-MOUSE MDRI-RAT MDR2-CRIGR MDR2-MOUSE MDR2-RAT MDR3-CAEEL MDR3-CRIGR MDR3-ENTHI MDR3-HUMAN MDR4-MOUSE MDR4-DROME MDR4-ENTHI MDR5-DROME MDR-LEITA MDR-PLAFF MESD-LEUME MGLA-ECOLI MGLA-HAEIN MGLA-MYCGE MGLA-MYCPN MGLA-SALTY MGLA-TREPA MKL-MYCLE MNTA-SYNY3 MOPC-AZOFI MOPC-ECOLI MOPC-MYCTU MOPC-RHOCA MOPC-ECOLI MRPI-HUMAN MRP2-HUMAN MRP2-RAT MRP3-HUMAN MRP3-RAT MRP4-HUMAN MRP5-HUMAN MRP5-MOUSE MRP5-RAT MRP6-HUMAN MRP6-RAT MSBA-ECOLI MSBA-HAEIN MSMK-STRMU MSTMX-BACSU MSRA-STAEF NASD-KLEPN NATA-BACSU NDVA-RHIME NIKE-ECOLI NIST-LACLA NOCP-AGRT5 NODI-BRAJA NODI-RHIGA NODI-RHILLO NODI-RHILT NODI-RHILV NODI-RHIME NODI-RHIS3 NRTC-SYNP7 NRTC-SYNY3 NRTD-SYNP7 OCCP-AGRTU OCCP-RHIME OPAA-BACSU OPBA-BACSU OPBA-BACSU OPDB-BACSU OPDD-ECOLI OPDD-HAEIN OPDD-LACLA OPDD-LACL OPDD-MYCGE OPDD-MYCPN OPDD-SALTY OPFF-BACSU OPFF-ECOLI OPFF-HAEIN OPFF-MYCGE OPFF-SALTY OPFF-STRMU OPFF-STRPY P29-MYCGE P29-MYCHR P29-MYCPN PDR5-YEAST PDRA-YEAST PDRB-YEAST PDRD-YEAST PDRF-YEAST PDRV-YEAST PEDD-PEDAC PHNC-ECOLI PHNK-ECOLI PHNL-ECOLI PMD1-SCHPO POTA-ECOLI POTA-ECOLI POTA-HAEIN POTA-MYCPN POTA-SALTY POTG-ECOLI PROV-ECOLI PRDTR-ERWCH PSTB-ECOLI PSTB-EDWTA PSTB-ENTCL PSTB-METJA PSTB-MYCGE PSTB-MYCTU PSTB-MYCPN PSTB-MYCTU PSTB-PASMU PSTB-RHILLO PSTB-SALTY PSTB-XYLEFA PXA1-YEAST PXA2-YEAST RBSA-BACSU RBSA-ECOLI RBSA-HAEIN RFBI-KLEPN RFB2-KLEPN RFBE-YEREN RTIB-ACFPI RT3B-ACFPL SAPD-ECOLI SAPD-HAEIN SAPD-SALTY SAPF-ECOLI SAPF-HAEIN SAPF-SALTY SCRT-DROME FBPC-SERMA SMOK-RHOSH SNO2-YEAST SPAT-BACSU SRTF-STRPY SSUB-BACSU SSUB-ECOLI ST6E-YEAST SVRD-PSESY TAGB-DICDI TAGC-DICDI TAGH-BACSU TAPI-HUMAN TAPI-MOUSE TAPI-RAT TAP2-HUMAN TAP2-MOUSE TAP2-RAT TAPB-ECOLI THIQ-ECOLI THIQ-HAEIN TLRK-STRFR TROB-TREPA UGPC-ECOLI UUP1-HAEIN UUP2-HAEIN UUP-BUCAI UUP-ECOLI UVRA-AQUAE UVRA-BACHD UVRA-BACSU UVRA-BORBU UVRA-CHLMU UVRA-CHLPN UVRA-CHLTR UVRA-DEIRA UVRA-ECOLI UVRA-HAEIN UVRA-HELPI UVRA-HAEIN UVRA-LACLA UVRA-METTH UVRA-MICLU UVRA-MYCGE UVRA-MYCPN UVRA-MYCTU UVRA-NEIGO UVRA-PARDE UVRA-PASMU UVRA-PROMI UVRA-PSELE UVRA-RHIME UVRA-RICPR UVRA-SALTY UVRA-SERMA UVRA-STRCO UVRA-SYNY3 UVRA-THEMA UVRA-THEHT UVRA-THETH UVRA-TREPA UVRA-VITST UVRA-ZYMO V296-BACSU WHIT-ANOAL WHIT-CERCA WHIT-DROME WHIT-LUCCU XYLG-ECOLI XYLG-HAEIN Y014-MYCGE Y014-MYCPN Y015-MYCGE Y015-MYCPN Y035-METJA Y035-TREPA Y036-HAEIN Y065-MYCGE Y065-MYCPN Y068-CHLTR Y075-SYNY3 Y089-METJA Y121-METJA Y124-THEMA Y179-MYCGE Y179-MYCPN Y180-MYCGE Y180-MYCPN Y182-SYNY3 Y187-MYCGE Y187-MYCPN Y303-MYCGE Y303-MYCPN Y304-MYCGE Y304-MYCPN Y318-BORBU Y339-CHLMU Y352-THEMA Y354-HAEIN Y361-HAEIN Y382-RHIME Y412-METJA Y415-SYNY3 Y416-CHLTR Y467-MYCGE Y467-MYCPN Y468-MYCGE Y468-MYCPN Y4FO-RHISN Y4GM-RHISN Y4MK-RHISN Y4OS-RHISN Y4TH-RHISN Y4TR-RHISN Y4YS-RHISN Y542-CHLPN Y663-HAEIN Y664-HAEIN Y697-CHLMU Y700-RICPR Y719-METJA Y796-METJA Y799-ANASP Y873-METJA Y888-HELPI Y888-HELPE Y986-MYCTU YA23-METJA YA51-HAEIN YA78-HAEIN YADG-ECOLI YATR-BACFI YAWB-SCHPO YBBA-ECOLI YBBL-ECOLI YBH-ECOLI YBIT-ECOLI YBT1-YEAST YBXA-BACSU YC72-HAEIN YC72-MYCTU YC73-MYCTU YC81-MYCTU YCBN-BACSU YCFI-YEAST YCFIV-ECOLI YCKI-BACSU YCXD-CYAPA YD34-MYCPN YD48-MYCTU YD49-MYCTU YD67-METJA YDCT-ECOLI YDDA-ECOLI YDDO-ECOLI YDDP-ECOLI YDIF-BACSU YE67-HAEIN YE70-HAEIN YE74-HAEIN YECC-ECOLI YEHC-ECOLI YEJF-ECOLI YEM6-YEAST YD01-SCHPO YFCS-YEAST YFEB-YERPE YFIB-BACSU YFIC-BACSU YNT9-SCHPO YG18-HAEIN YHBB-GAZO YHCB-ECOLI YHBB-HAEIN YHBB-KLEPN YHBB-PSEPU YHBB-THIFE YHCB-BACSU YHCH-BACSU YHDS-YEAST YHDC-ECOLI YHES-ECOLI YHES-HAEIN YHIH-ECOLI YHI9-MYCTU YJJK-ECOLI YK83-YEAST CED7-CAEEL YLIA-ECOLI YMEB-LACLA YN26-MYCTU YN99-YEAST YNID-ECOLI YOH5-YEAST YOJI-ECOLI YORI-YEAST YP64-MYCTU YPC3-CAEEL YPHE-ECOLI YQSC-CAEEL YQGI-BACSU YQKG-BACSU YQIZ-BACSU YRBF-ECOLI YRBF-HAEIN YSCL-STRGC YTRF-ECOLI YTRF-ECOLI YTMN-BACSU YTMN-BACSU YTRF-ECOLI YWIA-BACSU YXEO-BACSU YYBJ-BACSU ZNUC-BUCAI ZNUC-ECOLI ZNUC-HAEIN ZURA-LISIN ZURA-LISMO         </b> </p>

(continued)

TABLE 4. (Continued)

Family	Protein ID's
PS00301	CYSN-RHTR CYSN-XYLEA EF1A-ARCFU EF1A-DICDI EF1A-SULSO EF1S-PORPU EF2-CHICK EF2-MESAU EFTU-CHLTR EFTU-FERIS EFTU-GRALE EFTU-MYCPN EFTU-NEPOL EFTU-TOBAC EFTU-XYLEA LEPA-MYCHY LEPA-MYCLE LEPA-MYCPN TETQ-PREIN TYPA-SYNY3 CYSN-BUCAI CYSN-ECOLI CYSN-MYCTU CYSN-PSEAE CYSN-RHIME EF10-XENLA EF11-CRIGR EF11-DAUCA EF11-DROME EF11-EUPCR EF11-HORVU EF11-HUMAN EF11-MOUSE EF11-RHIRA EF11-SCHPO EF11-XENLA EF12-DAUCA EF12-DROME EF12-ABSQL EF12-EUPCR EF12-HORVU EF12-HUMAN EF12-MOUSE EF12-RHIRA EF12-SCHPO EF12-XENLA EF13-RHIRA EF12-SCHPO EF13-XENLA EF14-ABSQL EF14-AERPE EF14-AJECA EF14-APIME EF14-ARATH EF14-ARTSA EF14-ARXAD EF14-ASHGO EF14-AURPU EF14-BLAHO EF14-BOMMO EF14-BRARE EF14-CAEEL EF14-CANAL EF14-CHICK EF14-CRYNE EF14-CRYPV EF14-DESMO EF14-EIMBO EF14-ENTHI EF14-EUCGR EF14-GIALA EF14-HALHA EF14-HALMA EF1A-PLAFK EF1A-PODAN EF1A-LYCES EF1A-MAIZE EF1A-MANES EF1A-METIA EF1A-METTH EF1A-METWA EF1A-EUCGR EF1A-ONCVO EF1A-ORYSA EF1A-PEA EF1A-PLAFK EF1A-PODAN EF1A-PUCGR EF1A-THYEA EF1A-THYEA EF1A-TOBAC EF1A-TRIRE EF1A-TRYBB EF1A-VICFA EF1A-SCHCO EF1A-SORMSA EF1A-SOYBN EF1A-STYLE EF1A-SULAC EF1A-TETPY EF1A-THEAC EF1A-TOBAC EF1A-TOBAC EF1A-TRYBB EF1A-VICFA EF1A-WHEAT EF1A-YARLI EF1A-YEAST EF1C-PORPU EF2-AERPE EF2-ARCFU EF2-BETU EF2-BLAHO EF2-CAEEL EF2-CANAL EF2-CHLKE EF2-CRIGR EF2-CRYPV EF2-DESMO EF2-DICDI EF2-DROME EF2-ENTHI EF2-HALHA EF2-HUMAN EF2-METBU EF2-METIA EF2-METMT EF2-METTE EF2-METTH EF2-METVA EF2-MOUSE EF2-PYRAB EF2-PYRHO EF2-PYRPU EF2-RABIT EF2-RAT EF2-SCHPO EF2-SULAC EF2-SULSO EF2-THEAC EF2-YEAST EFG1-BORBU EFG1-STRO EFG1-SYNY3 EFG1-TREPA EFG1-YEAST EFG2-BORBU EFG2-STRO EFG2-SYNY3 EFG2-TREPA EFG2-YEAST EFGC-PEA EFGC-SOYBN EFGM-MYCTU EFGM-SYNY3 EFGM-RAT EFG-AGRTU EFG-APPPP EFG-AQUAE EFG-AQUPY EFG-BACHD EFG-BACST EFG-YEAST EFG-BUCAI EFG-CHLMU EFG-CHLPN EFG-CHLTP EFG-ECOLI EFG-HAEN EFG-HELPI EFG-HELPU EFG-MICLU EFG-MYCGE EFG-MYCLE EFG-MYCPN EFG-MYCTU EFG-NEIGO EFG-PASMU EFG-PLARO EFG-RICCN EFG-RICPR EFG-SALTY EFG-SPIPL EFG-STIAM EFG-STRYR EFG-STRRR EFG-SYNP6 EFG-THEMA EFG-THEAT EFG-THICU EFG-UREPA EFTI-SOYBN EFTI-STRO EFTI-STRUC EFTI-STRRR EFT2-SOYBN EFT2-STRRR EFT3-STRO EFT3-STRRR EFT4-PASMU EFT4-PASMU EFTU-AGRTU EFTU-APPPP EFTU-AQUAE EFTU-AQUY EFTU-ARATH EFTU-ASTLO EFTU-BACFR EFTU-BACHD EFTU-BACST EFTU-BACSU EFTU-BORBU EFTU-BOVIN EFTU-BRELN EFTU-BRYPL EFTU-BUCAI EFTU-BUCAP EFTU-BUCMH EFTU-BUCSC EFTU-BURCE EFTU-CAMJE EFTU-CHACO EFTU-CHLUA EFTU-CHLMU EFTU-CHLRE EFTU-CHLNU EFTU-CHLVU EFTU-CHLWU EFTU-CODFR EFTU-COLOB EFTU-CORGL EFTU-COSCS EFTU-CYAPA EFTU-CYME EFTU-CYTLU EFTU-DEIRA EFTU-DEISP EFTU-DERMA EFTU-ECOLI EFTU-EIKKO EFTU-EUGGR EFTU-FIBUS EFTU-FLAFE EFTU-FLESI EFTU-GLOSI EFTU-GLOVI EFTU-GONPE EFTU-GUITH EFTU-GYMSU EFTU-HAEN EFTU-HELPI EFTU-HELPU EFTU-HERAU EFTU-HUMAN EFTU-MANSQ EFTU-MESVI EFTU-MICLU EFTU-MYCGA EFTU-MYCGE EFTU-MYCHO EFTU-MYCLE EFTU-MYCTU EFTU-NEIGO EFTU-ODOSI EFTU-PANMO EFTU-PEA EFTU-PHOEC EFTU-PLARO EFTU-PLEBO EFTU-PORPU EFTU-PROHO EFTU-PSEAE EFTU-RECAM EFTU-RHILU EFTU-RICPR EFTU-SALTY EFTU-SALTY EFTU-SHEPU EFTU-SPIAU EFTU-SPIPL EFTU-STIAU EFTU-STRAU EFTU-STRCJ EFTU-STRLU EFTU-STRMU EFTU-STROK EFTU-STRPY EFTU-SYNP6 EFTU-SYNP7 EFTU-SYNY3 EFTU-TAXOC EFTU-THEAQ EFTU-THEMA EFTU-THETH EFTU-THICU EFTU-TREHY EFTU-TREPA EFTU-UREPA EFTU-WOLSU EFTU-YEAST EFTY-CANAL ERF2-PCIP ERF2-SCHPO ERF2-YEAST GSP1-HUMAN GUF1-YEAST HBS1-YEAST LEPA-AQUAE LEPA-BACHD LEPA-BACSU LEPA-BORBU LEPA-BORPE LEPA-BUCAI LEPA-CHLMU LEPA-CHLNU LEPA-CHLTP LEPA-ECOLI LEPA-HAEN LEPA-HELPI LEPA-HELPU LEPA-LACLA LEPA-MYCGE LEPA-MYCTU LEPA-MYCLE LEPA-PSEFL LEPA-RICPR LEPA-SALTY LEPA-STRO EFTU-SYNY3 LEPA-THEMA LEPA-TREPA NODQ-AZOBK NODQ-RHIME NODQ-RHIS3 NODQ-RHISB NODQ-RHINT OTRA-STRRM RFB-BACNO RFB3-BUCAI RFB3-ECOLI RFB3-HAEN RFB3-LACLA RFB3-PASMU RFB3-SALTY RFB3-STAAU RFB3-SYNY3 RFB3-DESBA SELB-ECOLI SELB-HUMAN SELB-METIA SELB-MOOTH SELB-MOUSE SN14-YEAST TET1-ENTFA TETS-ENTFA TET9-ENTFA TETM-NEIME TETM-STAAU TETM-STRLI TETM-STRN TPETM-UREUR TETO-CAMCO TETO-CAMJE TETO-STRMU TETO-STRN TPET-CTOPE TETQ-BACFR TETQ-BACTN TETQ-BACTN TETQ-PRRU TETS-LACLA TETS-LISMO TETW-BUTH TYPA-BACSU TYPA-BUCAI TYPA-ECOLI TYPA-HAEN TYPA-HELPI TYPA-HELPU USS1-HUMAN USS1-MOUSE YE14-SCHPO YNQ3-YEAST Y081-CAEEL

TABLE 5. DESCRIPTION OF THE PROSITE 2 DATASET

Family	Protein ID's
PS00070	DHAE-MACPR DHAX-HUMAN DHA1-BOVIN HPPC-ECOLI YH9Y-YEAST GABD-ECOLI MAOC-ECOLI DHA4-YEAST DHA3-BACSU DHA5-YEAST YLQ6-CAEEL DHAS-CHICK DHAM-BOVIN PUT2-HUMAN MMSA-CAEEL ALDA-ECOLI ALDB-ECOLI ASTD-ECOLI ASTD-PSEAE CALB-CAUCR CALB-PSEAE CALB-PSESP CROM-OCTDO CROM-OMMSL DHA1-BACSU DHA1-CHICK DHA1-ENTHI DHA1-HORSE DHA1-HUMAN DHA1-MOUSE DHA1-RAT DHA1-SHEEP DHA2-ALCEU DHA2-BACST DHA2-BACSU DHA2-HUMAN DHA2-MOUSE DHA2-RAT DHA2-YEAST DHA3-YEAST DHA4-HUMAN DHA4-MOUSE DHA4-RAT DHA5-BOVIN DHA5-HUMAN DHA6-HUMAN DHA6-YEAST DHA7-HUMAN DHA8-HUMAN DHA9-POLMI DHA8-AMAHF DHA8-ATRHO DHA8-BACSU DHA8-BETVU DHA8-ECOLI DHA8-GADCA DHA8-HORVU DHA8-ORYSA DHA8-RHIME DHA8-SPIOL DHAC-RAT DHAE-ELED DHAF-VIBHA DHAG-HUMAN DHAG-PIG DHAL-AGABI DHAL-ALITAL DHAL-ASPNQ DHAL-BACST DHAL-CLAHE DHAL-DEIRA DHAL-ECOLI DHAL-EMENI DHAL-ENCUB DHAL-MYCTU DHAL-PSEOL DHAL-PSESP DHAL-RHORU DHAL-STROCO DHAL-VIBCH DHAM-HORSE DHAM-HUMAN DHAM-LEITA DHAM-MESAU DHAM-MOUSE DHAM-RAT DHAN-MACPR DHAP-BOVIN DHAP-HUMAN DHAP-MOUSE DHAP-RAT DHAX-PEA DHAX-YEAST DHAY-YEAST DMPC-PSESP FEAB-ECOLI FTDH-HUMAN FTDH-RAT GABD-DEIRA GABD-RHISN GABD-SYNY3 GAPN-MAIZE GAPN-NICLE GAPN-PEA GAPN-STRMU MMSA-BACSU MMSA-BOVIN MMSA-HUMAN MMSA-PSEAE MMSA-RAT NAHF-PSESP PUT2-AGABI PUT2-YEAST PUTA-ECOLI PUTA-KLEAE PUTA-RHIME PUTA-SALTY ROCA-BACSU SSSD-HUMAN SSSD-RAT THCA-RHOER UGA5-YEAST XYC2-ACIGB XYLC-PSEPU XYLG-PSEPU Y4UC-RHISN YDCW-ECOLI YM00-YEAST YNEI-ECOLI
PS00077	COXI-THETH COXI-BACFI COXI-DIDMA COXI-ASCSU COXI-HORSE COXI-EPHEQ FIXN-AZOCA COXI-SYNYV COXI-CRION COXI-ALLMA AOXI-AERPE COXI-PEA COXI-RHOSH COXI-SOYBN COXI-PLABE CO13-THETH CO14-BRAJA COX1-ACACA COX1-ALBCO COX1-ALBTU COX1-AMICA COX1-ANAPL COX1-ANOGA COX1-ANOQU COX1-APLI COX1-APTAD COX1-ARATH COX1-ARTSF COX1-ASTPE COX1-BACP3 COX1-BACSU COX1-BALMU COX1-BALPH COX1-BETVU COX1-BOVIN COX1-BRAJA COX1-CAEEL COX1-CANFA COX1-CANSI COX1-CAPHI COX1-CARAU COX1-CASBE COX1-CERSI COX1-CHICK COX1-CHLRE COX1-CHOB1 COX1-CHOCR COX1-CHOFU COX1-CHOOC COX1-CHORO COX1-COTIA COX1-CROLA COX1-CYACA COX1-CYPCA COX1-DASNO COX1-DINSE COX1-DROME COX1-DRONO COX1-DROYA COX1-EMENI COX1-EQUAS COX1-FELCA COX1-GADMO COX1-GEOSD COX1-GOMVA COX1-HALGR COX1-HALHA COX1-HANWI COX1-HIPAM COX1-HUMAN COX1-KLULA COX1-LATCH COX1-LEITA COX1-LEPOC COX1-LEPSQ COX1-LOCM1 COX1-LUMTE COX1-MACRO COX1-MAIZE COX1-MARPO COX1-MEGAT COX1-METSE COX1-MOUSE COX1-MYCTU COX1-MYTED COX1-MYXGL COX1-NEUCR COX1-NOTPE COX1-OENBE COX1-ONCMY COX1-ORNAN COX1-ORYSA COX1-PANBU COX1-PAPHA COX1-PARLI COX1-PARTE COX1-PECCA COX1-PELSU COX1-PETMA COX1-PHOVI COX1-PHYME COX1-PIHYO COX1-PIG COX1-PISOC COX1-PLACH COX1-PLAFA COX1-PODAN COX1-POLOR COX1-POLSP COX1-POLSX COX1-PMNM COX1-PONPA COX1-PROWI COX1-RABIT COX1-RAT COX1-RHEAM COX1-RHILE COX1-RHISA COX1-RHUN COX1-RHOCA COX1-RICPR COX1-SACDO COX1-SALSA COX1-SALTR COX1-SCAPL COX1-SCHPO COX1-SCYCA COX1-SHEEP COX1-SORBI COX1-SQUAC COX1-STRCA COX1-STRUC COX1-SYNY3 COX1-TETPY COX1-TINMA COX1-TRIRU COX1-TRYBB COX1-WHEAT COX1-XENLA COX1-YEAST COXN-BRAJA CX1A-PARDE CX1B-PARDE CYOB-BUCAI CYOB-ECOLI CYOB-PSEPU FIXN-AGRT7 FIXN-BRAJA FIXN-RHIME NORB-PSEAE NORB-PSEST QOX1-ACEAC QOX1-BACSU QOX1-SULAC QOXM-SULAC
PS00118	PA21-NAJMO PA21-HORSE PA2H-BUNFA PA2E-PSEAU PA2C-CRODU PA2H-BOTIR PA2C-PSEAU PA2Z-HUMAN PA22-BUNMU PA23-NAJNG PA21-TRIGA PA21-ACAAN PA21-BOTPI PA2X-RAT PA2Z-PIG OC90-CAVPO OC90-HUMAN OC90-MOUSE PA20-BUNMU PA20-NOTSC PA20-PSEAU PA21-AGKHA PA21-AGKHP PA21-AGKPI PA21-BOTAS PA21-BOTJA PA21-BOTIR PA21-BOTMO PA21-BOVIN PA21-BUNMU PA21-CANFA PA21-CAVPO PA21-ERIMA PA21-HEMHA PA21-HUMAN PA21-LATSE PA21-MATBI PA21-MOUSE PA21-NAJME PA21-NAJOX PA21-NOTSC PA21-OXYSC PA21-PIG PA21-PSEAU PA21-RAT PA21-SHEEP PA21-TRIFL PA21-VIPAA PA21-VIPAZ PA22-ACAAN PA22-AGKHA PA22-AGKHP PA22-ASPC PA22-BITNA PA22-BOTAS PA22-BOTMO PA22-BOTPI PA22-CERGO PA22-ERIMA PA22-HELPU PA22-LATCO PA22-MATBI PA22-NAJKA PA22-NAJME PA22-NAJMO PA22-NOTSC PA22-OXYSC PA22-TRIGA PA22-TRIST PA22-VIPAZ PA23-AGKHP PA23-BOTAS PA23-BOTPI PA23-BUNMU PA23-HELPU PA23-HUMAN PA23-LATSE PA23-NAJKA PA23-NAJME PA23-NAJMO PA23-NOTSC PA23-OXYSC PA23-PSEAU PA23-TRIGA PA24-BUNMU PA24-DABRU PA24-LATSE PA24-TRIGA PA25-HUMAN PA25-MOUSE PA25-PSEAU PA25-RAT PA25-TRIGA PA25-TRIST PA26-BUNFA PA26-TRIGA PA27-DABRU PA27-TRIGA PA29-PSEAU PA2A-BUNFA PA2A-CRODU PA2A-HUMAN PA2A-MICNI PA2A-MOUSE PA2A-PSEAU PA2A-PSEPO PA2A-PSETE PA2A-RABIT PA2A-RAT PA2A-VIPAA PA2A-VIPDA PA2B-BUNFA PA2B-CRODU PA2B-MICNI PA2B-PSEPO PA2B-PSETE PA2B-TRIFL PA2B-TRIMU PA2B-VIPAA PA2C-MOUSE PA2C-PSETE PA2C-RAT PA2C-VIPAA PA2X-HUMAN PA2D-MOUSE PA2D-PSEAU PA2D-PSETE PA2E-HUMAN PA2E-MOUSE PA2F-HUMAN PA2F-MOUSE PA2G-PSEAU PA2H-AGKPI PA2H-ATRNM PA2H-LATCO PA2H-XENLA PA2I-VIPAA PA2I-VIPAZ PA2M-AGKCL PA2M-CAVPO PA2M-CROSS PA2N-BUNFA PA2N-CROSS PA2N-ECHCA PA2N-VIPAA PA2X-HUMAN PA2X-MOUSE PA2X-MOUSE PA2X-NOTSC PA2X-TRIFL PA2Y-HUMAN PA2Y-MOUSE PA2Y-TRIFL PA2Z-MOUSE PA2-APLA PA2-APIME PA2-BIFCA PA2-BITGA PA2-BOMTE PA2-CERCE PA2-CROAD PA2-CROAT PA2-DABRR PA2-ENHSC PA2-HELHU PA2-LATLA PA2-NAJAT PA2-NAJPA PA2-OPHHA PA2-RHONO PA2-TRIOK PA2-VIPBB
PS00180	GLNA-COLGL GLN4-PEA GLN2-DROME GLNA-HELPU GLNA-PANAR GLN3-RHILP GLN1-ARATH GLN5-MAIZE GLNA-PIG GLNA-PYRHO GLNA-THIFE GLNA-SALTY GLN3-PHAVU GLNA-NICPL GLN2-DAUCA GLN1-ALNGL GLN1-BRAJA GLN1-CHLRE GLN1-DAUCA GLN1-DROME GLN1-FRAAL GLN1-LOTIA GLN1-MAIZE GLN1-MEDA GLN1-MYCTU GLN1-ORYSA GLN1-PEA GLN1-PHAVU GLN1-RHILV GLN1-RHIME GLN1-SOYBN GLN1-STRPR GLN1-STRYR GLN1-VITVI GLN2-ARATH GLN2-BRAJA GLN2-CHLRE GLN2-FRAAL GLN2-HORVU GLN2-MAIZE GLN2-MEDA GLN2-MYCTU GLN2-ORYSA GLN2-PEA GLN2-PHAVU GLN2-RHILP GLN2-RHIME GLN2-SOYBN GLN2-STRHY GLN2-STRYR GLN2-VITVI GLN3-HORVU GLN3-LUPAN GLN3-MAIZE GLN3-MEDA GLN3-ORYSA GLN3-PEA GLN3-RHIME GLN4-MAIZE GLN4-PHAVU GLNA-AGABI GLNA-ANASP GLNA-AQUAE GLNA-ARCFU GLNA-AZOBK GLNA-AZOCA GLNA-AZOV1 GLNA-BACCE GLNA-BACFR GLNA-BACSU GLNA-BOVIN GLNA-BUTPI GLNA-CAEEL GLNA-CHICK GLNA-CLOSA GLNA-CRIOA GLNA-DUNSA GLNA-ECOLI GLNA-FREDI GLNA-HAEN GLNA-HALNA GLNA-HALVO GLNA-HELPI GLNA-HUMAN GLNA-LACDE GLNA-LACLA GLNA-LACSA GLNA-LUPLU GLNA-METCA GLNA-METIA GLNA-METMP GLNA-METTH GLNA-METVO GLNA-MOUSE GLNA-NEIGO GLNA-PASMU GLNA-PINSY GLNA-PROVU GLNA-PYRAB GLNA-PYRPU GLNA-PYRKO GLNA-PYRWO GLNA-RHOC GLNA-RHOCA GLNA-RHOSH GLNA-SCHPO GLNA-SQUAC GLNA-STAAU GLNA-STROCO GLNA-SULAC GLNA-SULSO GLNA-SYNP2 GLNA-SYNY3 GLNA-THEMA GLNA-TRITH GLNA-VIBAL GLNA-VIBCH GLNA-VIGAC GLNA-XENLA GLNA-YEAST GLNC-BRANA GLNC-MAIZE YCJR-ECOLI

(continued)

TABLE 5. (Continued)

Table with 2 columns: Family and Protein ID's. Rows include families PS00215, PS00217, and PS00338 with extensive lists of protein IDs.

TABLE 6. DESCRIPTION OF THE GPCR DATASET

Table with 2 columns: Subfamily and Protein ID's. Subfamily 'Amine' lists numerous protein IDs such as 5H1A-RAT, 5H1B-CAVPO, etc.

(continued)





TABLE 6. (Continued)

Subfamily	Protein ID's
Prostanoid	O00326 PD2R-MOUSE PE22-MOUSE PE23-BOVIN PE23-HUMAN PE23-RABIT PF2R-MOUSE P12R-BOVIN Q9R261 TA2R-BOVIN O00325 O15191 O35932 O46657 O75228 PD2R-HUMAN PE21-HUMAN PE21-MOUSE PE21-RAT PE22-CANFA PE22-HUMAN PE22-RAT PE23-MOUSE PE23-PIG PE23-RAT PE24-HUMAN PE24-MOUSE PE24-RABIT PE24-RAT PF2R-BOVIN PF2R-HUMAN PF2R-RAT PF2R-SHEEP P12R-HUMAN P12R-MOUSE P12R-RAT Q9BGL8 Q9D627 Q9TU16 TA2R-CERAE TA2R-HUMAN TA2R-MOUSE TA2R-RAT
Nucleotide-like	AA1R-BOVIN AA1R-RAT AA2A-RAT O57466 P2Y3-MELGA P2Y6-HUMAN P2YR-RAT Q99MT6 Q9ERK9 Q9HIC0 AA1R-CANFA AA1R-CAVPO AA1R-CHICK AA1R-HUMAN AA1R-RABIT AA2A-CANFA AA2A-CAVPO AA2A-HUMAN AA2A-MOUSE AA2B-CHICK AA2B-HUMAN AA2B-MOUSE AA2B-RAT AA3R-CANFA AA3R-HUMAN AA3R-RABIT AA3R-RAT AA3R-SHEEP GPRZ-HUMAN GPRZ-MOUSE O00398 O08766 O35811 P2UR-HUMAN P2UR-MOUSE P2UR-RAT P2Y3-CHICK P2Y4-HUMAN P2Y5-CHICK P2Y5-HUMAN P2Y6-RAT P2Y8-XENLA P2Y9-HUMAN P2YR-BOVIN P2YR-CHICK P2YR-HUMAN P2YR-MELGA P2YR-MOUSE Q9BXC5 Q9BXC1 Q9BYU4 Q9CFZ4 Q9DE05 Q9JIS7 Q9N1U0 Q9PU18 Q9R202 Q9W6C4

## ACKNOWLEDGMENT

The authors wish to thank the anonymous referees for the valuable comments that have improved the quality and the presentation of this work.

## REFERENCES

- Almeida, J.S., and Vingá, S. 2002. Universal sequence map (USM) of arbitrary discrete sequences. *BMC Bioinformatics* 3(6).
- Altshul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. 1990. Basic local alignment tool. *J. Mol. Biol.* 215, 403–410.
- Bailey, T.L., and Elkan, C. 1994. Fitting a mixture model by expectation maximization to discover motifs in biopolymers. *2nd Int. Conf. on Intelligent Systems for Molecular Biology*, 28–36.
- Bailey, T.L., and Gribskov, M. 1998. Combining evidence using p-values: Application to sequence homology searches. *Bioinformatics* 14, 48–54.
- Bishop, C.M. 1995. *Neural Networks for Pattern Recognition*, Oxford University Press, New York.
- Brázma, A., Jonasses, I., Eidhammer, I., and Gilbert, D. 1998. Approaches to the automatic discovery of patterns in biosequences. *J. Comp. Biol.* 5(2), 277–303.
- Dayhoff, M.O., Schwartz, R.M., and Orcutt, B.C. 1978. *Atlas of Protein Sequence and Structure*, Vol. 5, Natl. Biomed. Res. Found., Washington, DC.
- Dempster, A.P., Laird, N.M., and Rubin, D.B. 1977. Maximum likelihood from incomplete data via the EM algorithm. *J. Roy. Statist. Soc. B* 39, 1–38.
- Durbin, R., Eddy, S., Krough, A., and Mitchison, G. 1998. *Biological Sequence Analysis: Probabilistic Models of Protein and Nucleic Acid*, Cambridge University Press, New York.
- Foressé, F.D., and Hagan, M.T. 1997. Gauss–Newton approximation to Bayesian regularization. *Proc. 1997 Int. Joint Conf. on Neural Networks*, 1930–1935.
- Henikoff, S.S., and Henikoff, J.G. 1994. Protein family classification based on searching a database of blocks. *Genomics* 19, 97–107.
- Hertz, G.Z., and Stormo, G.D. 1999. Identifying DNA and protein patterns with statistically significant alignments of multiple sequences. *Bioinformatics* 15(7/8), 563–577.
- Hofmann, K., Bucher, P., Falquet, L., and Bairoch, A. 1999. The PROSITE database, its status in 1999. *Nucl. Acids Res.* 27, 215–219.
- Horn, F., Weare, J., Beukers, M.W., Hörsch, S., Bairoch, A., Chen, W., Edvardsen, O., Campagne, F., and Vriend, G. 1998. GPCRDB: An information system for G protein-coupled receptors. *Nucl. Acids Res.* 21(1), 227–281.
- Hughey, R., and Krogh, A. 1996. Hidden Markov models for sequence analysis: Extension and analysis of the basic method. *CABIOS* 12(2), 95–107.
- Jaakkola, T., Diekhans, M., and Haussler, D. 2000. A discriminative framework for detecting remote protein homologies. *J. Comp. Biol.* 7(1–2), 95–114.
- Karchin, R., Karplus, K., and Haussler, D. 2002. Classifying G-protein coupled receptors with support vector machines. *Bioinformatics* 18(1), 147–159.
- Karplus, K., Barrett, C., and Hughey, R. 1998. Hidden Markov models for detecting remote protein homologies. *Bioinformatics* 14(10), 846–856.
- Lawrence, C.E., Altschul, S.F., Boguski, M.S., Liu, J.S., Neuwland, A.F., and Wootton, J.C. 1993. Detecting subtle sequence signals: A Gibbs sampling strategy for multiple alignment. *Science* 226, 208–214.
- Logan, B., Moreno, P., Suzek, B., Weng, Z., and Kasif, S. 2001. A study of remote homology detection. Technical report CRL 2001/05, Cambridge Research Laboratory.

- Ma, Q., and Wang, J.T.L. 2000. Application of Bayesian neural networks to protein sequence classification. *ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, 305–309.
- MacKay, D.J.C. 1992. Bayesian interpolation. *Neural Comp.* 4, 415–447.
- Rigoutsos, I., Floratos, A., Parida, L., Gao, Y., and Platt, D. 2000. The emergency of pattern discovery techniques in computational biology. *Metabolic Eng.* 2, 159–177.
- Vapnik, V.N. 1979. *Estimation of Dependencies Based on Empirical Data*, Nauka, Birmingham, AL.
- Wang, J.T.L., Ma, Q., Shasha, D., and Wu, C.H. 2001. New techniques for extracting features from protein sequences. *IBM: Systems Journal* 40(2), 426–441.
- Wu, C.H., Zhap, S., Chen, H.L., Lo, C.J., and McLarty, J. 1996. Motif identification neural design for rapid and sensitive protein family search. *CABIOS* 12(2), 109–118.

Address correspondence to:

*Konstantinos Blekas*

*Department of Computer Science and Biomedical Research Institute—FORTH*

*University of Ioannina*

*GR-45110 Ioannina, Greece*

*E-mail: kblekas@cs.uoi.gr*