

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 (Br  zma *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

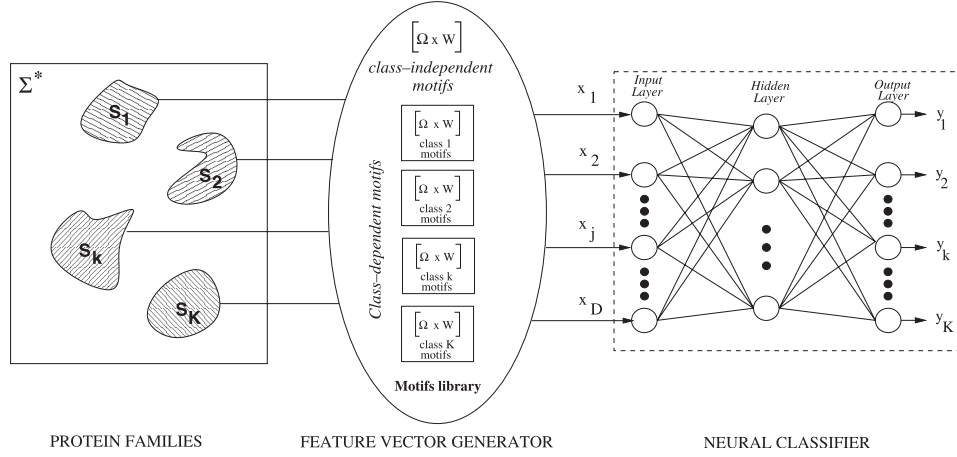


FIG. 1. The architecture of the proposed classification scheme.

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 family¹. 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.

¹The general relative frequencies of amino acids indicate the background information in a protein family and can be presented as a probabilistic vector ρ 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 matrix². 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.,

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

²The 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 θ ($\in [0, 1]$), 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 θ 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 $F(\mathbf{w})$:

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

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 α and β at the minimum point \mathbf{w}_{MP} can be computed as follows:

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

The quantity γ 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

Problem: PROSITE 1 ($K = 6$)			Problem: PROSITE 2 ($K = 7$)		
PROSITE family	Positive data	Training set (avg length of seqs)	PROSITE family	Positive data	Training set (avg length of seqs)
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 .³ 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 ROC₅₀ 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

³Experiments 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 ROC₅₀ 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>	<i>N_g 2-gram features</i>	<i>D motif-based features</i>
PROSITE 1	174	$5 \times 6 = 30$
PROSITE 2	285	$5 \times 7 = 35$
GPCR	152	$5 \times 7 = 35$

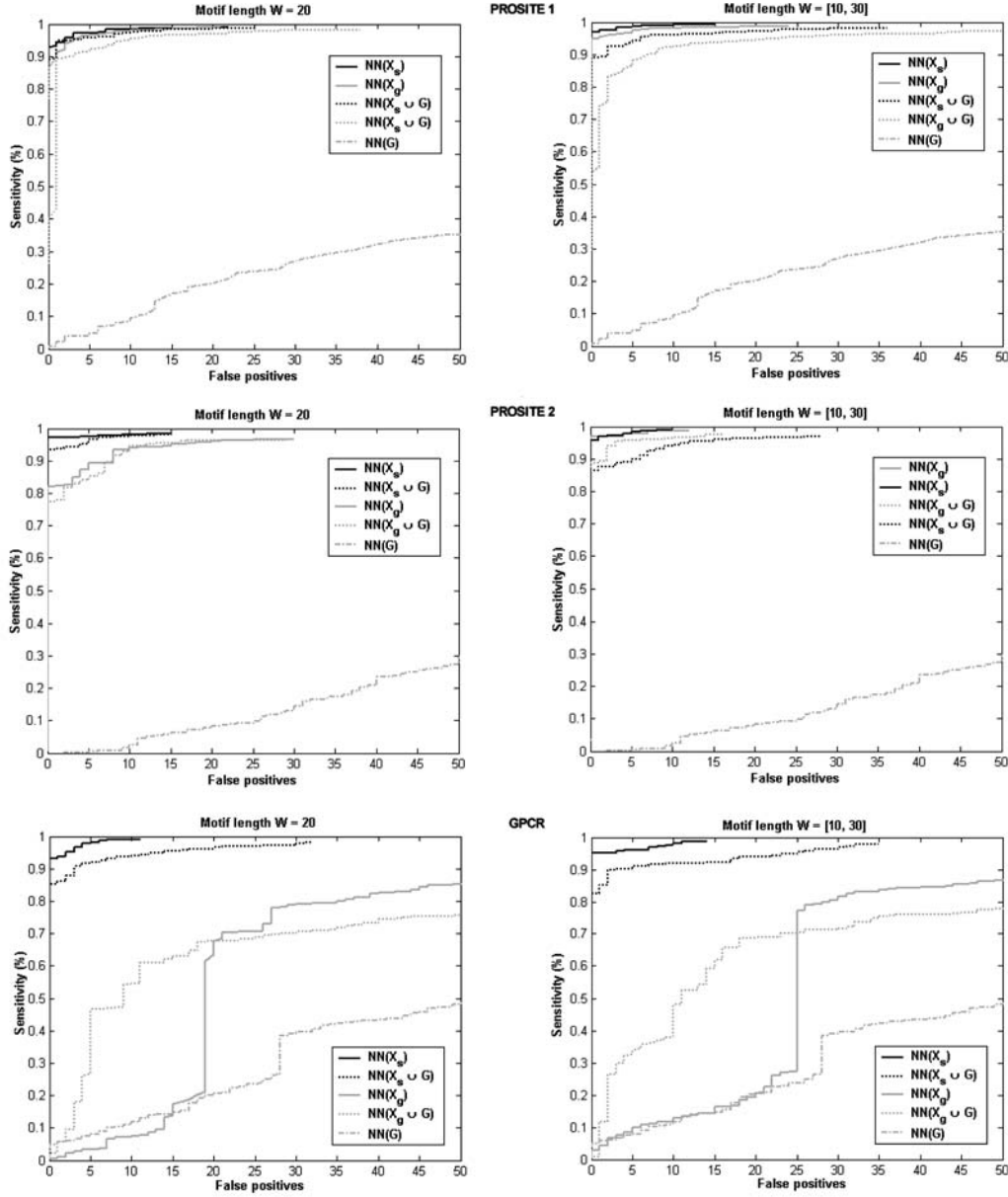


FIG. 2. ROC₅₀ curves illustrating the performance of the neural classifier on the three datasets using the five different feature vectors.

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 ROC₅₀ 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

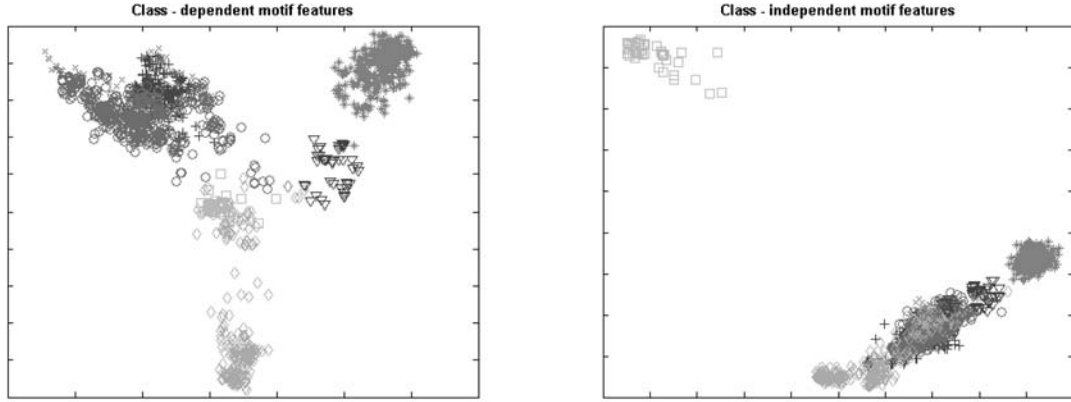


FIG. 3. The seven class regions in the GPCR dataset in the case of class-dependent and class-independent features. The data have been projected in two dimensions using PCA.

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.⁴ 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.⁵

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 (ROC₂₅). The performance of the neural classifier and MAST was given by two curves respectively⁶ concerning motifs of fixed ($W = 20$) and variable length ($W = [10, 30]$), while the last one corresponds to SAM performance.

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

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

⁶The curves for the neural classifier performance were the best plots from the corresponding ROC₅₀ diagrams in Fig. 2.

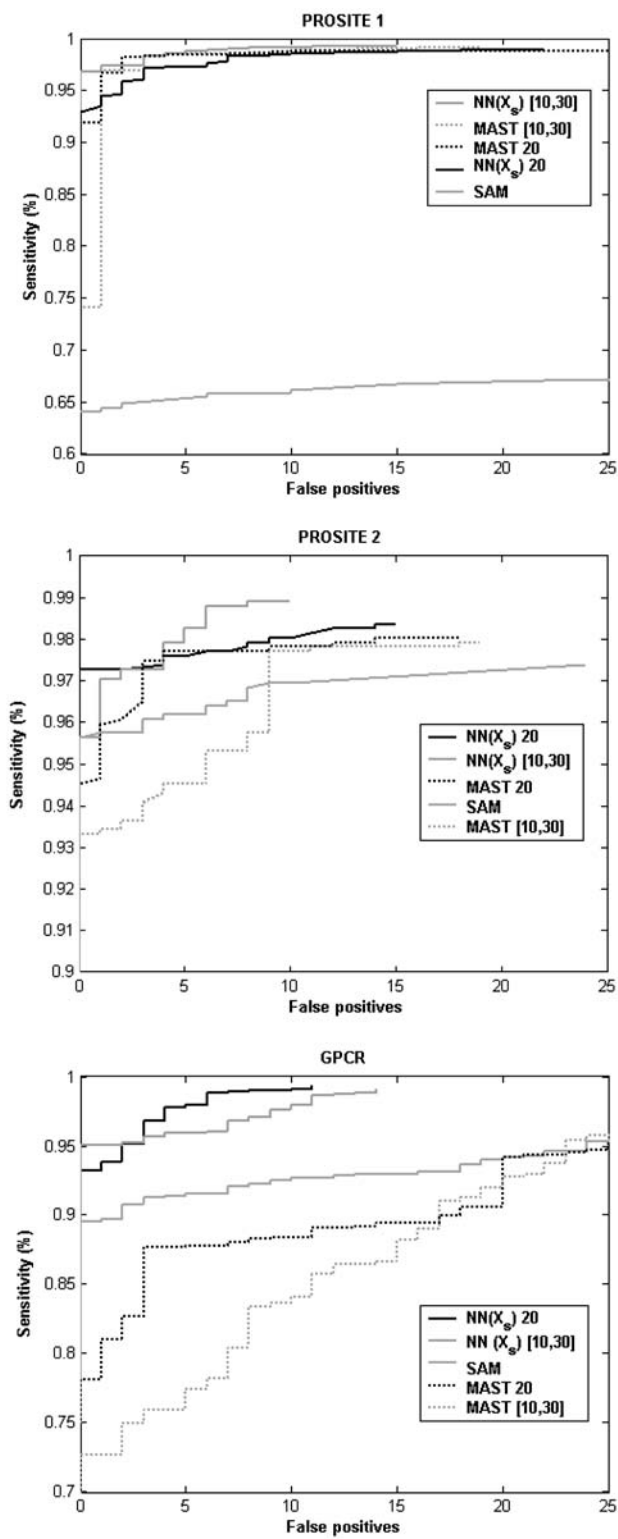


FIG. 4. ROC₂₅ curves for the three methods (neural (NN), MAST, and SAM) on the three datasets.

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

Family	Protein ID's
PS00030	CB20-HUMAN GAR2-SCHPO HRB1-YEAST HS49-YEAST IF34-MOUSE NAB4-YEAST PAB3-ARATH RB27-DROME RN15-YEAST ROA1-MOUSE ROC3-NICSY RU17-HUMAN RU17-YEAST RU1A-DROME RUIZB-HUMAN SFPQ-HUMAN U2AF-CAEEL U2AF-HUMAN U2AG-HUMAN K682-HUMAN A2BP-HUMAN A2BP-MOUSE ARP2-PLAFA ROAA-MOUSE CAZ-DROME CB20-XENLA CG79-HUMAN PM14-MOUSE CIRP-HUMAN CIRP-MOUSE CIRP-XENLA CPO-DROME CST2-HUMAN CSX1-SCHPO CTFI-SCHPO CUG1-HUMAN CUG1-MOUSE CWFS-SCHPO CYPE-DROME CYPE-HUMAN CYPE-MOUSE D111-ARATH ELAV-DROME ELAV-DROVI ELV1-HUMAN ELV1-MOUSE ELV2-HUMAN ELV2-MOUSE ELV3-HUMAN ELV3-MOUSE ELV4-HUMAN ELV4-RAT ELV4-RAT EWS-HUMAN EWS-MOUSE FCA-ARATH FUS-BOVIN FUS-HUMAN FUS-MOUSE G3B2-HUMAN G3B2-MOUSE G3BP-HUMAN G3BP-MOUSE G3BP-SCHPO GBP2-YEAST GR10-BRANA GRF1-HUMAN GRP1-SINAL GRP1-SORBI GRP2-SINAL GRP2-SORBI GRP7-ARATH GRP8-ARATH GRPA-MAIZE GRP-DAUCA IF34-CAEEL IF34-HUMAN IF34-SCHPO IF34-YEAST IF39-HUMAN IF39-SCHPO IF39-TOBAC IF39-YEAST IF4B-HUMAN IF4B-YEAST IF4H-HUMAN IF4H-MOUSE ISN1-YEAST LAA-XENLA LAB-XENLA LAHI1-SCHPO LAHI1-YEAST LA-AEDAL LA-BOVIN LA-DROME LA-HUMAN LA-MOUSE LA-RABIT LA-RAT MAT3-HUMAN MAT3-RAT MEI2-SCHPO MLO3-SCHPO MODU-DROME MSSP-HUMAN NAB3-YEAST NAM8-YEAST NGRI1-YEAST NONA-DROME NOP3-YEAST NOP4-YEAST NOP8-YEAST NOT4-YEAST NR54-HUMAN NRD1-SCHPO NRD1-YEAST NRPI-YEAST NSRI-YEAST NUCL-CHICK NUCL-HUMAN NUCL-MESAU NUCL-MOUSE NUCL-RAT NUCL-XENLA PAB1-HUMAN PAB1-MOUSE PAB2-ARATH PAB2-HUMAN PAB4-HUMAN PAB5-ARATH PABP-DROME PABP-SCHPO PAB1-XENLA PABP-YEAST PABX-ARATH PES4-YEAST PR71-PICAN PTB-HUMAN PTB-MOUSE PTB-PIG PTB-RAT PUB1-YEAST RB56-HUMAN RB87-DROME RB97-DROME RBM3-HUMAN RBM3-MOUSE RBM5-HUMAN RBM6-HUMAN RBM7-HUMAN RB8A-HUMAN RBM9-HUMAN RBMA-HUMAN RBMA-RAT RBMB-HUMAN RBMS-CHICK RBMS-HUMAN RBMS-MOUSE RBMS-XENLA RBPI-DROME RBPA-SYNY3 RDP-HUMAN RDP-MOUSE RN24-SCHPO RNPI-YEAST R021-XENLA R022-XENLA R031-XENLA R032-XENLA ROA0-HUMAN ROAI-BOVIN ROAI-DROME ROAI-HUMAN ROAI-MACMU ROAI-RAT ROAI-SCHAM ROAI-XENLA ROA2-HUMAN ROA2-MOUSE ROA3-HUMAN ROAB-ARTSA ROC1-ARATH ROC1-NICPL ROC1-NICSY ROC1-SPIOI ROC2-ARATH ROC2-NICPL ROC2-NICSY ROC3-ARATH ROC4-NICSY ROC5-NICSY ROC-HUMAN ROC-RAT ROH-HUMAN ROD-RAT ROF-HUMAN ROG-HUMAN ROG-MOUSE ROHI-HUMAN ROH2-HUMAN ROH-LHUMAN ROM-HUMAN ROR-HUMAN ROU2-HUMAN RS31-ARATH RS40-ARATH RS41-ARATH RT19-ARATH RU17-ARATH RU17-DROME RU17-MOUSE RU17-XENLA RU1A-HUMAN RU1A-XENLA RU1A-YEAST RX21-DROME S3B4-HUMAN SCE3-SCHPO SFK1-ARATH SFK1-HUMAN SFK2-CAEEL SFK2-CHICK SFK2-SFR2-HUMAN SFR2-MOUSE SFR3-HUMAN SFR4-HUMAN SFR5-HUMAN SFR5-MOUSE SFR5-RAT SFR6-HUMAN SFR6-RABIT SFR7-HUMAN SFR9-HUMAN SFRB-HUMAN SQD-DROME SR55-DROME SR44-HUMAN SR44-RAT SPR1-SCHPO SSB1-YEAST SXL-CERCA SXL-CHRRU SXL-DROME SXL-DROSU SXL-MEGSC SXL-MUSDO TIA1-HUMAN TIA1-MOUSE TIAR-HUMAN TIAR-MOUSE TR2A-HUMAN TRA2-DROME TRA2-DROVI U2AF-CAEBR U2AF-DROME U2AF-MOUSE U2AF-SCHPO U2AG-DROME U2AG-MOUSE U2AG-SCHPO U2RI-HUMAN U2RI-MOUSE U2RI-MOUSE U2RI-MOUSE WH13-YEAST XM52-DROME Y17-HUMAN YAC4-SCHPO YAG3-SCHPO YAS9-SCHPO YBF1-YEAST YBLC-SCHPO YQC9-SCHPO YD3-SCHPO YDB2-SCHPO YDC1-SCHPO CWF2-SCHPO YDR6-SCHPO YFK2-YEAST YGSB-YEAST YHC4-YEAST YHH5-YEAST YIS1-YEAST YIS5-YEAST YIS9-YEAST YKV4-YEAST YVL1-CAEEL YML7-YEAST YN26-YEAST YN8T-YEAST YNL0-YEAST YNRS-YEAST YP68-YEAST YP85-CAEEL YQ01-CAEEL YQ04-CAEEL YQQA-CAEEL YQOC-CAEEL YRA1-YEAST YS07-CAEEL YSX2-CAEEL
PS00038	AHR-RAT ARRS-MAIZE CBF1-YEAST DA-DROME ESM7-DROME HEN2-MOUSE HE33-MOUSE MAD4-MOUSE MITF-RAT MX11-BRARE MYC-HYLIA MYF6-HUMAN MYOD-BRARE NDHF1-MESAU SIM1-MOUSE TAL2-MOUSE TAL-HUMAN TF21-MOUSE TW51-HUMAN USF2-RAT AHR-HUMAN AHR-MOUSE ALRC-MAIZE ARN2-HUMAN ARN2-MOUSE ARNT-DROME ARNT-HUMAN ARNT-MOUSE ARNT-RABIT ARNT-RAT ASC1-HUMAN ASC1-MOUSE ASC1-RAT ASC1-XENLA ASC2-HUMAN ASC2-MOUSE ASC2-RAT AST3-DROME AST4-DROME AST5-DROME AST8-DROME ATH1-HUMAN ATH1-MOUSE NDF6-MOUSE NDF4-MOUSE NDF4-XENLA NGN2-MOUSE NGN3-MOUSE ATO-DROME BET3-MESAU BMAL1-HUMAN CBFI-KULUA CLOC-DROME CLOC-HUMAN CLOC-MOUSE CYCL-DROME DEI-DROME DPN-DROME EMC-DROME ESC1-SCHPO EMS3-DROME EMS5-DROME ESM8-DROME ESM8-DROME ESM8-DROME ESMC-DROME ESMC-DROME ESMC-DROME ESMC-DROME HAIR-DROVI HANI-CHICK HANI-HUMAN HANI-MOUSE HANI-RAT HANI-SHEEP HANI-XENLA HAN2-BRARE HAN2-CHICK HAN2-HUMAN HAN2-MOUSE HAN2-XENLA HEN1-HUMAN HEN1-MOUSE HEN2-HUMAN HESI-CHICK HESI-HUMAN HESI-MOUSE HESI-RAT HES2-HUMAN HES2-MOUSE HES2-RAT HES3-RAT HES5-MOUSE HESS-RAT HEY1-HUMAN HEY1-HUMAN HEY1-MOUSE HIFA-HUMAN HIFA-MOUSE HTF4-CHICK HTF4-HUMAN HTF4-MESAU HTF4-MOUSE HTF4-PAPHA HTF4-RAT HTF4-XENLA ID1-HUMAN ID1-MOUSE ID1-RAT ID2-HUMAN ID2-MOUSE ID2-RAT ID3-HUMAN ID3-MOUSE ID3-RAT ID4-HUMAN ID4-MOUSE INO2-YEAST INO4-YEAST ITF2-CHICK ITF2-CHICK ITF2-HUMAN ITF2-MOUSE ITF2-RAT LYL1-HUMAN LYL1-MOUSE MAD4-HUMAN MAD-HUMAN MAD-MOUSE MAX-BRARE MAC-CHICK MAX-HUMAN MAX-MOUSE MAX-RAT MAX-XENLA MP25-XENLA MP25-XENLA MITF-HUMAN MITF-RAT MLX-HUMAN MLX-MOUSE MNT-HUMAN MNT-MOUSE MUSC-HUMAN MUSC-MOUSE MX11-HUMAN MX11-MOUSE MX11-RAT MYC1-CYPACA MYC1-XENLA MYC2-CYPACA MYC2-MARMO MYC2-SPEBE MYC2-XENLA MYCB-RAT MYCL-HUMAN MYCL-MOUSE MYCL-XENLA MYCM-HUMAN MYCM-XENLA MYCN-CHICK MYCN-HUMAN MYCN-MARMO MYCN-MOUSE MYCN-RAT MYCN-SERCA MYCN-XENLA MYCS-RAT MYC-ASTVU MYC-AVM12 MYC-AVM12 MYC-AVM12 MYC-AVM12 MYC-AVM12 MYC-AVM12 MYC-BRARE MYC-CALJA MYC-CANFA MYC-CARAU MYC-CHICK MYC-FELCA MYC-FLVT1 MYC-HUMAN MYC-MARMO MYC-MOUSE MYC-MOUSE MYC-ONCMY MYC-PANTR MYC-PIG MYC-RAT MYC-SHEEP MYF5-BOVIN MYF5-CHICK MYF5-COTIA MYF5-HUMAN MYF5-MOUSE MYF5-NOTV1 MYF5-XENLA MYF6-CHICK MYF6-MOUSE MYF6-RAT MYF6-XENLA MYO1-ONCMY MYO2-ONCMY MYOD-CAEBR MYOD-CAEEL MYOD-CHICK MYOD-COTIA MYOD-DROME MYOD-HUMAN MYOD-MOUSE MYOD-PIG MYOD-RAT MYOD-SHEEP MYOD-XENLA MYOG-CHICK MYOG-COTIA MYOG-HUMAN MYOG-MOUSE MYOG-PIG MYOG-RAT NDF1-CHICK NDF1-HUMAN NDF1-MOUSE NDF1-RAT NDF1-XENLA NDF2-HUMAN NDF2-MOUSE NDF2-RAT NGN1-HUMAN NGN1-MOUSE NGN1-RAT NDFM-CHICK NPA1-HUMAN NPA1-MOUSE NPA2-HUMAN NPA2-MOUSE NCUI-NEUCR PAS1-HUMAN PAS1-MOUSE PHO4-YEAST OLG2-HUMAN OTG1-YEAST RTG3-YEAST SCL-CHICK SIM1-HUMAN SIM2-HUMAN SIM2-MOUSE SIMA-DROME SIM-DROME SRE1-CRIGR SRE1-HUMAN SRE1-MOUSE SRE1-PIG SRE1-RAT SRE2-CRIGR SRE2-HUMAN SUM1-LYTIVA SYFA-DROME TAL2-HUMAN TAL-MOUSE TAP4-HUMAN TAP4-DROME TF21-HUMAN TF2-HUMAN TF22-MESAU TF2-MOUSE TF2-RAT TF2-XENLA TF2-HUMAN TF3-MOUSE TFEB-HUMAN TFEB-MOUSE TRH-DROME TWST-DROME TWS1-MOUSE TWST-XENLA TYF7-YEAST USF1-HUMAN USF1-MOUSE USF1-RABIT USF1-XENBO USF2-HUMAN USF2-MOUSE USF-STRPV WS14-HUMAN WS14-MOUSE YAWC-SCHPO YLB7-CAEEL YMH7-CAEEL YNP2-CAEEL YRY3-CAEEL TWST-CAEEL
PS00061	2BH2-STREX ADH2-DROMN ADH-DROMA ADH-DROMN ADH-DROSU DECR-RAT DBH7-RAT DHGA-BACME DHG-BACSU DH2-RABIT DH2-SHEEP DHK1-STRVN ENTA-ECOLI MAS1-AGR19 PCDH-HUMAN Y019-THEMA YAE8-SCHPO YF43-MCTYU YW4-CAEEL OXIR-STRAT 25KD-SRPHB CHD2-EMTE ACT3-STR3 ADH1-CERCA ADH1-DROHY ADH1-DROMN ADH1-DROME ADH1-DROMT ADH1-DROMU ADH1-DRONA ADH1-CERCA ADH2-DROAR ADH2-DROBU ADH2-DROHY ADH2-DROMO ADH2-DROMU ADH-DROMY ADH2-DROWH ADHR-DROMO ADHR-DROGR ADHR-DROGU ADHR-DROIM ADHR-DROLE ADHR-DROMAD ADHR-DROMD ADHR-DROME ADHR-DROPE ADHR-DROPS ADHR-DROSD ADHR-DROTE ADHR-DROYA ADH-BACOL ADH-DROAD ADH-DROAF ADH-DROMA ADH-DROBO ADH-DROCK ADH-DRODI ADH-DROER ADH-DROFL ADH-DROGR ADH-DROGU ADH-DROHA ADH-DROHE ADH-DROIM ADH-DROLA ADH-DROLE ADH-DROMD ADH-DROME ADH-DRONI ADH-DROOR ADH-DROPE ADH-DROPI ADH-DROPL ADH-DROPS ADH-DROSI ADH-DROSD ADH-DROTE ADH-DROTS ADH-DROWI ADH-DROYA ADH-SCAAL ADH-ZAPTU ARDH-CANAL ARDH-CANTR ARDH-PCIST BA71-EUBSP BA72-EUBSP BDHA-ALCEU BDHA-RHIMO BDH-BOVIN BDH-HUMAN BDH-RAT BEND-ACISA BL14-NEUTR DTD2-PSEPU BPBH-PDESE BPBH-COMTE BPBH-PSEP BPBH-PSEPI BPBH-PSESI BPBH-RHIGO BUDC-KLEPN BUDC-KLETE CBR2-CAEEL CBR3-MOUSE CBR2-PIG CMTB-BSPY DECUR-HUMAN DBH1-HUMAN DBH1-MOUSE DBH1-RAT DBH2-HUMAN DBH2-MOUSE DBH2-RAT DBH3-HUMAN DBH3-MOUSE DBH3-SCHPO DBH4-HUMAN DBH4-MOUSE DBH4-RAT DBH7-HUMAN DBH7-MOUSE DBH8-CAELA DBH8-HUMAN DBH8-MOUSE DBHA-BACSU DBHK-MOUSE DBHV-CAEEL DBHW-CAEEL DBHX-ANPL DBHX-CAEEL DBHY-CAEEL DHC3-HUMAN DHCA-HUMAN DHCA-MOUSE DHCA-RABIT DHCA-RAT DHGI-BACME DHG2-BACME DHG2-BACSU DHG3-BACME DHG4-BACME DHGB-BACME DHGC-BACME DHII-HUMAN DHII-MOUSE DHII-RAT DHII-SAISC DHII-SHEEP DHIZ-BOVIN DHIZ-HUMAN DHIZ-MOUSE DHIZ-RAT DHKI-STRVN DHKR-STRCM DHMA-FLA1 DPHR-HUMAN DPHR-RAT DHSO-RHOSS DIHD-COMTE DIDH-PSEP DLTE-BACSU EPID-MYCTU FABG-ACTAC FABG-AQUAE FABG-ARATH FABG-BACSU FABG-BRANA FABG-BUCAI FABG-CHLMU FABG-CHLPM FABG-CHLTR FABG-CUPLA FABG-ECOLI FABG-HAEIN FABG-MYCTU FABG-MYCSM FABG-MYCTU FABG-PERAIE FABG-PSEAE FABG-RICPR FABG-SALTU FABG-VIBCH FABG-VIBHA FAGI-SYNY3 FAG2-SYNY3 FBP2-DROME FIXR-BRAJIA FOX2-CANTRX FOX2-NEUCR FOX2-YEAST FVT1-HUMAN GNO-GLUOX GS39-BACSU

(continued)

TABLE 4. (Continued)

Family	Protein ID's
PS00198	<p>DHSB-CYACA DHSB-PARDE DHSB-RICCN FER3-PLEBO FER-ALAC FER-CLOST FIXG-RHIME FIXX-BRAJA HMC6-DESVH MAUM-METEX NIFJ-ECOLI NUIC-ARATH NUIM-NEUCR PORD-METIA PSAC-ORYSA RNFB-PASMU RNFC-ECOS7 Y208-METIA YD49-METIA YFHL-ECOLI AEGA-ECOLI ASRA-SALTY ASRC-SALTY COOF-RHOU DCA1-METMA DCA2-METMA DCA-METIA DCA-METISO DCA-METTH DCMG-METTE DHSB-SCHPO DHSB-BACSU DHSB-CAEEL DHSB-CHOCR DHSB-COXB 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-HAEN DPYD-CAEEL DPYD-HUMAN DPYD-PIG DSRB-ARCFU DSVB-DESGI DSVB-DESVH FDHB-METFO FDHB-METIA FDHB-METTF FDHB-WOLSU FDNH-ECOLI FDOH-ECOLI FDXH-HAEN FDXN-ANASP FDXN-ANAVA FDXN-AZOC FDXN-BRAJA FDXN-RHILT FDXN-RHIME FDXN-RHISN FDXN-RHOCA FER1-AZOV1 FER1-CAUCR FER1-CHLL1 FER1-DESAF FER1-DESDN FER1-DESV1 FER1-METIA FER1-RHOA FER1-RHOU FER1-SULTO FER2-CHLL1 FER2-DESDN FER2-DESV1 FER2-METIA FER2-RHOCA FER2-RHOU FER2-SULTO FER2-THEAC FER3-ANASP FER3-ANAVA FER3-DESAF FER3-METIA FER3-RHISN FER3-RHOCA FER4-METIA FER5-METIA FER6-METIA FER7-METIA FER8-METIA FER9-AZOV1 FER9-AZOV1 FER-ACIAM FER-BACSC FER-BACST FER-BACSU FER-BACTH FER-BUTME FER-CHILL1 FER-CHIRV1 FER-CLOAC FER-CLOBU FER-CLOPA FER-CLOPE FER-CLOSP FER-CLOTH FER-CLOTS FER-DESGI FER-ENTHI FER-MEGL FER-METRA FER-METTE FER-METTL FER-MOOTH FER-MYCSM FER-MYCTU FER-MYCTU FER-PEPAS FER-PSEPK FER-PSEST FER-PYRAB FER-PYRIS FER-RICPR FER-SACER FER-STRGR FER-SULAC FER-THEAC FER-THELI FER-THEMA FER-THEH FIXX-AZOCA FIXX-AZOV1 FIXX-ECOLI FIXX-RHILE FIXX-RHILP FIXX-RHILT FIXX-RHIME FIXX-RHISN FPRB-MYCLE FPRB-MYCTU FRD1-AQUAE FRD2-AQUAE FRDB-ECOLI FRDB-HAEN FRDB-HELPI FRDB-HELPU FRDB-MYCTU FRDB-PROVU FRDB-WOLSU FRHG-METIA FRHG-METTH FRHG-METVO GLCF-ECOLI GLPC-ECOLI GLPC-HAEN HMC2-DESVH HYBA-ECOLI HYCB-ECOLI HYCF-ECOLI HYDN-ECOLI HYFA-ECOLI HYFH-ECOLI IORA-ARCFU IORA-METTH IORA-PYRAB IORA-PYRHO IORA-PYRKO MAUM-METFL MAUM-METME MAUM-PARDE MAUM-METEX MAUM-METFL MAUN-PARDE NAFI-ECOLI NAFI-HAEN NAFI-ECOLI NAFI-HAEN NAPH-ECOLI NAPH-HAEN NIFJ-ANASP NIFJ-ENTAG NIFJ-KLEPN NIFJ-RHOU NIFJ-SYNY3 NQ09-PARDE NQ09-THEH NRFC-ECOLI NRFC-HAEN NUGC-RHIME NUIC-MAIZE NUIC-MARPO NUIC-MESVI NUIC-ORYSA NUIC-PLEBO NUIC-SPHOL NUIC-SYNY3 NUIC-TOBAC 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-RICCN NUOI-RICPR PHFI-CLOPA PHFI-DESVH PHFI-DESVO PHFI-SALTY PORD-METTH PORD-PYRAB PORD-PYRFU PORD-PYRHO PORD-THEMA PSAC-ANASP PSAC-ANTSP PSAC-ARATH PSAC-CHLRE PSAC-CHLVU PSAC-CYACA PSAC-CYAPA PSAC-EUGGR PSAC-FREDI PSAC-GUTH 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 RNFB-BUCAI RNFB-ECOS7 RNFB-ECOLI RNFB-HAEN RNFB-PSEAE RNFB-RHOCA RNFB-VIBCH RNFC-BUCAI RNFC-ECOLI RNFC-HAEN RNFC-PASMU RNFC-PSEAE RNFC-RHOCA RNFC-VIBCH VORC-METTH VORD-PYRAB VORD-PYRFU VORD-PYRHO Y092-METIA Y264-METIA Y492-MYCTU Y578-METIA Y726-METIA Y870-METIA YA43-HAEN YCCM-ECOLI YCXI-PORPU YDIJ-ECOLI YDIJ-HAEN YDIT-ECOLI YEIA-ECOLI YFHL-HAEN YFRA-PROVU YG84-METTH YGFS-ECOLI YGFT-ECOLI YGL5-BACST YJES-ECOLI YJWJ-ECOLI YKGF-ECOLI YNFG-ECOLI YSAA-ECOLI</p>
PS00211	<p>ABC2-HUMAN APPD-BACSU FTSE-HAEN HISP-SALTY KST1-ECOLI LCCL-LACLA LMRA-LACLC LOLD-BUCAI MDLB-BUCAI MKL-MYCTU MODC-HAEN MRP2-RABIT NIKD-ECOLI NODI-AZOCA NODI-RHISN NOSF-PSEST NRTD-SYNY3 OPPF-LACLA OPPF-MYCPN POTA-MYCGE RFBB-MYXXA SUFC-ECOLI UVRA-BRUAB UVRA-STRMU VEXC-SALTI WHIT-ANOAL Y348-CHLPN YF08-METIA YJJK-HAEN YXDL-BACSU AAPF-RHILV AB11-HUMAN AB11-MOUSE AB11-RABIT AB11-RAT ABC1-HUMAN ABC1-MOUSE ABC1-SCHPO ABC2-MOUSE ABC3-HUMAN ABC6-HUMAN ABC7-HUMAN ABC7-MOUSE ABC8-HUMAN ABCA-AERSA ABCR-HUMAN ABCX-ANTSP ABCX-CYACA ABCX-CYAPA ABCX-GALSU ABCX-GUTH ABCX-ODOSI ABCX-PORPU ABCX-STRMU METN-ECOLI METN-HAEN 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-CRIGR ACC8-HUMAN ACC8-RAT ACC9-HUMAN ACC9-MOUSE ACC9-RABIT ACC9-RAT ADCC-STRPN ADPI-YEAST FBPC-ACITPL FBPC-ECOLI FBCL-HAEN AGLK-RHIME ALD-HUMAN ALD-MOUSE ALSA-ECOLI AMIE-STRPN AMIF-STRPN AOTF-PSEAE APPF-BACSU APRD-PSEAE ARAG-ECOLI ARTP-ECOLI ARTP-HAEN ATMI-YEAST BCRA-BACLI BEXA-HAEN BFER1-SCHPO BP71-YEAST BRAP-PSEAE BRAG-PSEAE BROW-DROME BROW-DROVI BTUD-ECOLI BZTD-RHOCA CBIO-SALTY CBRD-ERWCH CCMA-BRAJA CCMA-ECOLI CCMA-HAEN CCMA-PARDE CCMA-RHOCA CDR1-CANAL CDR2-CANAL CDR3-CANAL CDR4-CANAL CFTR-BOVIN CFTR-CAVPO CFTR-HUMAN CFTR-MACMU CFTR-MOUSE CFTR-RABIT CFTR-RAT CFTR-SHEEP CFTR-SQUAC CFTR-XENLA CHVA-AGRT5 CHVD-AGRTU COMA-STRPN CTBD-NEIMA CTBD-NEIMB CVAB-ECOLI CYAB-BORPE CYDC-BACSU CYDC-ECOLI CYDC-HAEN CYDD-BACSU CYDD-ECOLI CYDD-HAEN CYSA-CHLVU CYSA-ECOLI CYSA-MARPO CYSA-MESVI CYSA-SALTY CYSA-SYNY7 CYSA-SYNY3 DPPP-BACSU DPPP-ECOLI DPPP-HAEN DPPP-ECOLI DPPP-HAEN DRRA-STRPE ECSA-BACSU EF3A-YEAST EF3B-YEAST EF3-CANAL 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 HEPF-ANASP HEPF-CAUCR HISP-ECOLI FBC2-HAEN HLY2-ECOLI HLYB-ACTAC HLYB-ECOLI HLYB-PASHA HLYB-PASSP HLYB-PROVU HMTL-SCHPO HNUV-YERPE HST6-CANAL KST5-ECOLI LACK-AGRRD LCN3-LACLA LCNC-LACLA LIVF-ARCFU LIVF-ECOLI LIVF-METIA LIVF-SALTY LIVG-ARCFU LIVG-ECOLI LIVG-METIA LIVG-SALTY LOLD-BUCAP LOLD-ECOLI LOLD-HAEN LOLD-NEIMA LOLD-NEIMB LOLD-VIBCH LOLD-XYLEA MACB-ECOLI MALK-ECOLI MALK-ENTAE MALK-PHOLU MALK-SALTY MAM1-SCHPO MCHF-ECOLI MDL1-CANAL MDL1-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 MDR3-MOUSE MDR4-DROME MDR4-ENTHI MDR5-DROME MDR-LEITA MDR-PLAFF MESD-LEUME MGLA-ECOLI MGLA-HAEN MGLA-MYCGE MGLA-MYCPN MGLA-SALTY MGLA-TREPA MKL-MYCLE MNTA-SYNY3 MODC-AZOV1 MODC-ECOLI MODC-MYCTU MODC-RHOCA MODF-ECOLI MRP1-HUMAN MRP2-HUMAN MRP2-RAT MRP3-HUMAN MRP3-RAT MRP4-HUMAN MRP5-HUMAN MRP5-MOUSE MRP5-RAT MRP6-HUMAN MRP6-RAT MSBA-ECOLI MSBA-HAEN MSMK-STRMU MSMX-BACSU MSRA-STAEF NASD-KLEPN NATA-BACSU NDVA-RHIME NIKE-ECOLI NIST-LACLA NOCP-AGRT5 NODI-BRAIA NODI-RHIGA NODI-RHILLO NODI-RHILT NODI-RHILV NDI1-RHIME NODI-RHIS3 NRTC-SYNY7 NRTC-SYNY3 NRTD-SYNY7 OCCP-AGRTU OCCP-RHIME OPAA-BACSU OPBA-BACSU OPBA-BACSU OPDD-BACSU OPDD-ECOLI OPDD-HAEN OPDD-LACLA OPDD-LACLC OPDD-MYCGE OPDD-MYCPN OPDD-SALTY OPFF-BACSU OPFF-ECOLI OPFF-HAEN OPFF-MYCGE OPFF-MYCTU OPFF-STRMU OPFF-STRPY P29-MYCGE P29-MYCHR P29-MYCPN PDR5-YEAST PDRA-YEAST PDRB-YEAST PDRX-YEAST PDRF-YEAST PEBB-CAMIE PEDD-PEDAC PHNC-ECOLI PHNK-ECOLI PHNL-ECOLI PMD1-SCHPO POTA-ECOLI POTA-HAEN POTA-MYCPN POTA-SALTY POTG-ECOLI PROV-ECOLI PROV-SALTY PRTD-ERWCH PSTB-ECOLI PSTB-EDWTA PSTB-ENTCL PSTB-METIA PSTB-MYCGE PSTB-MYCTU PSTB-MYCTU PSTB-PASMU PSTB-RHILLO PSTB-SALTY PSTB-SALTY PSTB-XYLEFA PXA1-YEAST PXA2-YEAST RBBA-BACSU RBBA-ECOLI RBBA-HAEN RFB1-KLEPN RFB2-KLEPN RFB6-YEREN RT1B-ACITPL RT3B-ACITPL SAPD-ECOLI SAPD-HAEN SAPD-SALTY SAPF-ECOLI SAPF-HAEN SAPF-SALTY SCRT-DROME FBPC-SERMA SMOK-RHOSH SNQ2-YEAST SPAT-BACSU SRIT-STRPY SSUB-BACSU SSUB-ECOLI ST6E-YEAST SYRD-PSESY TAGB-DICDI TAGC-DICDI TAGH-BACSU TAPI-HUMAN TAPI-MOUSE TAPI-RAT TAP2-HUMAN TAP2-MOUSE TAP2-RAT TAUB-ECOLI THIQ-ECOLI THIQ-HAEN TLRC-STRFR TROB-TREPA UGPC-ECOLI UUPI-HAEN UUP2-HAEN UUP-BUCAI UUP-ECOLI UVRA-AQUAE UVRA-BACHD UVRA-BACSU UVRA-BORBU UVRA-CHLMU UVRA-CHLPN UVRA-CHLTR UVRA-DEIRA UVRA-ECOLI UVRA-HAEN UVRA-HELPI UVRA-HELPU 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-THEH UVRA-TREPA UVRA-VITST UVRA-ZYMMO Y296-BACSU WHIT-ANOVA WHIT-CERCA WHIT-DROME WHIT-LUCCU XYLG-ECOLI XYLH-HAEN Y014-MYCGE Y014-MYCPN Y015-MYCGE Y015-MYCPN Y035-METIA Y035-TREPA Y036-HAEN Y065-MYCGE Y065-MYCPN Y068-CHLTR Y075-SYNY3 Y089-METIA Y121-METIA 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-HAEN Y361-HAEN Y382-RHIME Y412-METIA Y415-SYNY3 Y416-CHLTR Y467-MYCGE Y467-MYCPN Y468-MYCGE Y468-MYCPN Y4FO-RHISN Y4GM-RHISN Y4MK-RHISN Y4QS-RHISN Y4TH-RHISN Y4TH-RHISN Y4TS-RHISN Y542-CHLPN Y663-HAEN Y664-HAEN Y697-CHLMU Y700-RICPR Y719-METIA Y796-METIA Y799-ANASP Y873-METIA Y888-HELPI Y888-HELPU Y986-MYCTU Y223-METIA YA51-HAEN YA78-HAEN YADG-ECOLI YATR-BACFI YAWB-SCHPO YBBA-ECOLI YBBL-ECOLI YBHF-ECOLI YBHT-ECOLI YBT1-YEAST YBXA-BACSU YC72-HAEN YC72-MYCTU YC73-MYCTU YC81-MYCTU YCBN-BACSU YCFI-YEAST YCIV-ECOLI YCKI-BACSU YCXD-CYAPA YD34-MYCPN YD48-MYCTU YD49-MYCTU YD67-METIA YDCT-ECOLI YDDA-ECOLI YDDO-ECOLI YDDP-ECOLI YDIF-BACSU YE67-HAEN YE70-HAEN YE74-HAEN YECC-ECOLI YEHX-ECOLI YEJF-ECOLI YEM6-YEAST YD01-SCHPO YFC8-YEAST YFEB-YERPE YFIB-BACSU YFIC-BACSU YNT9-SCHPO YG18-HAEN YHGB-AZOCA YHGB-ECOLI YHGB-HAEN YHGB-KLEPN YHGB-PSEPU YHGB-THIFE YHCG-BACSU YHCH-BACSU YHDS-YEAST YHYZ-ECOLI YHES-ECOLI YHES-HAEN YHIE-ECOLI YH19-METIA YJJK-ECOLI YK83-YEAST CED7-CAEEL YLIA-ECOLI YMEB-LACLA YN26-MYCTU YN99-YEAST YNID-ECOLI YOH5-YEAST YOIJ-ECOLI YORI-YEAST YP64-MYCTU YPC3-CAEEL YPHE-ECOLI YQSC-CAEEL YQGI-BACSU YQGK-BACSU YQIZ-BACSU YRBF-HAEN YSCI-STRGC YTRF-ECOLI MNTB-BACSU YTMN-BACSU YTRE-BACSU YWJA-BACSU YXEO-BACSU YYBJ-BACSU ZNUC-BUCAI ZNUC-ECOLI ZNUC-HAEN ZURA-LISIN ZURA-LISMO</p>

(continued)

TABLE 4. (Continued)

Family	Protein ID's
PS00301	<p> 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-EUPCR EF12-HORVU EF12-HUMAN EF12-MOUSE EF12-RHIRA EF12-SCHPO EF12-XENLA EF13-RHIRA EF12-SCHPO EF13-XENLA EF12-DAUCA EF12-ABSOL EF1A-AERPE EF1A-AJECA EF1A-APIME EF1A-ARATH EF1A-ARTSA EF1A-ARXAD EF1A-ASHGO EF1A-AURPU EF1A-BLAHO EF1A-BOMMO EF1A-BRARE EF1A-CAEEL EF1A-CANAL EF1A-CHICK EF1A-CRYNE EF1A-CRYPV EF1A-DESMO EF1A-EIMBO EF1A-ENTHI EF1A-EUCGR EF1A-GIALA EF1A-HALHA EF1A-HALMA EF1A-PLAFK EF1A-HYDAT EF1A-LYCES EF1A-MAIZE EF1A-MANES EF1A-METIA EF1A-METTH EF1A-METVA EF1A-NEUCR EF1A-ONCVO EF1A-ORYSA EF1A-PEA EF1A-PLAFK EF1A-PODAN EF1A-PODCU EF1A-PUCGR EF1A-PYRAB EF1A-PYRAE EF1A-PYRHO EF1A-PYRWO EF1A-RHYAM EF1A-SCHCO EF1A-SORMA EF1A-SOYBN EF1A-STYLE EF1A-SULAC EF1A-TETPY EF1A-THEAC EF1A-THECE EF1A-TOBAC EF1A-TRIIE EF1A-TRYBB EF1A-VICFA EF1A-WHEAT EF1A-YARLI EF1A-YEAST EF1C-PORPU EF2-AERPE EF2-ARCFU EF2-BETVU 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-PYRFU EF2-RABIT EF2-RAT EF2-SCHPO EF2-SULAC EF2-SULSO EF2-THEAC EF2-YEAST EFG1-BORBU EFG1-STRCO EFG1-SYNY3 EFG1-TREPA EFG1-YEAST EFG2-BORBU EFG2-STRCO EFG2-SYNY3 EFG2-TREPA EFG2-YEAST EFGC-PEA EFGC-SOYBN EFGM-MYCTU EFGM-SYNY3 EFGM-TREPA EFGM-YEAST EFGM-PAPP EFGM-AQUAE EFGM-AQUY EFGM-BACHD EFGM-BACST EFGM-BACSU EFGM-BUCAL EFGM-CHLMU EFGM-CHLPN EFGM-CHLTR EFGM-ECOLI EFGM-HAEN EFGM-HELPI EFGM-HELPU EFGM-MICLU EFGM-MYCGE EFGM-MYCLE EFGM-MYCPN EFGM-MYCTU EFGM-NEIGO EFGM-PASMU EFGM-PLARO EFGM-RICCN EFGM-RICPR EFGM-SALTY EFGM-SPIPL EFGM-STAMM EFGM-STRPY EFGM-STRRRA EFGM-SYNP6 EFGM-THEMA EFGM-THETH EFGM-THICU EFGM-UREPA EFT1-SOYBN EFT1-STRCO EFT1-STRCU EFT1-STRRRA EFT2-SOYBN EFT2-STRRRA EFT3-STRCO EFT3-STRRRA EFT3-PASMU EFT3-PASMU EFTU-ARCFU 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-CHLPE EFTU-CHLRE EFTU-CHLVI EFTU-CHLVU EFTU-CODFR EFTU-COLOB EFTU-CORGL EFTU-COSCS EFTU-CYAPA EFTU-CYCMC EFTU-CYTLT EFTU-DEIRA EFTU-DEISP EFTU-DEIRA EFTU-DEIRA EFTU-ECOLI EFTU-EIKKO EFTU-EUGGR EFTU-FIBSU EFTU-FLAIE EFTU-FLESI EFTU-GLOS1 EFTU-GLOVI EFTU-GONPE EFTU-GUITH EFTU-GYMET EFTU-HAEN EFTU-HELPI EFTU-HELPU EFTU-HERAU EFTU-HUMAN EFTU-MANSO 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-RHIO EFTU-RICPR EFTU-SALTY EFTU-SCHPO 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 EFT2-CANAL ERF2-PCIP1 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-CHLPI LEPA-CHLTR LEPA-ECOLI LEPA-HAEN LEPA-HELPI LEPA-HELPU LEPA-LACLA LEPA-MYCGE LEPA-MYCTU LEPA-PASMU LEPA-PSEFL LEPA-RICPR LEPA-SALTY LEPA-STRCO LEPA-SYNY3 LEPA-THEMA LEPA-TREPA NODQ-AZOB R NODQ-RHIME NODQ-RHIS3 NODQ-RHISB NODQ-RHIRT OTRA-STRRM RF3-BACNO RF3-BUCAI RF3-ECOLI RF3-HAEN RF3-LACLA RF3-PASMU RF3-SALTY RF3-STAU RF3-SYNY3 SELB-DESHA SELB-ECOLI SELB-HAEN SELB-HUMAN SELB-METIA SELB-MOOTH SELB-MOUSE SN14-YEAST TET1-ENTFA TET5-ENTFA TET9-ENTFA TETM-NEIME TETM-STAAU TETM-STRLI TETM-STRPN TETM-UREUR TETO-CAMCO TETO-CAMJE TETO-STRMU TETO-STRPN TETP-CLOPE TETQ-BACFR TETQ-BACTN TETQ-BACTN TETQ-PRERU TETS-LACLA TETS-LISMO TETW-BUTHI TYPA-BACSU TYPA-BUCAI TYPA-ECOLI TYPA-HAEN TYPA-HELPI TYPA-HELPU USS1-HUMAN USS1-MOUSE YE14-SCHPO YNQ3-YEAST Y081-CAEEL </p>

TABLE 5. DESCRIPTION OF THE PROSITE 2 DATASET

Family	Protein ID's
PS00070	<p> DHAE-MACPR DHAX-HUMAN DHA1-BOVIN HPCC-ECOLI YH9-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-OCCTO 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 DHA4-BOVIN DHA5-HUMAN DHA6-HUMAN DHA6-YEAST DHA7-HUMAN DHA8-HUMAN DHA9-POLMI DHA6-AMAHF DHA6-ATRHO DHA6-BACSU DHA6-BETVU DHA6-ECOLI DHA6-GADCA DHA6-HORVU DHA6-ORYSA DHA6-RHIME DHA6-SPIOL DHA6-RAT DHA6-ELEP DHA6-VIBHA DHA6-HUMAN DHA6-PIG DHA6-AGABI DHA6-ALITAL DHA6-PNG DHA6-BACST DHA6-CLAHE DHA6-DEIRA DHA6-ECOLI DHA6-EMENI DHA6-ENCPU DHA6-MYCTU DHA6-PSEOL DHA6-PSESP DHA6-RICOR DHA6-STRCO DHA6-VIBCH DHAM-HORSE DHAM-HUMAN DHAM-LEITA DHAM-MESAUI 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 SSDH-HUMAN SSDH-RAT THCA-RHOER UGA5-YEAST XYC2-ACIGB XYLC-PSEPU XYLG-PSEPU Y4UC-RHISN YDCW-ECOLI YM00-YEAST YNEI-ECOLI </p>
PS00077	<p> COX1-THETH COX1-BACFI COX1-DIDMA COX1-ASCSU COX1-HORSE COX1-EPHEQ FIXN-AZOCA COX1-SYNYV COX1-CRION COX1-ALLMA AOX1-AERPE COX1-PEA COX1-RHOSH COX1-SOYBN COX1-PLABE CO13-THETH CO14-BRAJA COX1-ACACA COX1-ALBCO COX1-ALBTU COX1-AMICA COX1-ANAPL COX1-ANOGA COX1-ANOUQ COX1-APILI COX1-APTAU 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-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-LESPS COX1-LOCM 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-PIG COX1-PIG COX1-PISOC COX1-PLACH COX1-PLAFA COX1-PODAN COX1-POLOR COX1-POLSP COX1-POLSK COX1-POMNI 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-STRCO COX1-STRPY 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 COX1-ACEAC QOX1-BACSU QOX1-SULAC QOXM-SULAC </p>
PS00118	<p> PA21-NAJMO PA21-HORSE PA2H-BUNFA PA2E-PSEAU PA2C-CRODU PA2H-BOTJR PA2C-PSEAU PA2Z-HUMAN PA22-BUNMU PA23-NAJNG PA21-TRIGA PA21-ACAN PA21-BOTPI PA2X-RAT PA22-PIG OC90-CAVPO OC90-HUMAN OC90-MOUSE PA20-BUNMU PA20-NOTSC PA20-PSEAU PA21-AGKHA PA21-AGKHP PA21-AGKPI PA21-BOTAS PA21-BOTJA PA21-BOTJR 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-ONYSC PA21-PIG PA21-PSEAU PA21-RAT PA21-SHEEP PA21-TRIFL PA21-VIPAA PA21-VIPAZ PA22-ACAN PA22-AGKHA PA22-AGKHP PA22-ASPSC PA22-BITNA PA22-BOTAS PA22-BOTMO PA22-BOTPI PA22-CERGO PA22-ERIMA PA22-HELSSU PA22-LATCO PA22-MATBI PA22-NAJKA PA22-NAJME PA22-NAJMO PA22-NOTSC PA22-ONYSC PA22-TRIGA PA22-TRIST PA22-VIPAZ PA23-AGKHP PA23-BOTAS PA23-BOTPI PA23-BUNMU PA23-HELSSU PA23-HUMAN PA23-LATSE PA23-NAJKA PA23-NAJME PA23-NAJMO PA23-NOTSC PA23-ONYSC 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-VIPPA 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-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 PA2L-VIPAA PA2M-AGKCL PA2M-CAVPO PA2M-CROSS PA2N-BUNFA PA2N-CROSS PA2N-ECHCA PA2N-VIPAA PA2X-HUMAN PA2X-HUMAN PA2X-MOUSE PA2X-NOTSC PA2X-TRIFL PA2Y-HUMAN PA2Y-MOUSE PA2Y-TRIFL PA2Z-MOUSE PA2Z-MOUSE PA2-AIPLA PA2-APIME PA2-BITCA PA2-BITGA PA2-BOMTE PA2-CERCE PA2-CROAD PA2-CROAT PA2-DABRU PA2-ENHSC PA2-HELHO PA2-LATLA PA2-NAJAT PA2-NAJPA PA2-OPHHA PA2-RHONO PA2-TRIOK PA2-VIPBB </p>
PS00180	<p> 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-PHAYU GLNA-NICPL GLN2-DAUCA GLN1-ALNGL GLN1-BRAJA GLN1-CHLRE GLN1-DAUCA GLN1-DROME GLN1-FRAAL GLN1-LOTIA GLN1-MAIZE GLN1-MEDSA GLN1-MYCTU GLN1-ORYSA GLN1-PEA GLN1-PHAYU GLN1-RHILV GLN1-RHIME GLN1-SOYBN GLN1-STRRP GLN1-STRRV GLN1-VITVI GLN2-ARATH GLN2-BRAJA GLN2-CHLRE GLN2-FRAAL GLN2-HORVU GLN2-MAIZE GLN2-MEDSA GLN2-MYCTU GLN2-ORYSA GLN2-PEA GLN2-PHAYU GLN2-RHILP GLN2-RHIME GLN2-SOYBN GLN2-STRRY GLN2-STRRV GLN2-VITVI GLN3-HORVU GLN3-LUPAN GLN3-MAIZE GLN3-MEDSA GLN3-ORYSA GLN3-PEA GLN3-RHIME GLN3-MAIZE GLN4-PHAYU GLNA-AGABI GLNA-ANASP GLNA-AQUAE GLNA-ARCFU GLNA-AZOB R GLNA-AZOCA GLNA-AZOV1 GLNA-BACCE GLNA-BACFR GLNA-BACSU GLNA-BOVIN GLNA-BUTPI GLNA-CAEEL GLNA-CHICK GLNA-CLOSA GLNA-CRIO GLNA-DUNSA GLNA-ECOLI GLNA-FREDI GLNA-HAEN GLNA-HALNI 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-PROVI GLNA-PYRAB GLNA-PYRFU GLNA-PYRKO GLNA-PYRMO GLNA-RHACA GLNA-RHOCA GLNA-RHOSH GLNA-SCHPO GLNA-SQUAC GLNA-STAAU GLNA-STRCO 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 YC1K-ECOLI </p>

(continued)

TABLE 5. (Continued)

<i>Family</i>	<i>Protein ID's</i>
PS00215	UCP5-HUMAN ARI3-NEUCR SAI8-MOUSE YIA6-YEAST SHM1-YEAST ADT1-BOVIN ADT2-WHEAT UCP3-BOVIN M20M-RAT YAD8-SCHPO UCP1-MOUSE TXTPT-HUMAN DNC-HUMAN ADT3-YEAST ADT3-HUMAN ADT1-ARATH ADT7-GOSHI ADT1-HUMAN ADT1-MAIZE ADT1-MOUSE ADT1-RAT ADT1-SOLTU ADT1-WHEAT ADT1-YEAST ADT2-ARATH ADT2-HUMAN ADT2-MAIZE ADT2-MOUSE ADT2-RAT ADT2-SOLTU ADT2-YEAST ADT2-BOVIN ADT2-ANOGA ADT2-CHLKE ADT2-CHLRE ADT2-DROME ADT-LKULLA ADT-NEUCR ADT-ORYSA ADT-SCHPO BT1-MAIZE CG69-HUMAN CMCI1-CAEEL CMCI1-DROME CMCI1-HUMAN CMCI1-YEAST CMCI2-CAEEL CMCI2-HUMAN CMCI2-MOUSE CMCI3-CAEEL DIC-HUMAN DIC-MOUSE ECHP-MOUSE FLX1-YEAST GDC-BOVIN GDC-HUMAN GDC-RAT LEU5-YEAST M20M-BOVIN M20M-HUMAN M20M-MOUSE MCAT-HUMAN MCAT-RAT MFT-HUMAN MPCC-BOVIN MPCC-CAEEL MPCC-CHOFU MPCC-HUMAN MPCC-RAT MPCC-YEAST MR53-YEAST MR54-YEAST ODC1-YEAST ODC2-YEAST ODC-HUMAN ORT1-HUMAN ORT1-YEAST ORT2-HUMAN P47A-CANBO P47B-CANBO PET8-YEAST PM34-HUMAN PM34-MOUSE PMT-YEAST MR12-YEAST SAI8-HUMAN SFCl-YEAST TXTPT-BOVIN TXTPT-CAEEL TXTPT-RAT TXTPT-YEAST UCP1-BOVIN UCP1-HUMAN UCP1-MESAU UCP1-RABIT UCP1-RAT UCP2-BRARE UCP2-CANFA UCP2-CYPCA UCP2-HUMAN UCP2-MOUSE UCP2-PIG UCP2-RAT UCP3-CANFA UCP3-HUMAN UCP3-MOUSE UCP3-PIG UCP3-RAT UCP4-HUMAN UCP4-YEAST YD1K-SCHPO YDE9-SCHPO YE08-SCHPO YEA6-YEAST YE03-YEAST YFL5-YEAST YG20-YEAST YGSF-YEAST YM39-YEAST YMC1-YEAST YMC2-YEAST YQ51-CAEEL
PS00217	GTR1-RAT IOLF-BACSU CSBC-BACSU GTR5-HUMAN KHT2-KULLA PH84-YEAST NANT-ECOLI GH3-SCHPO HUPI-CHLKE HGT1-CANAL GTR4-RAT GTRI-CHICK MMLH-ALCEU OUSA-ERWCH PHDK-NOSCK AGT1-YEAST ARAE-BACSU ARAE-ECOLI ARAE-KLEOX BENK-AICITA CTI1-ECOLI CTI1-KLEPN CTI1-SALTU GAL2-YEAST GAPL-ECOLI GH2-SCHPO GH2Y-SCHPO GH75-SCHPO GHT1-YEAST GLCP-SYNX GLFV-ZYMMO GT10-HUMAN GT11-HUMAN GTRI1-BOVIN GTRI1-HUMAN GTRI1-LEIDO GTRI1-MOUSE GTRI1-PIG GTRI1-RABIT GTRI1-SHEEP GTR2-BOVIN GTR2-CHICK GTR2-HUMAN GTR2-LEIDO GTR2-MOUSE GTR2-PIG GTR2-RAT GTR3-BOVIN GTR3-CANFA GTR3-CHICK GTR3-DROME GTR3-HUMAN GTR3-MOUSE GTR3-PIG GTR3-RABIT GTR3-RAT GTR3-SHEEP GTR4-BOVIN GTR4-CANFA GTR4-HUMAN GTR4-MOUSE GTR4-PIG GTR5-BOVIN GTR5-MOUSE GTR5-RABIT GTR5-RAT GTR6-HUMAN GTR8-BOVIN GTR8-HUMAN GTR8-MOUSE GTR8-RAT GTR9-HUMAN HEX6-RICCO HGTI-KULLA HUP2-CHLKE HUP3-CHLKE HXT0-YEAST HTX1-YEAST HXT2-YEAST HXT3-YEAST HXT4-YEAST HXT5-YEAST HXT6-YEAST HXT7-YEAST HXT8-YEAST HXT9-YEAST HXTA-YEAST HXTC-YEAST HXTD-YEAST HXTE-YEAST HXTF-YEAST HXTG-YEAST ITR1-SCHPO ITR1-YEAST ITR2-SCHPO ITR2-YEAST JENI-YEAST KGTP-ECOLI LACP-KULLA MA3T-YEAST MA6T-YEAST MAXT-YEAST MHPT-ECOLI MUCK-ACALIA MYCT-HUMAN PKAC-ACALIA PKAK-PSEPU PRIOT1-LEIEN PROP-ECOLI PROL-SALTU QYA-NEUCR QDT2-EMENI RAGI1-KULLA RC03-NEUCR RG27-YEAST SHIA-ECOLI SNF3-YEAST STA-RICCO STC-RICCO STL1-YEAST STP1-ARATH STP1-SPIOL TH11-TRYBB TH12-TRYBB TH23-TRYBB BX5-ECOLI XYLT-LACBR YH21-HAEIN Y418-HAEIN YAAU-ECOLI YAEC-SCHPO YH04-HAEIN YB91-YEAST YECI-BACSU YDFJ-ECOLI YDJIE-ECOLI YDIJK-ECOLI YFE0-YEAST YFIG-BACSU YGCS-ECOLI YGK4-YEAST YHJE-ECOLI YIR0-YEAST YJHB-ECOLI YOI1-CAEEL YYAJ-BACSU
PS00338	SOMA-TRIVU PRL-CHICK PRL-PAROL SOMA-MACMU PRL-MOUSE SOMA-ACALA PLL2-MESAU SOML-SIGGU SOMA-ESOLU SOM2-CARAU SOMA-CANFA PRL-SHEEP SOM2-HUMAN PRL-HORSE SOMA-PANTR GHRI1-RAT GHR3-RAT GHR4-RAT PLFI-MOUSE PLF2-MOUSE PLF3-MOUSE PLFR-MOUSE PLL1-BOVIN PLL1-MOUSE PLL1-RAT PLL2-BOVIN PLL2-MOUSE PLL2-RAT PLL-HUMAN PLL1-SHEEP PRL1-ALLMI PRL1-CRONO PRL1-ONCKE PRL1-OOREMO PRL2-ALLMI PRL2-CRONO PRL2-ONCKE PRL2-ONCTS PRL2-OOREMO PRL-ANGAN PRL-BALBO PRL-BOVIN PRL-BUFJA PRL-CAMDR PRL-CAPHI PRL-CARAU PRL-CHEMY PRL-CORAU PRL-CYPCA PRL-DICLA PRL-FELCA PRL-FELPA PRL-HYPNO PRL-HYPMO PRL-ICTPU PRL-LOXAF PRL-MACMU PRL-MELGA PRL-MESAU PRL-MONDO PRL-MUSVI PRL-ONCMY PRL-PIG PRL-PROAT PRL-RABIT PRL-RAT PRL-SALSA PRL-SPAUA PRL-TRIVU PRR1-BOVIN PRK2-BOVIN PRK3-BOVIN PRK4-BOVIN PRRA-RAT PRRB-RAT PRRC-RAT SOM1-ACIGU SOM1-CARAU SOMA-ONCKE SOM1-ONCNE SOM1-SPAUA SOM2-ACIGU SOM2-MACMU SOM2-ONCMY SOM2-ONCNE SOM2-PANTR SOM2-SPAUA SOMA-ACABU SOMA-ANAPL SOMA-ANGIA SOMA-BALBO SOMA-BOVIN SOMA-BUBBU SOMA-BUFMA SOMA-CALIA SOMA-CARDIE SOMA-CEREL SOMA-CHEMY SOMA-CHICK SOMA-CORAU SOMA-CORLV SOMA-CRONO SOMA-CTEID SOMA-CYPCA SOMA-DICLA SOMA-FELCA SOMA-FUGRU SOMA-GALSE SOMA-HETFO SOMA-HORSE SOMA-HUMAN SOMA-ICTPU SOMA-KATPE SOMA-LABRO SOMA-LAMPA SOMA-LATCA SOMA-LEPOS SOMA-LOXAF SOMA-MELGA SOMA-MESAU SOMA-MISMI SOMA-MONDO SOMA-MORSA SOMA-MOUSE SOMA-MUSVI SOMA-NYCYP SOMA-ODOAR SOMA-ONCKE SOMA-ONCKI SOMA-ONCMA SOMA-ONCTS SOMA-OOREMO SOMA-ORENT SOMA-PAGMA SOMA-PANPG SOMA-PAROL SOMA-PERFV SOMA-PIG SOMA-PRIGL SOMA-PROAN SOMA-PSERC SOMA-RABIT SOMA-RANCA SOMA-RAT SOMA-SAIBB SOMA-SALSA SOMA-SCIOC SOMA-SEBSO SOMA-SEROU SOMA-SHEEP SOMA-SIGGU SOMA-SOLSE SOMA-SPAUA SOMA-STRCA SOMA-THUAL SOMA-THUTH SOMA-TRITC SOMA-VERVA SOMA-VULVU SOMA-XENLA SOMA-XENLA SOMI-XENLA SOMI-ACITR SOMI-ANGAN SOMI-CARAU SOML-CYCLU SOML-GADMO SOML-HIPHI SOMI-LOCU SOML-ONCKE SOML-PAROL SOML-PROAN SOML-SCIOB SOML-SOLSE SOMI-SOML-TETMU

TABLE 6. DESCRIPTION OF THE GPCR DATASET

<i>Subfamily</i>	<i>Protein ID's</i>
Amine	<p> 5H1A-RAT 5H1B-CAVPO 5H1B-CRIGR 5H1B-HUMAN 5H1B-RABIT 5H1D-MOUSE 5H2A-CRIGR 5H2A-MOUSE 5HTB-DROME ACMI-DROME ACM3-PIG ACM4-MOUSE B2AR-MESAU DBDR-XENLA HH2R-MOUSE Q04198 O61232 OAR2-LOCMI 5H1A-FUGRU 5H1A-HUMAN 5H1A-MOUSE 5H1B-DIDMA 5H1B-FUGRU 5H1B-MOUSE 5H1B-RAT 5H1B-SPAEH 5H1C-CANFA 5H1C-CAVPO 5H1D-FUGRU 5H1D-HUMAN 5H1D-RABIT 5H1D-RAT 5H1E-HUMAN 5H1F-CAVPO 5H1F-HUMAN 5H1F-MOUSE 5H1F-RAT 5H2A-HUMAN 5H2A-MOUSE 5H2A-PACMO 5H2A-PIG 5H2A-RAT 5H2B-HUMAN 5H2B-MOUSE 5H2B-RAT 5H2C-HUMAN 5H2C-MOUSE 5H2C-RAT 5H4-CAVPO 5H4-HUMAN 5H4-MOUSE 5H4-RAT 5H5A-HUMAN 5H5A-MOUSE 5H5A-RAT 5H5B-MOUSE 5H5B-RAT 5H6-HUMAN 5H6-MOUSE 5H6-RAT 5H7-CAVPO 5H7-HUMAN 5H7-MOUSE 5H7-RAT 5H7-XENLA 5HT1-APLCA 5HT1-DROME 5HT2-APLCA 5HTA-DROME SHT-BOMMO SHT-HELVI SHT-LYMTST AIAA-BOVIN AIAA-CAVPO AIAA-HUMAN AIAA-MOUSE AI5A-ORYLA AI5A-RABIT AI5A-RAT AI5A-HUMAN AI5A-MESAU AI5A-MOUSE AI5A-RAT AI5A-HUMAN AI5A-MOUSE AI1A-RABIT AI1A-RAT A2AA-BOVIN A2AA-CAVPO A2AA-HUMAN A2AA-MOUSE A2AA-PIG A2AA-RAT A2AB-CAVPO A2AB-HUMAN A2AB-MOUSE A2AB-ORYAL A2AB-RAT A2AC-CAVPO A2AC-DIDMA A2AC-HUMAN A2AC-MOUSE A2AC-RAT A2AR-CAURAU A2AR-LABOS ACMI-HUMAN ACMI1-ACMCM ACMI1-MOUSE ACMI1-PIG ACMI1-RAT ACMI2-CHICK ACMI2-HUMAN ACMI2-MOUSE ACMI2-PIG ACMI2-RAT ACMI3-BOVIN ACMI3-CHICK ACMI3-GORGO ACMI3-HUMAN ACMI3-MOUSE ACMI3-PANTR ACMI3-PONPY ACMI3-RAT ACMI4-CHICK ACMI4-HUMAN ACMI4-RAT ACMI4-XENLA ACMI5-HUMAN ACMI5-MOUSE ACMI5-RAT BIAR-BOVIN BIAR-CANFA BIAR-FELCA BIAR-HUMAN BIAR-MACMU BIAR-MELGA BIAR-MOUSE BIAR-PIG BIAR-RAT BIAR-SHEEP BIAR-XENLA B2AR-BOVIN B2AR-CANFA B2AR-FELCA B2AR-HUMAN B2AR-MACMU B2AR-MOUSE B2AR-PIG B2AR-RAT B3AR-BOVIN B3AR-CANFA B3AR-CAPHI B3AR-FELCA B3AR-HUMAN B3AR-MACMU B3AR-MOUSE B3AR-RAT B3AR-SHEEP B4AR-MELGA D1DR-CAURAU D1DR-FUGRU D1DR-OREMO D2D1-XENLA D2DR-BOVIN D2DR-CERAE D2DR-FUGRU D2DR-HUMAN D2DR-MELGA D2DR-MOUSE D3DR-CERAE D3DR-HUMAN D3DR-MOUSE D3DR-RAT D4DR-HUMAN D4DR-MOUSE D4DR-RAT D5DR-FUGRU DADR-DIDMA DADR-HUMAN DADR-MACMU DADR-PIG DADR-RAT DADR-XENLA DBDR-HUMAN DBDR-RAT DCDR-XENLA DOP1-DROME DOP2-DROME GREI-BALAM GRE2-BALAM HHIR-BOVIN HHIR-HUMAN HHIR-MOUSE HHIR-RAT HH2R-CANFA HH2R-CAVPO HH2R-HUMAN HH2R-RAT HH3R-CAVPO HH3R-HUMAN HH3R-RAT HH4R-HUMAN O02146 O15969 O15970 O17470 O17496 O18512 O42315 O42316 O42317 O42322 O60451 O61730 O76267 O77254 O96716 O97171 OAR1-LOCMI OARI-LYMTST OAR2-LYMTST OAR-BOMMO OAR-DROME OAR-HELVI P90927 P91096 P97842 Q13167 Q13675 Q13729 Q24038 Q63004 Q923X3 Q923X6 Q923X7 Q923X8 Q923X9 Q923Y0 Q923Y1 Q923Y2 Q923Y3 Q923Y4 Q923Y5 Q923Y6 Q923Y7 Q923Y8 Q923Y9 Q96N94 Q96R19 Q96R91 Q96J10 Q98841 Q98842 Q98843 Q98844 Q98988 Q99MB0 Q9BM8A Q9BKZ0 Q9D282 Q9DBL0 Q9G1S6 Q9GT70 Q9GU11 Q9GK99 Q9GKA0 Q9GKL2 Q9GL56 Q9GL57 Q9GLP5 Q9QG54 Q9MYT5 Q9MZ00 Q9MZU2 Q9MZU3 Q9N623 Q9N626 Q9N927 Q9N928 Q9N980 Q9N981 Q9N982 Q9N983 Q9NGQ2 Q9NZR3 Q9PSA6 Q9PSA7 Q9PTFE Q9QW44 Q9QW71 Q9QW73 Q9QWS2 Q9QX37 Q9TSW7 Q9TTTM Q9US47 Q9UTD5 Q9UD63 Q9UN267 Q9UPA9 Q9VB33 Q9VS03 Q9VD16 Q9ND80 Q9YHA5 </p>

(continued)

TABLE 6. (Continued)

[illegible]

(continued)

TABLE 6. (Continued)

Subfamily	Protein ID's
Prostanoid	O00326 PD2R-MOUSE PE22-MOUSE PE23-BOVIN PE23-HUMAN PE23-RABIT PF2R-MOUSE PI2R-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 PI2R-HUMAN PI2R-MOUSE PI2R-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 Q9BXA5 Q9BXC1 Q9BYU4 Q9CFZ4 Q9DE05 Q9JJS7 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 Vinga, 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.
- Br  zma, 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.
- Foresse, 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., H  rsch, 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., Altshul, 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