

Hidden Markov Models in Computational Biology

Applications to Protein Modeling

Anders Krogh^{1†}, Michael Brown¹, I. Saira Mian²
Kimmen Sjölander¹ and David Haussler^{1‡}

¹Computer and Information Sciences

²Sinsheimer Laboratories

University of California, Santa Cruz,
CA 95064, U.S.A.

Hidden Markov Models (HMMs) are applied to the problems of statistical modeling, database searching and multiple sequence alignment of protein families and protein domains. These methods are demonstrated on the globin family, the protein kinase catalytic domain, and the EF-hand calcium binding motif. In each case the parameters of an HMM are estimated from a training set of unaligned sequences. After the HMM is built, it is used to obtain a multiple alignment of all the training sequences. It is also used to search the SWISS-PROT 22 database for other sequences that are members of the given protein family, or contain the given domain. The HMM produces multiple alignments of good quality that agree closely with the alignments produced by programs that incorporate three-dimensional structural information. When employed in discrimination tests (by examining how closely the sequences in a database fit the globin, kinase and EF-hand HMMs), the HMM is able to distinguish members of these families from non-members with a high degree of accuracy. Both the HMM and PROFILESEARCH (a technique used to search for relationships between a protein sequence and multiply aligned sequences) perform better in these tests than PROSITE (a dictionary of sites and patterns in proteins). The HMM appears to have a slight advantage over PROFILESEARCH in terms of lower rates of false negatives and false positives, even though the HMM is trained using only unaligned sequences, whereas PROFILESEARCH requires aligned training sequences. Our results suggest the presence of an EF-hand calcium binding motif in a highly conserved and evolutionary preserved putative intracellular region of 155 residues in the α -1 subunit of L-type calcium channels which play an important role in excitation-contraction coupling. This region has been suggested to contain the functional domains that are typical or essential for all L-type calcium channels regardless of whether they couple to ryanodine receptors, conduct ions or both.

Keywords: hidden Markov models; multiple sequence alignments; globin; kinase; EF-hand

1. Introduction

The rate of generation of sequence data in recent years provides abundant opportunities for the development of new approaches to problems in computational biology. In this paper, we apply

hidden Markov models (HMMs§) to the problems of statistical modeling, database searching, and multiple alignment of protein families and protein domains. To demonstrate the method, we examine three protein families. Each family consists of a set of proteins that have the same overall three-dimensional structure but widely divergent sequences. Features of the sequences that are determinants of folding, structure and function should be present as conserved elements in the family of sequences. We consider the globins, whole proteins ranging in length from 130 to 170 residues (with few exceptions) and two domains, the protein kinase catalytic domain (250 to 300 residues) and the EF-hand calcium-binding motif (29 residues). The same

† Present address: Electronics Institute, Build 349, Technical University of Denmark, 2800 Lyngby, Denmark.

‡ Author to whom all correspondence should be addressed.

§ Abbreviations used: HMM, hidden Markov models; EM, Expectation-Maximization; ML, maximum likelihood; MAP, maximum *a posteriori*; NLL-score, negative log likelihood score.

approach can be used to model families of nucleic acid sequences as well (Krogh *et al.*, 1993b).

A hidden Markov model (Rabiner, 1989) describes a series of observations by a "hidden" stochastic process, a Markov process. In speech recognition, where HMMs have been used extensively, the observations are sounds forming a word, and a model is one that by its hidden random process generates these sounds with high probability. Every possible sound sequence can be generated by the model with some probability. Thus, the model defines a probability distribution over possible sound sequences. A good word model would assign high probability to all sound sequences that are likely utterances of the word it models, and low probability to any other sequence. In this paper we propose an HMM similar to the ones used in speech recognition to model protein families such as globins and kinases. In speech recognition, the "alphabet" from which words are constructed could be the set of phonemes valid for a particular language; in protein modeling, the alphabet we use is the 20 amino acids from which protein molecules are constructed. Where the observations in speech recognition are words, or strings of phonemes, in protein modeling the observations are strings of amino acids forming the primary sequence of a protein. A model for a set of proteins is one that assigns high probability to the sequences in that particular set.

The HMM we build identifies a set of positions that describe the (more or less) conserved first-order structure in the sequences from a given family of proteins. In biological terms, this corresponds to identifying the core elements of homologous molecules. The model provides additional information, such as the probability of initiating an insertion at any position in the model and the probability of extending it. The structure of the model is similar to that of a profile (Waterman & Perlwitz, 1986; Barton & Sternberg, 1990; Gribskov *et al.*, 1990; Bowie *et al.*, 1991; Lüthy *et al.*, 1991), but slightly more general. Once we have built the model from unaligned sequences, we can generate a multiple alignment of the sequences using a dynamic programming method. By employing it for database searching, the model can be used to discriminate sequences that belong to a given family from non-members. Finally, we can study the model we have found directly, and see what it reveals about the common structure underlying the various sequences in the family.

Our method of multiple alignment differs quite markedly from conventional techniques, which are usually based on pairwise alignments generated by dynamic programming schemes (Waterman, 1989; Feng & Doolittle, 1987; Barton, 1990; Subbiah & Harrison, 1989). The alignments produced by these methods often depend strongly on the particular values of the parameters required by the method, in particular the gap penalties (Vingron & Argos, 1991). Furthermore, a given set of sequences is likely to possess both fairly conserved regions and

highly variable regions, yet conventional global methods assign identical penalties for all regions of the sequences. Substitutions, insertions, or deletions in a region of high conservation should ideally be penalized more than in a variable region, and some kinds of substitutions should be penalized differently in one position compared to another. That is one of the motivations for the present work. The statistical model we propose corresponds to multiple alignment with variable, position-dependent gap penalties. Furthermore, these penalties are in large part learned from the data itself. Essentially, we build a statistical model during the process of multiple alignment, rather than leaving this as a separate task to be done after the alignment is completed. We believe the model should guide the alignment as much as the alignment determines the model.

We are not the first group to employ hidden Markov models in computational biology. Lander & Green (1987) used hidden Markov models in the construction of genetic linkage maps. Other work employed HMMs to distinguish coding from non-coding regions in DNA (Churchill, 1989). Later, simple HMMs were used in conjunction with the EM algorithm to model certain protein-binding sites in DNA (Lawrence & Reilly, 1990; Cardon & Stormo, 1992) and, more recently, to model the N-caps and C-caps of alpha helices in proteins (D. Morris, unpublished results). These applications of HMMs and the EM (Expectation-Maximization) algorithm, including our own, presage a more widespread use of this technique in computational biology. During the time that we have been developing this approach, several related efforts have come to our attention. One is that of White, Stultz and Smith (White *et al.*, 1991; Stultz *et al.*, 1993), who use HMMs to model protein superfamilies. This work is more ambitious than our own, since superfamilies are harder to characterize than families. It is not yet clear how successful their work has been since no results are reported for sequences not in the training set. If there are weaknesses in their method, it is possible that these are due to the use of handcrafted models and reliance on prealigned data for parameter estimation. In contrast, our models have a simple regular structure, and we are able to estimate all the parameters of these models, including the size of the model directly from unaligned training sequences. Interestingly enough, they independently propose an alternate HMM state structure similar to ours† in section 6.3 of their paper (White *et al.*, 1991), where they discuss the relationship of their work to Bowie and co-workers (Bowie *et al.*, 1991), but they do not pursue this further. It is possible that the type of models we use may work better for characterizing superfamilies than those investigated by White *et al.* However, it is more likely that they are too simple, and that richer and more varied state

† Instead of using delete states, they have direct transitions between each pair of match states m_i and m_j with $i < j$.

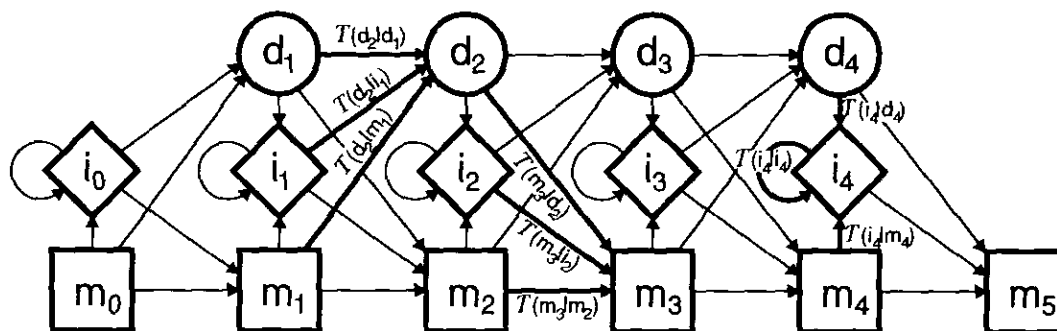


Figure 1. The model.

structure along the lines they propose is required for this problem. We recently found that Asai *et al.* (1993) have applied HMMs to the problem of predicting the secondary structure of proteins, obtaining prediction rates that are competitive with previous methods in some cases. In addition, Tanaka *et al.* (1993) also discuss the relationship between the HMM method for obtaining multiple alignments and previous methods. Finally, in work most closely related to our own, since the time we presented a preliminary report on this work (Haussler & Krogh, 1992; see also Haussler *et al.*, 1992), Baldi *et al.* (1993) have further demonstrated the usefulness of this technique by producing multiple alignments for immunoglobins and protease as well as globins and kinases†.

2. Methods

(a) HMM architecture

Consider a family of protein sequences that all have a common three-dimensional structure, for example the globins. The common structure in these sequences can be defined as a sequence of positions in space where amino acids occur. In the case of globins, whose structure contains principally α -helices, the 150 or so helical positions have been named A1, A2, ..., A16, B1, ... etc., where the letter denotes the α -helix, and the number indicates the location within that α -helix (see for example Bashford *et al.*, 1987). For each of these positions there is a (distinct) probability distribution over the 20 amino acids that measures the likelihood of each amino acid occurring in that position in a typical globin, as well as the probability that there is no amino acid in that position (i.e. that a sequence belonging to this family may have a gap at that position in a multiple alignment). These have been called profiles (Waterman & Perlwitz, 1986; Barton & Sternberg, 1990; Gribskov *et al.*, 1990; Bowie *et al.*, 1991; Lüthy *et al.*, 1991). A profile of globins can be thought of as a statistical model for the family of globins, in that for any sequence of amino acids, it defines a probability for that sequence, in such a way that globin sequences tend to have much higher probabilities than non-globin sequences.

The type of hidden Markov model we use as a statistical model for a protein family can be viewed as a generalized profile. However, instead of describing the HMM directly

in terms of the probability it assigns to each protein sequence, we find that it is easier to first think of an HMM as a structure that generates protein sequences by a random process. This structure and corresponding random process is illustrated in Figure 1 and can be described as follows.

The main line of the HMM contains a sequence of M states, which we call match states, corresponding to positions in a protein or columns in a multiple alignment (M equals 4 in Fig. 1). Each of these states can generate a letter x from the 20-letter amino acid alphabet according to a distribution $\mathcal{P}(x|m_k)$, $k=1 \dots M$. The notation $\mathcal{P}(x|m_k)$ means that each of the match states m_k , $1 \leq k \leq M$, have distinct distributions. For each match state m_k , there is a delete state d_k that does not produce any amino acid but is a "dummy" state used to skip m_k . Finally, there are a total of $M+1$ insert states to either side of the match states which generate amino acids in exactly the same way as the match states, but use probability distributions $\mathcal{P}(x|i_k)$. In Figure 1, match, delete and insert states are shown as boxes, circles and diamonds, respectively. For convenience, we have added a dummy "BEGIN" state and a dummy "END" state, denoted m_0 and m_{M+1} , respectively, which do not produce any amino acid.

From each state, there are three possible transitions to other states, also shown in Figure 1. Transitions into match or delete states always move forward in the model, whereas transitions into insert states do not. Note that multiple insertions between match states can occur, since the self-loop on the insert state allows a transition from the insert state to itself. The transition probability from state q to state r is called $\mathcal{F}(r|q)$. Our notation is summarized in Table 1.

A sequence can be generated by a "random walk" through the model as follows: Commencing at state m_0 (BEGIN), choose a transition to m_1 , d_1 , or i_0 randomly

Table 1
Notation

x	Amino acid
s	Sequence of amino acids ($s = x_1 \dots x_L$)
L	Length of sequence
q, r	State in HMM
path	A sequence of states, $q_1 \dots q_N$
N	Number of states in a path
M	Length of model
m, i, d	Match, insert and delete states
m_0, m_{M+1}	Begin and end states
$\mathcal{P}(x q)$	Probability distribution of amino acids in state q
$\mathcal{F}(r q)$	Probability of a transition from state q to r

† They have developed a variant of the method described here that employs a gradient descent training algorithm in place of the EM algorithm.

according to the probabilities $\mathcal{F}(m_1|m_0)$, $\mathcal{F}(d_1|m_0)$, and $\mathcal{F}(i_0|m_0)$. If m_1 is chosen, generate the first amino acid x_1 from the probability distribution $\mathcal{P}(x|m_1)$, and choose a transition to the next state according to probabilities $\mathcal{F}(\cdot|m_1)$, where \cdot indicates any possible next state. If this next state is the insert state i_1 , then generate amino acid x_2 from $\mathcal{P}(x|i_1)$ and select the next state from $\mathcal{F}(\cdot|i_1)$. If delete (d_2) is chosen next, generate no amino acid, and choose the next state from $\mathcal{F}(\cdot|d_2)$. Continue in this manner all the way to the END state, generating a sequence of amino acids $x_1, x_2 \dots x_L$ by following a path of states $q_0, q_1 \dots q_N, q_{N+1}$ through the model, where $q_0 = m_0$ (the BEGIN state) and $q_{N+1} = m_{M+1}$ (the END state). Because the delete states do not produce any amino acid, N is larger than or equal to L . If q_i is a match or insert state, we define $l(i)$ to be the index in the sequence $x_1 \dots x_L$ of the amino acid produced in state q_i . The probability of the event that the path $q_0 \dots q_{N+1}$ is taken and the sequence $x_1 \dots x_L$ is generated is

$$\text{Prob}(x_1 \dots x_L, q_0 \dots q_{N+1}|\text{model}) \\ = \mathcal{F}(m_{N+1}|q_N) \times \prod_{i=1}^N \mathcal{F}(q_i|q_{i-1}) \mathcal{P}(x_{l(i)}|q_i), \quad (1)$$

where we set $\mathcal{P}(x_{l(i)}|q_i) = 1$ if q_i is a delete state. The probability of any sequence $x_1 \dots x_L$ of amino acids is a sum over all possible paths that could produce that sequence, which we write as follows:

$$\text{Prob}(x_1 \dots x_L|\text{model}) \\ = \sum_{\text{paths } q_0 \dots q_{N+1}} \text{Prob}(x_1 \dots x_L, q_0 \dots q_{N+1}|\text{model}). \quad (2)$$

In this way a probability distribution on the space of sequences is defined. The goal is to find a model (i.e. a proper model length and probability parameters) that accurately describes a family of proteins by assigning large probabilities to sequences in that family.

This particular structure for the HMM was chosen because it is the simplest model that captures the structural intuition of a protein: (a) a sequence of positions, each with its own distribution over the amino acids; (b) the possibility for either skipping a position or inserting extra amino acids between consecutive positions; and (c) allowing for the possibility that continuing an insertion or deletion is more likely than starting one. This choice appears to have worked well for modeling the protein families that we have examined, but other types of HMMs may be better at other tasks (e.g. the more elaborate models for protein superfamilies used by White *et al.*, 1991; Stultz *et al.*, 1993). The important feature of the HMM method is its generality. One can choose any structure for the states and transitions that is appropriate for the problem at hand. Examples of more general HMM architectures are given in sections (d) and (e), below.

(b) *Estimating the parameters of an HMM from training sequences*

All the parameters in the HMM (i.e. the transition probabilities and the amino acid distributions) could in principle be chosen by hand from an existing alignment of protein sequences, as in Gribskov *et al.* (1990), White *et al.* (1991), Stultz *et al.* (1993), or from information about the three-dimensional structure of proteins, as in Bowie *et al.* (1991), White *et al.* (1991), Stultz *et al.* (1993). The novel approach we take is to "learn" the parameters entirely automatically from a set of unaligned primary sequences, using an EM algorithm. This approach can in

principle find the model that best describes a given set of sequences.

Given a set of training sequences $s(1), \dots, s(n)$, one can see how well a model fits them by calculating the probability that it generates them. This probability is simply a product of terms of the form given by equation (2), i.e.

$$\text{Prob}(\text{sequences}|\text{model}) = \prod_{j=1}^n \text{Prob}(s(j)|\text{model}), \quad (3)$$

where each term $\text{Prob}(s(j)|\text{model})$ is calculated by substituting $x_1 \dots x_L = s(j)$ in equation (2). This is called the likelihood of the model. One would like this value to be high. The maximum likelihood (ML) method of model estimation is to find the model that maximizes the likelihood (3).

An alternate approach to ML estimation is the maximum *a posteriori* (MAP) approach. Here, we assume a prior probability distribution over all possible parameters of the model embodying prior beliefs on what a model should be like. This can then be used to "penalize" models that are known to be bad or uninteresting. We discuss this further in Krogh *et al.* (1993a). In MAP estimation, we try to maximize the posterior probability of the model given the sequences. Using Bayes rule, the posterior probability can be calculated as

$$\text{Prob}(\text{model}|\text{sequences}) \\ = \frac{\text{Prob}(\text{sequences}|\text{model}) \text{Prob}(\text{model})}{\text{Prob}(\text{sequences})}. \quad (4)$$

Here $\text{Prob}(\text{model})$ is the prior probability distribution, and $\text{Prob}(\text{sequences})$ can be viewed as a normalizing constant. Since this normalizing constant is independent of the model, MAP estimation is equivalent to maximizing

$$\text{Prob}(\text{sequence}|\text{model}) \text{Prob}(\text{model}). \quad (5)$$

over all possible models. The MAP approach is closely related to minimum description length (Jurka & Milosavljevic, 1991) and minimum message length (Allison *et al.*, 1992) methods.

There is no known efficient way to directly calculate the best HMM model either in the ML or MAP sense. However, there are algorithms that given an arbitrary starting point find a local maximum by iteratively re-estimating the model in such a way that the likelihood (or the posterior probability) increases in each iteration. The most common one is the Baum-Welch or forward-backward algorithm (Rabiner, 1989; Lawrence & Reilly, 1990), which is a version of the general EM method often used in statistics (Dempster *et al.*, 1977). The process of the EM algorithm can be viewed as an iterative adaptation of the model to fit the training sequences. The steps in this process can be summarized as follows:

(1) An initial model is created by assigning values to the transition probability $\mathcal{F}(r|q)$ and the amino acid generation probability $\mathcal{P}(x|q)$ for each x, q and r , where x is one of the 20 amino acids and q and r are states in the HMM connected by a transition arc. If one already knows some features present in the sequences, or constraints on the sequences, these may sometimes be encoded in the initial model. The current model is set to this initial model.

(2) Using the current model, all possible paths for each training sequence are considered in order to get a new estimate $\hat{\mathcal{F}}(r|q)$ of the transition probability $\mathcal{F}(r|q)$ and a new estimate $\hat{\mathcal{P}}(x|q)$ of the amino acid generation probability $\mathcal{P}(x|q)$ for each x, q and r . The transition

probability estimate $\hat{\mathcal{F}}(r|q)$ is obtained by counting the number of times a transition is made from state q to r , for all paths of all training sequences, weighted by the probability of the path. The estimate $\hat{\mathcal{P}}(x|q)$ is made in a similar manner, by counting the number of times the amino acid x is aligned to the state q .

(3) In the next step of ML estimation, a new current model is created by simply replacing $\mathcal{F}(r|q)$ by $\hat{\mathcal{F}}(r|q)$ and $\mathcal{P}(x|q)$ by $\hat{\mathcal{P}}(x|q)$ for each x , q and r . In MAP EM estimation, the parameters $\hat{\mathcal{F}}(r|q)$ and $\hat{\mathcal{P}}(x|q)$ are further modified by considering the prior probability of the model before they are used to replace the old parameters.

(4) Steps (2) and (3) are repeated until the parameters of the current model change only insignificantly.

Since the quality of the current model (as measured by equations (3) or (5)) increases in each iteration, and no model is arbitrarily good, the process eventually terminates and produces a model that is, at least locally, the best model for the training sequences to within some specified precision of the parameters (Dempster *et al.*, 1977). Typically, this occurs very rapidly (e.g. in less than 10 iterations) even for large models and large sets of training sequences.

The main computational bottleneck in the algorithm is step (2), since individually examining each possible path for every training sequence would generally take time exponential in the length of the longest training sequence. However, it is possible to use a dynamic programming technique known as the forward-backward procedure to speed up this step. Using this method, the new parameter estimates can be calculated in time proportional to the number of states in the model multiplied by the total length of all the training sequences. Details are given in the excellent tutorial article on HMMs by Rabiner (1989).

The forward part of the forward-backward procedure can also be used to efficiently compute $-\log \text{Prob}(\text{sequence}|\text{model})$, the negative logarithm of the probability of a sequence given the model (as defined in equation (2)), without summing over all possible paths for the sequence (Rabiner, 1989). We call this the negative log likelihood (NLL)-score of the sequence. The average NLL-score of a training sequence is inversely related to the likelihood of the model, given by equation (3), and hence serves as a numerical measure of progress for each iteration of the EM procedure. The NLL-score can also be used to evaluate how well the model fits a novel "test" sequence not present in the training set, as described in section (c) below.

(c) *The Viterbi algorithm and multiple alignment from an HMM*

The forward-backward procedure is related to the dynamic programming technique used to align one sequence to another, or more generally to align a sequence to a profile. A variant of the forward-backward procedure known as the Viterbi algorithm is similar to the standard profile alignment algorithm (Waterman & Perlwitz, 1986; Barton & Sternberg, 1990; Gribskov *et al.*, 1990). Instead of calculating the NLL-score for a sequence, which implicitly involves all possible paths for that sequence through the model, the Viterbi algorithm computes the negative logarithm of the probability of the single most likely path for the sequence. We can write this as

$$-\log \max_{\text{paths}} \text{Prob}(s, \text{path}|\text{model}), \quad (6)$$

where $\text{Prob}(s, \text{path}|\text{model})$ is given in equation (1), with $s = x_1 \dots x_L$ and $\text{path} = q_0 \dots q_{N+1}$. Instead of first maxi-

mizing the probability of the path and then taking the negative logarithm, it is convenient (and equivalent) to simply minimize the negative logarithm of the probability over all paths. This minimum we will call the distance from the sequence to the model,

$$\begin{aligned} \text{dist}(s, \text{model}) &= \min_{\text{paths}} \{-\log \text{Prob}(s, \text{path}|\text{model})\} \\ &= \min_{\text{paths}} \sum_{i=1}^{N+1} [-\log \mathcal{F}(q_i|q_{i-1}) - \log \mathcal{P}(x_{i(i)}|q_i)] \end{aligned}$$

This distance from a sequence to a model is analogous to the standard "edit distance" from one sequence to another (with gap penalties), see e.g. Waterman (1989), but is perhaps more related to the distance from a sequence to a profile. The term $-\log \mathcal{P}(x_{i(i)}|q_i)$ represents a penalty for aligning the amino acid $x_{i(i)}$ to the position represented by state q_i in the model. The term $-\log \mathcal{F}(q_i|q_{i-1})$ corresponds to a penalty for using the transition from q_{i-1} to q_i in the model. If this is a transition from a match state to a delete state, then this represents a gap-initiation penalty; if it is from a delete state to a delete state it represents a gap-extension penalty; if it is from a match state to an insert state, it represents an insertion-initiation penalty; and if it is a transition from an insertion state to itself (a "self-loop"), then it represents an insertion extension penalty. One of the main features of this distance measure is that all these penalties depend on the position in the model, whereas they would be fixed in most standard pairwise alignment methods. Often the most likely path has a significantly higher probability than all other paths, and in that case the distance defined here will be approximately equal to the NLL-score defined earlier.

The computation time for the Viterbi algorithm is proportional to the number of states in the model multiplied by the length of the sequence being aligned, i.e. the same as the time for the forward-backward algorithm. In addition, with a simple extension to the algorithm, the most probable path itself can be found using the usual backtracking technique (Rabiner, 1989). This is the method we use to obtain our multiple alignments: each sequence is aligned to the model by the Viterbi algorithm, after which the mutual alignment of the sequences among themselves is then determined.†

(d) *Using the HMM to cluster sequences and discover subfamilies*

When a relatively large number of sequences are available, it is sometimes possible to obtain improved results by dividing these sequences into clusters of similar sequences and training a different HMM for each cluster/subfamily. The results of this are illustrated in more detail in Results section (a). Given a large set of unlabeled and unaligned sequences, a simple extension of the hidden Markov model enables us to use the EM training algorithm to automatically partition the sequences into clusters of similar sequences. By iteratively splitting clusters, this method might be useful for building phylogenetic trees in a "top-down" manner. However, when the clusters become too small there will be an insufficient number of sequences in each cluster to construct an accurate model, so some "bottom-up" processing may still be necessary.

In order to discover w clusters in the data, we make w copies of the HMM, one for each cluster. We call these

† We make no attempt to align portions of the sequence that use the insert states of the model.

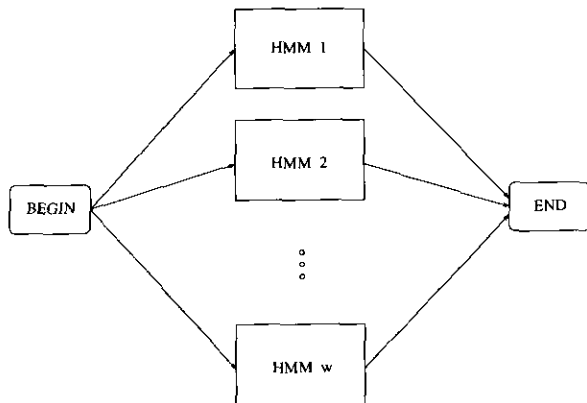


Figure 2. HMM architecture for discovering subfamilies.

components of the (composite) HMM. Presently, the number w of clusters and the initial lengths of the models for these clusters are determined empirically. We then add a new begin state with w outgoing transitions, one to each of the begin states of the component HMMs (see Fig. 2).

This new begin state is analogous to the other begin states in that it generates no amino acid. We then train this composite model with the EM algorithm as described in section (b), above. The EM re-estimation of a component model is the same as the re-estimation of a single model, except that the weight that a sequence has in the re-estimation of a component is proportional to the probability of the sequence given that component model. Thus, sequences that have better NLL-scores for a particular HMM component have greater influence in re-estimating the parameters of that component, and this causes the parameters of that component to change in such a way that the component further "specializes" in modeling those sequences. The "surgery" procedure described below in section (g) is used to adapt the length of that component to further specialize it. In this manner, the individual components evolve during training to represent clusters in the training sequences. This way of using EM is called mixture modeling in the statistics literature (Duda & Hart, 1973; Everitt & Hand, 1981), and is known as "(soft) competitive learning" in the neural network literature (Nowlan, 1990).

When the model is trained, the probability of a sequence given any of the submodels can be calculated, i.e. the probability that the sequence belongs to the corresponding cluster/subclass. The negative logarithm of this probability corresponds to the NLL-score calculated for a simple HMM. As with the standard HMM we use, this yields a quantitative measure of how well the model fits the data. The clusters found can also be compared to

known subfamilies of the sequences. Experiments with the clustering of globin sequences are described in Results section (a).

(e) Modeling protein domains with an HMM

There are many cases when one does not want to build a statistical model of a family of whole proteins like globins, but instead to build a model of a structural motif or domain that occurs as a subsequence in many different kinds of proteins, such as the EF-hand motif (Nakayama *et al.*, 1992) or the kinase catalytic domain (Hanks & Quinn, 1991). Here we expect our model to only match a relatively small subsequence of any given protein, with many other unmatched amino acids appearing before and after this subsequence. One approach to this problem is to alter the dynamic programming method used to align a sequence to a model so that it tries all possible ways of aligning each subsequence of the sequence to a model (Waterman, 1989). We use a simpler (but almost equivalent) method in which only the HMM model is altered, so that the same standard procedures (forward-backward and Viterbi) which we use for models of whole proteins can be used without modification for models of domains.

Consider a training set of many unaligned sequences consisting not of complete proteins, but of a specific domain. Our first step is to train an HMM for these sequences exactly as described earlier. As shown in Figure 1, this HMM will have initial and final "dummy" match states m_0 and m_{N+1} (where $N+1=5$ in Fig. 1) that do not match any amino acid. To alter the HMM to represent a protein domain, we create 2 new insert states i_B and i_E , adding i_B to the model before the state m_0 and i_E at the end of the model after m_{N+1} (see Fig. 3).

We then add a new dummy BEGIN state before i_B and a new dummy END state after i_E . Eight new transitions are also added to the model. The first 4 are from BEGIN to i_B , from m_{N+1} to i_E , and the self-loops from i_B to itself and from i_E to itself. These all have the same probability p , for some p between 0 and 1. The second 4 transitions are from i_B to m_0 , from BEGIN to m_0 , from i_E to END, and from m_{N+1} to END. These all have probability $1-p$. The new states added before and after the model, along with these transitions, form 2 new modules, 1 for matching the extra amino acids that occur in the sequence before the domain, and the other for matching the amino acids after the domain.

The choice of the parameter p does affect the way that the overall model aligns with a given sequence. To see how, it is convenient to think of the negative logarithm of the probability of a transition as a penalty for using that transition, as described in section (c), above. In the modified model, all sequences must suffer a penalty of $-\log(1-p)$ to enter and again to exit the domain part of the model, no matter which path they take. Hence this penalty is a fixed cost, which can be ignored when

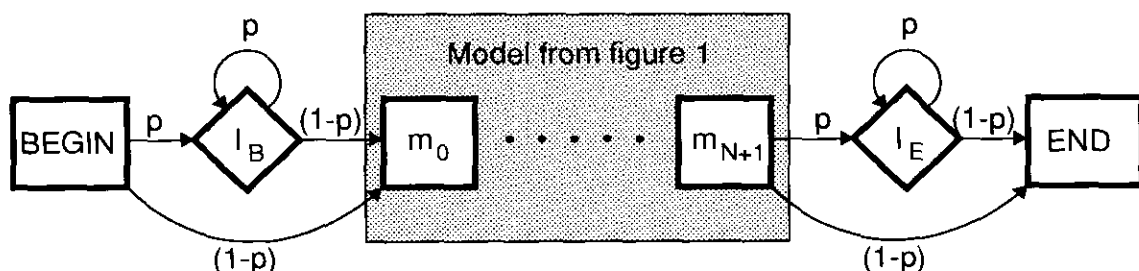


Figure 3. HMM architecture for modeling domains.

comparing the distances or NLL-scores of 2 sequences with respect to the model. In addition to this penalty, all sequences will suffer a penalty of $K(-\log p + \log 20)$, where $K \geq 0$ is the number of amino acids that are not matched to the original domain model, but are instead matched in the states i_B and i_E . The $-\log 20$ term arises because we set the probabilities of each amino acid to $1/20$ in the insertion states i_B and i_E (see Krogh *et al.*, 1993a). Thus p will determine the "pressure" on the sequence to align something to the domain model, i.e. if p is low it is advantageous to squeeze many amino acids into the domain model, using the insert states in this part of the model. If it is high, it is possible that most sequences would prefer to pass through the delete states in the domain model, aligning everything instead to the new modules before and after it. It is straightforward to estimate p the same way as all the other parameters, the only additional problem is that the same value must be used in all the transitions that use this value, "tying" these parameters to each other. Otherwise the model might become biased towards aligning the domain either near the beginning of the sequence or near the end of the sequence. We have not attempted to estimate p . Rather, we have used a fixed $p = 1$ with good results. (This should be thought of as a limit of p approaching 1, otherwise $-\log(1 - p)$ is infinite.)

Using this construction, it may also be possible to discover interesting domains by training on whole protein sequences, and letting EM determine which part of the proteins to model. Furthermore, if more than one occurrence of the same domain is expected in some sequences, then this model can be further modified to find all occurrences. This is accomplished by simply adding a transition from the END state back to the BEGIN state.

(f) Searching a database with an HMM

Once an HMM is built for a family of proteins, it can be used to search a database such as PIR or SWISS-PROT for other proteins in this family. Similarly, if an HMM is built for a protein domain or motif, then it can be used to search for occurrences of this domain or motif in the database, much like a PROSITE expression (Bairoch, 1992), a commonly used method for searching for patterns found in protein sequences. Like a profile (Waterman & Perlwitz, 1986; Barton & Sternberg, 1990; Gribskov *et al.*, 1990; Bowie *et al.*, 1991; Lüthy *et al.*, 1991), an HMM has an advantage over a PROSITE expression for database searches. It takes into account a large amount of statistical information in matching a sequence, and weighs this information appropriately, rather than relying on relatively rigid matching rules.

As described in section (b), above, the forward part of the forward-backward dynamic programming method calculates a NLL-score for any test sequence that measures how well it fits the model. This NLL-score is the negative logarithm of the probability of the sequence given the model. It turns out that this raw NLL-score is too dependent on the length of the test sequence to be used directly to decide if the sequence is in the family modeled by the HMM or not. However, we can overcome this problem by normalizing this NLL-score appropriately.

Whenever we build an HMM for a family of proteins or for a protein domain, we run all the proteins in a standard database (for instance, SWISS-PROT) through this HMM and compute the NLL-score for each sequence. A scatter plot of sequence length *versus* NLL-score for our kinase catalytic domain model is given in Figure 9.

Most proteins tend to lie on a fairly straight line (towards the top of the plot) indicating that the NLL-score for these proteins is proportional to their lengths. These proteins are the ones that do not contain the kinase catalytic domain and thus look like "random proteins" to the kinase catalytic domain model. In contrast, the proteins that do contain the kinase catalytic domain tend to have NLL-scores that are much lower than expected for proteins of their length, and hence appear below the linear band of non-kinase proteins.

We can quantify the difference between NLL-scores for proteins containing the kinase catalytic domain and NLL-scores for proteins not containing the domain by a simple statistical method, as follows. Using a local windowing technique,† we first calculate a smooth average curve for the roughly linear band of the NLL-score *versus* length plot. The standard deviation around this average curve is also calculated. Using this, we calculate the difference between the NLL-score of a sequence and the average NLL-score of typical sequences of that same length, measured in standard deviations. This number is called the Z-score for the sequence. We then choose a Z-score cut-off, either *a priori* or by looking at the histogram of Z-scores for sequences in the database (see Fig. 10), and use it to decide if a given sequence fits the model or not. We have found that a Z-score of approximately 5 appears a good choice in most cases we have examined, but we suggest carefully checking the histogram by eye before deciding on a cut-off. For example, for our HMM of the kinase catalytic domain, sequences with Z-scores below 5 are classified as not containing the kinase catalytic domain, and sequences with Z-scores above 5 are classified as containing the catalytic domain. If the Z-score of a sequence indicates that it contains the catalytic domain, we can align the sequence to the catalytic domain HMM to find out where this domain occurs in the sequence. The time it takes to do a database search is proportional to the number of residues in the database times the length of the model. For our globin model (length 147) we can search the SWISS-PROT database (about 8,375,000 residues) in approximately 2 CPU hours on a Sun Sparcstation 1. Using the shorter EF-hand model (length 29) it takes only 18 CPU seconds (11 user min) on a Sun Sparcstation 2. A parallel implementation of the search procedure (not yet implemented) will speed up these searches substantially, as it has the EM training procedure.

While the statistical techniques we have used to determine Z-scores are still quite crude, we have found that the HMMs are sufficiently good models that these techniques work well enough in practice. However, it may be that more sophisticated techniques are needed in certain cases.

† The average curve is calculated as follows. For each length i starting at $i = 1$, the length l_i is computed such that there are at least 500 proteins of lengths i to l_i and less than 500 proteins of lengths i to $l_i - 1$. The length interval i to l_i is called a window. The average curve is piecewise linear through the points corresponding to the average length and average NLL-score for each window. The first and last parts of the curve are calculated by linear regression in the first and last window, respectively. The standard deviation of the points from the smooth curve is also calculated for each window. The estimate of the average curve can be improved by eliminating outliers, i.e. NLL-scores that lie many standard deviations from the average. We iterate the process of removing outliers and re-estimating the average curve until no more outliers remain.

(g) *Initial model, local minima, and choice of model length*

As mentioned in section (b), above, when estimating the model from the training sequences, the EM algorithm does not guarantee convergence to the best model. It is basically a steepest-descent-type algorithm that climbs the nearest peak (local maximum) of the likelihood function (or the posterior probability in MAP estimation). Since finding the globally optimal model seems to be a difficult optimization problem in general (Abe & Warmuth, 1990), we have experimented with various heuristic methods to improve the performance of the method.

Probably the best method is to give the model a hint if something is already known about the sequences, which is often the case. A good starting point makes it much more likely that the nearest peak is at least close to optimal. This is done by setting the probabilities in the initial model to values reflecting that knowledge. If, for instance, an alignment of some of the sequences is available, it is straightforward to translate that into a model by simply calculating the relative frequency of the amino acids and the transition frequencies in each position, as in the profile method (Gribskov *et al.*, 1990).

It is of course even more interesting if the model can be found from a *tabula rasa*, i.e. using no knowledge about the sequences. For that we have used an initial model where all equivalent probabilities are the same, i.e. $\mathcal{P}(m_{k+1}|m_k)$ is independent of the position k in the model, and similarly for all other transition probabilities, and $\mathcal{P}(x|m_k)$ is also independent of k . To avoid the smaller local maxima, noise is added to the model during the iteration before each re-estimation. Initially quite a lot of noise is added, but over 10 iterations the noise is decreased linearly to zero. Since noise is added directly to the model, it is not like the usual implementation of simulated annealing, but the principle is the same. The "annealing schedule" is presently rather arbitrary, but it does seem to give reasonable results† if it is applied several times, and the best of the models found is used as the final model.

It is important that the best model be selected, since suboptimal models do produce inferior alignments in general. However, when studying alignments from suboptimal globin models, we noted that they tend to align some regions well, occasionally getting better alignments in those regions than the best overall model found, while in other regions they are completely incorrect. This leaves open the intriguing possibility of combining the best solutions found for different regions into a new overall best model. We have not yet explored this possibility.

The length of the model is also a crucial parameter that needs to be chosen *a priori*. However, we have developed a simple heuristic that selects a good model length, and even helps in the problem of local maxima. The heuristic is this: after learning, if more than a fraction‡ γ_{del} of the paths of the sequences choose d_k , the delete state at position k , that position is removed from the model. Similarly, if more than a fraction γ_{ins} make insertions at position k (in state i_k), a number of new positions equal to the average number of insertions made at that position are inserted into the model after position k . After these

changes in the model, it is retrained, and this cycle is repeated until no more changes are needed. We call this "model surgery".

(h) *Over-fitting and MAP estimation*

A model with too many free parameters cannot be estimated well from a relatively small data set of training sequences. If we try to estimate such a model, we run into the problem of overfitting, in which the model fits the training sequences very well, but gives a poor fit to related (test) sequences that were not included in the training set. We say that the model does not "generalize" well to test sequences. This phenomenon has been well documented in statistics and machine learning (see e.g. Geman *et al.*, 1992; Berger, 1985). One way to deal with this problem is to control the effective number of free parameters in the model by using prior information. This can be accomplished with MAP estimation. Parameters that we assume (*via* our prior distribution on models) can be well-estimated *a priori* in effect become less adaptive, because it takes a lot of data to override our prior beliefs about them, whereas those about which we have only weak prior knowledge are estimated in almost the same manner as in maximum likelihood estimation. In this way, the model can have a very large number of parameters, but a much smaller number of "effectively free" parameters. To make MAP estimation practical, we use Dirichlet distributions as priors. The details of the method are described elsewhere (Krogh *et al.*, 1993a; Brown *et al.*, 1993).

3. Results

(a) *Globin experiments*

The modeling was first tested on the globins, a large family of heme-containing proteins involved in the storage and transport of oxygen that have different oligomeric states and overall architecture (for a review see Dickerson & Geis (1983)). Hemoglobins are tetramers composed of two α chains and two other subunits (usually β , γ , δ or θ). Myoglobin is a single chain, some insect globins are present as dimers and some intracellular invertebrate globins occur in large complexes of many subunits.

Globin sequences were extracted from the SWISS-PROT database (release 19) by searching for the keyword "globin". Eliminating the false positives, resulted in 625 genuine globin sequences of average length 145 amino acids. We left three non-globins in the sample for illustrational purposes giving a total of 628 sequences. The sample of globins in the database is not the random sample a statistician would prefer, but is perhaps one of the best and largest collections of protein sequences from a homologous family. Searching for the words "alpha", "beta", "gamma", "delta", "theta", and "myoglobin" in the data file yielded 224 alpha, 199 beta, 16 gamma, 8 delta and 5 theta chains and 79 myoglobins, which adds up to 531 sequences. These should naturally be considered minimum numbers, but they give a good picture of how skewed the sample is.

To test our method, we trained an HMM using the method described in Methods sections (b) and

† An alternate method that also appears to give good results has been developed by Baldi *et al.* (Baldi *et al.*, 1993; Baldi & Chauvin, 1993). This method uses stochastic gradient descent in place of the EM method, which may help in avoiding local minima.

‡ Currently we choose γ_{del} and γ_{ins} each to be 1/2.


```

Helix      AAAAAAAAAAAAAAAAAA  BBBB BBBB BBBB BBBB CCCCCCCCCC  DDDDDDEE
HBA_HUMAN  -----VLSPADKTNVKAAWGKVG A--HAGEYGAEALERMFLSFPTTKTYFPHF--DLS-----HGSA
HBB_HUMAN  -----VHLTPEEKSAVTALWGKV---NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNP
MYG_PHYCA  -----VLSEGEWQLVLHVWAKVEA--DVAGHGQDILIRLFKSHPETLEKFDKFRFKHLKTEAEKASE
GLB3_CHITP -----LSADQISTVQASFDKVKG-----DPVGILYAVFKADPSIMAKFTQFAG--KDLESIKGTA
GLB5_PETMA PIVDTGSVAPLSAAEKTIRSAWAPVYS--TYETSGVDILVKFFTSTPAAQEFFPKFKGLTTADQLKSA
LGB2_LUPLU -----GALTESQAALVKSSWEEFNA--NIPKHTHRFFILVLEIAPAACKDLFS--FLK--GTSEVPQNNP
GLB1_GLYDI -----GLSAAQRQVIAATWKDIAGADNGAGVGKDCCLIKFLSAHPQMAAVFG--FSG----AS---DP

Helix      EEEEEEEEEEEEEEEEE  FFFFFFFFFF  FGGGGGGGGGGGGGGGGGGG
HBA_HUMAN  QVKGHGKKVADALTNVAHV---D--DMPNALSALSDLHAKL--RVDPVNFKLLSHCLLVTLAAHLP AE
HBB_HUMAN  KVKAHGKQVLFVSDGLAHL---D--NLKGTATLSELHCDKL--HVDPENFRLLGNVLVLCVLAHFFGKE
MYG_PHYCA  DLKKGVTVLTALGAILKK---K--GHHEAELKPLAQSHATKH--KIPIKYLEFISEAIIHVLHSRHPGD
GLB3_CHITP PFETHANRIVGFFSKIIGEL---P---NIEADVNTFVASHKPRG---VTHDQLNFRAGFVSYMKAHT--D
GLB5_PETMA DVRWHAERIINAVNDAVASM--DDTEKMSMKLRDLSGKHAJSF--QVDPQYFKVLA AVIADTV AAG----
LGB2_LUPLU ELQAHAGKVFVLYEAAIQLVTVVDTATLKNLGSVHVS KG---VADAHFPVVK EAILKTIKEVVGAK
GLB1_GLYDI GVAALGAKVLAQIGVAVSHL--GDEGKMVAQMKAVGVRHKGYGNKHIKAQYFEPLGASLLSMEHRIGGK

Helix      HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH
HBA_HUMAN  FTPAVHASLDKFLASVSTVLT SKYR-----
HBB_HUMAN  FTTPVQAAYQKV VAGVANALAHKYH-----
MYG_PHYCA  FGADAQGAMNKALELFRKDIAAKYKELGYQG
GLB3_CHITP FA-GAEAAWGATLD TFFGMIFSKM-----
GLB5_PETMA -----DAGFEKLSMICILLRSAY-----
LGB2_LUPLU WSEELNSAWTIA YDELAIVIKKEMNDAA---
GLB1_GLYDI MNAAAKDAWAAAYADISGALISGLQS-----

```

Figure 4. Seven representative globin sequences of known structure and their alignment taken from Bashford *et al.* (1987). The letters A to H in Helix denote the 8 different α -helices. Some regions, especially CD, D and FG, are not well defined. The sequences and their SWISS-PROT identifiers are Human α (HBA_HUMAN), human β (HBB_HUMAN), sperm whale myoglobin (MYG_PHYCA), larval *Chironomus thummi* globin (GLB3_CHITP), sea lamprey globin (GLB5_PETMA), *Lupinus luteus* leghemoglobin (LGB2_LUPLU), and bloodworm globin (GLB1_GLYDI). (In SWISS-PROT 19 a \$ is used instead of an “_” in the identifiers.)

(g). We used a homogeneous initial model that contained no knowledge about the globin family. Its probability parameters were derived from the prior, and were the same for all equivalent transitions (i.e. 9 different transition probabilities). All amino acid probabilities (the \mathcal{P} distributions) were set equal to the distribution of the amino acids given by Krogh *et al.* (1993a). In the insert states we used a probability of 1/20 for all amino acids. The only model parameters set by hand are the initial transition probabilities and corresponding regularization parameters (see Krogh *et al.*, 1993a). From our experience, the method does not seem to be very sensitive to the choice of these parameters, but it would require considerable further experimentation to verify this quantitatively.

For our training set, we picked 400 sequences at random from the 628 sequences. We withheld the remaining 228 sequences in order to test the model on data not used in the training process. The model was trained using noise and model surgery ($\gamma_{del} = \gamma_{ins} = 0.5$), as described in Methods section (g). This procedure was repeated about 20 times with model lengths chosen randomly between 145 and 170. The average run-time was around 60 CPU minutes on a Sun Sparcstation I. For each run we computed a

NLL-score for the model, which was the average of the NLL-scores for the training sequences, as defined in Methods section (b). The final NLL-scores varied considerably for these runs but the best was 210.7.

We then took this model, produced ten new models by adding noise, and optimized these. These models all generated approximately the same NLL-score and we picked the model with the best NLL-score, 210.3, having a length of 147. We validated this model† in two ways: from the alignments it produced, and by its ability to discriminate between globins and non-globins. The results are described below.

(i) Multiple sequence alignments

A multiple alignment of many globin sequences has been produced by Bashford *et al.* (1987) by including into the alignment procedure tertiary-structure information of seven globins (Fig. 4). This

† We stress that the final model was chosen according to an objective measure, namely the NLL-score on the training set, and not retroactively on the basis of how well it did in multiple alignment or database search tasks.

was achieved by aligning these seven sequences and then aligning the rest of the 226 studied to the closest of these seven. In contrast, generating multiple alignments with HMMs requires no prior knowledge of underlying structure. Using the globin HMM, we produced a multiple alignment of all the 625 globin sequences by the Viterbi algorithm as described in Methods section (c). Figure 5 shows this alignment for the seven sequences from Bashford *et al.* (1987).

The alignment found in this experiment agrees extremely well with the structurally derived alignment of Bashford *et al.* Our alignment differs in the region between the C and E helices. However, this is a highly variable area since only some globins possess a D helix. The difference in the F/G-helices is more pronounced, with the remaining discrepancies possibly representing an alternative alignment. Four of the insertions the model chose are in variable regions between or at the end of helices, i.e.

between secondary structure elements. The last two insertions appear in the F/G region.

(ii) *Database search: discriminating globins from non-globins*

The globin HMM model we found was also tested on all the 25,044 proteins in the SWISS-PROT database release 22.0 of length less than 5000 amino acids (which is all but 2). A NLL-score and a Z-score were computed for each of these sequences as described in Methods section (f). These are plotted in Figures 6 and 7 as a scatter plot and a histogram, respectively. For the histogram (but not the scatter plot), the data were filtered as follows:

All sequences with a Z-score >3.5 and either more than a total of 25, or more than 15% unknown residues were removed (a total of 23). Currently, we treat an unknown amino acid, X, as being the most probable amino acid at the position it is matched to,

```

Helix          AAAAAAAAAAAAAAAAAA  BBBB BBBB BBBB BBBB BCCCCCCCCCCC  DDDDDDEE
                *****
HBA_HUMAN  V.....LSPADKTNVKAAGKVGVA.. HAGEYGAELERMFLSFPTTKTYFPHF-DLSHGSAQ----
HBB_HUMAN  Vh.....LTPEEKSAVTALWGKV--. NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNP
MYG_PHYCA  V.....LSEGEWQLVHLVWAKVEA.. DVAGHGQDILIRLFKSHPETLEKFRDFKHLKTEAEMKASE
GLB3_CHITP -.....LSADQISTVQASFDKV--. KGDVPG--ILYAVFKADPSIMAKFTQF-AGKDLESIKGTA
GLB5_PETMA PivdtgsvapLSAAEKTIRSAWAPVYS.. TYETSGVDILVKFFTSTPAAQEFPKFKGLTTADQLKKSAA
LGB2_LUPLU Ga.....LTESQAALVKSSWEEFNA.. NIPKHTHRFFILVLEIAPAAKDLF-SFLKGTSEVPQ-NNP
GLB1_GLYDI G.....LSAAQRQVIAATWKDIAGadNGAGVGVKDCLIKFLSAHPQMAAVF-GF----SGASD---P

Helix          EEEEEEEEEEEEEEEEEEE  FFFFFFFF  FFFFFGGG  GGGGGGGGGGGGGGG
                *****
HBA_HUMAN  -VKGHGKVVADALTNVAHVDD...MPNALSALSDLHA...HKLRVDPV.NFKLLSHCLLVTLAAHLP
HBB_HUMAN  KVKAHGKVLGAFSDGLAHLDN...LKGTFATLSELHC...DKLHVDPE.NFRLLGNVLVCLVAHHFG
MYG_PHYCA  DLKKGHVTVLTALGAILKKKGH...HEAELKPLAQSHA...TK-HKIPIKYLEFISEAIIHVLHSRHP
GLB3_CHITP PFETHANRIVGFFSKIIGELPN...IEADVNTFVASHK...PR-GVTHD.QLNFRAGFVSYMKAH--
GLB5_PETMA DVRWHAERIINAVNDAVASMDDtek..MSMKLRDLSGKHA...KSFQVDPQ.YFKVLAAVIADTVAA---
LGB2_LUPLU ELQAHAGKVFKLVEAAIQLQVtgvvvTDATLKNLGSVHV...SK-GVADA.HFPVVKEAILKTIKEVVG
GLB1_GLYDI GVAALGAKVLAQIGVAVSHLGDegk..MVAQMKAVGVRHKgygNK-HIKAQ.YFEPLGASLLSAMEHRIG

Helix          HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH
                *****
HBA_HUMAN  AEFTPAVHASLDFKFLASVSTVLTISKY.....R
HBB_HUMAN  KEFTPPVQAAYQKVVAGVANALAHKY.....H
MYG_PHYCA  GDFGADAQGAMNKALELFRKDIAAKYkelgyqG
GLB3_CHITP TDF-AGAEAAWGATLDTFFGMIFSKM.....-
GLB5_PETMA GD-----AGFEKLMSMICILLRSAY.....-
LGB2_LUPLU AKWSEELNSAWTIIAYDELAIVIKEMnda...A
GLB1_GLYDI GKMNAAKDAWAAAYADISGALISGLq.....S

```

Figure 5. The alignment of the same 7 globins as in Fig. 4, as obtained from our model trained on 400 randomly chosen globin sequences. The capital letters represent amino acids aligned to the main line of the model, -, to deletions in the model, and lower-case letters to amino acids treated as insertions by the model. The . is used as a fill character to accommodate insertions. No attempt has been made to align the insertion regions. In the line above the alignments * indicates complete agreement of a column with the structural alignment (Fig. 4) and + denotes a minor deviation (the only accepted difference is a reasonable displacement of a gap). The regions between the helices are not checked in this way. The training set contained 5 of the 7 globins, not HBA_HUMAN and GLB5_PETMA.

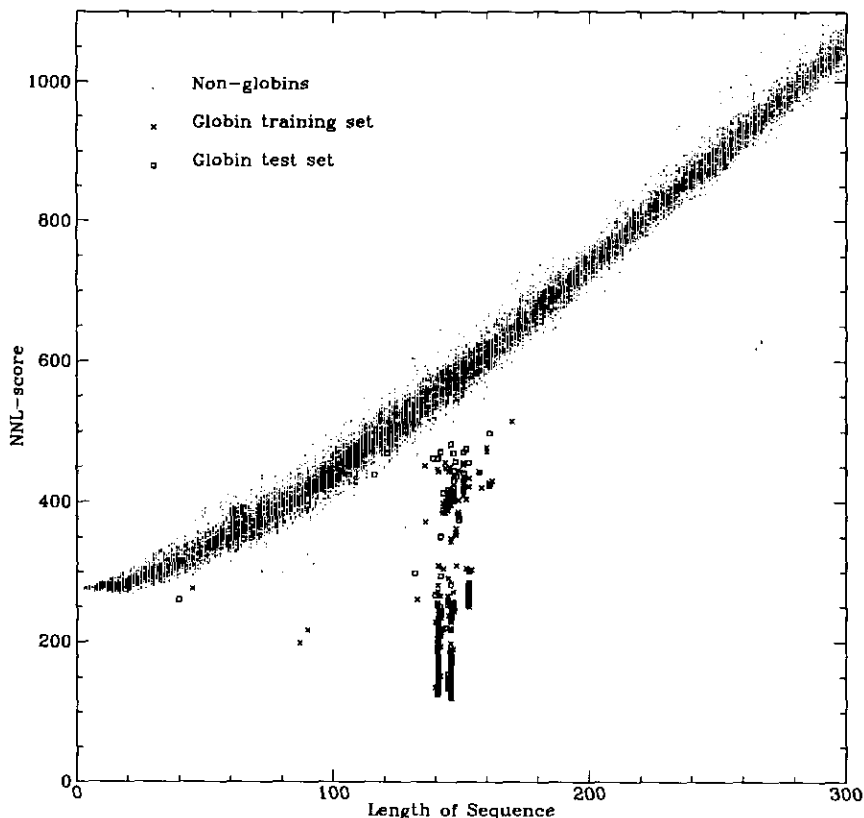


Figure 6. Plot of NLL-score versus sequence length for globins and non-globins. All sequences of length less than 300 from the SWISS-PROT 22 database are shown, including partial sequences and 3 false globins from the globin file, and sequences from the database containing many Xs.

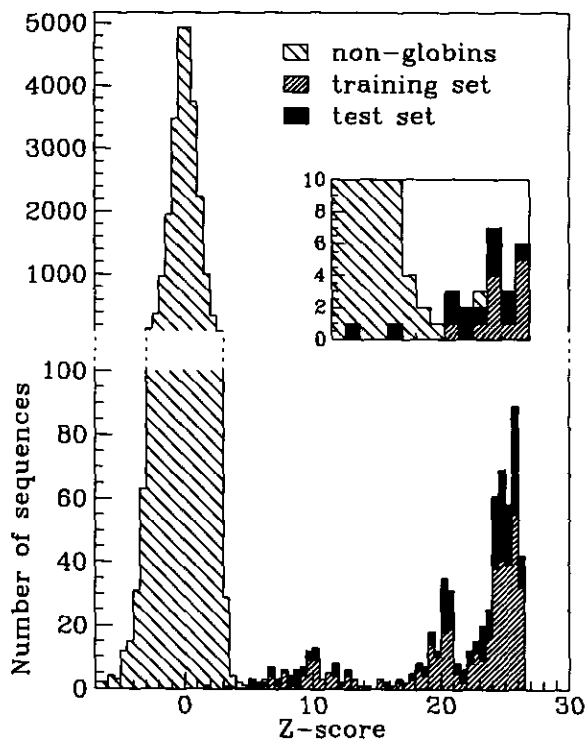


Figure 7. Histogram showing the number of sequences with a certain Z-score. The training set of 397 globins, the test set of 231 globins, and the rest of the sequences from SWISS-PROT 22 after "filtering" are shown. The insert shows expansion of the region around a Z-score of 5.

so sequences with many Xs spuriously match the model very well.

Since we searched a newer release of SWISS-PROT (release 22) than the one from which the globin training set was extracted (release 19), eight new globins were found and incorporated into the test set.

Five globin fragments of length 19 to 45 were removed from the data.

Three non-globin sequences in the globin file that were identified as outliers in Figure 6 were removed. One of these non-globins was left as part of the training set to illustrate the robustness of the method.

The model distinguishes extremely well between globins and non-globins. Choosing a Z-score cutoff of 5 we would miss 2 out of 628 globins† and get essentially no false positive globins. There is one "non-globin", a bacterial hemoglobin-like protein (SWISS-PROT id HMP_ECOLI), that may or may not be counted as a false positive. Only one sequence, the heme containing catalase of *Penicillium vitale* (CATA_PENVI, Z-score 4.7), has a Z-score between 4.2 and 5.1, so any cutoff in this range will essentially give the same separation. The two sequences falling between a Z-score of 1 and 4

† 628 in the original data set, plus 8 new, minus 3 spurious, minus 5 fragments. 397 were left from the training and the remaining 231 made up the test set.

(GLB_PARCA and GLB_TETPY) are protozoan, whereas the other globins are metazoan. The primary sequences of these globins are similar and have little similarity with other eukaryotic globins. Note also that both of these sequences are in the test set.

(iii) *Discovering subfamilies of globins*

We also performed an experiment to automatically discover subfamilies of globins using the method described in Methods section (d). An HMM with ten component HMMs was used. The initial lengths of the components were chosen randomly between 120 and 170, but were adjusted by model surgery during training. We trained this HMM on all 628 globins and then calculated the NLL-score for each sequence for each of the ten component HMMs. A sequence was classified as belonging to the cluster represented by the component HMM that gave the lowest NLL-score, i.e. the one giving the highest probability to that sequence.† Three of these clusters were empty and the remaining seven non-empty ones represented chains from known globin subfamilies:

Class 1. 233 sequences: principally all α , a few ζ (an α -type chain of mammalian embryonic hemoglobin), π/π' (the counterpart of the α chain in major early embryonic hemoglobin P), and θ -1 chains (early erythrocyte α -like).

Class 2. 232 sequences, almost all β , a few δ (β -like), ε (β -type found in early embryos), γ (comprise fetal hemoglobin F in combination with 2 α chains), ρ (major early embryonic β -type chain) and θ chains (embryonic β -type chain).

Class 3. 71 myoglobins.

Class 4. 58 sequences. The 13 highest scoring in this cluster are leghemoglobins. This class contains a variety of sequences including the three non-globins in original data set.

Class 5. 19 sequences. Midge globins.

Class 6. Eight sequences. Globins from agnatha (jawless fish).

Class 7. Seven sequences. Varied.

We have not repeated this experiment using different randomization to ascertain if better results can be obtained. However, we are encouraged by the results of this first experiment since it is able to classify correctly the major globin subfamilies (alpha, beta and myoglobin).

(iv) *The final globin model*

Examination of the model itself yields information on the structure of globins. Figure 8 shows the normalized frequency counts (the numbers used to re-estimate the parameters of the model) from some parts of the final model. The thickness of a line

indicates what fraction of the 400 training sequences made that transition or used that particular amino acid. A broken line indicates that less than 5% of the sequences used that transition. (The continued delete is mostly due to fragments that have to make many deletions.) The histogram in a match state shows the distribution of amino acids that were matched to that state. The number in an insert shows the average length of an insertion beginning at that position.

For the amino acids the ordering proposed by Taylor (1986) is used. Starting from the top, the amino acids are medium-sized and non-polar, small and medium polar (around G and P), medium sized and polar (around K), large medium-polar (around F and Y), and finally below they are medium-large and non-polar. There does seem to be some tendency for the distributions to peak around neighboring amino acids when using this ordering, as one would expect. When one looks at the whole model, regions that are highly conserved are also readily distinguished from the more variable regions, both as a function of the probability that a position is skipped, and the entropy of the distribution of amino acids at that position.

(b) *Kinase experiments*

Protein kinases are defined as enzymes that transfer a phosphate group from a phosphate donor onto an acceptor amino acid in a substrate protein (Hunter, 1991; Hanks *et al.*, 1988). Based upon the acceptor amino acid specificity, they have been classified into serine/threonine, tyrosine, histidine, cysteine, aspartyl and glutamyl kinases. Only enzymes in the first two categories have been well characterized and recent developments indicate that some can phosphorylate both alcohol (serine/threonine) and phenol (tyrosine) groups, the so-called dual-specificity protein kinases (Lindberg *et al.*, 1992). It is the region comprising the catalytic domain of these hydroxyamino acid phosphorylating enzymes that we model by an HMM and which we subsequently refer to as protein kinases or simply kinases. Despite the differences in size, substrate specificity, mechanism of activation, subunit composition and subcellular localization, all these kinases share a homologous catalytic core containing 12 conserved subdomains or regions (Hanks & Quinn, 1991; Hanks *et al.*, 1988).

Because the kinase catalytic domain is only a subsequence embedded in a larger protein, the kinase experiments differed from the globin experiments. The HMM used in the globin experiments modeled the entire protein rather than simply a segment of a protein as is the case for the kinase family. Modeling domains requires several modifications to our standard HMM training which are described in Methods section (e).

The training set for these experiments is a group of 193 sequences from the March 1992 release of the protein kinase catalytic domain database maintained by Hanks & Quinn (1991). This set is

† We can also calculate the posterior probability of each cluster by looking at the transition probabilities out of the global start state, and thereby obtaining a posterior distribution over the 10 clusters for each sequence. However, these posteriors are very sharply peaked, so this adds little to the analysis.

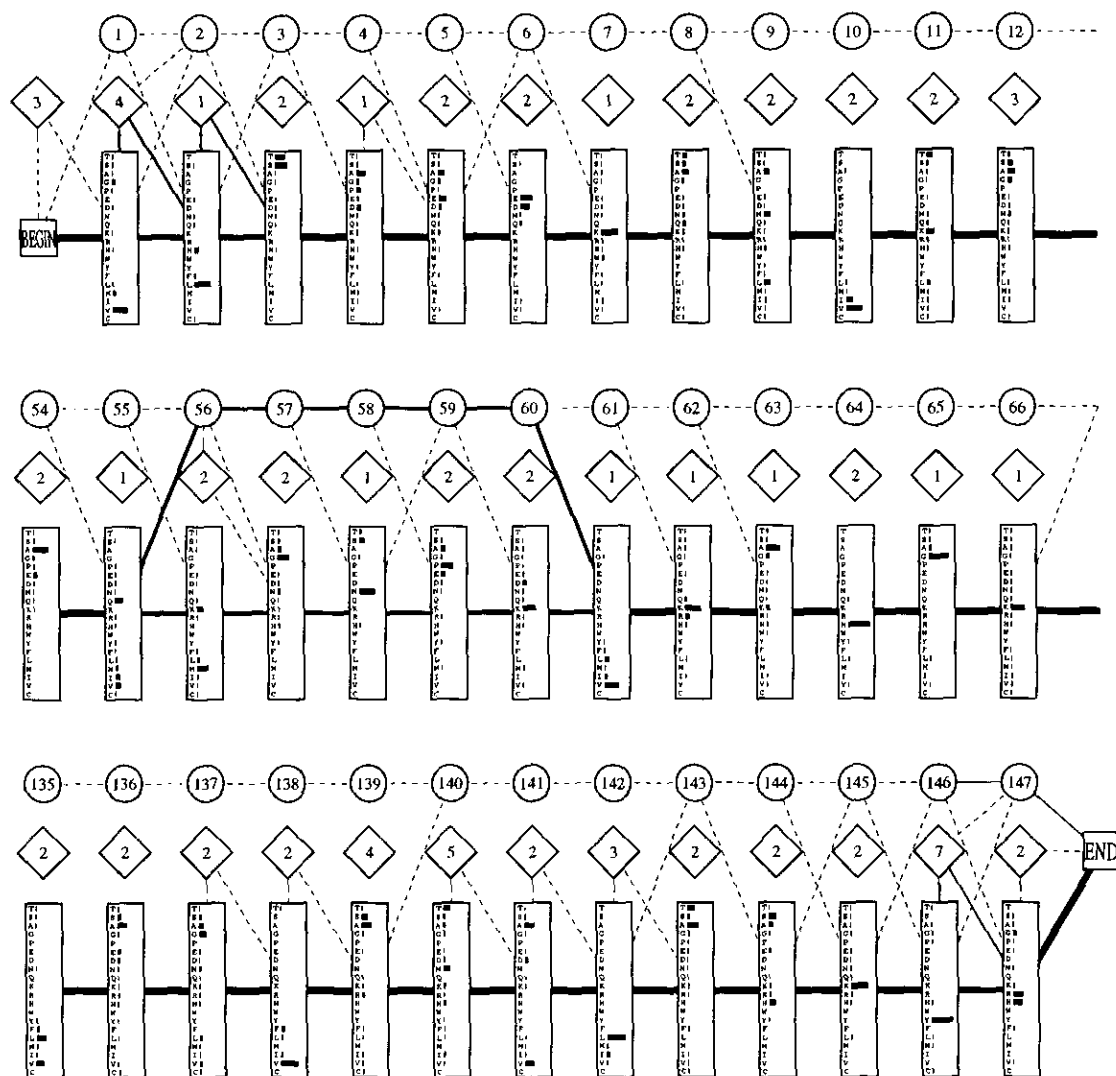


Figure 8. Parts of the final globin model. The position numbers are shown in the delete states.

composed of serine/threonine, tyrosine and dual-specificity kinases principally from vertebrates and higher eukaryotes but also includes some from lower eukaryotes and viruses.

We trained ten HMMs on all 193 (unaligned) sequences in this data set using the prior distributions described by Krogh *et al.* (1993a). No parameters of the modeling process were set manually and the initial model lengths ranged from 242 to 282 positions (this encompasses the average length of the sequences in our kinase catalytic domain training set). At the end of the ten training runs, the best kinase model had a NLL-score (the average $-\log P(\text{sequence}|\text{model})$ over the training set) of 588.39 and a length of 254. Modules were added at the beginning and end of this model as described in Methods section (e). We tested this model in the same manner as described earlier for the globin model.

Our main tests were discrimination tests, in which we utilized the model to search the SWISS-PROT version 22 database (25,044 sequences) for proteins containing the kinase catalytic domain.

As described in Methods section (f), a NLL-score was computed for each of the sequences in the database and this information was used to compute a sequence's deviation from the average curve as measured by a Z-score. The data were then filtered to remove all sequences with any unknown residues (353) and all sequences having length less than 200 (4230), since complete protein kinase catalytic domains range from 250 to 300 residues (Hanks *et al.*, 1988). This filtering removed a total of 4386 sequences. A scatter plot of NLL-score *versus* length for the SWISS-PROT sequences is given in Figure 9.

A cutoff of 6.0 was chosen because there are no sequences with Z-scores between 4.935 and 6.773. See Figure 10 for a histogram of the resulting Z-scores. Any sequence having a Z-score > 6.0 was therefore classified as containing the kinase catalytic domain while those with Z-score < 6.0 were classified as not possessing the domain. With this cutoff, 296 sequences were classified as containing the kinase catalytic domain. The remaining 20,357 sequences were rejected.

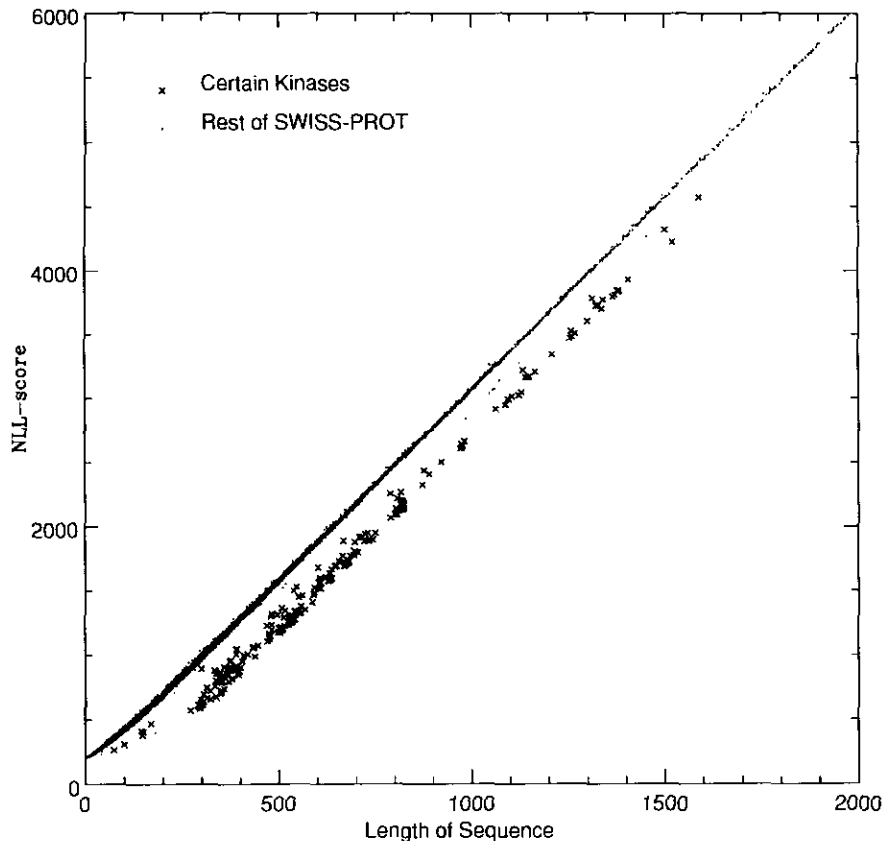


Figure 9. Scatter plot of NLL-score versus length for sequences in SWISS-PROT using the Kinase HMM.

The general issue of estimating the number of false negatives and false positives when distinguishing sequences belonging to a given family

from non-members is a complex one. In the case of the globins, it is "relatively" straightforward since it is possible to identify all the globins in the database by performing a keyword or title string search. The situation for the kinase domain or the EF-hand motif (see section (c) below) is less obvious and thus more problematic. For instance, while a given protein may possess the sequence characteristics for this motif or domain, functionally, the region may not bind calcium or possess kinase activity. We have attempted to address this complicated matter as best we can as described below. However, we stress that we do not feel able to give a definitive answer as to the number of true false negatives and true false positives in our kinase or EF-hand database discrimination tests.

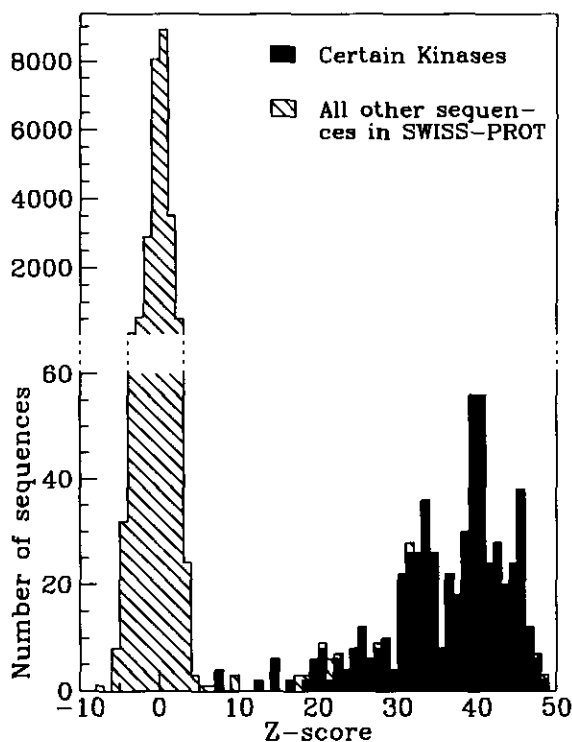


Figure 10. Histogram showing the number of sequences with a certain Z-score relative to the kinase model.

A list of potential protein kinases was created from the union of sequences designated as being kinases from four independent sources: our HMM, PROSITE (a dictionary of sites and patterns in proteins (Bairoch, 1992)), PROFILESEARCH (a technique used to search for relationships between a protein sequence and multiply aligned sequences (Gribskov *et al.*, 1990)) and a keyword search.

Two regions of the catalytic domain of eukaryotic protein kinases have been used to build PROSITE signature patterns. The first pattern corresponds to an area believed to be involved in ATP binding (PROSITE entry PROTEIN_KINASE_ATP, sequence motif [LIV]G.G.[FYM][SG].V). There are two signature patterns for the second region important for catalytic activity: one specific for serine/

Table with columns 1, 11, 21, 31, 41, 51, 71, 81, 91, 101, 111, 121, 131. It shows sequence alignments for various proteins including CAPK-ALPHA, WEE1+, TIK, SPR1, RSK1-N, PTT, PKC-ALPHA, PDGFR-B, P85, NIK1, NCK1, IKS.R, HSVK, ERK1, EGFR, ECK, DPK1, CLK, CDC2H5, CAK1P-ALPHA, C-SAC, CAF, KLSK_HUMAN, KLSK_MOUSE, ARK3_HUMAN, ARK3_BOVIN, BVR1_SCHPO, CTR_ARBP, ANPA_BAT, ANPA_HUMAN, ANPB_HUMAN, ANPB_MOUSE, ANPB_BAT, CYG5_STRPU, VPSF_YEAST, HSER_BAT, HSER_HUMAN, KR2_Y2VD, KR2_HSV11, RRI_HSV11, KR2_EBV, KR2_YACCV, KR2_YACCC, AK3_ECOLI, ASP_MOUSE, DHON_BACSV, FLIG_BACSV, CALQ_RABIT, NUIK_PODAN, migrlrg15LVLI, U1SR_HSV6U, KR7_YACCC, KR97_HCNVA, RKA6_JCIB1A, RKA6_ECOLI, KGPB_BOVIN, EGFR_CHICK, RKA1_ECOLI, RDKT_DROME, KPCC_HUMAN. Each row shows the sequence at positions 1, 11, 21, 31, 41, 51, 71, 81, 91, 101, 111, 121, 131.

A (cont)

Fig. 11.

threonine kinases (PROTEIN_KINASE_ST, [LIVMFYC].[HY].D[LIVMFY]K.2N[LIVMFYC]3) and the other for tyrosine kinases (PROTEIN_KINASE_TYR, [LIVMFYC].[HY].D[LIVMFY] [RSTA].2N[LIVMFYC]3). Since PROSITE expressions do not allow for flexible gapping or insertions, a profile of kinases was constructed from an alignment of seven kinases and employed for database discrimination tests (M. Gribskov, personal communication) using the program PROFILESEARCH (Gribskov et al., 1990). The seven kinases used to generate the profile are, bovine cAMP-dependent protein kinase (PIR code

OKBOG), bovine cGMP-dependent protein kinase (OKBO2C), bovine protein kinase C (KIBOC), human mos kinase-related transforming protein (TVHUF6), human ref-a kinase-related transforming protein (TVHUMS), mouse pim-1 kinase-related transforming protein (TMSP1), and human fes/fps kinase-related transforming protein (TVHUFF). The keyword search consisted of searching the descriptions of the sequences in SWISS-PROT for the following strings: "SERINE/THREONINE-PROTEIN KINASE, SER/THR-PROTEIN KINASE, PROTEIN-SERINE/THREONINE KINASE, PROTEIN-SER/THR

Subdomain	141	151	161	171	181	191	201	211	221	231	241	251	261	271
PROSITE	-----><--v----->-----<-vIa----->													
X-ray	BB	B	BBB	B	BB	BBB	BB	BBB	AAA	AA	AA	A	AA	AAA
X-ray	44	4	444	S	55	555	55	555	DDD	DD	DD	D	DD	EEE
1 CAPK-ALPHA	P	FL	V	KLR	F	SFKD	NSV	LY	XVN	EYVPGG	E	NFS	HL	RR
2 WEE1+	D	HL	V	ELR	D	SWEN	GGF	LY	MQV	ELCEGG	S	LDR	FL	EE
3 TIR	V	NI	V	QYHscwg	V	VDYD	PEHmddcary12LF	IQN	EFCDKG	T	LEQ	WN	RNR	N
4 SPK1	P	RI	V	RLR	G	EYED	TES	YY	XVN	EFVSGG	D	LWD	FV	AA
5 ASK1-N	P	FV	V	KLR	V	AFQT	EGR	LY	LIL	DPLKGG	D	LFT	KL	SA
6 PYT	D	HL	V	ELR	D	YEIT	DQY	LY	MCN	EYKGGG	S	LDR	FL	EE
7 PKC-ALPHA	P	FL	T	QLN	S	CFQT	VDR	LY	FVN	EYVGGG	D	LWY	RI	QQ
8 PDGFR-B	L	NY	V	RLN	G	ACTK	GGP	IY	IIT	EYKAGY	D	LVD	YL	RR
9 PBS2	P	VI	V	DFY	G	AFFI	EGA	VY	MCN	EYKGGG	S	LDR	FL	EE
10 NIK1	P	FV	V	MLV	N	VVST	MDN	IF	LQL	DYCEGG	D	LST	FL	SE
11 NCK1	P	NI	V	RLQ	V	FFTH	LSPqdk	VYqLAm	ECLP	E	T	LQI	EI	NRYvtM
12 INS-R	H	NY	V	ALL	G	VVSK	GGP	YL	VVN	ELKANG	D	LKS	YL	AS
13 HSYK	P	AI	L	PLL	D	LRVY	SGV	TC	LVL	PKYQA	D	LVT	YL	SR
14 ERK1	E	NV	I	GIR	D	ILMA	PTLeamrd	VY	IVQ	DLME-T	D	LYK	LL	KS
15 EGFR	P	HY	C	RLN	G	ICLT	S-T	VQ	LIT	QLKPPG	C	LID	VY	RE
16 ECK	H	NI	I	RLE	G	VISA	YKP	NN	IIT	EYKGGG	A	LQR	FL	RE
17 DPYK1	P	NV	V	QFL	G	ACTAggEDH		HC	IVT	EYVGGG	S	LQD	FL	TD
18 CLA	Phatz	RC	V	QML	E	VFEH	RGH	IC	IVF	ELL-GL	S	TYD	FI	RSa
19 CDC2HS	P	NI	V	SUQ	D	VLDQ	S-R	LY	LIF	EFLS-N	D	LKK	YL	DSip
20 CAX11-ALPHA	P	NI	V	RLH	D	SISE	EGH	HY	ILF	DLTGG	E	LFE	DI	VA
21 C-SRC	E	KL	V	QLY	A	VVSE	E-P	IV	IVT	EYKSGG	S	LQD	FL	KE
22 C-RAF	V	NI	L	LFN	G	YTKR	D-V	LA	IVT	QVCEDS	S	LYR	NL	HV
23 KLSK_HUMAN	Q	KL	V	RLY	A	VVTQ	E-P	IV	IIT	EYKGGG	S	LVD	FL	KTp
24 KLSK_MOUSE	P	RL	V	RLY	A	VVTQ	E-P	IV	IIT	EYKGGG	S	LVD	FL	KTp
25 ARRB_HUMAN	P	FI	V	CNS	Y	AFHT	PDR	LS	FIL	DLNMGG	D	LHY	HL	SQ
26 ARRB_BOVIN	P	FI	V	CNS	Y	AFHT	PDR	LS	FIL	DLNMGG	D	LHY	HL	SQ
27 BYR1_SCHPO	P	YI	V	GFY	G	AFQY	KHM	IS	LCM	EYKGGG	S	LDA	TL	RE
28 CYOR_ARBP	D	NI	C	PP1	G	ACID	APH	IC	ILN	HYCAGG	S	LQD	IL	EM
29 ANPA_RAT	E	HL	T	RFV	G	ACTD	PPM	IC	ILT	EYCPGG	S	LQD	IL	EM
30 ANPA_HUMAN	E	HL	T	RFV	G	ACTD	PPM	IC	ILT	EYCPGG	S	LQD	IL	EM
31 ANPB_HUMAN	N	HL	T	RPI	G	ACID	PPM	IC	IVT	EYCPGG	S	LQD	IL	EM
32 ANPB_MOUSE	E	HL	T	RFV	G	ACTD	PPM	IC	ILT	EYCPGG	S	LQD	IL	EM
33 ANPB_RAT	N	HL	T	RPI	G	ACID	PPM	IC	IVT	EYCPGG	S	LQD	IL	EM
34 CYGS_STAPU	D	NI	C	PP1	G	ACID	APH	IS	ILN	HYCAGG	S	LQD	IL	EM
35 VPSF_YEAST	P	NV	L	NYS	K	LLET	MRA	GY	ILR	QHLKM	N	LVD	ML	S
36 HSER_RAT	Dyy	HL	T	KFY	G	TVKL	DTK	IF	GVV	EYCEGG	S	LRE	VL	ND
37 HSER_HUMAN	V	HL	T	KFY	G	TVKL	DTK	IF	GVV	EYCEGG	S	LRE	VL	ND
38 KR2_VZVD	S	SI	V	CLL	G	FSLQ	TK	-Q	LIF	PAYD-N	D	NDE	YI	VR
39 KR2_HSV11	H	NI	I	APL	G	FSLQ	QR	-Q	IVF	PAYD-N	D	LGR	Y	YgqLAS
40 KR1_HSV11								-KS	AQ	NVt+eCCLYG	S	REG	YF	YR
41 KR2_EBV	K	AL	V	DYL	S	ACTSc	R-A	LF	R	qPFC	S	LQD	YG	HW
42 KR2_VACCV	dnvtr14L	AI	P	DLY	G	IGET	DDY	NF	FV	--IKN	--LG	RV	FA	P
43 KR2_VACCC	dnvtr14L	AI	P	DLY	G	IGET	DDY	NF	FV	--IKN	--LG	RV	FA	P
44 AK3_ECOLI	P	NV	I	reeieALL	E	NTIV	LAEaaalata	PA	LTD	ELVSHG	E	LKStllfveIL	RE	R
45 PSP_MOUSE														
46 DHON_BACSU														
47 FLIG_BACSU	Qkhd													
48 CALQ_RABIT														
49 NVA_PODAN	P	NP	V	GY	G	LLQAF	ADA	LK	LLLK	EYVA-P	T	QAM	II	LFS
50 NVA_ECOLI														
51 U15L_HSV6U														
52 KRF1_VACCC	Avineni	NV	I	NYP		HrmdhPFEH	EKRtneyerg	NI	ITF	PLALY	S	ADR	VD	TE
53 UL97_HCRVA	S	GL	I	RTR	A	AGEQ	QQP	PS	LV	--GTG	--V	S	R	G
54 KRAS_ACIBA	L	KV	P	ELI	N	TFDD	EQP	EF	NIT	KAIMA	--R	PI	SA	L
55 KRAS_ECOLI														
56 KGPB_BOVIN														
57 EGFR_CHICK	htel32													
58 KRA1_ECOLI														
59 XDTX_DROME	Vunkm													
60 KPCG_HUMAN	P													

A (cont)

Fig. 11.

KINASE, TYROSINE-PROTEIN KINASE, TYR-PROTEIN KINASE, PROTEIN-TYROSINE KINASE, PROTEIN-TYR KINASE, V-ABL, C-ABL, V-FGR, C-FGR, V-FMS, C-FMS, V-FPS/FES, V-FES/FPS, C-FPS/FES, C-FES/FPS, V-FYN, C-FYN, V-KIT, C-KIT, V-ROS, C-ROS, V-SEA, C-SEA, V-SRC, C-SRC, V-YES, C-YES, V-ERBB”.

Of the 296 SWISS-PROT 22 sequences that were above the Z-score cutoff of 6.0 and were thus classified as containing a kinase domain by our HMM, 278 were similarly classified by PROSITE, PROFILESEARCH and the keyword search. These 278

sequences may be considered to constitute “certain kinases”. Figure 11 shows the multiple sequence alignment generated by our HMM of some representative kinases from this set (sequences 1 to 22). Sequences 23 to 40 are the 18 sequences (296 minus 278) that were designated as kinases by the HMM and one or two of the three other methods. For PROSITE, we consider a sequence to be a kinase if it satisfies one or more of the three patterns PROTEIN_KINASE_ATP, PROTEIN_KINASE_ST or PROTEIN_KINASE_TYR as a true positive (“T” in Fig. 11B). PROSITE false negatives (“N”), potential hits (“P”) and false positives (“F”,

Subdomain	281	291	301	311	321	331	341	351	361	371	381	391	401	411
PROSITE	-----> <-VII- - - - -> <-VII- - - - ->													
X-ray	AAA	AAA	AAAAA	A	BB	BB	B	BB	B	BBB	B	BBB	B	BBB
X-ray	EEE	EEE	EEEEE	E	66	66	7	77	7	8	8	888	8	999
1 CAPA-ALPHA	YAA	Q	IYL	TFEYL	H	S	L	DL	TYADLKP	N	L	L	I	D
2 WEI1+	ILV	E	VAL	CLQFI	H	H	K	NY	VHDLKPA	N	V	M	I	T
3 TIK	LYE	Q	IVT	GVVEYI	H	S	K	GL	IHRDLKPG	N	I	F	L	V
4 SPR1	ISR	Q	ILT	ATKVI	H	S	N	GI	SHRDLKP	N	I	L	I	E
5 RSK1-N	YLA	E	LAL	GLDM	H	S	L	GI	TYRDLKPE	N	I	L	I	D
6 PTT	YWA	N	MLE	AVNTI	H	Q	R	GI	IHRDLKPG	N	I	F	L	V
7 PKC-ALPHA	YAA	E	ISI	GLFFL	H	K	R	GI	TYRDLKD	N	V	M	I	D
8 PDGFR-B	FSY	Q	VAN	GNEFL	A	S	K	WC	VHRDLAAR	N	V	L	I	C
9 P8S2	IAN	A	VIK	GLKEL	E	Q	R	NI	IHRDLVAPT	N	I	L	I	D
10 KIK1	NLP	Q	LTO	ALMFI	H	L	L	EP	VHLDVAPS	N	V	L	I	T
11 NCR1	YTV	Q	IAR	GMLY	H	G	L	GV	CHRDIPS	N	V	L	V	D
12 INE-R	NAA	E	IAD	GRAYL	H	A	K	KP	IHRDLKPA	N	I	L	I	D
13 HSK	VSR	Q	LLS	AVDYI	H	R	Q	CI	IHRDIAT	N	I	F	L	V
14 ERK1	FLY	Q	ILR	GLKVI	H	S	A	NV	IHRDLKPS	N	L	I	L	D
15 EGFR	WCV	Q	IAR	GMVY	E	D	R	KL	VHRDLAAR	N	V	L	I	D
16 ECK	XLR	G	IAR	GKYL	E	N	K	NY	VHRDLAAR	N	I	L	I	D
17 DPK1	LAL	D	IAR	GMVY	L	H	G	TP	IHRDLKPS	N	L	I	L	D
18 CLK	NAY	Q	ICK	GMVY	L	S	N	KL	THDLKPE	N	I	L	I	D
19 CD2HS	YLY	Q	ILQ	GIVCF	H	S	R	NY	IHRDLKPA	N	L	I	L	D
20 CAM1-ALPHA	CIQ	Q	ILE	AVLHC	H	Q	K	GV	VHDLKPA	N	L	I	L	D
21 G-SRC	NAA	Q	IAS	GRAYE	E	A	K	NY	VHRDLAAA	N	I	L	I	D
22 C-RAF	IAR	Q	TAQ	GNDY	L	A	K	NI	IHRDKSN	N	I	F	L	V
23 KLSK_HUMAN	NAA	Q	IAC	GRAFI	E	E	R	NY	IHRDLKPA	N	I	L	I	D
24 KLSK_MOUSE	NAA	Q	IAC	GRAFI	E	E	Q	NY	IHRDLKPA	N	I	L	I	D
25 ARAB_HUMAN	YAA	E	ITL	GLENN	H	N	R	FV	VYRDLKPA	N	I	L	I	D
26 ARAB_BOVIN	YAA	E	ITL	GLENN	H	N	R	FV	VYRDLKPA	N	I	L	I	D
27 BYR1_SCHPO	IIM	S	KVR	GLVLY	H	N	V	L	IHRDLKPS	N	V	V	N	S
28 CYG1_ARBPB	LIA	D	ILK	GLVY	H	S	S	EX	SHGLKSS	N	C	V	V	D
29 ANPA_RAT	LTM	D	IVR	GMLFL	H	G	S	GL	SHGLKSS	N	C	V	V	D
30 ANPA_HUMAN	LTM	D	IVR	GMLFL	H	G	S	GL	SHGLKSS	N	C	V	V	D
31 ANPB_HUMAN	LIM	D	IVR	GRAFL	H	N	S	IL	SHGLKSS	N	C	V	V	D
32 ANPA_MOUSE	LTM	D	IVR	GMLFL	H	G	S	GL	SHGLKSS	N	C	V	V	D
33 ANPB_RAT	LTM	D	IVR	GRAFL	H	N	S	IL	SHGLKSS	N	C	V	V	D
34 CYG2_STRPB	LIA	D	IVR	GMLFL	H	S	S	EX	SHGLKSS	N	C	V	V	D
35 YP5F_YEAST	IAP	Q	LLW	ALKDI	H	N	L	NI	VHGDIAT	N	I	L	I	D
36 HSER_RAT	VLM	D	IAR	GSYL	H	S	S	KI	VHGDLKST	N	C	V	V	D
37 HSER_HUMAN	VLY	D	IAR	GSYL	H	S	S	KI	VHGDLKST	N	C	V	V	D
38 KR2_VZVD	VFL	D	LAQ	ALFTL	H	R	C	GL	THLDVKG	N	I	L	I	D
39 KR2_HSV1J	CFT	E	LAR	AVVFL	H	T	C	GI	SHLDIACA	N	I	L	I	D
40 KR1_HSV1J	ILA	D	LTO	KLAL	I	R	K	GI	VHGLKASE	N	I	L	I	D
41 KR2_EBV	GFQ	G	LKD	AVYFL	H	R	H	GL	FHSDISPS	N	I	L	I	D
42 KR2_VACCV	-CV	T	KIK	TLEFI	H	S	Q	GF	THGKIERN	N	I	L	I	D
43 KR2_VACCC	-CV	T	KIK	TLEFI	H	S	R	GF	THGKIERN	N	I	L	I	D
44 AK3_ECOLI	VNR	ndr	grae	pa	IA	AL	ABL	A	A	L	Q	L	P	R
45 PSP_MOUSE	FTS	Q	ILW	GL	ILW	GL	ILW	GL	ILW	GL	ILW	GL	ILW	GL
46 QH02_BACSU	NH	G	IVN	G	IVN	G	IVN	G	IVN	G	IVN	G	IVN	G
47 FLIG_BACSU	YAR	Q	VLE	---	KAL	G	Ed	---	RAE	N	IL	N	L	T
48 CALQ_RABIT	LAA	Q	VLE	---	KAL	G	Ed	---	RAE	N	IL	N	L	T
49 NUN1_PODAN	YAY	Q	VLE	---	KAL	G	Ed	---	RAE	N	IL	N	L	T
50 RUYA_ECOLI	LAL	A	ILS	GNS	---	A	Q	QF	---	V	---	---	---	---
51 U1SR_HSV6U	VNW	---	yyev	cklad	AVVFL	H	L	K	NI	NHPDISP	N	I	L	I
52 RNF1_VACCC	YIK	---	IFL	qall	YI	ki	yp	cc	---	NF	I	H	AD	L
53 UL97_RCNVA	CT	---	LAD	---	A	I	K	FN	---	Q	C	R	V	---
54 KKA6_CICBA	IYK	---	ALN	---	LLMSI	---	A	i	---	A	i	---	A	i
55 KKA6_ECOLI	TAY	---	VLR	---	---	---	---	---	---	S	D	---	---	---
56 KGPB_BOVIN	---	---	---	---	---	---	---	---	---	---	---	---	---	---
57 EGFR_CHICK	---	---	---	---	---	---	---	---	---	---	---	---	---	---
58 KKA1_ECOLI	TAF	---	VLE	---	---	---	---	---	---	---	---	---	---	---
59 RDK1_DRONE	---	---	---	---	---	---	---	---	---	---	---	---	---	---
60 KPC6_HUMAN	---	---	---	---	---	---	---	---	---	---	---	---	---	---

A (cont)
Fig. 11.

sequences which do not belong to the set under consideration) are ignored. Among the 18 sequences classified as kinases by our HMM, eight (23 to 26, 35, 38 to 40) were also deemed to be kinases by the keyword search and PROSITE, and one (27) by PROFILESEARCH and PROSITE. The remainder (28 to 34, 36 to 37, those indicated by % in Fig. 11B) are particulate guanylyl cyclases and except for 36 to 37, PROFILESEARCH also defines them as possessing a kinase domain. These guanylyl cyclases contain a single transmembrane domain, a cyclase catalytic domain and an intracellular protein kinase-like

domain in which protein kinase activity has not been seen to date (reviewed by Garbers, 1992). Although these sequences are not kinases in terms of function, they possess all the conserved subdomains (subdomain I, the nucleotide binding loop is modified in some) and the majority of conserved residues present in certain kinases (see Subdomain of Fig. 11A and positions indicated by *). Sequences 41 to 50 are the top ten sequences in SWISS-PROT immediately below our cutoff of 6.0. Of these, the first three (41 to 43) were classified as kinases by two out of PROSITE, PROFILESEARCH and the keyword search. Our cutoff was

Subdomain	421	431	441	451	461	471	481	491	501	511	521	531	541	551
PROSITE>C-X.....>C-X.....													
λ-rayA.AAAAAAAAAAAAAA.....A.A													
λ-rayF.FFFFFFFFFFFFFFFF.....G.G													
1 CAPK-ALPHAWTLGGT.P.EY.LAPE.IIL.....SK.....G-YNR.A.VDWALGVLTVEAA.G.....YP.P.....F.....F.A.A.-DQP.....I.....Q													
2 WEE1+MEAECC.Q.EY.IAPE.VLA.....XH.....L-YDK.P.ADIFSLGIVTVEAAh1.....VL.Pdqgsgvq17.....P.....R.LSST.....D.....W													
3 TIKTARTGT.L.QY.NSPE.QLF.....LK.....H-YGK.E.VDIFALGLILAELL.....HT.C.....F.....T.E.-SEK.....I.....K													
4 SPR1KTFGGT.L.AY.VAPE.VIR.....GKtvsdpe12MEYSS.L.VDWVSGCLVYVILT.G.....HL.P.....F.....S.G.-STQ.....D.....Q													
5 ASK1-WYSFGGT.V.EY.KAPE.VVW.....RQ.....G-RHN.S.ADVWSGVLN.....S.....G.....G.KDRK.....E.....T													
6 PTTDSQVGT.V.WY.NPPE.AIKmsarengKSR.....SKISP.K.SDVWSGCLLYVYVT.G.....KT.P.....F.....Q.....Q.LINQI.....S.....K													
7 PAC-ALPHARTFGGT.P.DY.IAPE.IIA.....YQ.....P-YGK.S.VDWVSGVLLYEMLA.G.....QP.P.....F.....D.....G.-EDE.....D.....E													
8 PDGFR-BKGSTFL.P.LK.WNAP.ESI.....FW.....SLYTT.L.SDVWSGFIILWEIFTG.....GT.P.....Y.....P.E.LPNR.....E.....Q													
9 P8S2KTFGGT.Q.SY.KAPE.RIR.....SLppdr.....ATYTV.Q.SDVWSGCLILEMAL.G.....RY.P.....Y.....P.P.-ETY.....DnifeQ													
10 RIK1D-LEGD.R.VY.IAPE.ILA.....SH.....M-YGK.P.ADVYSLGSLKRIEATv.....VL.Pengvewr16L.....P.....M.LKDL.....L.....L													
11 RCK1ISYICG.R.FY.KAPE.LIL.....GC.....TQVTT.Q.TDVLGCGVNGMLI.G.....KA.T.....F.....Q.G.QEPL.....L.....Q													
12 IMS.KGGAGLL.P.VR.WNAP.ESL.....KD.....GVFTT.S.SDVWSGVLWEITSIA.....EQ.P.....Q.....G.LSWE.....L.....V													
13 H5VAYGIAGT.I.DT.KAPE.VLA.....GD.....P-YTT.T.VDWSAGLVIFETAVH.....AS.L.....FsaaprpAA.....G.-PCD.....S.....Q													
14 ERA1TEYVAT.R.VY.KAPE.IIL.....MS.....KGYTK.S.DVWSGCLIAEMLS.N.....RP.I.....F.....P.G.KHVL.....D.....Q													
15 DGFRE-GKV.P.IK.WNAL.ESI.....LN.....RIYTH.Q.SDVWSGVTWELNTHG.....SK.P.....Y.....D.....G.IPAS.....E.....I													
16 ECKTSGGKI.P.TR.WTAP.EA1.....SY.....RKFTS.A.SDVWSGIVRWEVYTG.....EN.P.....Y.....W.E.LSMH.....E.....V													
17 DPKY1TOSVGC.I.PY.KAPE.VFK.....GD.....S-NSE.K.SDVWSGCVLVEFLT.S.....DE.P.....Q.....D.....D.KAPN.....K.....K													
18 CLXSTLVST.R.WY.KAPE.VIL.....AL.....G-WSQ.P.CDVWSGCLILEYVL.G.....FT.V.....F.....S.....T.HDSR.....E.....H													
19 CDC2HSTHEVYT.L.WY.NSPE.VLL.....GS.....ARYST.P.VDWSIGTIFAELAT.K.....KP.L.....F.....W.G.DSEI.....D.....Q													
20 CAM11-ALPHAFGFAGT.P.GY.LSPE.VLR.....XD.....P-YGK.P.VDWSAGVLLYVILV.G.....YP.P.....F.....W.D.-EDQ.....H.....R													
21 C-SACRQGAKF.P.IK.WTAP.EAA.....LY.....GRFTI.K.SDVWSGFIILTEITKG.....RV.P.....Y.....P.....G.MVNR.....E.....V													
22 C-RAFEQPTGS.V.LW.KAPE.VIR.....HQdan.....P-FSF.Q.SDVWSGIVLVEINT.G.....EL.P.....Y.....S.H.IWNR.....D.....Q													
23 KLSA_HUMANREGAKF.P.IK.WTAP.EA1.....NY.....GTFYI.K.SDVWSGILLTEIVTHG.....RI.P.....Y.....P.G.-MTP.....E.....V													
24 KLSA_MOUSEREGAKF.P.IK.WTAP.EA1.....WY.....GTFYI.K.SDVWSGILLTEIVTHG.....RI.P.....Y.....P.G.-MTP.....E.....V													
25 ARKB_HUMANHASVGT.H.GY.KAPE.VLQ.....KG.....VAYDS.S.ADVWSGCLIFALLR.G.....HS.P.....F.....R.Q.HATK.....Dkh.E													
26 ARKB_BOVINHASVGT.H.GY.KAPE.VLQ.....KG.....VAYDS.S.ADVWSGCLIFALLR.G.....HS.P.....F.....R.Q.HATK.....Dkh.E													
27 BYR1_SCRPOQTFVGT.S.TY.NSPE.RIR.....GG.....K-YTV.K.SDVWSGILSIELAT.Q.....EL.Pws.....F.....S.M.I.DDSigil.....D.....I													
28 CYGLARBP0GENKLA.A.KKAWTAP.EHL.....REgknhp.....G-GTP.K.SDVWSGFIILTEYVS.R.....QE.P.....F.....Hen.D.-LELA.....D.....I													
29 ANPA_BATTLFAKR.L.--WTAP.ELlrmapp.....AR.....G--SQ.A.CDVWSGFIILQEIAL.S.....GV.P.....F.....Vreg.....L.....D.LSPR.....E.....I													
30 ANPA_HUMANY-ARKL.--WTAP.ELL.....NAappv.....K-GSQ.A.CDVWSGFIILQEIAL.R.....SG.V.....F.....Hvegid.LSPR.....E.....I													
31 ANPB_HUMANY-ARKL.L.WT.-APE.LLS.....GB.....PLPTTgmK.ADVWSGFIILQEIAL.R.....SG.P.....F.....Ylegid.LSPR.....E.....I													
32 ANPA_MOUSETLFAKR.L.--WTAP.ELlrmapp.....AR.....G--SQ.A.CDVWSGFIILQEIAL.S.....GV.F.....F.....Vreg.....L.....D.LSPR.....E.....I													
33 ANPB_BATY-ARKL.L.WT.-APE.LLS.....GW.....PLPTTgmK.ADVWSFIILQEIAL.R.....SG.....F.....Ylegid.LSPR.....E.....I													
34 CYGS_STRPOGDHAKL.R.GL.WTSP.EHL.....Kdegampta.....G-SP.Q.GDYVSAFIIITELYS.R.....QE.P.....F.....H.EneDLA.....D.....I													
35 VPSF_YEASTf12ytd..TSAKAT.--CY.LAPE.NFB.....SKlydgkann.GRLTK.....E.DVWSGCVIAEJFaeC.....RP.J.....F.....-ML.....S.....Q													
36 HSER_BATKDL-----WTAP.EHL.....RQ.....ATISQ.....K.GELYSFSIIAQEIL.R.....KE.T.....F.....Y.T.LSCR.....D.....Qnek													
37 HSER_HUMANKDL-----WTAP.EHL.....RQ.....ANISQ.....K.CDVWSGIIAQEIL.R.....KE.T.....F.....Y.T.LSCRdrn.....E.....K													
38 RR2_VZVDFALVLS.H.FT.WQPP.EIL.....LDyngz:glr15QRVGL.A.IDLVALGQALIEVILJG.....RL.Pgq1pavb119Y.....Y.....G.MALS.....P.....D													
39 RR2_HSV11qfc1qe28HTLVG.H.GY.WQPP.ELIvkhylane+15LK.....HDVGL.....A.VDLVYALGQTLLELYVv.....YV.Apalgypvtr.F.....P.....G.....													
40 RR1_HSV11WNPICG.E.AY.NSPE.NSR.....DRvprpdsal2GTNGA.....G.I-----RE.....P.....Hli.....K.....G.DCVR.....A.....H													
41 RR2_EBVXSSKGA.Q.LY.R.-L.YCQ.....KE.....P-PSI.A.KDYV.....AP.Lcl1skcy124.....G.AQT4.....L.....A													
42 KRB2_VACCVyeadmi19NHLGAT.V.SR.RGDL.EHL.....GY.....C-----NIEWFG.G.....KL.P.....W.....KME.....S.....S													
43 KRB2_VACCCyeadmi19NHLGAT.V.SR.RGDL.EHL.....GY.....C-----NIEWFG.G.....KL.P.....W.....KME.....S.....S													
44 AK3_ECOLIDPRVVS.A.AR.RDPE.IAF.....AE.....A-----AEM77Gakv1bpA7.....L.....P.....A.VKS.....D.....J													
45 PSP_MOUSECIDLT.L.VP.LAGE.ASL.....VL.....PFIGK.T.VD1-5VSLDLIMSL.S.I.....KT.....Naq1g1pev14.....SWT.....D.....K													
46 DHOM_BACSUDVEGLD.A.AR.KNA-.ILA.....RL.....G-FSN.M.VDLE.....dvkvhg1S.....Q.ITDE.....D.....I													
47 FL1G_BACSUiqqebpq.7-----NAL.TLSy.....LD.....PVQ.....AGQILSELR.....F.....V.....VQA.....E.....V													
48 CALQ_RANITKED-----E.VIE.....YD.....GEFSad.....TLVEFL.....C.....LDVL.....E.....D													
49 WU1_PDBAGSLAST.A.GL.ISYE.LVL.....SS.....A-----ILLVRLT.G.....SL.....Nlevnieqi4F.....P.....L.-LPV.....F.....I													
50 ANVA_ECOLIGALVXL.Pg13K.KTAE.RL1vekmdrfk11GD.....L-FTP.A.ADL-----VLTSpA.....SP.....A.....D.AEZE.....A.....V													
51 V15R_HSV6UYRDACC-----K.VLAehvLL.....GL.....L-----VLTSpA.....F.....Y.....R.DVV.....E.....I													
52 KRF1_VACCVS-----A.....LN.....D-PDF.S.....A-----LLVRLT.G.....SL.....Nlevnieqi4F.....P.....L.-LPV.....F.....I													
53 VL97_HCRVAFPVAGL.R.VY.NSPE.LSA.....LG.....WVLGF.....CLN-----WVLGF.....R.LLDR.....R.....G													
54 KRA6_ACIBADD-----D.....DDPT.....E.....L.....W.....C.....D.HKT.V.....LelnE.....													
55 KRA6_ECOLIr1hai18TTHAGL.P.ER.GSIE.AGVvdddfdk.....RE.....G-WTA.Eq-----VWEANR.....LL.P.....L.....A.P.DPVV.....T.....H													
56 KGP_BOVINKQAST.....LQ.....GE.....P-RTK.R.....E.....A.LSAE.....P.....T.....T.....													
57 EGFR_CHICKQGLAEL.P.KK.NLSE.IIL.....GG.....VKIS.....P.....V.....WPK1cmdT.....V.....													
58 KRA1_ECOLIDALAVF.L.RR.L-----RST.P.V-----cncP.....F.....N.....S.DRV.....F.....R													
59 KDK_DRONENGSGG.....AR.....SQ.....G.....GA.....P.....T.....S.....G.-SCP.....W.....Qbege1p													
60 RPOC_HUMANKOKTR1-----VR.....AT.....L.....RPPVn.....E.....T													

A (cont)

Fig. 11.

chosen from a visual inspection of a histogram of Z-scores which indicated that 60 lay in a large gap (see Fig. 10). If the Z-score cutoff is lowered to the next largest gap (from Z-score 3.9 to 4.8) between sequences 43 to 44, then these three viral sequences (41 to 43) would also be categorized as kinases by the HMM.

Of the eight sequences (41, 51 to 53, 56 to 57, 59 to 60) that were not classified as kinases by our HMM but were classified only by the keyword search and PROSITE, one (41) is the first sequence below our cutoff discussed above. Four (56 to 57, 59 to 60) are partial sequences where the kinase

domain is absent. Three (51 to 53) possess divergent forms of many of the conserved regions and like 41 to 43, although they are below our cutoff, the HMM is able to generate an alignment that correctly identifies divergent forms of conserved regions. Finally, there are three aminoglycoside 3'-phosphotransferase sequences (54 to 55, 58) which are only designated as kinases because they satisfy the PROSITE expression for the catalytic loop.

Inspection of Figure 11B permits an estimation of the accuracy of the various methods in distinguishing kinases from non-kinases in database discrimination tests. The HMM generates six false

ID	Length	NLL-score	Z-score	HMM	PROFILE-SEARCH	Keyword	PROSITE		
							A	B1	B2
23 KLSK_HUMAN	509	1188.032	48.056	+	-\$	+	T	-	T
24 KLSK_MOUSE	509	1193.879	47.376	+	-\$	+	T	-	T
25 ARKB_HUMAN	689	1826.919	31.781	+	-\$	+	*	*	-
26 ARKB_BOVIN	689	1827.514	31.720	+	-\$	+	*	*	-
27 BYR1_SCHPO	340	808.153	27.540	+	+	-	N	T	-
28 CYGR_ARBPU	986	2839.392	22.121	+	+	-	-%	-	-
29 ANPA_RAT	1057	3062.107	21.418	+	+	-	-%	-	-
30 ANPA_HUMAN	1061	3072.615	21.390	+	+	-	-%	-	-
31 NPB_HUMAN	1047	3033.232	21.220	+	+	-	-%	-	-
32 ANPA_MOUSE	1057	3065.181	21.042	+	+	-	-%	-	-
33 ANPB_RAT	1047	3038.053	20.633	+	+	-	-%	-	-
34 CYGS_STRPU	1125	3277.621	18.745	+	+	-	-%	-	-
35 VPSF_YEAST	1454	4263.173	17.896	+	-	+	N	T	-
36 HSER_RAT	1075	3143.529	17.681	+	-	-	-%	-	-
37 HSER_HUMAN	1073	3139.039	17.552	+	-	-	-%	-	-
38 KR2_VZVD	510	1521.597	9.615	+	-	+	N	T	-
39 KR2_HSV11	518	1548.949	9.042	+	-	+	N	-	-
40 KR1_HSV11	230	710.448	6.773	+	-\$	+	N	T	-
41 KR2_EBV	455	1393.761	4.935	-	-	+	T	-	T
42 KRB2_VACCV	283	880.650	4.848	-	+	+	N	N	-
43 KRB2_VACCC	283	880.753	4.838	-	+	+	N	N	-
44 AK3_ECOLI	449	1385.412	3.900	-	-	-	-	-	-
45 PSP_MOUSE	235	754.545	3.804	-	-	-	-	-	-
46 DHOM_BACSU	433	1340.413	3.706	-	-	-	-	-	-
47 FLIG_BACSU	338	1055.096	3.699	-	-	-	-	-	-
48 CALQ_RABIT	395	1229.120	3.487	-	-	-	-	-	-
49 NU1M_PODAN	368	1149.759	3.415	-	-	-	-	-	-
50 RUVA_ECOLI	203	667.519	3.413	-	-	-	-	-	-
51 U15R_HSV6U	562	1728.770	3.171	-	-\$	+	T	-	T
52 KRF1_VACCC	439	1366.011	2.900	-	-\$	+	N	T	-
53 UL97_HCMVA	707	2165.296	2.854	-	-\$	+	N	-	T
54 KKA6_ACIBA	259	838.469	2.370	-	-	-	-	-	T
55 KKA8_ECOLI	271	885.548	1.182	-	-	-	-	-	T
56 KGPB_BOVIN	293	953.735	0.684	-	-	+	P	P	-
57 EGFR_CHICK	703	2179.703	0.065	-	-	+	P	-	P
58 KKA1_ECOLI	271	902.461	-0.467	-	-	-	-	T	-
59 KDTK_DROME	753	2334.760	-0.523	-	-	+	N	-	N
60 KPCG_HUMAN	318	1051.016	-1.486	-	-	+	P	P	-

B

Figure 11. A, Multiple sequence alignment generated by our kinase HMM of some of the sequences used to train the HMM (1 to 22) and test sequences from the SWISS-PROT 22 database (23 to 60) (see Results section (b)). Numerals appearing in the alignments indicate the number of amino acids to be inserted at that point, otherwise the notation follows the convention of Fig. 5. In Subdomain, the Roman numerals and * refer to the subdomains and residues conserved across 75 serine/threonine kinases given by Hanks & Quinn (1991). A and B in PROSITE refer to the ATP binding and catalytic regions, respectively, used to create 2 different signature patterns for kinases. X-ray identifies the location of the α -helices AA-AI and β -strands B1-B9 (read vertically) derived from the 2.7 Å crystal structure of the catalytic subunit of cAMP-dependent protein kinase (sequence 1) (Knighton *et al.*, 1991). Sequences 1 to 22 are representative kinases taken from the March 1992 Protein Kinase Catalytic Domain Database (Hanks & Quinn, 1991). These are: CAPK-ALPHA, cAMP-dependent protein kinase catalytic subunit, α -form; WEE1+, reduced size at division mutant wild-type allele gene product; TIK, mouse serine/threonine kinase; SPK1, *S. cerevisiae* kinase cloned with anti-p-Tyr antibodies; RSK1-N, amino domain of type 1 ribosomal protein S6 kinase; PYT, putative serine/threonine kinase cloned with anti-p-Tyr antibodies; PKC-ALPHA, protein kinase C, α -form; PDGFR-B, platelet-derived growth factor receptor B type; PBS2, polymix in B antibiotic resistance gene product; MIK1, *S. pombe mik1* acts redundantly with *wee1+*; MCK1, *S. cerevisiae* protein kinase; INS.R, insulin receptor; HSVK, Herpes simplex virus-US3 gene product; ERK1, rat insulin-stimulated protein kinase; EGFR, epidermal growth factor receptor (cellular homolog of *v-erbB*); ECK, receptor-like tyrosine kinase detected in epithelial cells; DPYK1, developmentally regulated tyrosine kinase in *D. discoideum*; CLK, mouse serine/threonine/tyrosine kinase; CDC2HS, human functional homolog of yeast *cdc2+*/CDC28; CAMII-ALPHA, calcium/calmodulin-dependent protein kinase II, α -subunit; C-SRC, cellular homolog of *v-src*; and C-RAF, cellular homolog of *v-raf/mil*. Sequences 2 to 4, 6, 10, 11, 14, 17 and 18 are the candidate dual-specificity protein kinases as defined by Lindberg *et al.* (1992). Sequences 23 to 40 are the SWISS-PROT 22 sequences designated as kinases by our HMM (Z-score > 6.0) but not by all 3 other methods, PROSITE, PROFILESEARCH and the keyword search. Sequences 41 to 50 are the top 10 sequences below our cutoff of 6.0 and 41 to 43 and 51 to 60 are sequences that were not classified as kinases by the HMM but were so by one or more (but not all) of the 3 other methods. Note that sequences identified as kinases by all 4 methods are not shown. All sequences that are less than 200 residues in length

the number of sequences denoted as kinases only by all three other methods is evaluated, the number of false negatives for each of the techniques differ from the more detailed analysis: two for the HMM (42 to 43), seven for PROFILESEARCH (23 to 26, 35, 38, 40) and none for PROSITE (ignoring known false negatives as above). This general problem is further highlighted by the guanylyl cyclases (indicated by % in Fig. 11B). If the definition of a kinase is based upon function and not possession of particular sequence patterns, then the guanylyl cyclases are the only false positives for both the HMM and PROFILESEARCH. The PROSITE patterns PROTEIN_KINASE_ATP, PROTEIN_KINASE_ST and PROTEIN_KINASE_TYR produce eight, none and two false positives, respectively, giving some indication of the actual PROSITE performance.

Overall, both the HMM and PROFILESEARCH appear to perform generally better than PROSITE in the discrimination tests, with the HMM possibly having a slight advantage over PROFILESEARCH.

The HMM database search did not suggest any new putative kinases in SWISS-PROT 22. However, a comparative examination of the HMM produced multiple sequence alignment and the crystal structure of the catalytic subunit of cAMP-dependent protein kinase (Knighton *et al.*, 1991) (sequence 1), a template for the protein kinase family, yields insights into the conserved regions and their functions in kinases of unknown structure. Figure 11A displays the location of secondary structure elements obtained from this crystal structure. An invariant Asp in subdomain VIb (Asp166 in Knighton *et al.*, 1991) that is proposed to be the catalytic base is known to diverge in guanylyl cyclases (28 to 34, 36 to 37) even though the immediate region is highly conserved (Garbers, 1992). Our results indicate that other invariant residues appear to be replaced as well. In the sea urchin spermatozoan cell-surface receptor for the chemotactic peptide "resact" (sequences 28 and 34), a Lys

in subdomain II (Lys72) that forms part of the ATP α - and β -phosphate binding site is changed to His. The heat-stable enterotoxin receptor of rat (36) replaces an Asp in subdomain IX (Asp200) that contributes directly to stabilization of the catalytic loop by Glu. Yeast VPS15 (sequence 35), a probable serine/threonine kinase that is autophosphorylated, lacks many of the residues in subdomain I. In addition, a conserved ion-pair that stabilizes ATP (Glu91-Lys72) would be disrupted in VPS15 because the Glu in subdomain III is altered to Arg resulting in the apposition of two positively charged residues. In the putative B12 kinases of two strains of vaccinia virus (42 to 43), the proposed Asp catalytic base is replaced by Lys (cf. guanylyl cyclases). This is accompanied by a further change in the "general" sequence of the catalytic loop: the normally positively charged residue at $n + 2$ has been altered to Glu. In general, all the sequences below our cutoff and the last one above it (40 to 60) appear to lack α -helix F (see X-ray in Fig. 11A). The functional and or structural consequences of these modifications on any kinase activity are not clear.

(c) EF-hand experiments

For these experiments we used the June 1992 database of EF-hand sequences maintained by Kretsinger and co-workers (Nakayama *et al.*, 1992). Sequences in this database are proteins containing one or more copies of the EF-hand motif, a 29-residue structure present in cytosolic calcium-modulated proteins (Nakayama *et al.*, 1992; Persechini *et al.*, 1989; Moncrief *et al.*, 1990). These proteins bind the second messenger calcium and in their active form function as enzymes or regulate other enzymes and structural proteins. The motif consists of an α -helix, a loop binding a Ca^{2+} followed by a second helix. Although a number of proteins possess the EF-hand motif, some of these regions have lost their calcium-binding property.

For our training set, we extracted the EF-hand structures from each of the 242 sequences in the

have been removed. B, Details on sequences 23 to 60 shown in the alignment (arranged in order of decreasing Z-score). NLI-score and Z-score are measures of how well the kinase HMM fits these SWISS-PROT 22 test sequence that were not present in the training set (see Results section (b) for more details). In HMM, PROFILESEARCH and Keyword, + denotes sequences that are classified as containing a kinase domain and - those that do not. For PROFILESEARCH, -\$ identifies sequences that do not appear in the results obtained from searching SWISS-PROT 25 (not 22 as in HMM, Keyword and PROSITE) provided to us by M. Gribskov (personal communication). Two PROSITE signature patterns for eucaryotic protein kinases have been derived and these are labeled A and B in the alignment. A is the region believed to be involved in ATP binding (PROSITE entry PROTEIN_KINASE_ATP) while B1 and B2 indicate the area important for catalytic activity in serine/threonine kinases (PROTEIN_KINASE_ST) and tyrosine kinases (PROTEIN_KINASE_TYR), respectively. In all instances, T signifies a true positive; N a false negative (a sequence which belongs to the set under consideration but which is not picked up by the pattern); P a "potential" hit (a sequence that belongs to the set but which is not picked up because the region that contains the pattern is not yet available in the data bank, i.e. a partial sequence); and ? an unknown (a sequence which possibly could belong to the set). * Indicates SWISS-PROT files which contain a cross reference to the specified PROSITE pattern, but these PROSITE entries do not contain a corresponding pointer to the SWISS-PROT file. - Signifies sequences that do not satisfy the kinase patterns and % denotes particulate forms of guanylyl cyclase receptors which contain an intracellular protein kinase-like domain but which have not been shown to possess kinase activity to date (reviewed by Garbers, 1992).

database, obtaining 885 EF-hand motifs having an average length of 29. For our first experiment we trained five HMMs on all 885 EF-hand motifs, using the standard techniques described earlier. (In subsequent experiments, described below, we trained on smaller subsets of these 885 sequences.) The best model had a final length of 29, and a NLL-score (the average $-\log P(\text{sequence}|\text{model})$) of 61.41.

As described in Methods section (e), we modified the final model to enable it to search the SWISS-PROT database for sequences containing the EF-hand motif. We computed Z -scores for all sequences as described in section (f) and Figure 12 shows the resulting histogram. In contrast to the kinases, a visual inspection of the histogram of Z -scores did not indicate the presence of a distinct gap thus making the selection of a cutoff more difficult. After choosing by eye a cutoff of 4.75 and excluding all sequences with unknown residues (Xs), the model classified 232 sequences as containing the EF-hand sequence motif.

As with the kinase experiments in the previous section, false positives and false negatives were identified in the following manner. A list of "certain EF-hands" was created from the union of sequences determined to be containing the EF-hand motif by three independent sources: PROSITE, a keyword search, and the results of Michael Gribskov's PROFILESEARCH. Details of the PROSITE and keyword searches are given by Krogh *et al.* (1993a). Two different PROFILESEARCH experiments were conducted for us by M. Gribskov (personal communication). The first employed a profile generated using the multiple sequence alignment of sequences classified as EF-hands by our HMM and the second was constructed using an alignment of the following four sequences: *Escherichia coli* galactose binding protein (JGECG, 1 EF-hand motif), rabbit parvalbumin (PVRB, 2), human troponin (TPHUCS, 4) and human calmodulin (MCHU, 4).

Although a sequence may possess multiple copies of the EF-hand (or any other) motif, only the one which most closely resembles that described by the HMM is identified. Of the 232 SWISS-PROT 22 sequences that were above the cutoff (Z -score >4.75) and were thus classified as containing an EF-hand motif by our HMM, 163 were similarly classified by PROSITE, both PROFILESEARCH experiments and the keyword search (if only one of the PROFILESEARCH experiments is considered, then there are an additional 14 sequences making a total of 177). These may be considered to constitute certain EF-hands and Figure 13 shows the multiple sequence alignment generated by our HMM of some representative EF-hands from this set (sequences 1 to 27). Of the 69 (232 minus 163) or 55 (232 minus 177) sequences above the cutoff and not categorized as EF-hands by all three other methods, 33 possess the motif but do not bind calcium (indicated by % in Fig. 13B) and six (64, 72, 88, 89, 91, 94) were classed as EF-hands by only one other method.

The identification of certain EF-hands as compared to certain kinases is not as straight-

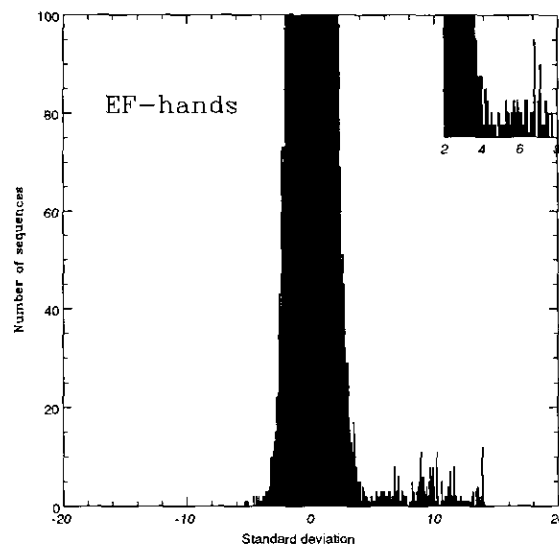


Figure 12. Histogram showing the number of sequences with a certain Z -score relative to the EF-hand model.

forward, making it difficult to ascertain the precise number of classification errors made by each technique. This problem arises partly because of the absence of a pronounced gap in the histogram of Z -scores and the resultant uncertainty in assigning an exact cutoff (Figs 10 and 12). The mnemonic developed to identify EF-hand homologs and distinguish them from analogs (Nakayama *et al.*, 1992) is known to generate errors and is unable to detect 8 of the 27 sequences known to be EF-hands (sequences 1 to 27 in Fig. 13). Therefore, the sensitivity and specificity of the EF-hand database discrimination tests is unlikely to be comparable to the kinases. Using Figure 13B, an estimate of the false negative rate for each method was determined by using the simple notion of evaluating the number of sequences classified as EF-hands by all methods other than the one being considered. (Those which possess the motif but do not bind calcium, denoted by % in Fig. 13B, are not considered.) Using this criterion, the number of false negatives are 1 for the HMM (101), 20 for PROFILESEARCH using four sequences (28, 47, 56 to 57, 59, 67, 74 to 82, 84 to 85, 92 to 93, 96), seven for PROFILESEARCH using our EF-hand alignment (28, 57, 74, 79 to 80, 92 to 93), one for the keyword search (58) and two for PROSITE (60, 70). A similar analysis of false positives produces six for the HMM (52, 71, 83, 86, 90, 94) nine for PROFILESEARCH using four sequences (97, 99, 111 to 112, 121 to 122, 129, 132 to 133), eight for the keyword (123, 126, 130-131, 134 to 137) and one for PROSITE (120). It should be noted however, that a search of SWISS-PROT 22 using the PROSITE pattern EF-HAND produces different results: three false negatives and 24 false positives (compared with 2 and 1 using the simple criterion). A total of 26 sequences were not desig-

	1	11	21	31	41	51	61	71
StructureH..H.HHRRH.H.H.....H.....LL...LLL.LLLL...H.HH.HHHHHH.....							
PROSITEE..F..K.EAFS.L.F.....D....KD...GGG.TITT...K.EL.GTVNRSL-							
Ca-bindingE..F..RASFW.H.F.....D....RK...ATG.MKDC...E.DF.RACLISW-							
1 CANHSE..F..K.EAFS.L.F.....D....KD...GGG.TITT...K.EL.GTVNRSL-							
2 aACTGGE..F..RASFW.H.F.....D....RK...ATG.MKDC...E.DF.RACLISW-							
3 VISINWE..L..S.WYF.G.F.....Qr...QC...SDG.AIAC...D.EF.PAIYGF-							
4 TFP24CFG..L..ARFR.R.L.....D....RD...RSR.SLDS...R.EL.QRGLAEL-							
5 TPHVCSE..F..KAAPD.H.F.....D....AD...GGG.DISV...K.EL.GTVNRSL-							
6 TPAP1A..L..QKAFD.S.F.....D....TD...SKG.FITP...E.TV.GILRHW-							
7 TCRP25V..A..RAIFE.H.Y.....D....RG...RKG.RIEN...T.DC.VPNIETA-							
8 SPEG2AL..F..KSSF.R.S.E.....D....TD...GGG.KITS...E.EL.RAAFCSI-							
9 SCBPBL1K..I..KPTFD.F.F1.....D....YM...KDG.SIQW...E.DF.EENIKRY-							
10 QUIDLWE..I..KDAFD.H.F.....D....ID...GGG.QITS...K.EL.RSVNKS-							
11 MORSERE..F..KEAFT.I.N.....D....QN...RDG.FIDK...N.DL.RDTFAAL-							
12 MORSa1E..F..KEAFL.L.F.....D....ST...GDS.KITL...S.QV.GDVLRAL-							
13 LPS1AA..L..KQEF.DnY.....D....TM...KDG.TVSC...A.EL.VKLKWT-							
14 LAV1A..L..VADFR.K.I.....D....TM...SWG.TLSR...K.EF.RHFVRL-							
15 EFN5E..L..AEGFR.V.L.....D....SN...GQK.TISIp...K.EV.SALNASV-							
16 CVPE..C..KRIFD.I.F.....D....RM...AEN.IAPV...S.DT.KDNLTKL-							
17 CRGHSQ..L..HYFK.H.H.....D....YD...GWN.LLDG...L.EL.STAITHV-							
18 CRSER..L..KRFR.R.V.....D....FD...GKG.ALER...A.DF.EKEAQHI-							
19 CDFRG..L..KELFR.K.I.....D....TD...NSG.TITF...D.EL.RDGLRRV-							
20 CDC31E..I..YEAFS.L.F.....D....KM...WDG.FLDY...H.EL.RVANAL-							
21 CALPLHST..C..RSVA.V.M.....D....SD...TTG.KLGF...E.EF.KYLWNI-							
22 CALCTBR..L..GRFR.K.L.....D....LD...NSG.SLSV...E.EF.WS-LPEL-							
23 CALBNGQ..F..FEIVH.H.Y.....D....SD...GNG.VMDG...K.EL.QMFIQEL-							
24 CALICEE..F..REAFW.H.F.....D....KD...GNG.TIST...K.EL.GIARRSL-							
25 BGHSA..L..IDVFH.Q.Y.....Sg...RE...GDKRLKK...S.EL.KELINNE-							
26 AEDV1R..R..KHFV.F.L.....D....VM...HNG.KISL...D.EN.VYKASDI-							
27 IF0R..R..IELFR.K.F.....D....KM...ETG.KLCY...D.EV.HSGCLEV-							
28 CALNASPMI	adslteveeE..F..K.EAFS.L.F.....D....KD...GGG.QITT...K.EL.GTVNRSL-gqnpaa109							
29 NLE1_HUMAN	apkhdkv129D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gekmkee35							
30 NLE1_RABIT	apkhdkv127D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gekmkee35							
31 NLEV_HUMAN	apkkpep130D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gerltee35							
32 NLEK_CHICK	ppkkpep129D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gerltee35							
33 NLEV_RAT	apkkpep135D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gerltee35							
34 NLE1_CHICK	ppkhdkv126D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gekmtee35							
35 NLE1_RAT	apkhdkv124D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gekmtee35							
36 NLE1_MOUSE	apkhdkv123D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gekmtee35							
37 NLEF_HUMAN	apkkpep132D..F..VEGLR.V.F.....D....KE...SWG.TVNG...A.EL.RHVLATL-gekmtee35							
38 NLEF_RAT	ppkkpep128D..F..VEGLR.V.F.....D....KE...SWG.TVNG...A.EL.RHVLATL-gekmtee35							
39 NLEF_MOUSE	ppkkpep128D..F..VEGLR.V.F.....D....KE...SWG.TVNG...A.EL.RHVLATL-gekmtee35							
40 NLEK_CHICK	mpkhdp121D..F..VEGLR.V.F.....D....KE...GNG.LVNG...A.EL.RHVLATL-gekmtee35							
41 NLE3_HUMAN	sfadqia185D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gekmtee35							
42 NLEY_HUMAN	mpkhdkv144D..Y.LEGFR.V.F.....D....KE...GNG.KVNG...A.EL.RHVLATL-gekmtee35							
43 NLE3_RABIT	sfadqia185D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gekmtee35							
44 NLE3_RAT	sfadqia185D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gekmtee35							
45 NLE3_MOUSE	sfadqia185D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gekmtee35							
46 NLE3_CHICK	sfadqia185D..F..VEGLR.V.F.....D....KE...GNG.TVNG...A.EL.RHVLATL-gekmtee35							
47 AACT_HUMAN	mdhyoaq149E..F..RASFW.H.F.....D....RD...HSG.TLGP...E.EF.KACLISL-gydlgd114							
48 NLE_HALLO	adrsdr186D..F..VEGLR.V.F.....D....KE...NMG.KING...A.EL.RHVLATL-gekmtee35							
49 NLES_HUMAN	mcdfteqtaE..F..KEAFQ.L.F.....D....RT...GGG.KITLY...S.QC.GDVNRAL-gqnpaa112							
50 NLEK_HUMAN	mcdfteqtaE..F..KEAFQ.L.F.....D....RT...GGG.KITLY...S.QC.GDVNRAL-gqnpaa112							
51 NLEM_CHICK	cdfsseqtaE..F..KEAFQ.L.F.....D....RT...GGG.KITLY...S.QC.GDVNRAL-gqnpaa112							
52 NLEM_CHICK	cdfsseqtaE..F..KEAFQ.L.F.....D....RT...GGG.KITLY...S.QC.GDVNRAL-gqnpaa112							
53 NLEK_HUMAN	eeemvnl30D..F..VEGLR.V.F.....D....KE...SWG.TVNG...A.EL.RHVLATL-gekmtee35							
54 NLE_PATYE	pklsqde184D..Y..NEAFK.T.F.....D....RE...GGG.FISG...A.EL.RHVLATL-gerltee43							
55 NLE_AEQIR	pklsqde184D..Y..NEAFK.T.F.....D....RE...GGG.FISG...A.EL.RHVLATL-gerltee43							
56 AACT_DRONE	mmeng1752E..F..KSSF.H.F.....D....KW...ATG.RLSP...E.EF.KSCLVSL-gysigk114							
57 RECO_CHICK	mgnezv100K..L..EWAFS.L.F.....D....VD...RNG.EVSR...S.EV.LEIITAI-2kmpiee83							
58 NLE_DICDI	measdaq178E..N..LDAFK.A.L.....D....KE...GNG.TIQQ...A.EL.RQLLTL-gdylatae59							
59 SPCA_DRONE	menfip2268E..F..SNFX.H.F.....D....KD...KSG.KLNH...Q.EF.KSCLRAL-gydlpmv118							
60 NLR_DICDI	mastrkr123E..L..KEAPE.L.F.....D....KD...ATG.FIRK...D.AL.RITCRQ-gyivme109							
61 NLE_TODPA	sqtkdei185E..F..NEAFK.T.F.....D....KE...GGG.LJSS...A.EI.RMWLKL-gerltee45							
62 SPCM_CHICK	mdpavv2331E..F..SNFX.H.F.....D....KD...KSG.KLNH...Q.EF.KSCLRAL-gydlpmv117							
63 CL1L_MOUSE	pqmeham49A..V..DKLNK.D.L.....D....QC...RDG.KVGF...Q.SF.LSLVAGL-tiacudyf18							
64 AACCS_CHICK	mmamqi1795E..F..ARINS.L.V.....D....PM...GGG.VVTF...Q.SF.IDPNTRE-tadtdtae73							
65 CL1L_RAT	pqmeham49A..V..DKLNK.D.L.....D....QC...RDG.KVGF...Q.SF.LSLVAGL-tiacudyf18							
66 LAV1_PHYPD	mayqea220A..L..VADFR.K.I.....D....TX...SWG.TLSR...K.EF.RHFVRL-gtdkksv106							
67 CAP3_RAT	mpvips9695S..C..ASIA.L.N.....D....TD...GSG.KLNL...Q.EF.HHLWRI-kauqik1k97							
69 NLEP_DRONE	mvdvpkre83D..F..IECLK.L.Y.....D....KE...ENG.TMLL...A.EL.QHALLAL-ge1lddeq43							
69 NLEL_DRONE	mvdvpkre83D..F..IECLK.L.Y.....D....KE...ENG.TMLL...A.EL.QHALLAL-ge1lddeq43							
70 SP2D_STRPU	maan1f111K..J..KEIE.K.A.....D....FP...NDG.KCSL...E.EF.VNKNVF-cs1111111							
71 CL1L_BOVIN	pqmeham49A..V..DKLNK.D.L.....D....QC...RDG.KVGF...Q.SF.FSLIAGL-tiacudyf18							
72 EHF5_TRYBB	mkkdkpy122E..N..KCAFL.H.V.....D....KQ...KTC.FVTK...K.QF.TELFATG-gecatpee41							
73 CL1L_PIG	pqmeham49A..V..DKLNK.D.L.....D....QC...RDG.KVGF...Q.SF.FSLIAGL-tiacudyf18							
74 FCAB_TRYBB	mgcagk170D..A..TVFM.E.I.....D....TM...GSG.VVTF...D.EF.SCWAVTA-k1qvgvdp34							
75 SCPI_LASTPO	ayewdzv59L..W..KEIAE.L.A.....D....FN...KDG.EVTI...D.EF.KKAVQV-cvghafaf04							
76 CAP2_RABIT	qtlirir297T..C..KIVD.H.L.....D....SD...GTC.KLGL...K.EF.YVLWRI-qkyqk1k96							
77 CAP3_HUMAN	ea1iam652S..C..ASIA.L.N.....D....TD...GSG.KLNL...Q.EF.HHLWRI-kauqik1k97							
78 CAP3_HUMAN	mflvna142T..C..SNVA.V.N.....D....SD...TTG.KLGF...E.EF.KYLWNI-krvqi1k97							

A
Fig. 13.

79	CAP2_HUMAN	magiaak875T..C.KINVD..N.L.....D.....SD...GSG.KLGL...K.EF.YILWTAI-qkyqkiyr96
80	KOGL_PIG	mksrgk157I..L.QENKK.E.I.....D.....VD...GSG.SVSL...A.EW.LRAGATT-vp1l13440
81	SCPA_PENSP	aysvdrf103F..I.ANQFK.A.I.....D.....VM...GDC.KVGL...D.EYrLDCITAS-afaevek159
82	SCPB_PENSP	aysvdrv59L..V.MEIAE.L.A.....D.....FM...KDC.EVTV...D.EF.KQAVQRM-ckghaf104
83	IPYR_ARATH	moeikde169E..I.RRFFE.D.V.....D.....K...KN...EMK.XVDV...E.AF.LPAQAAI-daiikdm65
84	SCP1_BRALA	glndfqk105R..I.PFLFK.G.N.....D.....VS...GDC.IVDL...E.EF.QNYCKMF-qlqcadvp51
85	SCP2_BRALA	glndfqk105R..I.PFLFK.G.N.....D.....VS...GDC.IVDL...E.EF.QNYCKMF-qlqcadvp51
86	PIF3_RAT	mdsgrdf143W..I.HSCLR.K.A.....D.....KN...KDN.RNMF...K.EL.KDFLRSE-niqvddg584
87	AACT_CHICK	mdhbydp786E..F.ARINS.I.V.....D.....PN...KMG.VYTF...Q.AF.LDFNSAE-tadtatd73
88	CAB_MOUSE	marpleea53A..F.QKVS.W.L.....D.....SM...RDX.EVDF...Q.EY.CVFLSCT-ammcue119
89	TEGU_SCHRA	matetk1e1E..F.IRAF.L.E.I.....D.....AD...SME.NIDN...Q.EL.IKYCRVY-r1dmk1150
90	CAB_RAT	marpleea53A..F.QKLN.N.L.....D.....SM...RDN.EVDF...Q.EY.CVFLSCT-ammcue119
91	G19P_HUMAN	ml1p11211E..L.AADAFK.E.L.....D.....DD...KDC.TYSV...T.EL.QTH-PEL-dtdgga287
92	TCH2_ARATHD.....KN...GDC.KISV...D.EL.KEVIRAL-epataspee25
93	KDGL_HUMAN	mksrgk150I..L.QENKK.E.I.....D.....VD...GSG.SVSL...A.EW.VRAGATT-vp1l13440
94	PIF3_BOVIN	peoqf182W..I.HSCLR.K.A.....D.....KN...KDN.RNMF...K.EL.QNFKLE-niqvddg584
95	CALN_LYTP1	kknkdtdeeE..I.REAFR.V.F.....D.....KD...GWC.FI-----EL-a.....
96	CAP1_HUMAN	mseeif588S..G.RSNVN.L.N.....D.....RD...GCG.KLGL...V.EF.YILWRII-ruyliftr97
97	CIC1_CYPCA	meegag142IE..F.KRIVA.E.Y.....D.....PE...ATG.RIKH...L.DV.VTLRRI-pp1gfg102
98	GUNF_CLOTN	mk1laf699E..H.QRFIA.A.A.....D.....VD...GCG.RINS...T.DL.YVL--RR-yilk1iet15
99	CIC1_KABIT	mepspp190EE..F.KAIVA.E.Y.....D.....PE...ARG.RIKH...L.DV.VTLRRI-pp1gfg144
100	V57A_BPT4	mseqtveq40T..L.AEIAH.A.V.....D.....G...IT...GD...TIAV...E.EI.VCAVNL-taesada12
101	CALG_CHICK	ellavfq43V..V.DRNNK.R.L.....D.....IM...SDG.QLDF...Q.EF.....
102	K1FH_MOSCO	mdhyvprd56E..L.EELLI.E.F.....D.....G...ILES...D.EW.TANLVGR-tateap.....
103	ARF1_DROME	mgvlyt85A..I.IYVVD.S.A.....D.....RD...RIG--ISK...D.EL.LYMLRE-elage1v67
104	ARDA_LLEPE	mealt1q151R..L.RGGFI.G.G.....D.....IE...VDC.SVSS...Q.FL.TALLNAS-plapadt247
105	RELI_HUMAN	mpr1f1b96E..L.KAALS.E.Raps1pel11P...AL...KDS.WLGF...E.EF.KKLRIRK-qseadan49
106	K11_BOVIN	setapaap41L..I.TRAVA.Aa.....D.....RE...RSG--VSL...A.AL.KAALAA-gydvkkn36
107	YCSL_CHLPPY	malenil155K..L.IEFLD.N.Y.....D.....K...VE...KAK.SITL...Q.QL.QSVLQRI-klmsakq26
108	DP3L_ECOLI	mayq1421S..L.NDALS.L.T.....D.....QAIesGDC.QVST...Q.AV.SANLRL-ddqqa1399
109	ARDA_SALTY	mealt1q151R..L.RGGFI.G.G.....D.....IE...VDC.SVSS...Q.FL.TALLNTA-plapadt247
110	ANK1_CAVCU	mnmvaf102H..L.EEVVL.A.L.....D.....L...KY...PA--QLDA...D.EL.KAANKGL-gtdet1216
111	CICC_RAT	mirafal524E..F.KRIVA.E.Y.....D.....PE...ARG.RIKH...L.DV.VTLRRI-pp1gfg16
112	CICC_KABIT	mlralv1525E..F.KRIVA.E.Y.....D.....PE...ARG.RIKH...L.DV.VTLRRI-pp1gfg16
113	LACA_LACLA	maivvgad21L..V.EEGFE.V.I.....D.....VT...KDC.Q-DF...V.DV.TLAVASE-vnkdeqal92
114	ARDA_BORPE	maglay1d31L..L.AAALA.E.G.....D.....S...TE...ITG.LLDS...D.DT.RVNLAL-rqlgvsv382
115	ARDA_SALTY	mealt1q151R..L.RGGFI.G.G.....D.....IE...VDC.SVSS...Q.FL.TALLNTA-plapadt247
116	ARDA_SALGL	mealt1q151R..L.RGGFI.G.G.....D.....IE...VDC.SVSS...Q.FL.TALLNTA-plapadt247
117	CAP1_CHICK	mpggia608S..W.LTIFR.Q.Y.....D.....LD...KSG.TMSS...Y.EN.RNALESA-gfklank167
118	PR10_CAVPOD.....AD...ENG.YIEG...K.FN.QRY--DK-oadqhvga68
119	SC1_RAT	mkav111593EhcI.TRFPE.E.C.....D.....PN...KDK.HITL...K.EW.GHCFGIK-eedidn11f
120	QR1_COTJA	mkcv111635EhcI.TRFPE.E.C.....D.....GD...QDK.L.ITL...K.EW.CRCFAIK-eedidn11f
121	RS37_MEUCR	gkzrkky22K..L-AVLA.Y.Y.....D.....K...VD...SDG.KIER...---LRRE-rcpnc2ga34
122	YTR1_SPLAU	mvmdbin49P..T.KFVAS.I.A.....D.....D.....G.RVTF...S.FF.VPLGLRL-daktplav66
123	SPCB_HUMAN	kfedflg225E..L.GELFA.Q.V.....D.....PsmgeEG...GDA.DLSI...ERF.LDLLEPL-grtkql15
124	OTNC_MOUSE	mraviff261EhcT.TRFPE.T.C.....D.....LD...KDR.YIAL...D.EW.AGCFGIK-ekidk1v1
125	CALC_KABIT	dgh.....S.V.....D.....TLSK...T.EF.LSFWTE-lactkdp16
126	SPCA_MOUSE	rvcdgde199A..V.QWVLD.T.A.....D.....E...SL...RDKAVGA...E.EI.QENLAQF-vqhvkk124
127	OTNC_HUMAN	mraviff262EhcT.TRFPE.T.C.....D.....LD...KDR.YIAL...D.EW.AGCFGIK-ekidk1v1
128	OTNC_BOVIN	mraviff263EhcT.TRFPE.T.C.....D.....LD...KDR.YIAL...D.EW.AGCFGIK-ekidk1v1
129	Y493_BPT4	mi.....E..L.NEQI1EL.G.....D.....Dg...TE...GDL.EYKL...Y.EY.MIWLARA-egidvsv69
130	KDGL_ECOLI	anzetgt48D..V.DA1TR.V.L.....D.....D.....LISS...V.NL.VN1VEIL-mnaizeav50
131	SPCA_HUMAN	etvves128E..L.ARLWD.L.L.....D.....L.ELTL...E.KG.DQLLRAL-kfqqvq443
132	INMC_ECOLI	nglk1bh15E..F.KG--G.E.Y.....D.....S...KDI.GDD...---GSVIESL-gmp1kdn49
133	DGAL_ECOLI	mmkhv1283Q..A.KATFD.L.A.....D.....Knl..AD...GKG.AADG...T.NW.KID--NA-vrvrvpyg20
134	SPCB_MOUSE	lqat1q189A..L.RRWVE.S.R.....D.....G...NT...LTO.CLGF...Q.EF.QDKARQA-eailanq18
135	SP10_YEAST	mkf1ev1s1T..T.TTLFW.St.....D.....S...TLNI...T.QL.YQIATGV-nqtlqee251
136	SRCH_HUMAN	mgbhrrp192R..R.EEAGG.A.S.....D.....S...EE...ESG.EDT6pqaQ.EY.GNYQPCS-1cgyest74
137	SRCH_KABIT	mgcrp144D..L.AEHGS.H.Ghgbee.....D.....ED...ED--VISS...E.RP.RNVLRRR-prghggee76

A (cont)

Fig. 13.

nated as EF-hands by the HMM but were classified so by PROFILESEARCH, PROSITE or the keyword search. Of these, 19 were classified as such by only one of these methods. This includes five fragments where the EF-hand motif is missing: human and murine spectrin alpha- and beta-chains (123, 126, 131, 134) and rabbit calgizzarin (125).

Inspection of the HMM produced alignment and examination of the putative calcium-binding ligands (Fig. 13) for the 20 sequences immediately below the cutoff (97 to 116) and the false negatives and positives suggests that many possess potential EF-hand motifs. This includes six sequences whose Z-scores lie above our cutoff but are not classed as EF-hands by any other method: chicken myosin light chain alkali, smooth muscle (52); bovine calpactin I light chain (71); *Arabidopsis thaliana*

inorganic pyrophosphatase (83); rat placental calcium-binding protein (90) (note however that the mouse protein, sequence 88, is designated an EF-hand by the keyword search); and rat and bovine 1-phosphatidylinositol-4,5-bisphosphate phosphodiesterase III (86 and 94). A notable example among the false negatives is the α -1 sub-unit of L-type calcium channels from carp and rabbit skeletal muscle (97, 99) and rat and rabbit cardiac muscle (111, 112). These proteins play an important role in excitation-contraction coupling and carry the calcium antagonist binding domains (reviewed by Grabner *et al.*, 1991). They possess a highly conserved and evolutionarily preserved putative intracellular region of 155 residues near the carboxyl terminus immediately following the fourth internal repeat. This region has been suggested to

ID	Length	NLL-score	Z-score	HMM	PROFILESEARCH		Keyword	Prosite
					Gribskov	HMM		
28 CALM_ASPNI	148	398.961	12.975	+	-	-	+	T
29 MLE1_HUMAN	193	542.924	11.662	+	+	+	-	%
30 MLE1_RABIT	191	537.011	11.661	+	+	+	-	%
31 MLEV_HUMAN	194	546.027	11.631	+	+	+	-	%
32 MLEC_CHICK	193	543.095	11.605	+	+	+	-	%
33 MLEV_RAT	199	561.007	11.561	+	+	+	-	%
34 MLE1_CHICK	190	534.042	11.516	+	+	+	-	%
35 MLE1_RAT	188	528.051	11.262	+	+	+	-	%
36 MLE1_MOUSE	187	525.056	11.224	+	+	+	-	%
37 MLEF_HUMAN	196	554.316	11.005	+	+	+	-	%
38 MLEF_RAT	192	542.332	10.892	+	+	+	-	%
39 MLEF_MOUSE	192	542.332	10.892	+	+	+	-	%
40 MLEX_CHICK	185	521.797	10.342	+	+	+	-	%
41 MLE3_HUMAN	149	411.100	10.201	+	+	+	-	%
42 MLEY_HUMAN	208	588.847	10.194	+	+	+	-	%
43 MLE3_RABIT	149	411.179	10.177	+	+	+	-	%
44 MLE3_RAT	149	411.207	10.169	+	+	+	-	%
45 MLE3_MOUSE	149	411.208	10.169	+	+	+	-	%
46 MLE3_CHICK	149	411.206	10.169	+	+	+	-	%
47 AACT_HUMAN	892	2642.237	9.957	+	-	+	+	T
48 MLE_HALRO	151	418.497	9.918	+	+	+	-	%
49 MLES_HUMAN	151	418.627	9.879	+	+	+	-	%
50 MLEN_HUMAN	151	418.627	9.879	+	+	+	-	%
51 MLEN_CHICK	150	415.631	9.798	+	+	+	-	%
52 MLEM_CHICK	150	415.631	9.798	+	-	-	-	%
53 MLEG_HUMAN	94	248.725	9.735	+	+	+	-	%
54 MLE_PATYE	156	433.703	9.629	+	+	+	-	%
55 MLE_AEQIR	156	433.703	9.629	+	+	+	-	%
56 AACT_DROME	895	2653.286	9.130	+	-	+	+	T
57 RECO_CHICK	192	548.396	8.848	+	-	-	+	T
58 MLE_DICDI	166	465.170	8.834	+	+	+	-	T
59 SPCA_DROME	2415	7205.568	8.787	+	-	+	+	T
60 MLR_DICDI	161	451.967	8.678	+	+	+	+	-
61 MLE_TODPA	159	446.406	8.616	+	+	+	-	%
62 SPEN_CHICK	2477	7392.895	8.157	+	-	+	-	T
63 CL1L_MOUSE	96	263.095	7.516	+	+	+	-	%
64 AACS_CHICK	897	2663.548	7.446	+	-	-	+	-
65 CL1L_RAT	94	257.103	7.423	+	+	+	-	%
66 LAV1_PHYPO	355	1039.236	7.298	+	-	+	-	T
67 CAP3_RAT	821	2436.445	7.150	+	-	+	+	T
68 MLEP_DROME	155	439.713	7.053	+	+	+	-	%
69 MLEL_DROME	155	439.713	7.053	+	+	+	-	%
70 SP2D_STRPU	141	397.689	6.990	+	+	+	+	-
71 CL1L_BOVIN	96	265.582	6.819	+	-	-	-	%
72 EHP5_TRYBB	192	554.482	6.797	+	+	+	-	-
73 CL1L_PIG	95	262.586	6.763	+	+	+	-	%
74 FCAB_TRYBB	233	676.012	6.684	+	-	-	+	T
75 SCP1_LASTPO	192	554.824	6.681	+	-	+	+	T
76 CAP2_RABIT	422	1242.278	6.589	+	-	+	+	T
77 CAP3_HUMAN	778	2307.499	6.577	+	-	+	+	T
78 CAPS_HUMAN	268	782.852	6.383	+	-	+	+	T
79 CAP2_HUMAN	700	2074.486	6.305	+	-	-	+	T
80 KPGL_PIG	734	2176.760	6.160	+	-	-	+	T

B

Fig. 13.

contain functional domains that are typical or essential for all L-type calcium channels regardless of whether they couple to ryanodine receptors, conduct ions or both (Grabner *et al.*, 1991). The

inferred EF-hands for these proteins occur within this conserved 155-residue segment.

The above results were for an HMM trained on all 885 EF-hand motifs from the Kretsinger database.

ID	Length	NLL-score	Z-score	HMM	PROFILESEARCH		Keyword	Prosite
					Gribskov	HMM		
81 SCPA.PENSP	192	556.636	6.071	+	-	+	+	T
82 SCPB.PENSP	192	557.071	5.924	+	-	+	+	T
83 IPYR.ARATH	263	769.241	5.909	+	-	-	-	-
84 SCP1.BRALA	185	535.787	5.827	+	-	+	+	T
85 SCP2.BRALA	185	535.816	5.818	+	-	+	+	T
86 PIP3.RAT	756	2244.255	5.713	+	-	-	-	?
87 AACT.CHICK	888	2641.411	5.684	+	-	-	+	N
88 CAB.MOUSE	101	284.695	5.589	+	-	-	+	-
89 TEGU.SCHMA	190	552.242	5.469	+	-	+	-	?
90 CAB.RAT	101	285.488	5.369	+	-	-	-	-
91 GI9P.HUMAN	527	1560.198	5.330	+	-	-	-	T
92 TCH2.ARATH	45	116.235	5.321	+	-	-	+	T
93 KDGL.HUMAN	735	2182.343	5.301	+	-	-	+	T
94 PIP3.BOVIN	695	2063.206	5.034	+	-	-	-	?
95 CALM.LYTP1	30	67.341	4.942	+	-	-	+	P
96 CAP1.HUMAN	714	2120.342	4.924	+	-	+	+	T
97 C1C1.CYPCA	1852	5530.321	4.714	-	+	-	-	-
98 GUNF.CLOTM	739	2196.618	4.602	-	-	-	-	?
99 C1C1.RABIT	1873	5593.640	4.550	-	+	-	-	-
100 V57A.BPT4	80	224.359	4.470	-	-	-	-	-
101 CALG.CHICK	65	178.908	4.438	-	+	+	+	T
102 N1PH.NOSCO	86	243.556	4.347	-	-	-	-	-
103 ARFL.DROME	180	524.609	4.300	-	-	-	-	-
104 AROA.KLEPN	427	1264.280	4.296	-	-	-	-	-
105 REL1.HUMAN	185	540.676	4.249	-	-	-	-	-
106 H11.BOVIN	104	298.227	4.240	-	-	-	-	-
107 YGSX.CHLPY	110	316.022	4.210	-	-	-	-	-
108 DP3X.ECOLI	643	1910.667	4.186	-	-	-	-	-
109 AROA.SALTY	427	1264.760	4.130	-	-	-	-	-
110 ANX1.CAVCU	346	1022.514	4.043	-	-	-	-	-
111 C1CC.RAT	2169	6481.468	4.011	-	+	-	-	-
112 C1CC.RABIT	2171	6487.460	4.010	-	+	-	-	-
113 LACA.LACLA	141	407.967	3.986	-	-	-	-	-
114 AROA.BORPE	442	1310.475	3.985	-	-	-	-	-
115 AROA.SALTI	427	1265.295	3.945	-	-	-	-	-
116 AROA.SALGL	427	1265.295	3.945	-	-	-	-	-
117 CAP1.CHICK	704	2093.590	3.888	-	-	-	+	T
118 PR10.CAVPO	92	267.751	2.866	-	+	+	+	P
119 SC1.RAT	634	1888.351	2.662	-	-	-	-	T
120 QR1.COTJA	676	2015.770	1.941	-	-	-	-	T
121 RS37.NEUCR	78	229.363	1.766	-	+	-	-	-
122 YTR1.SPIAU	140	412.844	1.753	-	+	-	-	-
123 SPCB.HUMAN	274	814.811	1.610	-	-	-	+	-
124 OTNC.MOUSE	302	899.470	1.146	-	-	-	+	T
125 CALG.RABIT	35	106.946	1.126	-	+	+	+	P
126 SPCA.MOUSE	253	753.490	1.101	-	-	-	+	-
127 OTNC.HUMAN	303	902.914	0.988	-	-	-	+	T
128 OTNC.BOVIN	304	905.856	0.983	-	-	-	+	T
129 Y493.BPT4	102	305.597	0.603	-	+	-	-	-
130 KDGLECOLI	121	362.137	0.547	-	-	-	+	-
131 SPCA.HUMAN	595	1779.087	0.039	-	-	-	+	-
132 IMMC.ECOLI	85	257.069	0.025	-	+	-	-	-
133 DGAL.ECOLI	332	992.734	-0.028	-	+	-	-	-
134 SPCB.MOUSE	236	706.853	-0.161	-	-	-	+	-
135 SP10.YEAST	326	978.184	-1.203	-	-	-	+	-
136 SRCH.HUMAN	699	2098.086	-2.613	-	-	-	+	-
137 SRCH.RABIT	852	2556.715	-3.145	-	-	-	+	-

B

Figure 13. A, Multiple sequence alignment generated by our EF-hand HMM of some of the sequences used to train the HMM (1 to 27) and test sequences from the SWISS-PROT 22 database (28 to 137) (see Results section (c)). In Structure, H and L denote residues in an α -helical or loop conformation based upon EF-hands of known structure (Nakayama *et al.*,

There is considerable overlap between this training set and the EF-hand motifs found in SWISS-PROT 22, so in order to provide some clearer cross validation of our results we also did another series of experiments. In these experiments, models were estimated using training sets consisting of different numbers of randomly chosen EF-hand sequences from the database of 885 EF-hand sequences. For training sets consisting of 5, 10, and 20 random EF-hand sequences, 15 models were estimated, each using a different randomly chosen training set. For training sets consisting of 40, 80, 100, 200, and 400 random EF-hand sequences, five models were estimated. In all, 70 models were estimated. A model's performance after training was gauged on how well it performed on a test set which consisted of motifs from the database of 885 sequences that were not used in the training set. Thus for each model, two NLL-scores were computed (see Methods section (f)), one for the training set and one for the test set. These NLL-scores serve as a quantitative measure of how well the model is representing the sequence data. Figure 14 shows that for small training set sizes, the model overfits the training data. This is shown by low training NLL-scores but very high testing NLL-scores. This effect largely disappears when the training set size reaches about 100 sequences.

A model's performance was also gauged on how well it searches a database for sequences containing the EF-hand motif. For each training set size, one model was randomly chosen to search SWISS-PROT 22. A histogram of the resultant Z-scores was plotted and a cutoff was chosen by eye. The number of false positives was computed, as

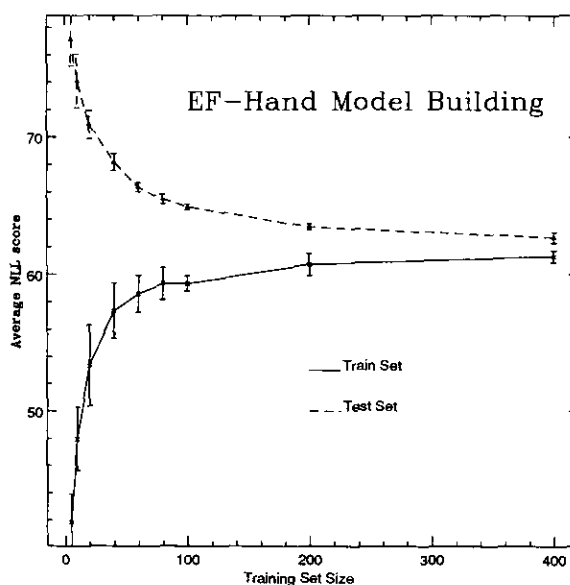


Figure 14. Average NLL scores for test and train sets for models with training sets of size 5, 10, 20, 40, 80, 100, 200 and 400. Error bars represent one standard deviation.

described earlier in this section, by taking a list of certain EF-hands (i.e. determined to contain the EF-hand motif by the 3 independent sources) and counting the number of sequences above the Z-score cutoff in the HMM database search that were not in the certain EF-hand list. Figure 15 shows that models built from small training sets have large numbers of false positives. Again, this effect dis-

1992). PROSITE denotes the positions used to generate the pattern EF-HAND. Ca-binding identifies the 6 residues involved in octahedrally coordinating the calcium ion (denoted by X, Y, Z, x, z and y). The oxygen atom at position y comes from the main-chain and so can be supplied by any amino acid. Sequences 1 to 27 are representatives of the various EF-hand subgroups in the June 1992 database of EF-hand sequences maintained by Kretsinger and co-workers (Nakayama *et al.*, 1992). These sequences are: CAMHS, *Homo sapiens* calmodulin; aACTGG, *Gallus gallus* α -actinin; VISININ, *G. gallus* visinin; TPP24CF, *Canis familiaris* p24 thyroïd protein; TPHUCS, *H. sapiens* skeletal troponin-C; TPAP1, *Astacus pontasticus* troponin-C-1; TCBP25, *Tetrahymena thermophila* TCBP-25; SPEC2A, *Strongylocentrotus purpuratus* spec2a; SCBPBL1, *Branchiostoma lanceolatum* SARC1; QUIDLN, *Loligo pealei* squidulin; MOHSCR, *H. sapiens* myosin (RLC-ventricle); MOHSA1, *H. sapiens* myosin (ELC-L1-skeletal); LPSIA, *Lytechinus pictus* α -Lps1; LAV1, *Physarum polycephalum* LAV1-2; EFH5, *Trypanosoma brucei* putative calcium binding protein; CVP, *B. lanceolatum* calcium vector protein; CRGHS, *H. sapiens* calmodulin-related gene; CMSE, *Saccharopolyspora erythraea* bacterial-CAM; CDPK, *Glycine max* calcium dependent protein kinase; CDC31, *Saccharomyces cerevisiae* cell division control protein 31; CALPLHS, *H. sapiens* calpain (light); CALCIB, *Bos taurus* calcineurin-B; CALBNGG, *G. gallus* calbindin; CALICE, *Caenorhabditis elegans* cal 1 gene; BCHS, *H. sapiens* β S-100 protein; AEQAV1, *Aequorea victoria* aequorin-1; and IF8, *Trypanosoma cruzi* flagellar calcium binding protein. 28 to 96 are the SWISS-PROT 22 sequences designated as EF-hands by our HMM (Z-score > 4.75) but not by all 3 other methods, PROSITE, PROFILESEARCH and the keyword search. Note that sequences identified as EF-hands by all 4 methods are not shown. 97 to 116 are the top 20 sequences below our cutoff of 4.75; 117 to 137 are sequences that were not classified as EF-hands by the HMM but were so by one or more (but not all) of the 3 other methods. B, Details on sequences 28 to 137 shown in the alignment (arranged in order of decreasing Z-score). NLL-score and Z-score are measures of how well the EF-hand HMM fits these database test sequence that were not present in the training set (see Results section (c) for more details). In HMM, PROFILESEARCH, Keyword and PROSITE + and - denote sequences that are and are not, respectively, classified as containing an EF-hand motif by the 4 specified methods. For PROFILESEARCH, Gribskov and HMM indicate results based upon profiles generated from four EF-hand sequences and our HMM alignments. T, N, P and ? in PROSITE have the same meaning as in Fig. 11. % indicates sequences which possess an EF-hand motif but do not bind calcium.

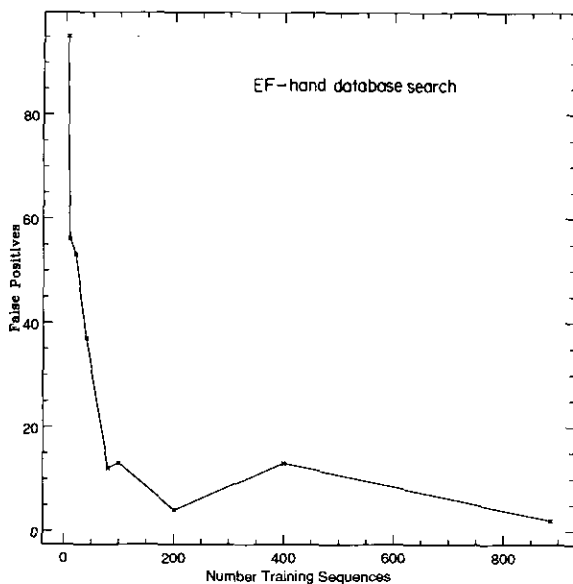


Figure 15. EF-hand database search false positives for models trained with 5, 10, 20, 40, 80, 100, 200, 400 and 885 sequences.

appears substantially when the training set size reaches about 100 sequences.

4. Discussion

A new method to model protein families using hidden Markov models has been introduced. The method is capable of tapping into the tremendous amount of statistical information contained in many unaligned sequences from the same family. For the cases of globins, kinases and EF-hands, the results have shown that by using this method, it is possible to obtain multiple alignments that mirror structural alignments, having only the unaligned primary sequences as input. The results have also shown that the model can be used successfully in database searches for putative analogs of sequences in a given protein family or domain. Finally, we believe that the model itself is a valuable tool for representing the family or domain.

The HMM method we have proposed requires that many sequences be available from the family or domain one wants to model. Since the number of sequences in the protein databases is growing rapidly, this may be less of a problem in the future, but it will always be a serious issue. Currently, only a relatively small number of sequences are available for most protein families and domains. For the globin family, we found that 400 sequences is certainly sufficient. Preliminary results indicate that 200 is enough, and even as few as 70 may suffice if they are chosen carefully from our database of 628 (70 chosen at random will be nearly all α - and β -chains). Our experiments using smaller numbers of EF-hand sequences for training, as described in Results section (c), show a similar trend. Using careful regularization, these numbers

might even be lowered further. However, there will be a limit on how small the number of available sequences can be if one hopes to obtain a reasonable model starting from a *tabula rasa*.

We believe that the answer to the problem of small training sets is to add more prior knowledge into the training process. One way to do this is by starting with a better initial model. We have performed several experiments in which we have started with a model obtained from a small set of aligned sequences, and then trained the model further using a larger set of unaligned sequences. These will be reported in a future paper. We find that this technique can often give better results. This also suggests that one application of HMMs may be in maintaining multiple alignments as the number of sequences in the alignment grows. Each time new sequences are added to a dataset of homologous sequences, we can begin with the HMM based on the alignment of the previous set of sequences, train it with the larger dataset that includes the new sequences, and then create a new multiple alignment for the larger dataset from this HMM. Not only will the new sequences be included in the new alignment, but the alignment of the old sequences may be improved by utilizing the statistical information present in the larger dataset.†

Another way to add more prior knowledge into the training process is to use a more sophisticated Bayesian prior. We are currently exploring the use of a prior on the probability distribution over the amino acids in a match state of the model consisting of a mixture of Dirichlet priors (Brown *et al.*, 1993). Using such a prior is like "soft-tying" the distributions in the states of the HMM. By soft-tying we mean a combination of the idea of tying states (see e.g. Rabiner, 1989), in which the number of free parameters is reduced by having groups of states all sharing the same distribution on the output alphabet (the 20 amino acids in this case), and the idea of soft weight sharing from Nowlan & Hinton (1992), in which the regularizer (in this case the prior for the distribution of amino acids) is also adaptively modified during learning. We have shown that this method can be used to estimate good EF-hand models using substantially fewer training sequences. Other types of more sophisticated priors can be obtained by switching from the alphabet of the primary sequences to a different representation based more on the structural or chemical properties of the amino acids in the sequence. We plan to explore these as well.

It is interesting to note that we have obtained quite good results in multiple alignment and database searching without using any special weighting schemes to make up for the statistical bias in our training sets (see e.g. Sibbald & Argos, 1990), or employing Dayhoff's matrix or any of its analogs (see e.g. Waterman, 1989) to take explicit mutation

† This point was suggested to us by an anonymous referee of one of our previous reports.

probabilities between amino acids into account. It also remains to be seen whether or not incorporating any of these extensions into the HMM approach will yield even better results.

We also believe that some of the errors made by our HMM models are due to the fact that these models are suboptimal, in the sense that their NLL-scores are not as low as they could be. This is because the EM procedure is not guaranteed to find the globally optimal model for a given training set. In other experiments, reported by Haussler *et al.* (1993), we trained an HMM for globins beginning with a model derived from the Bashford *et al.* (1987) alignment, and obtained a slightly lower NLL-score than any model from our experiments using EM on unaligned training sequences (208 compared to 210.3). Hence, we know that EM is not locating the globally optimal model in this case. An important open problem is to find a reliable way to prevent EM from getting stuck and returning a suboptimal solution.

Another issue is the adequacy of the hidden Markov model itself as a statistical model of the sequence variation within a protein family. Clearly an HMM provides at best a "first order" model of sequence variation. There are many kinds of interactions in proteins that are not easily modeled by HMMs, for example, pairwise correlations between amino acid distributions in positions that are widely separated in the primary sequence, but close in the three-dimensional structure (see e.g. Klinger & Brutlag (1993)). It would be very valuable to have more general models that incorporate such interactions while still remaining computationally tractable. We are currently exploring the potential of one model class of this type to capture the base-pairing in RNA families (Sakakibara *et al.*, 1993), and hope eventually to incorporate some of the features of these models into our protein models.

Finally, we are encouraged by the quality of the multiple sequence alignments generated by the HMMs and the accuracy of the database searches. For example, the kinase HMM is able to align correctly class III receptor tyrosine kinases which possess a domain that differs from other receptor tyrosine kinases by the insertion of a stretch of 70 to 100 residues (see the insertion between the D and E helices in sequence 8, the β -chain of the platelet-derived growth factor receptor, in Fig. 11A). With respect to the database discrimination tests, we would eventually like to see HMMs built for all the domains and families currently indexed by PROSITE expressions. In many cases, HMMs for subfamilies could be constructed automatically using the method described in Methods section (d). Once this is done, this might then lead to the construction of a simple "protein language parser" using HMMs. This parser could be constructed by connecting all these individual HMMs in parallel into a single large HMM with a global BEGIN and END state, and a transition from the END state back to the BEGIN state. In principle, this parser should be capable of finding all occurrences of each

of the PROSITE-indexed domains in a single long protein, using the Viterbi algorithm. The remaining portions of the sequence could be marked as "unknown". While this would not constitute a complete parse of the sequence, it would be very useful in providing some automatic annotation of new sequences, which is of critical importance as the rate of growth of the protein databases continues to accelerate. A related approach to protein annotation is given by Stultz *et al.* (1993), and a related HMM-based DNA parser for *E. coli* is described by Krogh *et al.* (1993b).

A comparative examination of the HMM produced kinase multiple sequence alignment and the crystal structure of the catalytic subunit of cAMP-dependent protein kinase (Knighton *et al.*, 1991) indicates a number of conserved residues in kinases of unknown structure that may be suitable for further experimental study (see Results section (b)). Results from our database discrimination tests suggest the presence of an EF-hand calcium-binding motif in a highly conserved and evolutionary preserved putative intracellular region of 155 residues in the α -1 subunit of L-type calcium channels which play an important role in excitation-contraction coupling (see Results section (c)). This region has been suggested to contain the functional domains that are typical or essential for all L-type calcium channels regardless of whether they couple to ryanodine receptors, conduct ions or both. Our EF-hand HMM indicates the following proteins may also possess this motif: chicken myosin light chain alkali (smooth muscle), bovine calpactin I light chain, *Arabidopsis thaliana* inorganic pyrophosphatase, rat placental calcium-binding protein and rat and bovine 1-phosphatidylinositol-4,5-bisphosphate phosphodiesterase III.

Although there are many experiments left to be done, based on our experience, we believe that HMMs and the EM algorithm have tremendous potential in the area of statistical modeling of biological macromolecules. Currently, most of this potential remains to be realized.

We thank Peter Brown, Søren Brunak, Richard Durbin, Harry Noller, Martin Vingron, Don Morris, and Michael Zuker for valuable comments on this work. Very special thanks to Richard Hughey for implementing our software on a MASPAP parallel machine and to MASPAP for providing computer time on their machine for some of these experiments, and very special thanks to Michael Gribskov for running his PROFILESEARCH program for the kinases and EF-hands so that we could compare the results to those found with the HMM. This work was supported by NSF grants CDA-9115268 and IRI-9123692, ONR grant N00014-91-J-1162, NIH grant number GM17129, and a grant from the Danish Natural Science Research Council. The full alignments and Z-score tables described in this paper are available in electronic form, and can be obtained by anonymous ftp from ftp.cse.ucsc.edu. Our HMM building program and other tools (written in C) will also be made available from the same ftp site.

References

- Abe, N. & Warmuth, M. (1990). On the computational complexity of approximating distributions by probabilistic automata. In *Proceedings of the 3rd Workshop on Computational Learning Theory*, pp. 52–66, Morgan Kaufmann, Rochester, NY.
- Allison, L., Wallace, C. S. & Yee, C. N. (1992). Finite-state models in the alignment of macromolecules. *J. Mol. Evol.* **35**, 77–89.
- Asai, K., Hayamizu, S. & Onizuka, K. (1993). HMM with protein structure grammar. In *Proceedings of the Hawaii International Conference on System Sciences*, pp. 783–791, IEEE Computer Society Press, Los Alamitos, CA.
- Bairoch, A. (1992). Prosite: a dictionary of sites and patterns in proteins. *Nucl. Acids Res.* **20**, 2013–2018.
- Baldi, P. & Chauvin, Y. (1993). A smooth learning algorithm for hidden Markov models. *Neural Computation*, in the press.
- Baldi, P., Chauvin, Y., Hunkapiller, T. & McClure, M. A. (1993). Hidden Markov models in molecular biology: new algorithms and applications. In *Advances in Neural Information Processing Systems 5* (Hanson, Cowan & Giles, eds), pp. 747–754, Morgan Kaufmann Publishers, San Mateo, CA.
- Barton, G. J. (1990). Protein multiple sequence alignment and flexible pattern matching. *Methods Enzymol.* **183**, 403–428.
- Barton, G. J. & Sternberg, M. J. (1990). Flexible protein sequence patterns: a sensitive method to detect weak structural similarities. *J. Mol. Biol.* **212** (2), 389–402.
- Bashford, D., Chothia, C. & Lesk, A. M. (1987). Determinants of a protein fold: unique features of the globin amino sequence. *J. Mol. Biol.* **196**, 199–216.
- Berger, J. (1985). *Statistical Decision Theory and Bayesian Analysis*. Springer-Verlag, New York.
- Bowie, J. U., Lüthy, R. & Eisenberg, D. (1991). A method to identify protein sequences that fold into a known three-dimensional structure. *Science*, **253**, 164–170.
- Brown, M. P., Hughey, R., Krogh, A., Mian, I. S., Sjölander, K. & Haussler, D. (1993). Using Dirichlet mixture priors to derive hidden Markov models for protein families. In *Proc. First Int. Conf. on Intelligent Systems for Molecular Biology* (Hunter, L., Searls, D., Shavik, J., eds), pp. 47–55, AAAI Press, Wash. D.C.
- Cardon, L. R. & Stormo, G. D. (1992). Expectation maximization algorithm for identifying protein-binding sites with variable lengths from unaligned DNA fragments. *J. Mol. Biol.* **223**, 159–170.
- Churchill, G. A. (1989). Stochastic models for heterogeneous DNA sequences. *Bull. Math. Biol.* **51**, 79–94.
- Dempster, A. P., Laird, N. M. & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *J. Roy. Statist. Soc. B*, **39**, 1–38.
- Dickerson, R. E. & Geis, I. (1983). *Hemoglobin: Structure, Function, Evolution and Pathology*, Benjamin/Cummings Pub. Co, Menlo Park, CA.
- Duda, R. O. & Hart, P. E. (1973). *Pattern Classification and Scene Analysis*, Wiley, New York.
- Everitt, B. S. & Hand, D. J. (1981). *Finite Mixture Distributions*, Chapman and Hall, London.
- Feng, D. F. & Doolittle, R. F. (1987). Progressive sequence alignment as a prerequisite to correct phylogenetic trees. *J. Mol. Evol.* **25**, 351–360.
- Garbers, D. L. (1992). Guanylyl cyclase receptors and their endocrine, paracrine and autocrine ligands. *Cell*, **71**, 1–4.
- Geman, S., Bienenstock, E. & Doursat, R. (1992). Neural networks and the bias/variance dilemma. *Neural Computation*, **4**, 1–58.
- Grabner, M., Friedrich, K., Knaus, H.-G., Striessnig, J., Scheffauer, F., Staudinger, R., Koch, W. J., Schwartz, A. & Glossmann, H. (1991). Calcium channels from *Cyprinus carpio* skeletal muscle. *Proc. Nat. Acad. Sci., U.S.A.* **88**, 727–731.
- Gribnikov, M., Lüthy, R. & Eisenberg, D. (1990). Profile analysis. *Methods Enzymol.* **183**, 146–159.
- Hanks, S. K. & Quinn, A. M. (1991). Protein kinase catalytic domain sequence database: identification of conserved features of primary structure and classification of family members. *Methods Enzymol.* **200**, 38–62.
- Hanks, S. K., Quinn, A. M. & Hunter, T. (1988). The protein kinase family: conserved features and deduced phylogeny of the catalytic domain. *Science*, **241**, 42–52.
- Haussler, D. & Krogh, A. (1992). Protein alignment and clustering. Presented at the conference Neural Networks for Computing.
- Haussler, D., Krogh, A., Mian, I. S. & Sjölander, K. (1992). Protein modeling using hidden Markov models: analysis of globins. Technical Report UCSC-CRL-92-23 University of California at Santa Cruz, Computer Science Dept., Santa Cruz, CA 95064.
- Haussler, D., Krogh, A., Mian, I. S. & Sjölander, K. (1993). Protein modeling using hidden Markov models: analysis of globins. In *Proceedings of the Hawaii International Conference on System Sciences*, IEEE Computer Society Press, Los Alamitos, CA.
- Hunter, T. (1991). Protein kinase classification. *Methods Enzymol.* **200**, 3–37.
- Jurka, J. & Milosavljevic, A. (1991). Reconstruction and analysis of human Alu genes. *J. Mol. Evol.* **32**, 105–121.
- Klinger, T. & Brutlag, D. (1993). Detection of correlations in tRNA sequences with structural implications. In *First International Conference on Intelligent Systems for Molecular Biology* (Hunter, L., Searls, D. & Shavlik, J., eds), AAAI Press, Menlo Park.
- Knighton, D. R., Zheng, J., Eyck, L. F. T., Ashford, V. A., Xuong, N.-H., Taylor, S. S. & Sowadski, J. M. (1991). Crystal structure of the catalytic subunit of cyclic adenosine monophosphate-dependent protein kinase. *Science*, **253**, 407–414.
- Krogh, A., Brown, M., Mian, I. S., Sjölander, K. & Haussler, D. (1993a). Hidden Markov models in computational biology: applications to protein modeling. Technical Report UCSC-CRL-93-32 University of California at Santa Cruz, Computer Science Dept., Santa Cruz, CA 95064.
- Krogh, A., Mian, I. S. & Haussler, D. (1993b). A hidden Markov model that finds genes in *E. coli* DNA. Technical Report UCSC-CRL-93-33 University of California at Santa Cruz, Computer Science Dept., Santa Cruz, CA 95064.
- Lander, E. S. & Green, P. (1987). Construction of multilocus genetic linkage maps in humans. *Proc. Nat. Acad. Sci., U.S.A.* **84**, 2363–2367.
- Lawrence, C. E. & Reilly, A. A. (1990). An expectation maximization (EM) algorithm for the identification and characterization of common sites in unaligned biopolymer sequences. *Proteins*, **7**, 41–51.
- Lindberg, R. A., Quinn, A. M. & Hunter, T. (1992).

- Dual-specificity protein kinases: will any hydroxyl do? *Trends Biochem. Sci.* **17**, 114–119.
- Lüthy, R., McLachlan, A. D. & Eisenberg, D. (1991). Secondary structure-based profiles: use of structure-conserving scoring table in searching protein sequence databases for structural similarities. *Proteins: Struct. Funct. Genet.* **10**, 229–239.
- Moncrief, N. D., Kretsinger, R. H. & Goodman, M. (1990). Evolution of EF-hand calcium-modulated proteins. I. Relationships based on amino acid sequences. *J. Mol. Evol.* **30**, 522–562.
- Nakayama, S., Moncrief, N. D. & Kretsinger, R. H. (1992). Evolution of EF-hand calcium-modulated proteins. II. Domains of several subfamilies have diverse evolutionary histories. *J. Mol. Evol.* **34**, 416–448.
- Nowlan, S. (1990). Maximum likelihood competitive learning. In *Advances in Neural Information Processing Systems* (Touretsky, D., ed), vol. 2, pp. 574–582, Morgan Kaufmann, San Mateo, CA.
- Nowlan, S. J. & Hinton, G. E. (1992). Soft weight-sharing. In *Advances in Neural Information Processing Systems 4* (Moody, Hanson & Lippmann, eds), Morgan Kaufmann Publishers, San Mateo, CA.
- Persechini, A., Moncrief, N. D. & Kretsinger, R. H. (1989). The EF-hand family of calcium-modulated proteins. *Trends Neurosci.* **12** (11), 462–467.
- Rabiner, L. R. (1989). A tutorial on hidden Markov models and selected applications in speech recognition. *Proc IEEE*, **77** (2), 257–286.
- Sakakibara, Y., Brown, M., Underwood, R., Mian, I. S. & Haussler, D. (1993). Stochastic context-free grammars for modeling RNA. Technical Report UCSC-CRL-93-16 University of California at Santa Cruz, Computer Science Dept., Santa Cruz, CA 95064.
- Sibbald, P. & Argos, P. (1990). Weighting aligned protein or nucleic acid sequences to correct for unequal representation. *J. Mol. Biol.* **216**, 813–818.
- Stultz, C. M., White, J. V. & Smith, T. F. (1993). Structural analysis based on state-space modeling. *Protein Sci.* **2**, 305–315.
- Subbiah, S. & Harrison, S. C. (1989). A method for multiple sequence alignment with gaps. *J. Mol. Biol.* **209**, 539–548.
- Tanaka, H., Ishikawa, M., Asai, K. & Konagaya, A. (1993). Hidden Markov models and iterative aligners. In *First International Conference on Intelligent Systems for Molecular Biology*, AAAI Press, Menlo Park.
- Taylor, W. R. (1986). The classification of amino acid conservation. *J. Theoret. Biol.* **119**, 205–218.
- Vingron, M. & Argos, P. (1991). Motif recognition and alignment for many sequences by comparison of dot-matrices. *J. Mol. Biol.* **218**, 33–43.
- Waterman, M. S. (1989). Sequence alignments. In *Mathematical Methods for DNA Sequences* (Waterman, M. S., ed.). CRC Press, Boca Raton, FL.
- Waterman, M. S. & Perlwitz, M. D. (1986). Line geometries for sequence comparisons. *Bull. Math. Biol.* **46**, 567–577.
- White, J. V., Stultz, C. M. & Smith, T. F. (1991). Protein classification by nonlinear optimal filtering of amino acid sequences. Unpublished manuscript.

Edited by F. Cohen

(Received 4 January 1993; accepted 23 September 1993)