# Hidden Markov Models in Computational Biology

## Applications to Protein Modeling

Anders Krogh[1][†], Michael Brown[1], I. Saira Mian[2]
Kimmen Sjölander[1] and David Haussler[1][‡]

[1]*Computer and Information Sciences*
[2]*Sinsheimer 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

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.

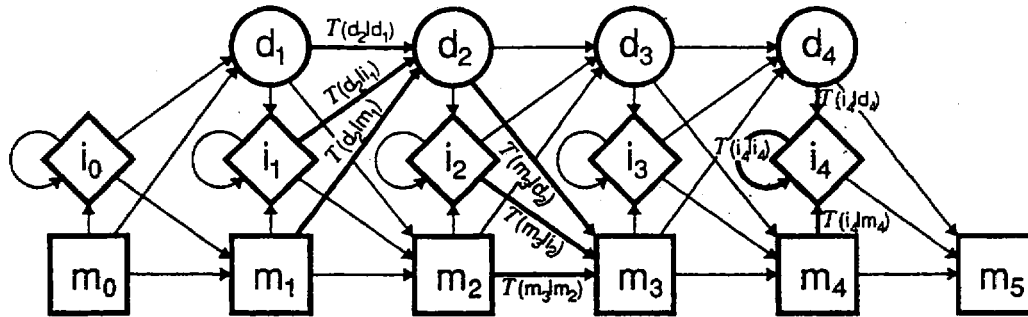structure along the lines they propose is required for this problem. We recently found that Asai *et al.* (1993) have applied HMMs to the problem of predicting the secondary structure of proteins.

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

according to the probabilities $\mathcal{T}(m_1|m_0)$, $\mathcal{T}(d_1|m_0)$, and $\mathcal{T}(i_0|m_0)$. 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 $\mathcal{T}(\cdot|m_1)$, where $\cdot$ indicates any possible next state. If this

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

Given a set of training sequences $s(1), \ldots, s(n)$, one can see how well a model fits them by calculating the probability that it generates them. This probability is simply a

probability estimate $\hat{\mathcal{T}}(r|q)$ 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 $\mathcal{T}(r|q)$ by $\hat{\mathcal{T}}(r|q)$ and $\mathcal{P}(x|q)$ by $\hat{\mathcal{P}}(x|q)$ for each $x$, $q$ and $r$. In MAP EM

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,

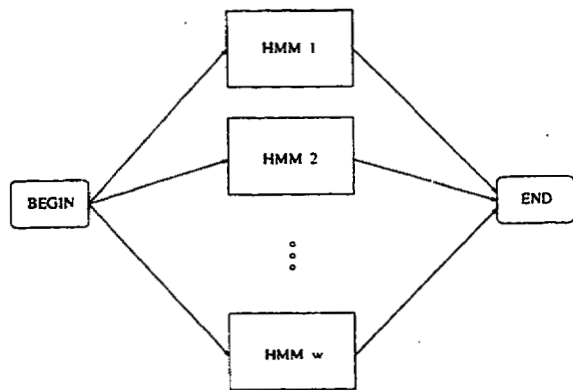$$\text{dist}(s, \text{model}) = \min_{\text{paths}} \{-\log \text{Prob}(s, \text{path}|\text{model})\}$$

**Figure 2.** HMM architecture for discovering subfamilies.

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

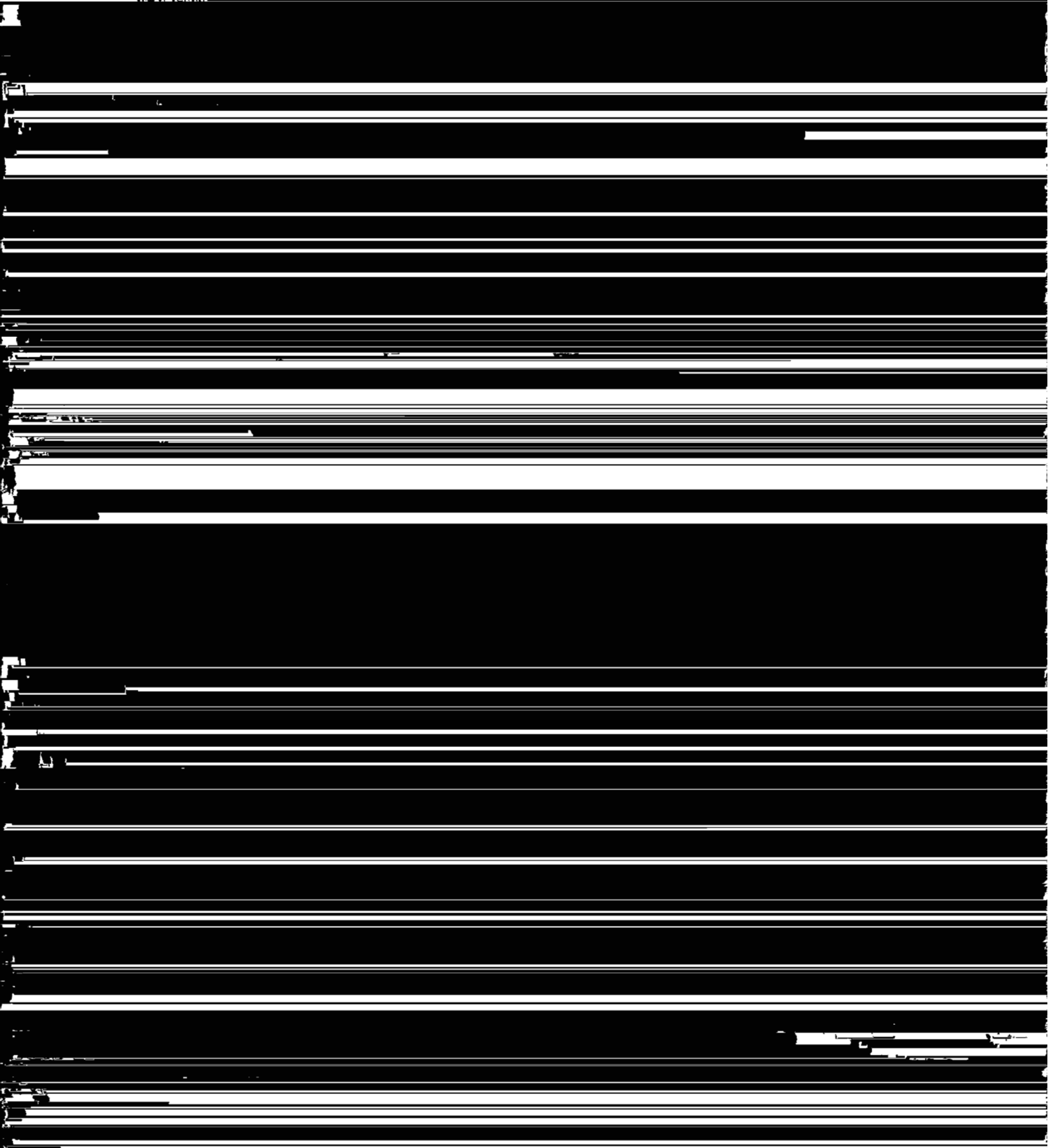### (e) *Modeling protein domains with an HMM*

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

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

comparing the distances or NLL-scores of 2 sequences

Most proteins tend to lie on a fairly straight line

### (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. $\mathcal{T}(m_{k+1}|m_k)$ 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‡ $\gamma_{del}$ 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 $\gamma_{ins}$ 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 $\alpha$ chains and two other subunits (usually $\beta$, $\gamma$, $\delta$ or $\theta$). 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 $\gamma_{del}$ and $\gamma_{ins}$ each to be 1/2.

```
Helix               AAAAAAAAAAAAAAAAA    BBBBBBBBBBBBBBBBBBCCCCCCCCCCC    DDDDDDDEE
HBA_HUMAN   ---------VLSPADKTNVKAAWGKVGA--HAGEYGAEALERMFLSFPTTKTYFPHF-DLS-----HGSA
HBB_HUMAN   --------VHLTPEEKSAVTALWGKV----NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNP
MYG_PHYCA   ---------VLSEGEWQLVLHVWAKVEA--DVAGHGQDILIRLFKSHPETLEKFDRFKHLKTEAEMKASE
GLB3_CHITP  ----------LSADQISTVQASFDKVKG------DPVGILYAVFKADPSIMAKFTQFAG-KDLESIKGTA
GLB5_PETMA  PIVDTGSVAPLSAAEKTKIRSAWAPVYS--TYETSGVDILVKFFTSTPAAQEFFPKFKGLTTADQLKKSA
LGB2_LUPLU  --------GALTESQAALVKSSWEEFNA--NIPKHTHRFFILVLEIAPAAKDLFS-FLK-GTSEVPQNNP
GLB1_GLYDI  ---------GLSAAQRQVIAATWKDIAGADNGAGVGKDCLIKFLSAHPQMAAVFG-FSG----AS---DP


Helix       EEEEEEEEEEEEEEEEEEE         FFFFFFFFFFFF  FFGGGGGGGGGGGGGGGGGGGGG
HBA_HUMAN   QVKGHGKKVADALTNAVAHV---D--DMPNALSALSDLHAHKL--RVDPVNFKLLSHCLLVTLAAHLPAE
HBB_HUMAN   KVKAHGKKVLGAFSDGLAHL---D--NLKGTFATLSELHCDKL--HVDPENFRLLGNVLVCVLAHHFGKE
MYG_PHYCA   DLKKHGVTVLTALGAILKK----K-GHHEAELKPLAQSHATKH--KIPIKYLEFISEAIIHVLHSRHPGD
GLB3_CHITP  PFETHANRIVGFFSKIIGEL--P---NIEADVNTFVASHKPRG---VTHDQLNNFRAGFVSYMKAHT--D
GLB5_PETMA  DVRWHAERIINAVNDAVASM--DDTEKMSMKLRDLSGKHAKSF--QVDPQYFKVLAAVIADTVAAG----
LGB2_LUPLU  ELQAHAGKVFKLVYEAAIQLQVTGVVVTDATLKNLGSVHVSKG---VADAHFPVVKEAILKTIKEVVGAK
GLB1_GLYDI  GVAALGAKVLAQIGVAVSHL--GDEGKMVAQMKAVGVRHKGYGNKHIKAQYFEPLGASLLSAMEHRIGGK


Helix       HHHHHHHHHHHHHHHHHHHHHHHHHHH
HBA_HUMAN   FTPAVHASLDKFLASVSTVLTSKYR------
HBB_HUMAN   FTPPVQAAYQKVVAGVANALAHKYH------
MYG_PHYCA   FGADAQGAMNKALELFRKDIAAKYKELGYQG
GLB3_CHITP  FA-GAEAAWGATLDTFFGMIFSKM-------
GLB5_PETMA  -----DAGFEKLMSMICILLRSAY-------
LGB2_LUPLU  WSEELNSAWTIAYDELAIVIKKEMNDAA---
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 *Chironomous 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 $\mathscr{P}$ 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

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

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

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 $>3\cdot5$ and either

**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.



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

(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 $\alpha$, a few $\zeta$ (an $\alpha$-type chain of mammalian embryonic hemoglobin), $\pi/\pi'$ (the counterpart of the $\alpha$ chain in major early embryonic hemoglobin P), and $\theta$-1 chains (early erythrocyte $\alpha$-like).

*Class 2.* 232 sequences, almost all $\beta$, a few $\delta$

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

**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



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

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.

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/

A (cont)

Fig. 11.

threonine kinases (PROTEIN KINASE ST... OKBOG) bovine cGMP-dependent protein kinase

| | 141 | 151 | 161 | 171 | 181 | 191 | 201 | 211 | 221 | 231 | 241 | 251 | 261 | 271 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Subdomain | ....................................................................................><-V-......................................................................>........<-V1a----------- | | | | | | | | | | | | | |
| PROSITE | ................................................................................................................................................................................................. | | | | | | | | | | | | | |
| X-ray | .......................................BB....B...BBB.....B.......BB..BBB................AAA....AA..AA..A.........................AA..A....AAA. | | | | | | | | | | | | | |
| X-ray | .......................................44....4...444.....5.......55..555................DDD....DD...DD...D........................EE...E....EEE. | | | | | | | | | | | | | |
| 1 CAPK-ALPHA | .......P......FL...V....KLE....F...SFKD..KSK.........LY..KVK..EYVPGG......E.KFS....KL...KR...I.....G..-..K...........FS.EP...N....AKF. | | | | | | | | | | | | | |
| 2 WEE1+ | .......D......NI...V....KLK....D...SVKH..GGF.........LY..KQV..ELCKKG......S.LDK.....FL..EE...Q....GqlS..K.........LD.EF...K....VVK. | | | | | | | | | | | | | |
| 3 TIK | .......Y......NI...V....QYKscwG...VDYD..PKNomedtsry12LF..IQK..EFCDKG......T.LEQ.....VK..KKr..N.....Q..S..K.........VD.KA...L....ILD. | | | | | | | | | | | | | |
| 4 SPK1 | .......P......KI...V....KLK....G...FYKD..TKS.........YY..KVK..EFVSGG......D.LKD.....FV..AA...K.....G..-..K.........VG.KD...A....GRK. | | | | | | | | | | | | | |
| 5 KSK1-K | .......P......FV...V....KLK....Y...AFQT..KGK.........LY..LIL..DFLAGG......D.LFT.....KL..SK...E.....V..-..K.........FT.EE...D....VKF. | | | | | | | | | | | | | |
| 6 PYT | .......D......KI...I....KLY....D...YKIT..DQY.........IY..KVK..ECGK-I......D.LKS.....VL..KK...K.....K..-..S.........ID.PN...E....KKS. | | | | | | | | | | | | | |
| 7 PKC-ALPKA | .......P......FL...T....QLK....S...CFQT..VDK.........LY..FVK..ETVKGG......D.LKY.....KI..QQ...V.....G..-..K.........FK.EP...Q....AVF. | | | | | | | | | | | | | |
| 8 PDGFR-B | .......L......KV...V....KLL....G...ACTK..GGP.........IY..IIT..EYCKYG......D.LVD.....YL..KK...K.....K..K..Tflqhhsdk9KlS.YK...D....LVG. | | | | | | | | | | | | | |
| 9 PBS2 | .......P......VI...V....DFY....G...AFFI..KGA.........YY..KCK..EYKDGG......S.LDK.....IY..DE...Ssei..G..G..-.........ID.EP...Q....LAF. | | | | | | | | | | | | | |
| 10 KIK1 | .......P......FV...V....KLV....K...VVSY..KDK.........IF..LQL..DYCKKG......D.LSL.....FL..SE...L....G11Q..V.........KD.PF...K....VVK. | | | | | | | | | | | | | |
| 11 KCK1 | .......P......KI...V....KLQ...Y....FFTK..LSPqdak......VVqALAK..ECLP-E......T.LQI.....EI...KKyrtK....K..L..E.........KP.LK...V....IKL. | | | | | | | | | | | | | |
| 12 IKS.R | .......N......KV...V....KLL....G...VVSK..CQP.........TL..FVK..ELKAKG......D.LKS.....YL..KS...L.....K..P..Kasnapgrpp.PT.LQ...E....KIQ. | | | | | | | | | | | | | |
| 13 KSYK | .......P......AI...L....PLL....D...LKVV..SGV.........TC..LVL..PKYQA-......D.LYT.....YL..SK...K.....L..K..P.........LG.KP...Q....IAA. | | | | | | | | | | | | | |
| 14 EKK1 | .......E......KV...I....GIK....D...ILKA..PTLsanrd....VY..IVQ..DLKE-T......D.LVK.....LL..KS...Q.....Q..-..-.........LS.KD...K....ICY. | | | | | | | | | | | | | |
| 15 EGFK | .......P......KV...C....KLL....G...ICLT..S-T.........VQ..LIT..QLKPFG......C.LLD.....VV..KE...K.....K..D..K.........IG.SQ...Y....LLK. | | | | | | | | | | | | | |
| 16 KCK | .......K......KI...I....KLE....G...VISK..YKP.........KK..IIT..ETKEKG......A.LDK.....FL..KE...K.....D..G..E.........FS.VL...Q....LVG. | | | | | | | | | | | | | |
| 17 DPYK1 | .......P......KV...V....QFK....G...ACTAttFDK.........KG..IVT..EYKGGG......S.LKD.....FL..TD...K.....F..-..K.........LL.KDqbI....KLK. | | | | | | | | | | | | | |

|  | 281 | 291 | 301 | 311 | 321 | 331 | 341 | 351 | 361 | 371 | 381 | 391 | 401 | 411 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Subdomain | ----------------------------------------------------->.<-VIb---------------------------------------------->.<-VII--------+--------+---------------------------------------->.<-V |
| PROSITE | .................................................B...BBBBBBBB.B.BB..B............................................................................. |
| X-ray | AAA..........A......AAA..AAAAA.A.........A..BB...BB.......B.BB..B.B....B.........B.....BBB.........B...BBB.........B.............. |
| X-ray | EEE..........E......EEE..EEEEE.E.........E..66...66.......7.77..7.7....8..........8....BBB.......:9....999.........9.............. |
| 1 CAPK-ALPHA | YAA..........Q.......IVL..TFEYL.H.........S.L.DL...IYRDLKPE.N.LL...I.D...QQG...........Y......IQVT.D..FGF.A...KKV.........KG...-...--R.T.. |
| 2 WEE1+ | ILV..........E.......VAL..GLQFI.H.........H.R.HY..VHLDLKPA.N.VH..I.T..FEG............T......LKIG.D..FGL.A...S--.........VW...P...VPR.G.. |
| 3 TIR | LYE..........Q.......IVT..GVEYI.H.........S.R.GL...IHRDLKPG.N.IF..L.V...DER...........H......IKIG.D..FGL.A...TAL.........EK...D...GRS.R.. |
| 4 SPK1 | ISR..........Q.......ILT..AIRVI.H.........S.R.GI...SHRDLKPD.N.IL..I.E...QDDpv........L......VKIT.D..FGL.A...K-V.........QG...N...GSF.N.. |
| 5 KSR1-F | YLA..........E.......LAL..GLDKL.H.........S.L.GI...IYRDLAPE.N.IL..L.D...EEG...........H......IKLT.D..FGL.S...REA.........ID...H...EKR.A.. |
| 6 PYT | YVK..........N.......KLE..AVNTI.H.........Q.N.GI...VHSDLKPA.N.FL..I.V...-DG...........N......LKLI.D..FGI.A...NQR.........QP..Dtt.SVV.A.. |
| 7 PKC-ALPHA | YAA..........E.......ISI..GLFFL.H.........R.R.GI...IYRDLKLD.N.VH..L.D...SEG...........H......IKIA.D..FGN.C...REN.........NN...D...GVT.T.. |
| 8 PDGFR-B | FSY..........Q.......VAN..GNEFL.A.........S.K.NC..VHRDLAAR.N.VL..I.C...EGA...........L......VKIC.D..FGL.A...RDI.........NN...Ds.NYI.S.. |
| 9 PBS2 | IAY..........A.......VIN..GLRELAE.........Q.N.NI..IHRDVAPT.N.IL..C.Sa..NQG...........T......VKLC.D..FGV.S...GNL.........VA..-...-SL.A.. |
| 10 NIK1 | KLF..........Q.......LTQ..ALNFI.H.........L.L.EF..VHLDVKPS.N.VL..I.T...RDG...........N......LKLG.D..FGL.A...TSL.........PY..-...SSR.V.. |
| 11 NCK1 | YTY..........Q.......IAR..GRLYL.H.........G.L.GV..CHRDIKPS.N.VL..V.Dp..ETG...........V......LKIC.D..FGS.A...KKL.........EH...N...QP-.S.. |

```
              421      431      441      451      461      471      481      491      501      511      521      531      541      551
Subdomain   III------------------------><-II------------------------------------------------------------------------>...<-X------------------
PROSITE     ....................................................................................................................................
X-ray       .........................................................A.AAAAAAAAAAAAAAA.........................................................
X-ray       .........................................................................F.FFFFFFFFFFFFFFF........................................G....A.......
X-ray       ....................................................................................................................................G....G.......
 1 CAPK-ALPHA   .......VTLCGT.P...EY.LAPE..IIL.........SK.........G-YEK...A.VDWVALGVLIYEMAA.G......YP...P..........F.......F....A..-DQP.....I...Q.......
 2 VEE1*        .......MEREGD.C...EY.IAPE..VLA.........MR.........L-YDK...P.ADIFSLGITVFEAAAnI...VL...Pdsgqsvqh17*......P....R..LSGT....D....M.......
 3 TIK          .......TRATGT.L...QY.MSPE..QLF.........LK.........M-YGK...E.VDIFALGLILAELL-.-.....MT...G...........F.......T....E..-SEK....I....K.......
 4 SPRJ         .......RTFCGT.L...AY.VAPE..VIR.........GRdtsvsdpe12MEYSS...L.VDMVSMGCLVYVILT.G......ML...P..........F.......S....G..-STQ....D....Q.......
 5 KSK1-M       .......TSFCGT.V...EY.MAPE..VVR.........RQ.........G-MTM...S.ADMVSYGVLR-----.-.....--...-..........-.....-..G..KDRK....E....T.......
 6 PYT          .......SSQVGT.V...MY.MPPE..AIKdmaserenghSK.........SRISP...R.SDVMSLGCILYYMTY.G......KT...P..........F.......Q....Q..IIMQi...S....K.......
 7 PKC-ALPHA    .......RTFCGT.P...DY.IAPE..IIA.........YQ.........P-YGK...S.VDMVAVGVLLYEMLA.G......QP...P..........F.......D....G..-EDE....E....Q.......
 8 PDGFR-B      .......RGSTFL.P...LK.WRAP..ESI.........FM.........SLYTT...L.SDVMSFGILLMEIFT1G.....GT...P..........Y.......P....E..LPMM....E....Q.......
 9 PBS2         .......XTWIGC.Q...SY.MAPE..RIR.........SLmpdr......ATYTV.<.Q.SDIVSLGLSILEMAL.G.....RY...P..........Y.......P....-ETY.....Dnifa0.......
10 RIK?         .......D-LEGD.R...YY.IAPE..ILA.........SM.........M-YCK...P.ADVYSLGLSMIEAATaV....VL...Psngvewqr16L........P....M..LKDL....L....L.......
11 MCK1         .......IGYYCS.R...FY.RAPE..LII.........GC.........TQYTT...Q.IDIWGLGCVMGEALI.G.....RA...I..........F'......Q....G..QEPL....L....Q.......
12 CMS.R        .......GGRGLL.P...VR.WRAP..EGL.........RD.........GVFTT...S.SDRMSFGVVLMEITS1A....EQ...P..........P.......Q....G..LSKE....Q....V.......
13 MSVR         .......YGIAGT.I...DT.MAPE..VLA.........GD.........P-YTT...T.VDIMSAGLVIFETAVaM....AS...L..........FsaprgphR....G..-PCD....S....Q.......
14 ERA1         .......TEYVAT.R...MY.MAPE..IRL.........MS.........RGYTA...S.IDIVSVGCILAEMLS.K.....RP...I..........F.......P....G..RMYL....D....Q.......
15 EGFR         .......E-GGKV.P...IK.WMAL..ESI.........LK.........RIYTH...Q.SDVMSYGVTVVELMTtG....SR...P..........Y.......D....G..IPAS....E....I.......
16 ECK          .......TSGGKI.P...IR.VTAP..EAI.........SY.........RKFTS...A.SDVMSFGIVRMVEVRTyG...ER...P..........Y.......M....E..LSMK....E....V.......
17 DPYK1        .......TGSVGC.I...PY.RAPE..VFK.........GD.........S-MSE...R.SDVYSVGRVLFELLT.S.....DE...P..........Q.......Q....D..MRPK....R....M.......
18 CLK          .......GTLVST.R...MY.RAPE..VIL.........AL.........G-MSQ...P.CDMVSIGCILIEYYL.G.....FT...V..........F.......G....T..MDSR....E....M.......
19 CDC28S       .......TREVYT.L...MY.MSPE..VLL.........CS.........ARYST...P.VDIVSIGTIFAELAT.K.....RP...L..........F.......M....G..DSEI....D....Q.......
20 CAMII-ALPHA  .......FGFAGT.P...GY.LSPE..VLR.........RD.........P-YGK...P.VDLMACGVILYILLV.G.....YP...P..........F.......M....M..-EDQ....M....R.......
21 C-SRC        .......RQGAKF.P...IK.WTAP..EAA.........LY.........GRFTI...K.SDVMSFGILLTELTTaG....RV...P..........Y.......P....G..MVPR....E....V.......
22 C-RAF        .......EQPTGS.V...LW.MAPE..VIR.........RQdan......P-FSF...Q.SDVYSVGIVLYELMT.G.....EL...P..........Y.......S....K..IMMR....D....Q.......
23 KLSR_HUMAN   .......REGAKF.P...IK.VTAP..EAI.........MY.........GTFTI...K.SDVMSFGILLTEIVThG....AI...P..........Y.......P....G..MTMP....E....V.......
24 KLSR_MOUSE   .......REGAKF.P...IK.VTAP..EAI.........MY.........GTFTI...K.SDVMSFGILLTEIVThG....AI...P..........Y.......P....G..MTMP....E....V.......
25 ARAR_HUMAN   .......RASVGT.M...GY.MAPE..VLQ.........RG.........VAYDS...S.ADMFSLGCMLFKLLR.G.....MS...P..........F.......R....Q..MKTK....DRR..E.......
26 ARAR_BOVIN   .......RASVGT.M...GY.MAPE..VLQ.........RG.........VAYDS...S.ADMFSLGCMLFKLLR.G.....MS...P..........F.......R....Q..MKTK....DRR..E.......
27 BYR1.SCHPO   .......QTFVGT.S...TY.MSPE..RIR.........GG.........R-YTV...K.SDIMSLGISIIELAT.Q....EL...Pus.........F.......S....M..IDDSigiI.D....L.......
28 CTGR.ARBPU   .......GERAKL.A...RKIWTAP..ERL.........REgkamhp...G-GTP...R.GDIYSFSIILTERYS.R....DE...P..........F.......Mem..D..LELA....D....I.......
29 AMPR.RAT     .......TLFARR.L...--.VTAP..ELLrmaspp...AR.........C--GSQ..A.GDVYSFGIILQEIALrS....GV...F.......Yveg..L....D..LSPK....E....I.......
30 AMPA.HUMAN   .......Y-ARRL.-...--.VTAP..ELL.........RRaappv....R-GSQ...A.GDVYSFGITLQEIAL.R....SG...V..........F.......Hveg1D..LSPK....E....I.......
31 AMPR.HUMAN   .......Y-ARR.L...VT.-APE..LLS.........GM.........PLPTTgmqR.ADVYSFGIILQEIAL.R....SG...P..........F.......Yleg1D..LSPK....E....I.......
32 AMPA.MOUSE   .......TLFARR.L...--.VTAP..ELLrmaspp...AR.........C--GSQ..A.GDVYSFGIILQEIALrS....GV...F.......Yveg..L....D..LSPK....E....I.......
33 AMPR.RAT     .......Y--ARR.L...VT.-APE..LLS.........GM.........PLPTTgmqR.ADVYSFAIILQEIAL.R....SG...P..........F.......Yleg1D..LSPK....E....I.......
34 CYGS.STRPU   .......GDMAKLaR..QL.VTSP..ERL.........RQegompta...C--SP...Q.GDIYSFAIILTELYS.R....OE...P..........F.......M....EmeKDLA...D....I.......
35 VPSF.YEAST   flIytd..TSRRAT.-..CY.LAPE..RFM.........SKlyqdgkmn.GRLTK...E.MDIFSLGCVIAEIFAeG...RP...I..........F.......-...-..-ML....S....Q.......
36 MSER.RAT     .......RDL---.-...--.VTAP..ERL.........RQ.........ATISQ...K.GELYSFSIIAQEIIL.R....RE...T..........F.......Y....T..LSCR....D....Qmeh....
37 MSER.HUMAN   .......RDL---.-...--.VTAP..ERL.........RQ.........AMISQ...K.GDVVSYGIIAQEIIL.R....KE...T..........F.......Y....T..LSCRdrm..E.......
38 KR2.VZVD     .......FRLVLS.M...GY.MQPP..EIL.........LDyingtglt1SQRVGL..A.IDLYALGQALLEVILiG...RL...Pgqlpisvh14Y.....Y....G..MRLS....P....D.......
39 KR2.HSV11    qtslqe20MTLVG-.M..GY.MQPP..ELLwkylmmer15LK.........MDVGL...A.VQLYALGQTLLELVVaV...YV...Apslgvpvtr.F.....P....G..----....-...-.-.......
40 KR1.HSV11    .......VMPIGT.E...AY.ASPE..RSR.........DRvpdrpdsa12GTMGA..G.I----------.-.....RE...P..........M1i...R....G..OGYR....A....H.......
41 KR2.EBV      .......KSSKGR.Q...LY.R--L..YCQ.........RE.........P-FSI...K.EDTY----------.-.....KP...Lcllskcyi24-.....-....G..AGTA....L....R.......
42 KRR2.VACCV   ynedmil9MKLGAT.V..SR.AGDL..ERL.........GY.........C----..-.--------KIEMFG.G....KL...P..........V.......-...-..KME....S....G.......
43 KRR2.VACCC   ynedmil9MKLGAT.V..SR.AGDL..ERL.........GY.........C----..-.--------KIEMFG.G....KL...P..........V.......-...-..KME....S....G.......
44 AK1.ECOL1    .......DPRVVS.A...AK.RIDE..IAF.........AE.........A----..-.--------AEKATsGRkvLhpAT..L......L......P....A..-VRS....D....I.......
45 PSP.MOUSE    ........-GIDL.T...VP.LAGE..ASL.........VL.........PFIGK...T.VD1-SVSLDLIMSLS.I....KT...MaqcglpevI4-.....-...-..-SMT....D....R.......
46 DHOM.BACSU   .......DVEGLD.A...AR.RRA-..ILA.........RL.........G-FSR...K.VDLE---------.-.....--...-..........-.dvkvbgiS...Q..ITDE....S....I.......
47 FLIG.BACSU   lqqebpq.T----.-...-..--MAL..ILSy.......LD.........PVQ-..-.-.----AGQILSELR-.-.....--...-..........-...-..P....E..-VQA....I....V.......
48 CALQ.RABIT   .......RED---.-...--.--E..VIE.........YD.........GEFSAd...-.--------TLVEFL.-.....--...-..........-...-..-...-..LDVL....E....D.......
49 MVIK.PODAK   .......GSLRST.A...QL.ISTE..LVL.........SS.........A----..-.--------ILLVIRLT.G....SL...MlsvnieeqI4F......P....L..-SPV....F....I.......
50 RVVA.ECOL1   .......CALVRL.PgiGK.KTAE..RLIvemkdrft11GD.........L-FTP...A.ADL-------VLTSpA.....SP...A..........I.......D....D..AEQE....A....V.......
51 ?1SR.MSV6M   .......TRDACC.-..mK.VIAFnuTL1.........GI.........I----..-.---------------.-.....--...-..........F.......Y....R..-DVV....F....I.......
```

```
52 KRF1.VACCC   .......S-----.-...--.-----.--A.........LM.........D-FDF...S.--------------.-.....--...-..........-..-..-...-..-QVA....G....I.......
53 GL9T.MCHVA   .......FFWAGL.R..RY.CMSE..LSA.........LG.........MVLGF...-.-------CLK----.-.....--...-..........-...-..R..LLDR....R....G.......
54 RKA6.ACIBA   .......DD----.-...--.--I....DQ.........DDFDT...E.---------------.-.....--...-.......W.....G...D..MKTY....LslvnE.......
55 RKA6.ECOL1   rrthailoTTMAGL.P..ER.GSIE..ACVvdvddfdke.RE.........G-VTA...Eq-------------VGEARM.R....LL...P..........L......A..P..OPVV....T....M.......
56 KGPR.BOVIE   .......RQSAST.-...--.----.-LQ.........GE.........P-RTK...R.--------------.-.....--...-..........-...-..Q..A.ISAE....P....T.......
57 FGBR.CHICK   .......GGLRT.R..MY.RISE..ILM.........GG.........VVISp...-.---------------.-.....--...-..........-...-..-...MRRhlsmoIY...V.......
```

```
         561   571   581   591   601   611   621   631   641   651   661   671
Subdomain ----------------------------------><-XI------------------------------------>.........
PROSITE   ..................................................................................
X-ray     A.AAAAA...A...........................AA.....AA....A..AAAAA............AAAAAA............
X-ray     G.GGGGG...G.............................HH.....HH....H..HHHHH............IIIIII.>.........
X-ray     .......................................................................IIIIII.>.........
```

| # | Label |
|---|-------|
| 1 | CAPK-ALPHA |
| 2 | WEE1+ |
| 3 | TIK |
| 4 | SPK1 |
| 5 | RSK1-N |
| 6 | PYT |
| 7 | PKC-ALPHA |
| 8 | PDGFR-B |
| 9 | PBS2 |
| 10 | NIK1 |
| 11 | NCK1 |
| 12 | INS.R |
| 13 | HSVK |
| 14 | EXK1 |
| 15 | EGFR |
| 16 | ECK |
| 17 | DPYK1 |
| 18 | CLK |
| 19 | CDC2HS |
| 20 | CAMII-ALPHA |
| 21 | C-SRC |
| 22 | C-RAF |
| 23 | KLSR_HUMAN |
| 24 | KLSR_MOUSE |
| 25 | AKRB_HUMAN |
| 26 | AKRB_BOVIN |
| 27 | BYR1_SCHPO |
| 28 | CTGR_AMPPU |
| 29 | AMPA_RAT |
| 30 | AMPA_HUMAN |
| 31 | AMPB_HUMAN |
| 32 | AMPA_MOUSE |
| 33 | AMPB_RAT |
| 34 | CTGS_STRPU |
| 35 | VPSF_YEAST |
| 36 | WSER_RAT |
| 37 | WSER_HUMAN |
| 38 | RA2_VZVD |
| 39 | RA2_HSV11 |
| 40 | RR1_HSV11 |
| 41 | RR2_EBV |
| 42 | RRB2_VACCV |
| 43 | RRB2_VACCC |
| 44 | AK3_ECOLI |
| 45 | PSP_MOUSE |
| 46 | DHOM_BACSU |
| 47 | FLIG_BACSU |
| 48 | CALQ_RABIT |
| 49 | NUIR_PODAN |
| 50 | RDVA_ECOLI |
| 51 | U15R_NSV6U |
| 52 | RRF1_VACCC |
| 53 | UL97_HCMVA |
| 54 | RRA6_ACIBA |
| 55 | RRA6_ECOLI |
| 56 | XGPB_BOVIN |
| 57 | EGFR_CHICK |
| 58 | KRA1_ECOLI |
| 59 | KDTK_DROME |
| 60 | KPCG_HUMAN |

A *(cont)*

**Fig. 11.**

negatives (41 to 43, 51 to 53) of which the first three fall immediately below our kinase cutoff. For PROFILESEARCH, there are 12 false negatives (23 to 26, 35, 38 to 41, 51 to 53) but it should be recalled that eight of these (those indicated by $ in Fig. 11B) do not appear in the results obtained from searching SWISS-PROT 25 provided to us by M. Gribskov (personal communication). We suspect that at least four (23 to 26) would be correctly classified as kinases by PROFILESEARCH leaving an estimate of three to eight false negatives. In the case of PROSITE, using our assumption of a kinase to be a true positive (T) sequence for any one of the three patterns, there are three false negatives (39, 42 to 43). However, the actual performance of the PROSITE patterns themselves is much worse; scans of SWISS-PROT 22 with each of the patterns PROTEIN_KINASE_ATP, PROTEIN_KINASE_ST and PROTEIN_KINASE_TYR individually yield 40, 2 and 3 false negatives, respectively.

The difficulty in quantifying the precise number of false positives and false negatives produced by the database discrimination tests may be illustrated by employing an alternative mechanism for assessing the number of false negatives. If simply

| 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_HSVI1 | 518 | 1548.949 | 9.042 | + | - | + | N | - | - |
| 40 KR1_HSVI1 | 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 PROFILE-SEARCH.

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 entertoxin 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). NLL-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

```
                1      11     21      31      41     51      61      71
Structure   ........H..H.HHHHH.H.H.........H....LL...LLL.LLLL...H.HH.HHHHHH.........
PROSITE     ..................*....**...***.****...*.**...........
Ca-binding  .................X....Y....Z..y.x.......z..............

1  CAMS      .........E..F.REAFS.L.F.........D....KD...GDG.TITT....K.EL.GTVNRSL-........:.
2  aACTGG    .........E..F.RASFN.N.F.........D....RR...RTG.NRDC....E.DF.RACLISR-.........
3  VISININ   .........E..L.SRMYE.G.F.........Qr...QC...SDG.RIRC....D.EF.ERIYGHF-.........
4  TPP24CF   .........G..L.ARFFR.R.L.........D....RD...RSR.SLDS....R.EL.QAGLAEL-.........
5  TPNUCS    .........E..F.KAAFD.N.F.........D....AD...GGG.DISV....K.EL.GTVNRNL-.........
6  TPAFI     .........A..L.QRAFD.S.F.........D....TD...SKG.FITP....E.TV.CIILRNH-.........
7  TCBP25    .........V..A.RRIFE.N.Y.........D....KG...RKG.RIEN....T.DC.VPRITEA-.........
8  SPEC2A    .........L..F.KSSFR.S.E.........D....TD...CDG.RITS....E.EL.RAAFRSI-.........
9  SCBPBL1   .........R..I.KFTFD.F.FI........D....YN...KDG.SIQW....E.DF.EENIRRY-.........
10 QOIOLN    .........E..I.KDAFD.N.F.........D....ID...GDG.QTTS....N.EL.RSVRKSL-.........
11 NOHSCR    .........E..F.REAFT.T.N.........D....QN...RDG.FIDR....N.DL.RDTFAAL-.........
12 NOHSA1    .........E..F.REAFL.L.F.........D....ST...GDS.KIIL....S.QY.GDVLRAL-.........
13 LPSJA    *.........A..L.KQEFR.DaY........D....TN...KDG.TVSC....A.EL.VKLNNVT*.........
14 LAVI      .........A..L.VADFR.R.I.........D....TN...SNG.TLSR....R.EF.RENFVRL-.........
```

```
18 CHSE      ..........R..L.RKRFD.R.M.........D....FD...GNG.ALER....A.DF.EREAQHI-.........
19 CDPR      ..........G..L.RELFR.N.I.........D....TD...NSG.TITF....D.EL.RDGLRRV-.........
20 CDC31     ..........E..I.VEAFS.L.F.........D....RR...NDG.FLOY....N.EL.RVANRAL-.........
21 CALPLRS   .........T..C.RSNVA.V.F.........D....SD...TTG.KLGF....E.EF.RYLWHHI-.........
22 CALCIB    ..........R..L.GRRFR.R.L.........D....LD...NSG.SLSV....E.EF.NS-LPEL-.........
23 CALDNGG   .........Q..F.FEIVN.N.Y.........D....SD...GNG.YNDG....K.EL.QNFIQEL-.........
24 CALICE    .........E..F.REAFR.N.F.........D....RD...GNG.TIST....K.EL.GIARRSL-.........
25 BCHS      .........A..L.IDVFN.Q.V........Sg...RE...GDRRRLKK....S.EL.RELIVNE-.........
26 AEQAVI    ..........R..N.RNRFN.F.L.........D....VN...NNG.RISL....D.EN.VYRASDI-.........
27 JF8       ..........R..R.IELFR.N.F.........D....RN...ETG.RLCY....D.EV.NSGCLEV-.........
28 CALM_ASPNI adaiteeqveE..Y.REAFS.L.F.........D....RD...GDG.QITT....K.EL.GTVNRSL-gqmpeee109
29 NLEJ_HUMAN apkkdvk129D..F.VEGLR.V.F.........D....RE...GNG.TVNG....A.EL.RNVLATL-gekmkeee35
30 NLEJ_RABIT apkkdvk127D..F.VEGLR.V.F.........D....RE...GNG.TVNG....A.EL.RNVLATL-gekmkeee35
31 NLEV_HUMAN apkkpep130D..F.VEGLR.V.F.........D....RE...GNG.TVNG....A.EL.RHVLATL-gerltede35
32 NLEC_CHICK ppkkpep129D..F.VEGLR.V.F.........D....RE...GNG.TVNG....A.EL.RHVLRTL-gerltede35
33 NLEV_RAT   apkkpep135D..F.VEGLR.V.F.........D....RE...GNG.TVNG....A.EL.RHVLRTL-gerltede35
34 NLEL_CHICK pkkdvkk126D..F.VEGLR.V.F.........D....RE...GNG.TVNG....A.EL.RHVLRTL-gekmkeee35
35 NLEJ_RAT   apkkdvk124D..F.VEGLR.V.F.........D....RE...GNG.TVNG....A.EL.RNVLATL-gekmkeee35
36 NLEJ_MOUSE apkkdvk123D..F.VEGLR.V.F.........D....RE...GNG.TVNG....A.EL.RNVLATL-gekmkeee35
37 NLEF_HUMAN apkkpep132D..F.VEGLR.V.F.........D....RE...SNG.TVNG....A.EL.RNVLATL-gekmteee35
38 NLEF_RAT   ppkkpep120D..F.VEGLR.V.F.........D....RE...SNG.TVNG....A.EL.RNVLATL-gekmeeee35
39 NLEF_MOUSE ppkkpep120D..F.VEGLR.V.F.........D....RE...SNG.TVNG....A.EL.RNVLATL-gekmeeae35
40 NLEX_CHICK mplkkpd121D..F.VEGLR.V.F.........D....RE...GNG.LVRG....A.EL.RNVLVTL-gekmteee35
41 NLES_HUMAN sfaadqia85D..F.VEGLR.V.F.........D....RE...GNG.TVNG....A.EL.RNVLATL-gekmkeee35
42 NLEY_HUMAN mppkkdv144D..Y.LECFR.V.F.........D....RE...GNG.KVNG....A.EL.RNVLTTL-gekmkeee35
43 NLES_RABIT sfaadqia85D..F.VEGLR.V.F.........D....RE...GNG.TVNG....A.EL.RNVLATL-gekmkeee35
44 NLES_RAT   sfaadqia85D..F.VEGLR.V.F.........D....RE...GNG.TVNG....A.EL.RNVLATL-gekmkeee35
45 NLES_MOUSE sfaadqia85D..F.VEGLR.V.F.........D....RE...GNG.TVNG....A.EL.RNVLATL-gekmkeee35
46 NLES_CHICK sfapdqid85D..F.VEGLR.V.F.........D....RE...GNG.TVAG....A.EL.RNVLATL-gekmtee35
47 AACT_HUMAN mdhydsq749E..F.RASFN.N.F.........D....RD...NSG.TLGP....E.EF.RACLISL-gydigpd314
```

```
79 CAP2_HUMAN   magiaakS75T..C.KIRVD.R.L..........D....SD...GSG.KLGL....K.EF.YILVTKI-qkyqkiyr96
80 ADGL_PIG     makergi1571..L.QERBR.E.I..........D....YD...GSG.SY5L....A.EM.LRAGATI-vpllv11548
81 SCPA_PENSP   ayewdar103F..I.ABQFK.A.I..........D....YB...GDG.RYGL....D.EYrLDCITRS-afaevkeiS9
82 SCPB_PENSP   ayewdarv59L..V.HEIAE.L.A..........D....FB...KDG.EYTV....D.EF.KQAVQKR-ckgkafaiO4
83 IPTR_ARATH   maeikde169E..I.RRFFE.D.T..........K....RK...ERK.KVDV....E.AF.LPAQAAI-daikdemd65
84 SCP3_BRALA   gindfqk105K..I.PFLFK.G.K..........D....VS...GDG.IYDL....E.EF.QRYCKRF-qlqradvp51
85 SCP2_BRALA   gindfqk105K..I.PFLFK.G.K..........D....VS...GDG.IYDL....E.EF.QRYCKRF-qlqradvp53
86 PIP3_RAT     mdagrdf143V..I.HSCLR.K.A..........D....KK...RDH.KRRF....K.EL.KDFLREL-aiqvddg504
87 AACT_CHICK   mdhhydp786E..F.ARIKS.I.V..........D....PF...RKG.VYTF....Q.AF.IDFRSRE-tadtdtad73
88 CAB_MOUSE    marpleea53A..F.QRVRS.R.L..........D....SH...RDR.EVDF....Q.EY.CYFLSCI-amacaetf19
```

| 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 SPCN_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 |

| 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 G19P_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_LYTPI | 30 | 67.341 | 4.942 | + | . | . | + | P |
| 96 CAP1_HUMAN | 714 | 2120.342 | 4.924 | + | . | + | + | T |
| 97 CIC1_CYPCA | 1852 | 5530.321 | 4.714 | . | + | . | . | . |
| 98 GUNF_CLOTM | 739 | 2196.618 | 4.602 | -. | . | . | . | ? |
| 99 CIC1_RABIT | 1873 | 5593.640 | 4.550 | . | + | . | . | . |
| 100 V57A_BPT4 | 80 | 224.359 | 4.470 | . | . | . | . | . |
| 101 CALG_CHICK | 65 | 178.908 | 4.438 | . | + | + | + | T |
| 102 NIFH_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 YCSX_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 CICC_RAT | 2169 | 6481.468 | 4.011 | . | + | . | . | . |
| 112 CICC_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 |

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
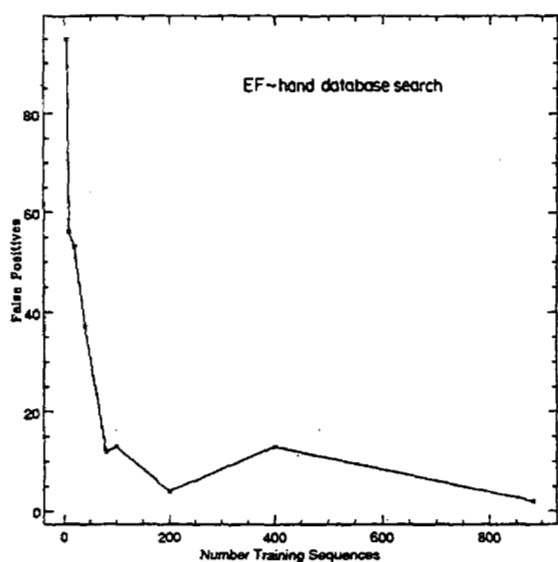


EF—Hand Model Building

——— Train Set

— — — Test Set

**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

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 distribu-

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-

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 $\alpha$-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 calpactain I light chain. *Arabidopsis thaliana* inorganic pyrophosphatase. rat placental calcium-binding protein and rat and bovine l-phosphatidylinositol-4,5-bisphosphate

## 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

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.

Gribskov. 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.

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 Kauffmann 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.