



Clustering protein sequences—structure prediction by transitive homology

Eva Bolten³, Alexander Schliep^{2,*}, Sebastian Schneckener³,
Dietmar Schomburg¹ and Rainer Schrader²

¹Institut für Biochemie, ²ZAIK/ZPR, Universität zu Köln, Weyertal 80, D-50937 Köln, Germany and ³Science Factory, Köln, Germany

Received on April 20, 2001; revised and accepted on July 9, 2001

ABSTRACT

Motivation: It is widely believed that for two proteins *A* and *B* a sequence identity above some threshold implies structural similarity due to a common evolutionary ancestor. Since this is only a sufficient, but not a necessary condition for structural similarity, the question remains what other criteria can be used to identify remote homologues.

Transitivity refers to the concept of deducing a structural similarity between proteins *A* and *C* from the existence of a third protein *B*, such that *A* and *B* as well as *B* and *C* are homologues, as ascertained if the sequence identity between *A* and *B* as well as that between *B* and *C* is above the aforementioned threshold. It is not fully understood if transitivity always holds and whether transitivity can be extended *ad infinitum*.

Results: We developed a graph-based clustering approach, where transitivity plays a crucial role. We determined all pair-wise similarities for the sequences in the SwissProt database using the Smith–Waterman local alignment algorithm. This data was transformed into a directed graph, where protein sequences constitute vertices. A directed edge was drawn from vertex *A* to vertex *B* if the sequences *A* and *B* showed similarity, scaled with respect to the self-similarity of *A*, above a fixed threshold. Transitivity was important in the clustering process, as intermediate sequences were used, limited though by the requirement of having directed paths in both directions between proteins linked over such sequences. The length dependency—implied by the self-similarity—of the scaling of the alignment scores appears to be an effective criterion to avoid clustering errors due to multi-domain proteins.

To deal with the resulting large graphs we have developed an efficient library. Methods include the novel graph-based clustering algorithm capable of handling multi-domain proteins and cluster comparison algorithms. Structural Classification of Proteins (SCOP) was used as

an evaluation data set for our method, yielding a 24% improvement over pair-wise comparisons in terms of detecting remote homologues.

Availability: The software is available to academic users on request from the authors.

Contact: e.bolten@science-factory.com; schliep@zpr.uni-koeln.de; s.schneckener@science-factory.com; d.schomburg@uni-koeln.de; schrader@zpr.uni-koeln.de

Supplementary information: <http://www.zaik.uni-koeln.de/~schliep/ProtClust.html>

INTRODUCTION

Finding the three-dimensional structure of proteins is one of the fundamental problems in molecular biology today. The improvements in throughput of classical methods for determining the structure—e.g. using x-ray diffraction analysis or NMR—cannot keep up with the ever-increasing speed at which proteins are sequenced. It is thus desirable to have methods allowing structure prediction solely from sequence data, either *ab initio*, modeling the molecular folding process, or homology based, using protein sequences with known structures as a template.

The main idea in the latter case is based on the fact that sequence similarity allows detection of homology, i.e. the existence of a common evolutionary predecessor, and thus to infer similar structure and even function virtue of this shared history (Pearson, 1997; Yona *et al.*, 1998). Note, that the same structure or function does not imply a common ancestor; likewise, a common ancestor does not imply a common function, but probably a shared fold.

The relation of sequence similarity—as obtained by pair-wise alignments—to structural or functional properties has been the goal of a number of publications (Brenner *et al.*, 1998; Pearson, 1995, 1997). An widely accepted rule-of-thumb is that 30% identity over aligned regions (Chothia and Lesk, 1986) suffices. More recent studies (Sander and Schneider, 1991; Rost, 1999) qualified this rule. We will call a sequence similarity above this threshold *significant*.

*To whom correspondence should be addressed.

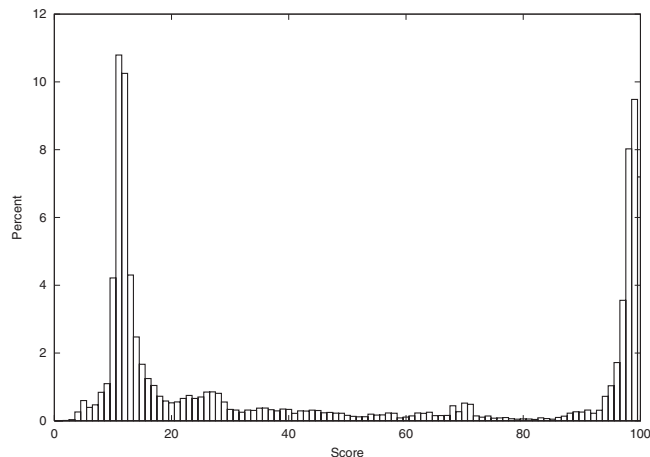


Fig. 1. Histogram of pair-wise alignment scores for all pairs from the same super-family in the SCOP1 data set. Note, the large proportion well below the 30% mark which shows that it is impossible to distinguish super-family and unrelated pairs only by their pair-wise alignment score.

There are lots of examples—e.g. pairs of SCOP super-family sequences with low similarity scores, cf. Figure 1—of homologue proteins with a sequence similarity below any reasonable threshold. Detecting those *distant homologues*, bringing light into the so-called *twilight zone* of low similarity, has been investigated with a number of different approaches (Abagyan and Batalov, 1997; Park *et al.*, 1997; Pearson, 1997; Gerstein, 1998; Krause and Vingron, 1998; Salamov *et al.*, 1999; Arvestad *et al.*, 2000). Fundamental in several of those approaches was the concept of *transitivity* of homology.

Transitivity is a property of (mathematical) relations. In the context of homology, the relation between proteins is defined as having a shared ancestor, which can be interpreted as follows. If for three given proteins *A*, *B* and *C*, *A* and *B* as well as *B* and *C* have a common ancestor, then also *A* and *C* have a common ancestor. This can be used to detect remote homologues, when the sequence similarity between *A* and *C* is too low to infer homology with a sufficient degree of confidence. If the level of sequence similarity between *A* and *B* as well as that between *B* and *C* allows to infer homology directly, we can then use transitivity to infer homology between *A* and *C* indirectly, using *B* as an *intermediate sequence* (cf. Figure 2). The question remains, if transitivity extends to arbitrary numbers of intermediate sequences, and whether it holds on data sets as large as SwissProt (Bairoch and Boeckman, 1992).

We will show later that multi-domain proteins constitute a problem in the use of transitivity for inferring remote homologues. We have developed a method capable of dealing with this problem, while still employing transi-

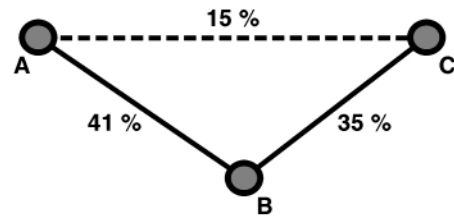


Fig. 2. An example of transitivity of homology. Note the low similarity score of *A* and *C*. The link between *A* and *C* via *B* is called *indirect* compared to the *direct* link between *A* and *B*, respectively *B* and *C*.

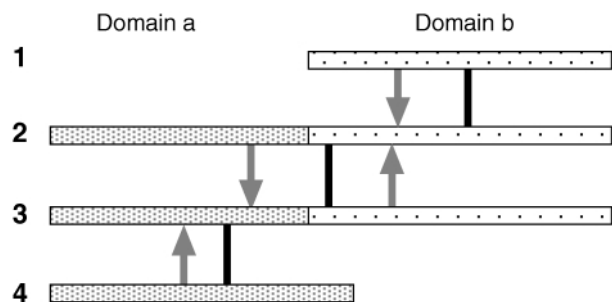


Fig. 3. The problem arising from multi-domain proteins is illustrated. If we use an undirected graph, the solid black edges provide a path from protein #1–#4. In the directed case, the grey edges avoid this possible problem.

tivity to a large degree. In the following sections, we will give a detailed account on this novel graph-based clustering algorithm, discuss the choices in the implementation and the data-sets used, present the results of our method and conclude with a discussion and an outlook on further extensions to our method.

ALGORITHM

We developed a novel graph-based clustering algorithm for structure prediction by transitive homology. Graphs are a natural model for objects and relations between those objects and have been successfully used for clustering in other contexts. See for example Roberts (1984), for an introduction to graph theory and Jain and Dubes (1988) for graph-based clustering. We identified protein sequences with nodes of our cluster graph and will from now on use the terms interchangeably.

Computing pair-wise similarities

At first, a complete undirected graph *G* was computed, where each edge between proteins *P* and *Q* was weighted with their raw Smith–Waterman (Smith and Waterman, 1981) local alignment score (See Section **Directing the**

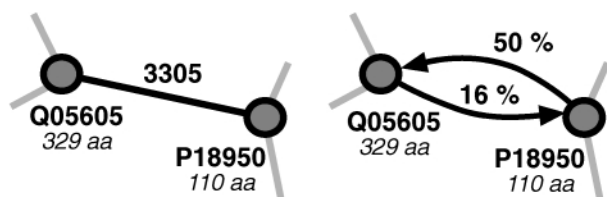


Fig. 4. A single edge (left) is replaced by two *directed* edges in opposite directions (right). Note the distinct similarity scores on the edges.

edges for implementation details), denoted by $\text{raw}(P, Q)$. Note, that an arbitrary similarity measure can be used as input for the clustering.

One concern in clustering protein-sequences are multi-domain proteins which form unwanted ‘bridges’, connecting clearly unrelated proteins. As Figure 3 shows, this is caused by the relation being symmetric, i.e. by not distinguishing between (A, B) and (B, A) . If protein A is more or less a domain of protein B , instead of the whole of protein A being similar to the whole of protein B , false positives will result during clustering (Park *et al.*, 1997).

A computationally inexpensive method for reliable prediction of domains solely from sequence information would be highly desirable to accurately establish such a domain relation between protein sequences. Unfortunately, to our knowledge, no such method exists.

Directing the edges

We developed the following simple criterion to approximate the knowledge we would obtain from a domain-prediction method. Noting that there has to be a difference in length between sequences if multi-domain proteins cause a problem, we decided to *direct* the edges in the graph. Each undirected edge was replaced by two *directed* edges, where the weight of the edge from P to Q , $w(P, Q)$, was computed as

$$w(P, Q) = \frac{\text{raw}(P, Q) * 100}{\text{raw}(P, P)},$$

and similarly for $w(Q, P)$, resulting in a similarity score between 0 and 100% scaled by the raw score of an alignment of P with itself. Note, that the raw self similarity score $\text{raw}(P, P)$ is approximately proportional to the length of P . In the case of two sequences P and Q of distinct length this results in distinct similarity scores for the edges (P, Q) and (Q, P) ; see Figure 4. The resulting graph is denoted by G_d .

Clustering in a threshold graph

The next step in the clustering procedure is to proceed to a threshold graph. That is, all edges from G_d with a

similarity score of less or equal than some fixed threshold τ are removed, resulting in the graph $G_d(\tau)$. All similarity values below this threshold are assumed to be produced by chance and not to be an indicator of true structural homology. Since using a single-link clustering algorithm would only solve part of the problem with multi-domain proteins, we made use of a standard concept in graph theory, the so-called *strongly connected component*.

DEFINITION 1. In a directed graph G , a Strongly Connected Component (SCC) is a maximal set C of nodes of G , such that for every pair of nodes p and q in C there is one directed path in G from p to q and one from q to p .

The nodes of a directed graph are partitioned into SCCs, all of which can be computed efficiently with complexity $O(n + e)$ (Sedgewick, 1990) for a graph with n nodes and e edges. Some care has to be taken in an efficient implementation of the algorithm and the underlying graph data structures to allow interactive handling of graphs of close to 100 000 nodes and 20 000 000 edges.

As at least two distinct paths, one in each direction, have to exist, using SCCs as clusters is comparable to using a more conservative algorithm such as average linking. Note, that in Figure 3 only proteins number two and three are in an SCC. Thus, using an SCC as a cluster does not make use of a large amount of information. Nevertheless, we chose to evaluate the performance of our algorithm on the basis of the SCCs alone to establish the validity of our approach.

IMPLEMENTATION AND EVALUATION

The algorithms presented here have been implemented in the C++ package ProtClust. They have been tested and used on various Sun Ultra computers (Ultra 5 up to Sun Enterprise 10 000), running Solaris 7 and earlier versions, using the GNU g++ compiler, as well as on a Compaq ES40 running Tru64 Unix V5.1, employing Compaq’s cxx compiler, version 6.20.

We used our own implementation of the Smith–Waterman local alignment algorithm (Smith and Waterman, 1981) for computing sequence similarity. The choice of the alignment algorithm was motivated by the superior sensitivity (Brenner *et al.*, 1998) for low-scoring alignments as compared with other, faster algorithms such as BLAST (Altschul *et al.*, 1990). The substitution matrix, an integerized version of BLOSUM80 (Blocks Substitution Matrix; Henikoff and Henikoff, 1993), was chosen based on the results of the investigations of one of the authors (Schneckener, 1998).

We chose the gap opening (*gop*) and gap extension penalties (*gep*) used in the alignment algorithm to be $gop = 90$ and $gep = 9$. This was decided after trying out a wide range of choices with extensive computations on

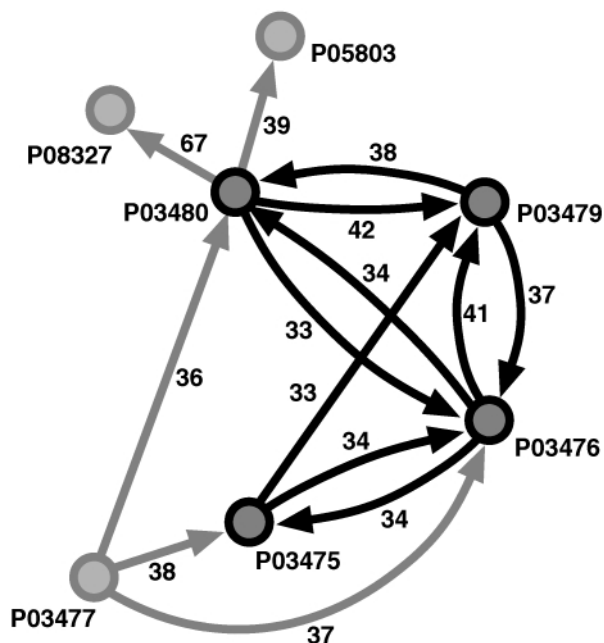


Fig. 5. An example of a SCC in SwissProt. The grey nodes are not part of the SCC, but are clearly related. Note the sub-threshold scores (no edge present) between nodes P03480 and P03475. The threshold used in this example was 32%.

the SCOP1 data set (see Section Data). After computing a pair-wise sequence alignment for all sequences against all, we performed single-link-clustering for thresholds from 5 to 95. By assigning a SCOP super-family to each cluster by maximal intersection we implied a measure of error by the number of sequences belonging to families distinct from the assigned. The gap penalties appeared to have a limited influence on the performance.

Data

We used SwissProt (Bairoch and Boeckman, 1992) version 39 as of June 2000, excluding all sequences with less than 40 amino acids (a.a.), resulting in a set of 86 494 protein sequences after removal of identical sequences. The total running time for the pair-wise Smith–Waterman alignment was on the order of 1400 CPU-days. We used a simple tool (Schliep, 1998) to distribute the work over about 40 computers with a total of 55 CPUs available.

For the evaluation, the SCOP database (Murzin *et al.*, 1995) was employed. It provides a high quality hand-crafted partition of protein sequences at different levels. For our application, the relevant levels are *family*—sets of sequences with more than 30% sequence identity and possible functional identity—*super-family*—sets of sequences likely to have a common ancestor, low sequence identity, but structural and functional similarity—and

fold—sets of sequences having structural similarity. All data sets defined in the following were based on version 1.37 of SCOP as obtained from <http://scop.mrc-lmb.cam.ac.uk/scop/pdbd.html>.

SCOP1. This data set of 2692 sequences contains all non-identical sequences from SCOP having at least 40 amino acids excluding sequences belonging to classes 8, 9 or 10. There are 65 464 pairs of homologue sequences; i.e. pairs where both sequences are in the same super-family and 3556 622 pairs where the sequences are in distinct super-families. On the SCOP-fold level the corresponding numbers are 76 660 and 3545 426 pairs, respectively.

SCOP1 + SP. To avoid having to uniquely identify sequences from the data set SCOP1 in the SwissProt data we used, we simply pooled the two sets of sequences, yielding a total of 85 961 sequences.

SCOP2. For comparison with Arvestad *et al.* (2000) we also used a subset of 609 randomly chosen sequences from SCOP. It contains sequences shorter than 40 amino acids, but no sequences from classes 8, 9 or 10.

Performance measure

A natural quality measure for detecting remote homologues is counting the number of truly homologous pairs of sequences and relating this to the number of errors made. A pair is taken to be truly homologous in this context, if both sequences are for example in the same SCOP super-family. We will call correctly identified true homologue pairs *True Positives* (TP), not identified true homologues *False Negatives* (FN), non-homologue pairs predicted to be homologue *False Positives* (FP) and non-homologue pairs correctly identified *True Negatives* (TN).

Sensitivity specifies the proportion of identified homologue pairs

$$\text{sens} = \frac{\#TP}{\#TP + \#FN}$$

and *specificity* the proportion of errors among the pairs predicted to be homologues

$$\text{spec} = \frac{\#TP}{\#FP + \#TP}$$

Note, that $\text{sens} = \text{spec} = 1$ would be the most highly desired performance, since it implies that neither false positive nor false negative errors are made.

DISCUSSION

It should be noted that, for a very large proportion of pairs of sequences from the same SCOP super-family, an alignment score in the twilight zone (cf. Figure 1) well

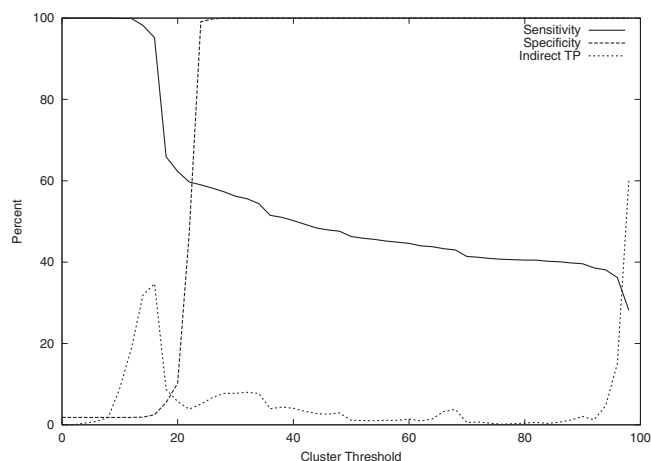


Fig. 6. Sensitivity, specificity and the percentage of indirectly linked true positives versus clustering threshold for the SCOP1 data set.

below reasonable thresholds for a pair-wise comparison is obtained. This supports the need for methods employing intermediate linking.

At first we evaluated the method on the SCOP1 data set. As Figure 6 shows, the algorithm has maximum specificity over a very wide range of thresholds. For a conservative threshold choice of 32, we obtain a sensitivity of 55.6%. One factor contributing to the relatively small proportion of homologues found due to intermediate linking, about 8.0% at the chosen threshold, is that SCOP1 is a small data set, providing not enough sequences for intermediate linking. Also, it cannot be expected to recover all SCOP relations, since the super-families contain sequences which are only *likely* to have a common ancestor. As a matter of fact, we are not aware of any method which fully recovers even SCOP families. The rapid change in sensitivity and specificity at a threshold of 23 is due to lifting off the noise floor; a similar behavior is observed for pair-wise sequence comparisons.

The evaluation on the SCOP1 + SP data set allows us to investigate two different aspects of our method. First, we should be able to benefit from the presence of more sequences in the clustering, as we are likely to find more indirect links. Secondly, there might be a large loss of specificity caused by unbounded transitivity. We computed the clustering for the whole data set, containing SCOP1 and SwissProt sequences. Therefore, sequences from SCOP1 might be joined by intermediate sequences from SwissProt not present in SCOP1 itself. The evaluation itself was only performed on pairs of SCOP1 sequences. Figure 7 shows a very favorable performance of the method on this substantially larger data set. Again, we achieve a 100% specificity over a wide range of thresholds while additional sensitivity is

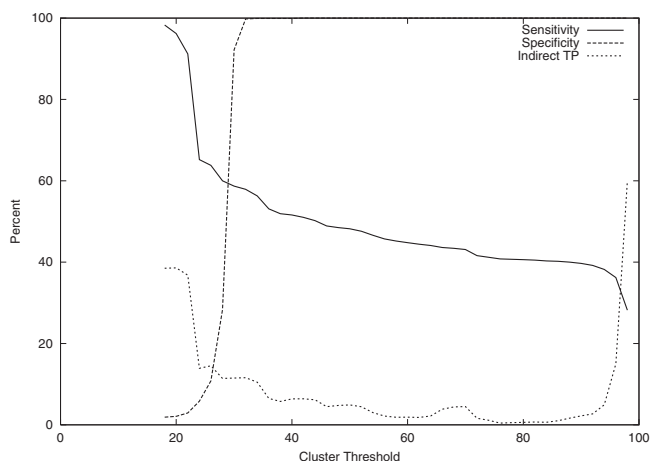


Fig. 7. Sensitivity, specificity and the percentage of indirectly linked true positives versus clustering threshold for the SCOP1 + SP data set.

gained. At a threshold choice of 32, a sensitivity of 57.9% and specificity of 99.8% is obtained. This is an absolute increase in sensitivity of 2.3, and a relative of 4.1%. The proportion of true positives found due to intermediate linking is increased by 3.6–11.6%. The noise floor is higher than in the previous case, as can be seen from the shift of the point of sudden decrease in specificity to the right in Figure 7. However, a more than 30-fold increase in the number of sequences only requires a minor adjustment of the threshold to again obtain near optimal specificity.

To investigate possible reasons for the somewhat surprising magnitude of the increase in sensitivity due to the additional sequences from SwissProt, we compared histograms of cluster sizes (cf. Table 1). The SCC-based clustering is of a very conservative nature, as can be observed in the number and sizes of clusters obtained. Few large and a large number of single element clusters are produced. Note, that larger clusters in general do not merge to one ‘super cluster’ as the threshold decreases, indicating a good separation of clusters and thus the validity of our approach. This observation can be readily explained by the fact that terminal sequences, connected only with an unidirectional path, are excluded from the SCC (cf. Figure 5).

To stress the importance of intermediate linking algorithms we compared our algorithm with Arvestad *et al.* (2000). While it employs only pair-wise sequence comparisons, their approach uses a more involved scoring method, optimized substitution matrices, and gap penalties, to achieve a substantial improvement over straight-forward pair-wise sequence comparisons. From Figure 8, noting that the graphs are nearly parallel to the x -axis and that a high sensitivity is obtained at nearly maximum specificity, we can deduce the quality and the

Table 1. Total number of SCC clusters, and SCC clusters of sizes 1, 2–5, 6–10 etc. for varying thresholds from 25 to 50% are displayed for the SCOP1 dataset (top) and the SCOP1 + SP data set (bottom). The bottom-most row are the corresponding numbers for SCOP super-families

Threshold	Total	1	2–5	6–10	11–20	21–50	51–100	>100
25	688	348	237	50	37	12	3	1
30	743	388	256	53	30	12	3	1
35	814	450	271	46	31	13	2	1
40	867	485	289	51	29	10	2	1
45	932	535	305	58	22	9	2	1
50	989	587	315	54	24	6	2	1
25	17 964	11 419	4566	915	560	384	83	37
30	23 395	15 079	5886	1123	692	433	120	62
35	27 154	17 933	6685	1231	709	432	109	55
40	30 748	20 791	7465	1266	710	379	95	42
45	34 323	23 814	8148	1263	654	330	79	35
50	38 033	27 138	8669	1262	585	298	51	30
—	384	153	126	43	32	24	5	1

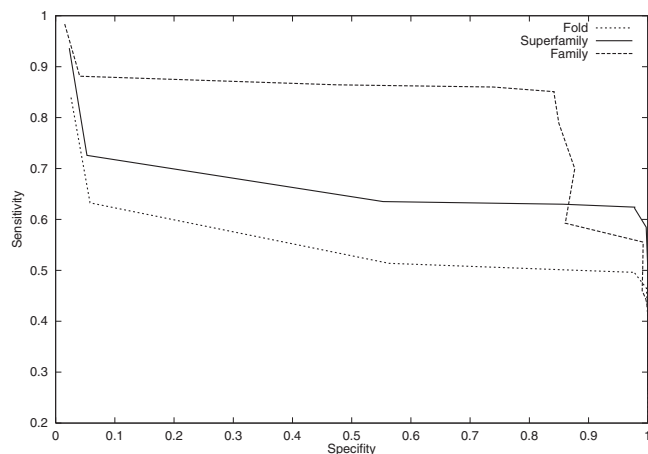


Fig. 8. Sensitivity versus specificity for the SCOP2 data set on the fold, super-family and family level.

conservative nature of our method on the SCOP2 test set. The SCCs based on the Smith–Waterman score scaled by self-similarity thus realize a 24% better sensitivity at virtually equal specificity.

The handling of the problem with multi-domain proteins was investigated by manual inspection of pair-wise alignments along paths connecting false positives (not shown) present at low clustering thresholds. The analysis of the alignments and a comparison with scores not scaled by self-similarity showed the expected behaviour. Unfortunately, a comprehensive analysis was not feasible due to the limited amount of domain data available.

The clustering algorithm, queries on the clusters etc. are low both in theoretical complexity and practical

running-time. Thus, they are even suitable for interactive work. As our experience shows, even the computation of the alignment matrix can be performed without need for specialized hardware or high-performance computers. Nevertheless, using substantially larger databases such as TrEMBL (Bairoch and Apweiler, 1999) might not be feasible without additional computing hardware.

We have demonstrated that the SCCs are conservative and highly specific clusters of protein sequences at the super-family level. They can be used as *seeds* for further advances in the clustering, ultimately obtaining an automated, biologically meaningful partition of protein space.

In the next stage, the information from nodes connected to but not part of an SCC will be taken into account. This should greatly improve the sensitivity of the method. As a further step, statistical significance, e.g. classical length dependent extremal value distribution fitting (Levitt and Gerstein, 1998), *E*-values, can be used to improve the quality of the threshold graph and thus the quality of the resulting clusters. Also, length dependent threshold functions (Arvestad *et al.*, 2000) should be investigated.

ACKNOWLEDGEMENTS

We would like to thank Oliver Leven, Sebastian Meller, Frank Nübel, Peter Pipenbacher, Barthel Steckemetz and Olav Zimmermann for helpful discussions and their support.

REFERENCES

- Abagyan, R.A. and Batalov, S. (1997) Do aligned sequences share the same fold? *J. Mol. Biol.*, **273**, 355–368.
- Altschul, S., Gish, W., Miller, W., Myers, E. and Lipman, D. (1990) Basic local alignment search tool. *J. Mol. Biol.*, **215**, 403–410.
- Arvestad, L., Ivansson, L., Lagergren, J. and Elofsson, A. (2000) What is the best method to determine if two proteins are related? A study on the structural and evolutionary significance of pairwise protein sequence alignments, submitted.
- Bairoch, A. and Apweiler, R. (1999) The Swiss-Prot protein sequence data bank and its supplement TrEMBL in 1999. *Nucleic Acids Res.*, **27**, 49–54.
- Bairoch, A. and Boeckman, B. (1992) The SwissProt protein sequence data bank. *Nucleic Acids Res.*, **20**, 2019–2022.
- Brenner, S.E., Chothia, C. and Hubbard, T.J. (1998) Assessing sequence comparison methods with reliable structurally identified distant evolutionary relationships. *Proc. Natl Acad. Sci. USA*, **95**, 6073–6078.
- Chothia, C. and Lesk, A.M. (1986) The relation between the divergence of sequence and structure in proteins. *EMBO J.*, **5**, 823–826.
- Gerstein, M. (1998) Measurement of the effectiveness of transitive sequence comparison, through a third ‘intermediate’ sequence. *Bioinformatics*, **14**, 707–714.
- Henikoff, S. and Henikoff, J. (1993) Performance evaluation of amino acid substitution matrices. *Proteins*, **17**, 49–61.

- Jain,A.K. and Dubes,R.C. (1988) *Algorithms for Clustering Data*. Prentice-Hall, Englewood Cliffs, NJ.
- Krause,A. and Vingron,M. (1998) A set-theoretic approach to database searching and clustering. *Bioinformatics*, **14**, 430–438.
- Levitt,M. and Gerstein,M. (1998) A unified statistical framework for sequence comparison and structure comparison. *Proc. Natl Acad. Sci. USA*, **95**, 5913–5920.
- Murzin,A.G., Brenner,S.E., Hubbard,T. and Chothia,C. (1995) Scop: a structural classification of proteins database for the investigation of sequences and structures. *J. Mol. Biol.*, **247**, 536–540.
- Park,J., Teichmann,S.A., Hubbard,T. and Chothia,C. (1997) Intermediate sequences increase the detection of homology between sequences. *J. Mol. Biol.*, **273**, 349–354.
- Pearson,W.R. (1995) Comparison of methods for searching protein sequence databases. *Protein Sci.*, **4**, 1145–1160.
- Pearson,W.R. (1997) Identifying distantly related protein sequences. *Comput. Appl. Biosci.*, **13**, 325–332.
- Roberts,F.S. (1984) *Applied Combinatorics*. Prentice-Hall, Englewood Cliffs, NJ.
- Rost,B. (1999) Twilight zone of protein sequence alignments. *Protein Eng.*, **12**, 85–94.
- Salamov,A.A., Suwa,M., Orengo,C.A. and Swindells,M.B. (1999) Combining sensitive database searches with multiple intermediates to detect distant homologues [in process citation]. *Protein Eng.*, **12**, 95–100.
- Sander,C. and Schneider,R. (1991) Database of homology-derived protein structures and the structural meaning of sequence alignment. *Proteins*, **9**, 56–68.
- Schliep,A. (1998) PBQ—the python batch queue. Available from <http://www.zpr.uni-koeln.de/~schliep/PBQ>.
- Schneckener,S. (1998) *Positionsgenaues Alignment von Proteinsequenzen*, PhD Thesis, Universität zu Köln.
- Sedgewick,R. (1990) *Algorithms in C*. Addison-Wesley, Reading, MA.
- Smith,T.F. and Waterman,M.S. (1981) Identification of common molecular subsequences. *J. Mol. Biol.*, **147**, 195–197.
- Yona,G., Linial,N., Tishby,N. and Linial,M. (1998) A map of the protein space—an automatic hierarchical classification of all protein sequences. *Ismb*, **6**, 212–221.