

Structure Comparison and Structure Patterns

INGVAR EIDHAMMER,¹ INGE JONASSEN,¹ and WILLIAM R. TAYLOR^{1,2}

ABSTRACT

This article investigates aspects of pairwise and multiple structure comparison, and the problem of automatically discover common patterns in a set of structures. Descriptions and representation of structures and patterns are described, as well as scoring and algorithms for comparison and discovery. A framework and nomenclature is developed for classifying different methods, and many of these are reviewed and placed into this framework.

Key words: protein structure, structure comparison, structure motifs.

1. INTRODUCTION

AS WE ENTER A POST GENOMIC AGE, considerable effort is being expended on structural-genomics programs with the aim of determining the structure of as many gene products as possible. These data, combined with the considerable volume of current structural data, will lead to the emergence of protein structure comparison as a critical technique—to aid not only in the understanding of the relationships between proteins in detail but also in the classification of the variety of structures into meaningful categories.

As with most comparison problems, complications arise as there is neither one best way to make the comparison nor to evaluate the answer. This situation pertains in sequence comparison where, although there is an optimal alignment algorithm for pairs of sequences, its results depend on a model of sequence relatedness that is based on much less certain ground. In structure comparison, we do not even have an algorithm that guarantees an optimal answer for pairs of structures and, with the added complexity of structure, relative to sequence, the models of relatedness are correspondingly more varied.

We have compiled this review with the aim of lending some order to the bewildering variety of methods that have been devised both to compare structures and to extract their essential components in the form of patterns. The result is neither wholly consistent nor complete as many of the more recent methods use a combination of many approaches. We hope it should, at least, provide a framework which will help in classifying the simpler methods and aid in breaking the more complex into their components.

Sequence and structure patterns can be used for characterizing families of proteins which are defined to be sets of functionally or structurally related proteins. The discovery of such patterns can help in the understanding of relationships between sequence, structure, and function of proteins and in a bigger context help to understand the working of living organisms. Sequence patterns can be used to predict, from sequence information only, the structure and function of new proteins. Once the sequence of a new

¹Department of Informatics, University of Bergen, Høyteknologisentret, N-5020 Bergen, Norway.

²Division of Mathematical Biology, National Institute for Medical Research, The Ridgeway, Mill Hill, London NW7 1A, U.K.

protein has been found, a database similarity search can identify proteins with similar sequences, and if the similarity is sufficiently strong, one can assume that the new protein has the same structure and sometimes function as the protein in the database. Sets of proteins having the same structure and or function can be analyzed, and one can find which residues are most informative about the structure/function, which are allowed to vary more freely, etc. Generally, one can define sequence patterns characteristic of each family. Then, comparing a new sequence to a collection of patterns (one for each structure family) is more efficient and also provides a more sensitive and specific test of family membership.

However, even if the two proteins have similar structures, they can have different functions, as a small number of amino acid changes may be responsible for changing the function of the protein. Families of proteins having the same function (represented by the same or different structures) may be formed, and it may be possible to find a signature (pattern) which is characteristic of the function so that the new protein can be matched against, and the outcome of the matching may be used to suggest presence or absence of the particular function.

In this way, patterns can play a major role in understanding relationships between different structures (to find recurring structural patterns and rules) and relationships between structure and function.

Sequence patterns can be used for characterizing proteins with structural similarity, for example, the PROSITE entry for the SH3 domain which is a small protein domain of about 60 amino acid residues in two tightly packed antiparallel beta sheets occurring in a variety of intracellular and membrane-associated proteins. The PROSITE entry gives a profile for this family. The function of the domain is not well understood, but if one finds that a new protein matches the profile with a sufficiently high score, one can be reasonably sure that the protein contains a SH3 domain. Patterns can also be used in relation to functional properties. For example, PROSITE gives a pattern characteristic for the sequences of a certain class of aminoacyl-tRNA synthetases (a group of enzymes which activate amino acids and transfer them to specific tRNA molecules as the first step in protein biosynthesis). If a new sequence matches this pattern, one can hypothesize that it is an enzyme having this function.

As the number of known protein structures increases (13049 in the PDB as of September 2000) (Bernstein *et al.*, 1977), there is need for methods describing and revealing common functionally important units in (related) structures. Such units could be described by structure patterns, analogous to the sequence patterns. These patterns could then be used to classify protein structures into structure families, for example by considering the occurrence of common arrangements of secondary structure elements in the core of proteins, as described by Koch *et al.* (1996). Structure patterns can also be used to make a hypothesis of the function of a protein, for example the “coordinate templates” for finding Ser-His-Asp catalytic triads in the serine proteinases and lipases as reported in Wallace *et al.* (1996). See Figure 1.

The terminology used in the literature for describing common similarities between sequences or structures is quite confusing. The words “pattern,” “motif,” “fingerprint,” “template,” “fragment,” “core,” and “site” are all used. By a pattern we mean here any description of sequence or structure properties for which one can either (1) decide whether a protein matches it or not or (2) assign a value to the match between the protein and the pattern. In Case (1) it will be called a *deterministic pattern* and in Case (2) a *probabilistic pattern*. A pattern may or may not match any sequence or structure and it may or may not have a biological meaning. Different formalisms for describing patterns exist. Following the definition above, for any pattern formalism there is a mechanism for deciding whether or not a sequence or a structure matches a pattern. In most cases, a pattern is allowed to match part of a sequence or a structure. The part of a sequence or a structure (e.g., a subsequence or a substructure) that matches a pattern is called an *occurrence* of (or a match to) the pattern.

For some pattern classes, one can define a partial ordering of the patterns with respect to their generality. If it can be proven that any structure that matches a pattern P_1 will also match pattern P_2 , then P_2 is a *generalization* of P_1 (or subsumes P_1) which we can denote by $P_1 \prec P_2$ (or $L(P_1) \subseteq L(P_2)$). This can also be extended to include structure descriptions, so that if pattern P matches structure S , then $S \prec P$ (or $S \in L(P)$). This is easiest to do formally for deterministic patterns (see, e.g., Shinohara and Arikawa (1995) and Conklin (1995)). For probabilistic patterns, it is not that straightforward, but for example for alignments, a subalignment (consecutive columns) can be seen as a generalization of an alignment.

By a *motif* we mean a pattern which has a biological meaning—i.e., it can be used to predict functional or structural properties of the protein, or it describes (nontrivial) features common to biologically (structurally or functionally) related proteins. Having found a pattern, one wants to evaluate whether it is a motif, i.e.,

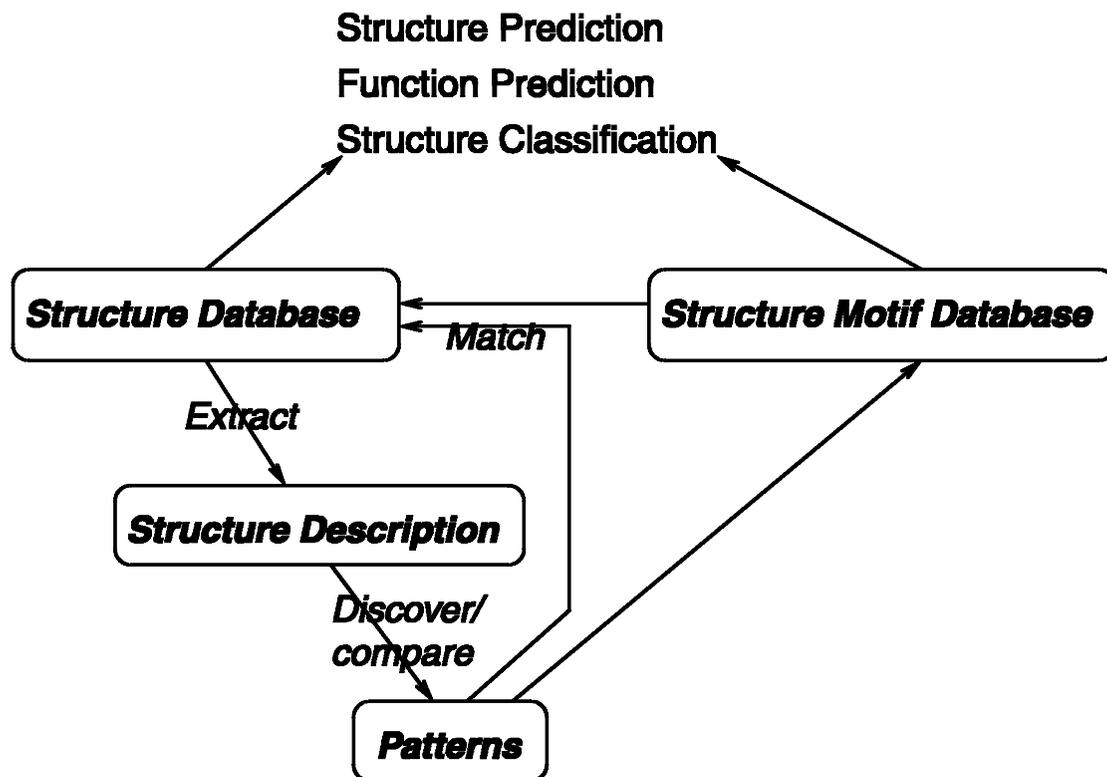


FIG. 1. Protein comparison and motif discovery in a context.

if it has a biological meaning. Because in automatic methods there is no way of assessing the biological meaning of a pattern directly, alternative mathematical methods for pattern evaluation have been developed. We call them *fitness* or *scoring* functions. The term motif can include sites (a small part of the structure having a specific functional or structural role, e.g., the active site in enzymes or metal binding sites), cores (an often bigger part of the structure in the interior of the protein, e.g., the hydrophobic core), secondary structures, and supersecondary motifs (constituted of secondary structures).

The PDB (Protein Data Bank, <http://www.rcsb.org/pdb/>) (Bernstein *et al.*, 1977) at Brookhaven National Laboratory is the international structural database. It is an archive of experimentally determined three-dimensional structures of biological macromolecules, together with an extensive annotation. A lot of other structural databases exist (see Baxeavanis, 2000).

In the following sections, we investigate different aspects regarding pairwise and multiple structure comparison, and the problem of automatic discovery of common patterns in a set of structures. Descriptions and representation of structures and patterns, as well as scoring and algorithms for comparison, are investigated. A framework and nomenclature is developed, and a number of methods are reviewed and placed into this framework. The current article partly builds on and extends two previous surveys: (Brown *et al.* (1996), which discusses structure comparison, and Brazma *et al.* (1998), which deals with the discovery of sequence patterns. Some method papers also include small reviews that have been useful in this work (Gibrat *et al.*, 1996; Godzik, 1996).

2. FRAMEWORK

Figure 2 shows the overall components of discovering structure motifs and matching new structures against them. The different steps include those following.

1. *Feature extraction*: In this step, the features to be used in the comparison of the structures or in the pattern discovery method are extracted. This might include comprehensive computing, e.g., assigning

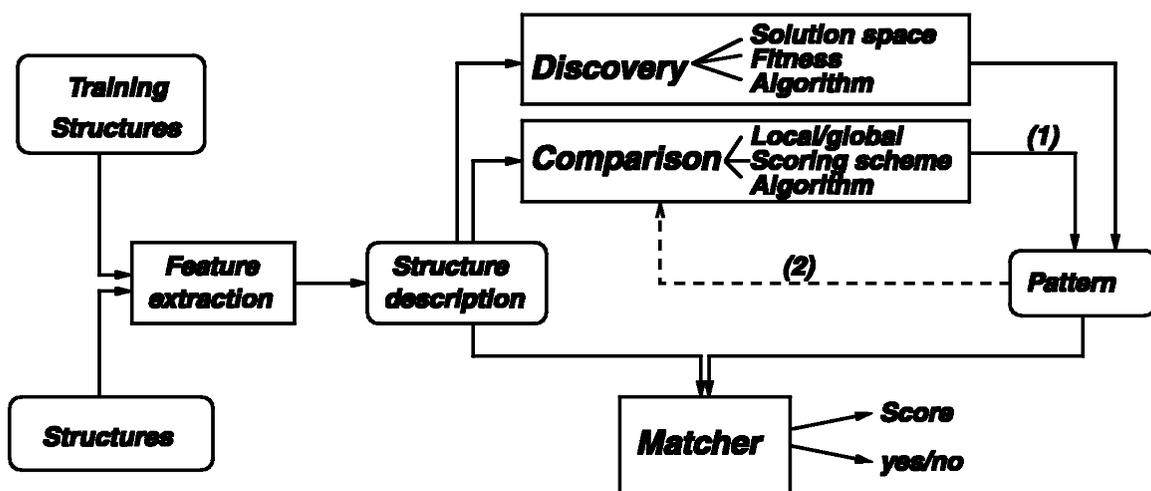


FIG. 2. Overview of the process of discovering patterns, either by direct discovery method or using a comparison algorithm, and pattern matching. The edge marked (1) indicates that the result of a comparison method may need some further processing in order to be expressed as one or several patterns. Furthermore, if the result of comparing two structure descriptions can be represented as a pattern and then compared with another structure (or another pattern), the edge (2) can be used to find patterns matching progressively bigger sets of structures. Note that “Matcher” and “Comparison” are very closely related and may in fact be implemented using the same method.

secondary structures to the residues. Not all methods perform this step explicitly, but all take into account a well-defined set of structural features and conceptually this can be seen as if they first extract the relevant features and then subsequently work on the resulting structure descriptions.

2. *Comparison*: This step takes as input a pair of structure descriptions (or a description/pattern pair) and finds (local or global) similarities between the two, optimizing a similarity measure and outputting a score. The similarity may also be represented as a pattern.
3. *Discovery*: patterns matching many or all of the input structures are found. The patterns are chosen from a solution space so that their fitness with respect to the input structures are as high as possible.
4. *Matcher*: This step takes as input one pattern and one structure and evaluates the match between the two; the output is “yes” or “no” if the pattern is deterministic or is a score if the pattern is probabilistic. Note the similarity between the matching and comparison operations; they may sometimes be implemented using the same method.

The steps involved in these operations will be discussed in more detail in the following sections. Matching is not discussed separately; for an example matching algorithm, see Gilbert *et al.* (1999).

3. DESCRIPTIONS AND REPRESENTATIONS

3.1. Feature extraction

This step is used to extract the features from the structure file which are relevant for pattern discovery, learning, matching, and comparison.

When deciding which features should be extracted and represented, certain aspects must be considered. First, one needs to decide on which structure level similarities are sought (e.g., atom group, residue, secondary structures). Also, one needs to decide whether the similarities should require sequence order to be preserved, a reasonable requirement if we assume that the proteins are evolutionarily related, but not if common features have arisen independently. The structure description to be used as input to the comparison or pattern discovery algorithms should contain only the features which we would like to compare and/or

to describe as patterns. A structure description will consist of geometry, topology, and properties, and for each of these we need to decide which features to include:

1. Geometry—e.g., coordinates or relative positions of atoms, residues, fragments, or SSEs.
2. Topology—e.g., the elements' order along the backbone.
3. Properties—e.g., physico-chemical properties of the elements (e.g., residues)

In our framework, patterns are to be found from structure descriptions so that they represent features common to a set of such descriptions. Patterns will therefore be generalizations of structure descriptions and are limited to features included in these.

In order to provide the comparison (pattern-discovery) algorithms with a good starting point, the structure descriptions should ideally satisfy the following properties:

1. *Invariant to trivial changes*—such as translation and rotation.
2. *Robust*—the description should not change drastically due to minor changes in the structure (e.g., experimental errors, hydrogen-bonding definition).
3. *Similar structures*—should get similar descriptions.
4. *Different structures*—should get different descriptions.

Loosely, the definition of “similar” and “different” will depend on what aspects of protein structure one wants to compare or capture in a pattern. For example, if the comparison is to be done at the level of packing of secondary structure element, the descriptions of structures with similar packing should be similar.

3.2. Structure descriptions

A natural way to describe a complex object like a protein structure is to break it into pieces (units) and to describe each unit separately and (most often) the relationship between the units. For structure comparison, the units are chosen in two ways:

1. *Element based*—natural structure elements (atoms, residues, fragments, secondary structures (SSEs)) are used as the basic units.
2. *Space based*—the space in which the structure is located is divided into (possibly overlapping) geometrically defined cells, e.g., by using a grid, or shells around center points.

3.2.1. Element based descriptions. Element-based descriptions are by far the most used type. Brown *et al.* (1996) give an overview of structure comparison methods using this approach and present a nomenclature for use in the comparison of structures. Following this nomenclature, three concepts can be used for description: *element class*, *properties*, and *relation* (Brown *et al.* use the word “feature”; we use “property” to avoid confusion with our use of “feature” at the higher level).

- *Element class*—the level of the description varies: atom (group), residue, backbone fragment, and secondary structure element.
- *Property*—is used for specifying the properties of each element, such as three-dimensional coordinates, physico-chemical properties, amino acid type, secondary structure type, curvature, and torsion.
- *Relation*—is used for describing the relation between the elements. In practice, the relations are binary, such as geometrical distances, difference in orientation, and bonding.

Depending on what element class is used, one might say that the description is on a coarse or a fine level. Secondary structure is on a coarse level, and the common similarities in this case are often common folds. Atom (group) and residue level are fine, while backbone fragment level might in some cases be considered as low, in other cases as high.

The two levels have different goals; the high level aims to classify the whole overall structure, but the low level aims to be used more directly in determining/finding active/binding sites.

3.2.2. Space based descriptions. In space-based descriptions, the 3D space in which a structure is located is divided into (possibly overlapping) cells. The part of the structure falling in each cell is described, for example, by the number of residues of each type or the number of residues with certain property values. Examples of space-based descriptions are given by Bagley and Altman (1995) and by Kastenmüller *et al.* (1998). Figure 3(a) shows cells as shells around points in the structure space.

3.2.3. Substructure descriptions. Some methods divide the structures into substructures, describe each substructure, and then find (local) patterns common to substructures from different structures. This initial step is often followed by a step where the identified local patterns are combined to form larger patterns. For example, Escalier *et al.* use this approach and element-based structure descriptions. They define a substructure to consist of elements where the distance between any two elements is below a given threshold (Escalier *et al.*, 1998).

The approach can also be used together with space-based structure descriptions. For example, Bagley and Altman (1995) define a substructure as the part of the structure that falls within a certain distance from a specified point. Similar sites constitute a pattern, and no combination of (local) patterns is made. Figure 3 illustrates substructures.

3.3. Structure representation

For each protein, the selected features are represented in a data structure. Typically each unit (element or cell) is represented separately and combined to form a representation of the whole structure. For example, each residue can be represented by its coordinates and amino acid, and the whole protein structure can then be represented as a list (ordered) or set (unordered) of residue representations.

We have grouped the different representations into five groups:

1. *Strings*—It is possible to represent structural features as one or several strings, (see, for instance, Matsuda *et al.* (1997) where residue i is represented by a letter showing the C_{α} 's position relative

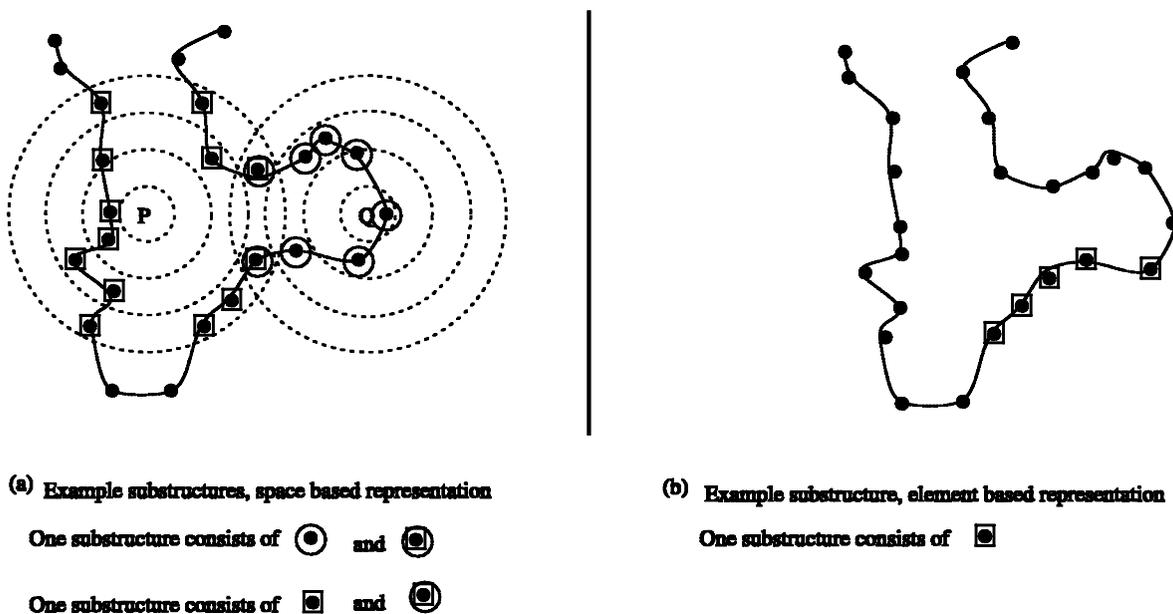


FIG. 3. Substructures illustrated in 2D. The structures are shown by solid curves; the dots represent elements. (a) Two substructures in a space-based description. One is the part of the structure falling in the shells around P and the other is the part falling in the shells around Q . (b) A substructure in an element-based description, where a substructure is a fragment (consequent residues).

to the positions of the C_α of residue $i - 2, i - 1, i + 1$). This representation is appropriate only for matching of (consecutive) substrings.

2. *List of unit descriptions*—For example, extract one or several coordinates per residue: the coordinates of the C_α atom (Taylor and Orengo, 1989), mean coordinates of side chain (Jonassen *et al.*, 1999), two pseudo atoms (Artymiuk *et al.*, 1994), or unit vectors (the vectors between succeeding C_α atoms (Chew *et al.*, 1999)). In addition to the coordinates, each unit can have associated additional properties, such as amino acid type, physico-chemical properties, degree of burial, SSE status, etc.
3. *Set of unit descriptions*—Same as list, but the units are unordered. Alexandrov *et al.* (1992) in their SARF program use a set of backbone fragments.
4. *Graphs*—A labeled graph can be used to represent all element-based feature descriptions, with nodes representing the elements and the edges representing the relations. For instance, one can identify secondary structure elements and label the edges with the distance and angular relationships between the SSEs of the nodes (Grindley *et al.*, 1993) or label them with the type of parallelism (Koch *et al.*, 1996).
5. *Feature arrays*—Features of the structure can be represented by fixed-length arrays. This representation can be used for both space-based and element-based description. Bagley and Altman (1995) have a space-based description where each cell is described by an array summing up the properties of its residues.

3.3.1. External/internal representations. For efficiency, most of the methods use internal representation of the structures, hidden from the user. Internal representations can, for example, facilitate more direct comparison of spatial similarity. In the “double dynamic programming” method by Taylor and Orengo (1989), a local coordinate system is constructed for each residue into which the remaining residues’ coordinates are transformed. Then, assuming a particular residue pair is aligned, the spatial similarity of the remaining residues can be assessed and scores assigned to be used in a (lower level) dynamic programming alignment step (see Section 4.4).

The same approach is used in geometric hashing (Fischer *et al.*, 1994; Wallace *et al.*, 1996; Alesker *et al.*, 1996); a local coordinate system is constructed for each residue, but in geometric hashing normally the coordinates of only the neighboring residues are transformed. Then, for a pair of residues (one from each structure), their spatial neighborhoods can be compared simply by counting the number of neighboring residues with similar position (in the respective local coordinate systems). Special hash tables are used to speed up the computations (see Section 4.6).

Another approach is taken by Jonassen *et al.* (1999) who represent the spatial neighborhood of each residue by a string of the spatially close residues (having spatial distance below some threshold). The strings do not contain information about the exact geometry of the neighbors, but do give the sequence order (along the backbone) of the neighboring residues. Patterns shared by sets of these strings can be discovered using efficient sequence pattern discovery methods (see Section 5.5).

The external representations can be generalized to be useful for pattern descriptions, which is not the case for the internal representations. When algorithms using an internal representation are meant to present the result to the user as a pattern, it must be described in another form. For example, each pattern found by the Spratt method (Jonassen *et al.*, 1999) is presented as a list of residues, for each residue is given the allowed amino acids and the 3D coordinates.

4. PAIRWISE STRUCTURE COMPARISON

In this section, we discuss methods for comparing pairs of structures, focusing on element-based structure descriptions. The most natural way to compare two objects, each represented by a collection of elements, is to try to find element correspondences between the two. More formally, for two objects A and B having elements a_1, a_2, \dots, a_m , and b_1, b_2, \dots, b_n , respectively, we define an *equivalence* as a set of pairs $L(A, B) = (a_{i_1}, b_{j_1}), (a_{i_2}, b_{j_2}), \dots, (a_{i_l}, b_{j_l})$. The equivalence is called an *alignment* if the elements of A and B are ordered and if the pairs in $L(A, B)$ are colinear, i.e., if $i_1 < i_2 < \dots < i_l$ and $j_1 < j_2 < \dots < j_l$. Many different equivalences exist and a scoring function is needed to rank them and to discriminate good equivalences from bad ones. We first briefly review different scoring schemes proposed

for pairwise structure comparison. Then we give an overview of different algorithms, which given two structure descriptions and a scoring scheme, aims to find the equivalences giving the highest score. A schematic overview of the steps involved is given in Figure 4.

The general problem is NP-hard (Lathrop, 1994); therefore, some simplifications have to be made either in the search or in the scoring.

4.1. Scoring equivalences

Assuming that one wants to assign high values to “good” equivalences, the comparison problem is one of finding an equivalence with a score as high as possible (maximization problem). The score of an equivalence will be high if the pairs are between elements with similar properties (for coordinates, if they can be superpositioned well—see below) and if relations between pairs of paired elements are similar.

When comparing two sequences, the traditional method is to align the sequences to pair up identical or similar residues. The score of the alignment is assigned using the similarity of the aligned (equivalenced) residues (as measured by a scoring matrix, e.g., Dayhoff (1978) and Henikoff and Henikoff (1992)) and as a function of the gaps inserted in order to obtain the alignment (Sankoff and Kruskal, 1983). Different

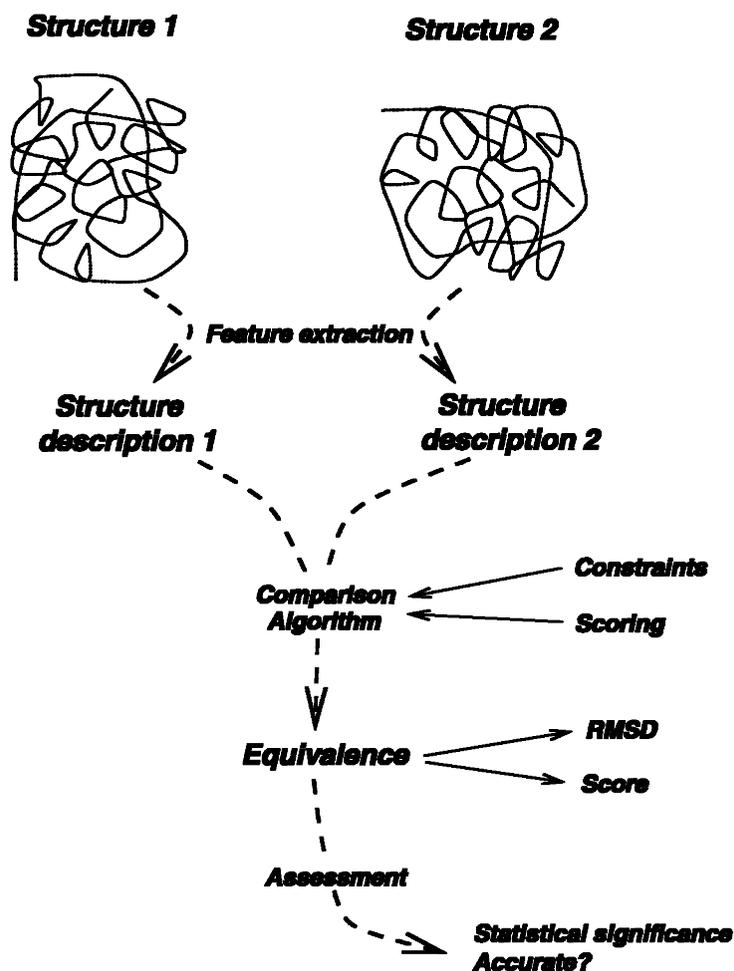


FIG. 4. Framework for pairwise structure comparison. First the relevant features are extracted and represented in structure descriptions. These are input to a comparison algorithm which finds an equivalence obeying certain constraints with scores high as possible. The equivalence is output together with the RMSD value. It can be assessed with respect to whether it is accurate (if a standard of truth is available) and as to whether it is statistically significant—i.e., if the similarity is stronger than could be expected by chance.

scoring schemes apply, depending on whether the alignment is to be local or global. Pairwise sequence comparison is very widely used, in particular in database similarity searches. For each similarity found, the alignment and its score should be given along with the statistical significance of the hit. Most often, this is given as an estimate of the probability of finding a hit at least this strong (having the same score or higher) by chance (Lipman and Pearson, 1985, Altschul *et al.*, 1990).

When comparing structures, an intuitive method is to try to put the structures on top of each other so that the equivalenced elements come as close as possible. The obtained distances can be used to quantify the similarity and to score the equivalence. This is called *superposition* of structures, and if the geometry of the structures are not changed in the process, it is referred to as rigid-body superposition.

Algorithms exist for superposing structure *A* on structure *B* by finding the superposition (*translation* of 3 distances and *rotation* of 3 angles) to minimize the *coordinate* root mean square deviation ($RMSD_C$) given by

$$RMSD_C = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i^A - x_i^B)^2} \quad (1)$$

where $(x_1^A, x_1^B), \dots, (x_N^A, x_N^B)$ are the coordinates (after superpositioning) of the equivalenced elements. For residue level structure descriptions, there may be one coordinate set per residue (e.g., C_α atom) or a number of coordinate sets per residue. An alternative measure is *distance* RMSD ($RMSD_D$). This alleviates the need for finding a translation and rotation of one of the structures and is given by

$$RMSD_D = \frac{1}{N} \sqrt{\sum_{i=1}^N \sum_{j=1}^N (d_{ij}^A - d_{ij}^B)^2} \quad (2)$$

where each d_{ij}^T is the spatial distance between elements *i* and *j* in structure *T*. The two measures are linearly related except that $RMSD_D$ is invariant under reflection.

It is observed (Brown *et al.*, 1996), that the translation is effected by relocating the origin of the coordinate system of each protein to the center of mass of its equivalenced elements. Finding the best superposition is easily done by standard pairwise least-square fitting algorithms (e.g., finding eigenvalues in a matrix constructed from the two coordinate sets (Kabsch, 1978)). Some methods will find a superposition with a mirror image of one of the structures (if this is best) and the complexities of this “feature” has been analyzed by Crippen and Mairov (see below).

Most scoring schemes for evaluating equivalences between structure descriptions contain factors related to the $RMSD_C$ or $RMSD_D$. Related measures include distance map similarity (Holm and Sander, 1993b) or contact map overlap (Godzik and Skolnick, 1994). Many structure comparison programs give as output an equivalence and the resulting $RMSD$ (or a closely related measure, such as weighted $RMSD$ or maximum $RMSD$) even if they do not use $RMSD$ internally to score equivalences (e.g., Taylor and Orengo (1989)). Internal measures can use information about, for example, amino acid types, physico-chemical properties, exposure/buriedness, secondary structure elements, and chemical bonds. Not all methods explicitly specify a scoring function to be optimized but can be based instead on a combination of (often heuristic) steps (Taylor and Orengo, 1989).

Falicov and Cohen (1996) use a scoring based on a minimal surface metric. Given a superposition of two structures, a surface between them can be described by a list of triangles, of which there are two types. Let a_i be the C_α of *i*'th residue in one of the structures, and b_j of the other. Then the triangles are either (a_i, a_{i+1}, b_j) or (a_i, b_j, b_{j+1}) . An atom might be the vertex in several consecutive triangles, thus allowing for insertion/deletion in a corresponding alignment. The triangulation corresponding to the minimal surface can be found by a dynamic programming procedure, and this gives the score for the transformation. A similar idea was proposed by Schulz and Schirmer (1979).

Another alternative scoring is URMS introduced by Chew *et al.* (1999). A structure is represented by a list of unit vectors $\{v_i\}$ with a common start point (origin). Each v_i is the unit vector along the vector $C_{\alpha_i}, C_{\alpha_{i+1}}$. They define URMS to be the minimum RMS distance for an equivalence between the two unit-vector representations. They claim that this measure is effective for finding small substructures shared by two structures and allows rapid computation using Fast Fourier Transform.

Assessing the significance of RMSD values obtained for candidate equivalence sets, one first needs to consider how many elements were equivalenced, as for random comparisons the expected RMS value is proportional to the square root of the number of equivalenced residues. Rules of thumb are often used (see Levitt and Gerstein (1998), Gerstein and Levitt (1998), and Section 6).

4.2. Comparison algorithms

A large number of methods exist for comparing pairs of sequences, pairs of structures, and also pairs of sequence and structure. The methods are used not only for comparing a given pair of objects, but also in database searches where the objects most similar to a query object are to be found as quickly as possible. A classification of the methods can be done along different axes, for example, the level of representation (atom, residue, frames, or secondary structure element), the constraints on topology (content, sequential, nonsequential (Brown *et al.*, 1996)), or techniques (superposition, dynamic programming, graph theory, geometric hashing, etc).

In this survey, we focus on the *algorithms* and make a classification of *approaches*. Our classification of the pairwise comparison methods comprises alignment-based methods, search methods, geometric hashing, statistical approaches, and clustering.

For each approach, we explain the fundamental principles and also explain (on a coarse level) some of the published *methods* in which the approach is used, focusing on the most-recently published methods. Figure 5 shows a schematic overview and the connections of the different approaches.

Pairwise structure alignment includes methods for aligning sequences where structural information is used and the colinearity is maintained. The most fundamental techniques for comparing structures are superposition and dynamic programming.

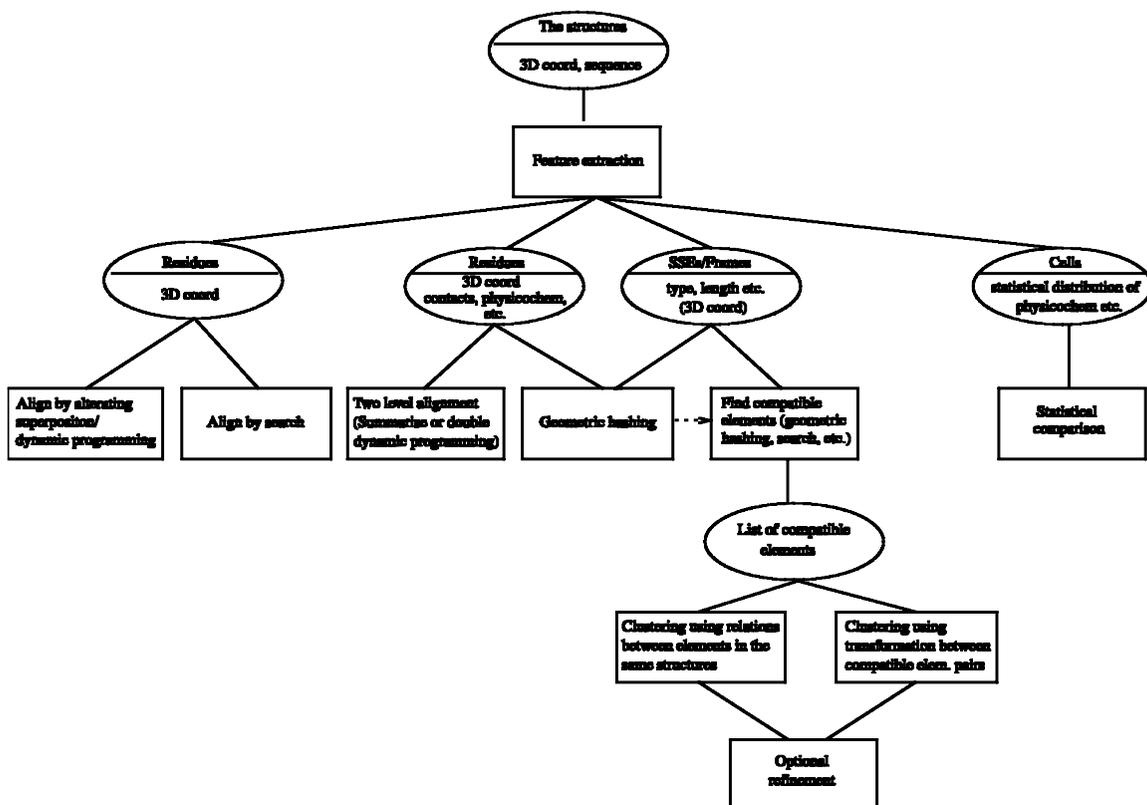


FIG. 5. Schematic diagram of the different approaches for pairwise comparisons. The ellipses representing information, the rectangles processing. The dotted line indicates that one process might be used by the other.

Dynamic programming (DP) (Needleman and Wunsch, 1970) for finding the best alignment of two sequences or structures uses a scoring matrix with a score for each pair of residues. One normally starts from the beginning of each structure and extends the alignment toward the end, filling in a DP matrix. An important assumption is that the optimality of earlier choices is not affected by decisions made later on. This assumption holds for sequence alignment using typical scoring schemes, e.g., edit distance or sum-of-pairs, but typically not for structure alignment. In the words of Gerstein and Levitt (1998), "structural alignment fails to converge globally because the possible matches for different segments are tightly linked as they are part of the same rigid 3D structure."

4.3. Alternating superposition alignment iteration

Given an equivalence of atoms from each structure, E_0 , a superposition algorithm can be used to find a transformation T_0 minimizing an RMSD measure. When the whole structures are superposed using this transformation, the distances between all pairs of atoms (residues) from the two structures can be used to define a new scoring matrix. This matrix can be used for a new alignment, and the equivalences (pair of residues) with least distances can be chosen for a new superposition. This iteration continues until convergence (the alignment is not changed) or some maximum number of iterations is done. Figure 6 illustrates the approach.

Rao and Rossmann (1973) and Rossmann and Argos (1975, 1976) first used this approach to align protein structures. Initial equivalences are found by searching for pairs of similar short fragments (3–4 residues). In the alignment phase, a constraint is used to assure colinearity (dynamic programming is not used). The scoring is determined by a probability function of the distances between the residues, and a measure of their conformational similarity is determined by use of the coordinates of the two (sequence) neighboring residues. In a related method, Cohen and coworkers (Satow *et al.*, 1986; Cohen, 1997) use dynamic programming in the alignment step. This is also used by Russell and Barton (1992) in the program STAMP, where the scoring matrix is the probability function of Rossmann and Argos. An analog iteration is used in the program CONTRUST by Ding *et al.* (1994), where the scoring between two residues consists of three components: sequence, local structure, and global structure components. The local structure component depends on the difference between the curvature and torsion calculated for the residues, and the global structure component depends on the distance between the two residues after superposition is done.

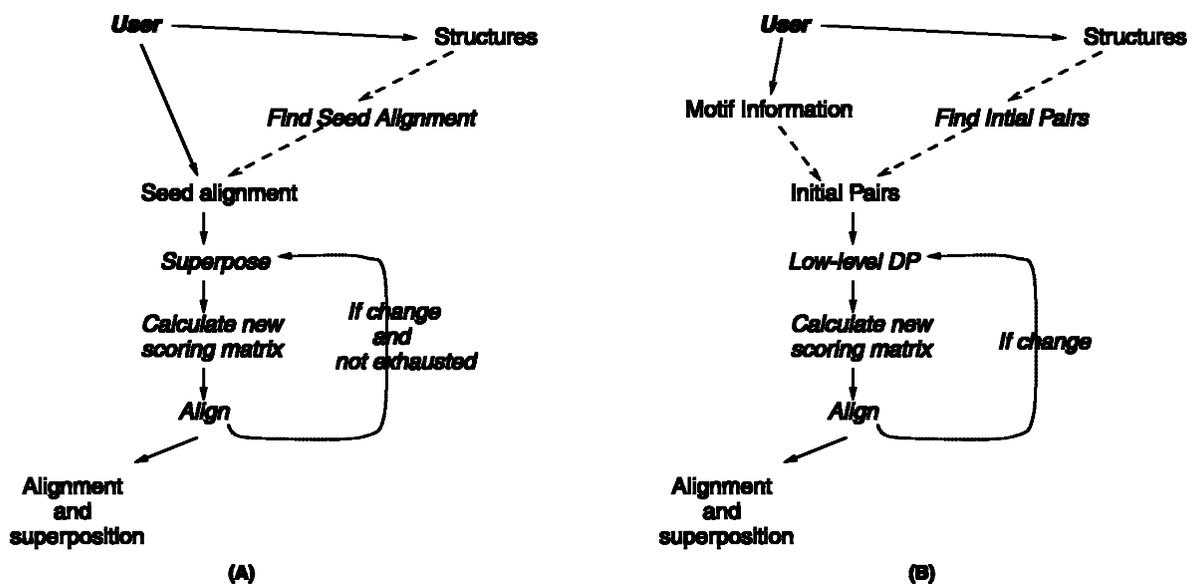


FIG. 6. (A) Outline of algorithm alternating between alignment and superpositioning. (B) For comparison, the outline of the SAP method is given.

The idea of Rossman *et al.* of automatically finding an initial equivalence is also used by Zu-Kang and Sippl (1996). Recognizing that the initial seed alignment is critical to the final result, they try different possible seed alignments and conclude that, depending on the seeds chosen, the alignments can be quite different (not share many residue pairs) while still giving similar quality measures (RMSD and number of aligned residues).

Petitjean (1998) uses iteration to find the “best” common subset of C_α coordinates of two structures, measured by RMSD. The initial alignment is chosen randomly, and as in the above methods, iterations with alternating superpositioning and alignment are performed to maximize the number of C_α coordinates aligned. An analog iteration is also used by Gerstein and Levitt (1998) in an extension of the ALIGN program (Cohen, 1997).

The method of performing alternating superpositioning and alignment has also been used for refining (postprocessing) the results found by other methods (e.g., Holm and Sander (1995, Section 4.6)).

4.4. Two-level alignment

Traditional dynamic programming (DP) guarantees to find an optimal alignment of sequences or structures if the scorings of matched pairs of elements are independent of each other and if inserted gaps can be penalized independently. When one assumes that the structures are already superposed, reasonable scoring schemes can be devised which allow the structures to be aligned optimally using DP (see Section 4.3). However, ideally, one might wish to simultaneously align and superpose the structures to optimize a score depending on how well aligned substructures superpose. Since any choice to align two substructures affects the scoring of the alignment of the complete structures, the independency requirement is violated and DP no longer can guarantee an optimal solution. Several heuristic algorithms which try to extend DP to solve this problem have been proposed.

A heuristic method was presented by Sali and Blundell (1990) in which several alignments are made, one for each type of relation. Each alignment is obtained using stochastic optimization (simulated annealing). The results are summarized into a residue-by-residue scoring matrix U where a residue pair is assigned a high score if the two residues are aligned in many of the relationship-based alignments. Finally, U is combined with property information to obtain a new scoring matrix, which is used in an alignment using dynamic programming.

Another method based on DP has been used in the program SSAP by Taylor and Orengo (1989). The method is called double dynamic programming since DP is used at two different levels. At the lower level, a series of DP matrices are calculated. For each, the alignment of one particular pair of residues is fixed. The highest scoring path from each lower-level DP matrix is propagated to the higher-level “summary” DP matrix in which DP is done to find the overall best alignment.

We will describe this algorithm in a bit more detail. Let i, k be residues in the two structures A, B , respectively, and let ${}^{ik}R$ be a lower-level score matrix. The element ${}^{ik}R_{jl}$ is a score showing the goodness of aligning A_j, B_l , given that A_i, B_k is aligned. For finding the score, local reference systems are defined at A_i and B_k , and the coordinates of the remaining residues in each structure are transformed into these coordinate systems (see Section 3.3.1). The score of aligning A_j, B_l depends on the distance between A_j and B_l in the respective coordinate systems defined by A_i and B_k , respectively.

The idea of Taylor and Orengo is similar to that used in the methods alternating between superpositioning and alignment (Section 4.3). While the latter methods find the DP matrix using an optimal superposition for the current equivalence, Taylor and Orengo do not need to decide (assume) one exact alignment to calculate the higher-level DP matrix. Instead, the residue pairs that participate in high-scoring lower-level alignments receive high values in the higher-level DP matrix and are likely to be included in the final alignment.

An iterative version of the SSAP algorithm was developed (Orengo and Taylor, 1990) which has been developed more recently into a program called SAP (Taylor, 1997). This is even more closely related to the methods of Section 4.3 (see Figure 6). In this version, lower-level DP matrices are only calculated for a subset of the possible residue pairs. In the first iteration, an initial set of residue pairs is used (randomly chosen or using local seed as secondary structures, burial, or motifs (Jonassen *et al.*, 1999)), and the set to be used in each subsequent iteration is calculated using the result of the last (high-level) alignment. The

procedure can be described as below, where $DP(R)$ is the result of aligning A, B using the matrix R , where R_{ik} is the score for aligning A_i, B_k . Therefore, $DP(R)$ returns a list of indices.

```

Q := 0                                Set high-level score matrix to zero
I := list of seed pairs
iter
  S := {0}m × n                        Set accumulated low-level score matrix to zero
  for each pair (ik) ∈ I do
    L :=  $DP^{(ik)}(R)$                 Low level dynamic programming
    for each pair (r, s) ∈ L do Sr,s := Sr,s + ikRr,s end
  end
  Q :=  $\frac{1}{2}Q + \log(1 + S/10)$       Update high-level score matrix
  L :=  $DP(Q)$                         High-level dynamic programming
  I := select new pairs
until convergence

```

4.5. Search

Several methods use general techniques for search in a well-defined search space. For example, May and Johnson (1995) use a genetic algorithm to find the superposition giving the best alignment. The search space is a set of transformations. Each transformation is scored by the score of the best alignment of the two structures under this transformation. The alignment is computed using dynamic programming where the score for a residue pair (one residue from each structure) is defined depending on the interresidue distance under the transformation.

A related method presented by Diederichs (1995) also performs a search for the best transformation. The search space is defined by a “Lattmann” angle space as a set of rotations associating equal volumes of angle space. For each rotation, a “translation grid” with 1Å spacing is defined, and for each grid point a count is made of how many C_α pairs (one from each structure) superpose well. This is used as one component of the scoring of the transformation. Optionally, the similarity of aligned residues and topological matching can be included in the scoring. Each good superposition is checked for local or global sequentiality. Local sequentiality is a measure of how frequently contiguous (backbone) fragments are superposed onto residues having the same sequence order. Global sequentiality is measured by the correlation between the indices i, j of matching residues. Different scoring schemes are used depending on whether, in addition to similarity in secondary structure architecture, one also wants topological similarity and whether direction of SSEs are taken into account.

Falicov and Cohen (1996), also search in the space of transformations. Each transformation is scored by the minimal surface metric (see Section 4.1). Three techniques were used for the search: Downhill Simplex, Powell’s Method, and Conjugate Gradient.

Holm and Sander (1996) describe a method based on a threading algorithm for sequence–structure comparison by Lathrop and Smith (1996). The sequence of their method is replaced by the structure (residue centers). The method is a branch-and-bound search method, where the solution space consists of all possible placements of the residues in a structure B relative to the SSE-segments of residues of another structure A . The structures are represented by distance matrices, and the upper bounds are calculated using distance matrices.

4.6. Geometric hashing

The aim of this technique is to discover common substructures, i.e., a set of elements from the two structures with the same mutual spatial relations. A highly redundant representation of the structures is used, which is *independent of rotation, translation, and sequence order*. This is normally obtained through the definition of local coordinate (reference) systems into which the coordinates of all (or a selected subset of) the elements (atom groups/residues/SSEs) are transformed. The methods use hash tables for storing and comparing local geometrical information (hence the name “geometric hashing”). The local reference systems are defined by use of 3D coordinate frames, and three points are needed for each frame (xyz).

The principle for the technique will be explained by a simple example using residue-based representation, where each residue is represented by the 3D coordinates of an atom (e.g., C_α).

Let A (the target) and B (the query) be two structures with m and n residues, respectively. For each $i, i = 2, \dots, m - 1$, define a reference system (R_i) using the neighboring residues ($i - 1, i, i + 1$) in A , and let $T(i, j)$ be a transformed representation of the position of residue j in reference system R_i . Further, let H be a hash table (with one dimension for simplicity) indexed by $T(i, j)$, such that an entry in H contains a list of reference systems. The hash table is filled by using all reference systems from A , by the following procedure:

Preprocess:

```

for  $i = 2, \dots, m - 1$  do
  calculate the reference system  $R_i$ 
  for  $j = 1..m, j \neq i$  do calculate  $T(i, j)$ ;  $H(T(i, j)) := H(T(i, j)) \cup (R_i)$  end
end

```

For each pair of elements i in A and k in B , we count the number of element pairs (j, l) such that j has the same spatial relation to i as l has to k and retain those for which this number is above a threshold (t). Then the pair of elements (i, k) is likely to be part of a common structure. The following procedure can be used (for illustration, see Figure 7):

Recognize:

```

for  $k = 2, \dots, n - 1$  do
  calculate the reference system  $R_k$ 
  set  $S(i)$  equal 0 for all  $i, i = 2, \dots, m - 1$  For counting
  for  $l = 1..n, l \neq k$  do
    calculate  $T(k, l)$ ; for each  $(R_i) \in H(T(k, l))$  do ++  $S(i)$  end
  end
   $C = C \cup \{(R_i, R_k) | S(i) > t\}$ 
end

```

The pair of elements (i, k) in C are likely to be part of a common substructure. To find the actual substructures, clustering can be performed (see Section 4.7).

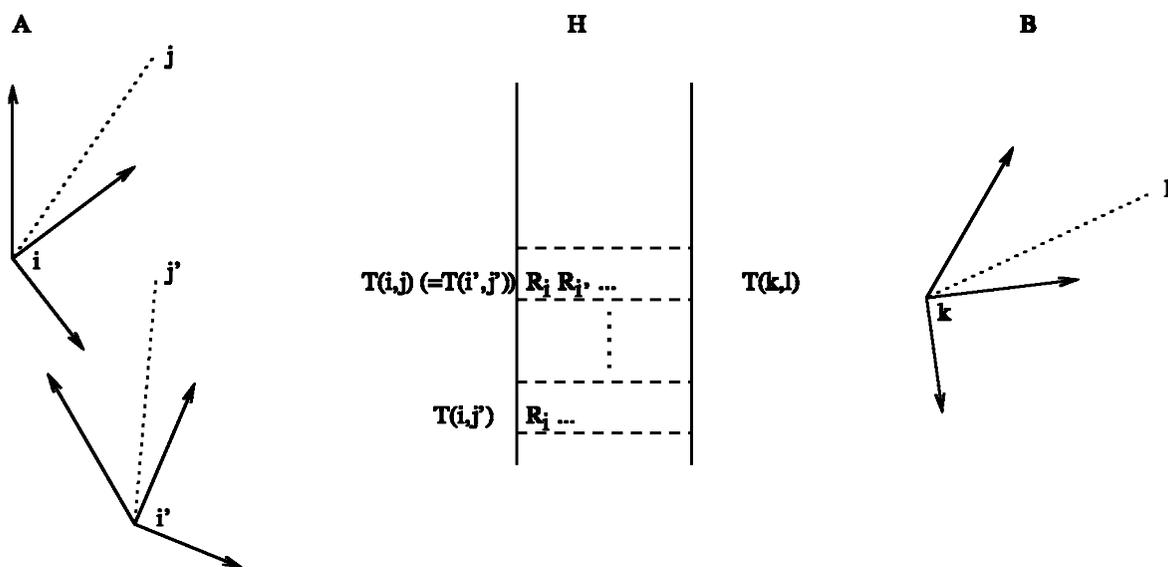


FIG. 7. A simple model for geometric hashing. The hash table is filled from structure A using the preprocess algorithm. When looking at the pair (k, l) when B is processed by the recognize algorithm, both $S(i)$ and $S(i')$ are increased by one.

All triples of points can be used for defining reference systems. In the above example, triples of residues were used. Since the number of reference systems can be very large, it may be desirable to use heuristics to reduce the number of triples to be considered. Also, labels can be assigned to the elements (Nussinov and Wolfson, 1991), and corresponding elements can be required to have equal labels (implemented by including the label in hashing $T(i, j, label)$). It is used also to simultaneously compare a structure to a set of structures A_1, A_2, \dots, A_n . The hash table is filled with data from each of the structures, but each datum is marked with the corresponding structure (Nussinov and Wolfson, 1991).

Holm and Sander (1995) have proposed a method where geometric hashing is used with SSEs as elements. The SSEs are represented as vectors, and right-handed coordinate systems are defined by using ordered pairs of SSEs, assigning the origin to be at the midpoint of the first vector, and directing the y-axis along that vector. By definition, the midpoint of the second vector lies in the z-positive yz halfplane. (A similar reference frame was used by Orengo *et al.* (1992).) The SSE vectors are stored at the location of their midpoint coordinates, for every coordinate system. Each cell in the hash table contains (a pointer to) a list holding the explicitly transformed coordinates of the SSE vectors, the elements constituting the coordinate system, and sequential number and type of the stored element. When deciding a match (i.e., element q in the query structure has the same relationship to two other elements (k, l) as an element p in the target has to (i, j)), the position of the midpoint of q and p , and the directions, must be similar to a specified degree; additionally, the involved SSEs must be of the same type and come in the same order along the proteins' backbones. The hash table contains additional information to be used under the comparison, such that finding the common substructures can be done without clustering. A refinement procedure analog to the iteration described in Section 4.3 is then performed.

Geometric hashing is often used for finding compatible elements in methods using clustering (see below).

4.7. Clustering

The result of a pairwise structure comparison is often a set or list of element pairs. Each pair contains one element from each structure and may have a score associated with it. The elements in each pair must be *compatible*. The requirement for compatibility varies. Residues may be defined to be compatible if they have the same amino acid or have amino acids in the same group. Secondary structure elements may be compatible if they are of the same type or have equal internal distances. Two or several pairs of compatible elements are *consistent* if the substructures consisting of the elements satisfy some constraints (e.g., the substructures are considered as similar).

The methods using clustering follow this scheme: (1) find pairs of compatible elements, (2) cluster using consistent compatible pairs to find the (k) substructure(s) with highest score(s), (3) optional refinement.

Let (A_i, B_k) and (A_j, B_l) be two compatible pairs. They can be joined (are consistent) if the substructure consisting of (A_i, A_j) is "similar" to the substructure consisting of (B_k, B_l) . To decide this, either *relations* between elements of the same structure or *transformations* between compatible elements from different structures are used. Figure 8 illustrates this, with examples. Transformation, in the most cases, means the transformation for an optimal superposition; consistency exists if two transformations are similar enough.

The clustering is then done based on either transformation or relation, or in some cases on both. In addition, an overall score might be used. The clustering problem is NP-hard, and heuristics are often used in practice. Several clustering techniques for use in comparing structures have been proposed in the literature. Each aims to group element pairs together so that the groups are as big or as high scoring as possible similar substructures. The methods differ in the following aspects:

- how the joining is done (whether if only a cluster and a pair can be joined, or two clusters can also be joined).
- how to test for consistency (is the new pair tested against the cluster as a whole, or against each pair in the cluster). Let a cluster $C = ((A_{i_1}, B_{k_1})(A_{i_2}, B_{k_2}) \dots (A_{i_m}, B_{k_m}))$, and a pair $P = (A_i, B_k)$. If the consistency is transitive, it is enough to test P for consistency against one of the pairs in C ; if not it must be tested against all.
- how is the next pair for clustering chosen
- can *parts* of elements be joined to clusters
- are the resulting clusters disjoint or overlapping

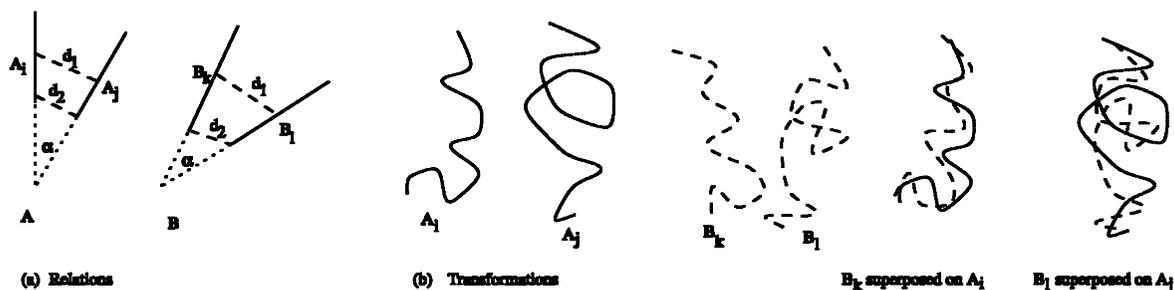


FIG. 8. Examples of relation and transformation. A_i, A_j are two elements from a structure A , and B_k, B_l are from the same for structure B . Assume (A_i, B_k) and (A_j, B_l) each are compatible. **(a)** Relations, the elements represented as vectors. The relation is between elements from the same structure, here represented by two distances and an angle (found by projection on a plane). The substructure (A_i, A_j) is consistent with (B_k, B_l) if the relation $(d_1, d_2, \alpha)_A$ is similar to the relation $(d_1, d_2, \alpha)_B$. **(b)** Transformation, the elements represented by the 3D coordinate for each residue. The transformation is between elements from different structures; in the example it is the transformation for the best superposition. If the transformation for the best superposition of B_k onto A_i is similar to the one for B_l onto A_j , then (A_i, A_j) is consistent with (B_k, B_l) .

The clustering methods generally allow for insertions and deletions and topological permutations (but there exist exceptions).

4.7.1. Clustering by use of transformations. The methods discussed in this section follow the scheme: (1) find compatible pair of elements, (2) find the optimal transformation between the compatible pairs, and (3) cluster compatible pairs using similar transformations, and (4) possibly perform final refinements.

The most widely used method for finding the compatible elements is geometric hashing. For each pair of compatible elements, a transformation describes how to transform the first element onto the second such that the "distance" between them becomes minimal. Often, *incremental* clustering is performed where it is tested for one pair at a time whether it can be included in the cluster by comparing its transformation to that of the first pair in the cluster.

Using geometric hashing (see Section 4.6), Nussinov and Wolfson (1991) first find pairs of compatible triplets. For each pair of compatible triplets, a transformation is calculated from one of the reference systems to the other, and in a clustering step, pairs with similar transformations are grouped together. The best superposition for the joined element sets is found.

Vriend and Sander (1991) cluster pairs of fragments. The fragments have at least some minimum length, and only fragments of equal length can be compatible. Compatibility is decided by first using distance geometry (internal C_α -distances). Fragments having similar internal geometry are superposed, and two fragments are compatible if the RMSD and the longest C_α - C_α distance between equivalent alpha-carbon are below specified upper limits. The compatible fragments might be elongated by adjacent residues. Then an incremental clustering is done, where a new fragment is compared to the first member of the cluster. To decide if the pair (A_i, B_k) is consistent with (A_j, B_l) , first the distance between the center of mass of A_i and A_j is compared to that between B_k and B_l . If this is satisfactory, the rotations of the best superposition of (A_i, B_k) and (A_j, B_l) are compared. After clustering, a final adjustment is done. In principle, each fragment pair should be used to start a new clustering process.

The SARF program by Alexandrov *et al.* (1992) also uses fragments as elements. Compatibility is decided by superposition (of fragments of equal length), and a compatible pair is given a score depending on the RMSD and the length. The clustering procedure is "best fit": a list of temporary clusters exists; initially every compatible fragment pair is a cluster. At each iteration in the clustering, all pairs of current clusters are temporarily grouped, and a score is calculated for each group. The group with the highest score is kept permanently. This continues as long as a clustering increases the score. In this manner each fragment pair does belong to only one cluster.

Another example of using geometric hashing to find pairs of compatible elements is the method of Fischer *et al.* (1994). They define spheres around each C_α atom. Geometric hashing (Nussinov and Wolfson, 1991) is used to find the number of C_α atoms inside two spheres (one from each structure) with "similar"

coordinates, and the two C_α atoms are defined to be compatible if there are at least 20. Pairs of compatible spheres having similar transformations are clustered by similarity of transformation. Later, the method was modified to also use secondary structure constraints (Fischer *et al.*, 1995). Spheres are constructed only around residues assigned α -helices or β -strands, and corresponding residues must have the same type.

Geometric hashing followed by clustering has also been used for analyzing structures described at secondary structure level. For example, Alesker *et al.* (1996) use geometric hashing for finding pairs of compatible secondary structure elements (SSEs). Then, for each pair, the six parameters describing the best superposition transformation are computed. For pairs of (compatible) pairs, a distance measure is defined depending on the similarity of the transformations. Clustering is performed in an iterative manner. In each iteration, the compatible pair with the minimal distance from one of the clusters is joined to this cluster. Pennec and Ayache (1998) also use SSEs and geometric hashing. They define an *information measure* for the transformations using different types of informations and cluster using "similar" transformations with decreasing information. When no more transformations can be joined to the current cluster, a new clustering is started on the rest transformations. The remaining transformations are then clustered using the same algorithm.

Verbitsky *et al.* (1997), use geometric hashing when comparing two structures, where one (the model) is allowed to contain a hinge point and the other is the target. The hash table is filled by use of the model, only saving the transformation between the coordinate frame at the hinge and the one at each of the other residues (C_α). The task then is to find the best candidates for the position (C_α) of the target corresponding to the hinge, together with two transformations, one for each part of the model. This is done using hash tables and clustering. Finally the candidates are evaluated by RMSD. The complexity is of $O(n^2)$, where n is the number of C_α atoms compared.

Clustering by transformation is also used by Chew *et al.* (1999). Let $R_{i,j}$ be the rotation matrix that best superposes two succeeding unit vectors ($i, i+1$) from structure A with two ($j, j+1$) from B (see Section 4.1). Two compatible connected substructures (elements) of length $k+1$ are found if $R_{i,j} \approx R_{i+1,j+1}, \dots, \approx R_{i+k,j+k}$ where $R \approx S$ if the two are approximately equal. Ordinary transformations are then used for clustering of compatible substructures.

When clustering transformations, it is also possible to define the similarity of two transformations by comparing the equivalences associated with each transformation. For example, Leibowitz *et al.* (1999) identify k residue equivalences and cluster these equivalences by the number of shared residue pairs (see Section 5.4).

Depending on the clustering procedure, a final adjustment or refinement might be worth while (Vriend and Sander, 1991; Holm and Sander, 1995). This is either an extension by searching for additional matching pairs of elements (Fischer *et al.*, 1994), or an iterative procedure, similar to the one described in Section 4.3.

4.7.2. Clustering by use of relations. Clustering methods compare element-based structure descriptions and build up shared substructures by first finding pairs of compatible elements (one element from each structure) and then progressively building up larger correspondences. In the previous section, we discussed methods where the consistency between substructure correspondences was checked by comparing transformations. In this section, we discuss methods where the consistency is assured by requiring similar relations between corresponding elements in the two structures.

Analyzing structures A and B , a list of all compatible elements (A_i, B_k) is generated. Each pair of compatible pairs is then tested for consistency. Let $\rho(A_i, A_j)$ denote the relation between any two elements (from the same structure). The pair (A_i, B_k) is consistent with (A_j, B_l) if and only if $\rho(A_i, A_j)$ and $\rho(B_k, B_l)$ satisfy some requirements and $\rho(A_i, A_j)$ is similar to $\rho(B_k, B_l)$. The problem then is to find the set of compatible pairs with highest score, such that any two of the compatible pairs are consistent. Note that this implies that testing a pair for joining a cluster requires testing for consistency against all the pairs in the cluster; this is due to nontransitivity.

In most methods, the elements are fragments or SSEs, represented as vectors. The relations considered when checking consistency include distances and angles (projected on specific planes) between the two vectors.

The problem can be formulated with help of graph theory in the following way: the structures are represented as graphs, where the nodes are the elements, and there is an edge between nodes A_i, A_j if the relation $\rho(A_i, A_j)$ satisfies some constraints. The nodes are labeled with element properties, and the

edges with relations. The problem is then, given two graphs, find a subgraph of each so that the subgraphs are similar and have maximum score. When the scoring is the number of compatible pairs, the problem can be easily solved by finding the maximal clique in a product graph (or connection graph) (Grindley *et al.*, 1993), which can, for example, be done by an algorithm by Bron and Kerbosch (1973). The (node) product graph is constructed by creating a node for each compatible node pair and an edge from node (A_i, B_k) to node (A_j, B_l) if they are consistent.

For example, Grindley *et al.* (1993) use SSEs, represented as vectors. SSEs of the same type are compatible, and the relation for deciding consistency is described by three parameters, distances and torsion angle between the elements, and the algorithm by Bron and Kerbosch is used. The same algorithm is used in the VAST protein structure comparison method (Madej *et al.*, 1995). Rufino and Blundel (1994) also use the same algorithm as Grindley, but the constraint for being compatible and also the relations are a bit stronger. A final dynamic programming step is performed for refinement.

Koch *et al.* (1996) also use graph theory and SSEs, but for making an edge between A_i, A_j they require that some of the van der Waals volumes around the atoms of A_i overlap with some of the van der Waals volumes around the atoms of A_j . The relations (edge labels) are “parallel,” “antiparallel,” or “mixed” (additional information can be added easily). They argue that it is sufficient to find *connected* subgraphs, thus reducing the problem. The algorithm is very similar to that of Grindley *et al.* but they define an *edge* product graph, in contrast to the *node* product graph. A node in an edge product graph consists of two edges (one from each structure). The difference is illustrated in Figure 9.

The method of Grindley *et al.* has been generalized by Mizuguchi and Gō (1995) who assign a *continuous number* to the edges of the product graph (thus generalizing from consistent/not consistent to giving a score of the consistency). In that way, each cluster is given a score $-\sum_{i \neq j} d((A_i, B_i), (A_j, B_j))$, where d is a measure of dissimilarity. They then search for the highest scoring clusters using a *parallel* iterative clustering procedure, where (the best) m new clusters are generated at each iteration, with the combination of two existing clusters.

A version of the SARF program, SARF2 (Alexandrov and Fischer, 1996), uses SSEs instead of backbone fragments. Compatibility is decided by SSEs of same type. The relation consists of five parameters describing distances and angle, and upper limits on these are used to decide if there will be an edge A_i, A_j . A recursive incremental clustering procedure is used. After finding the clusters, they apply an iterative procedure for refinement.

The DALI program of Holm and Sander (1993b) starts by using overlapping hexapeptide fragments. The decision of compatibility is done by use of distance matrices (C_α - C_α distances), and the relations also are distance matrices. A list of *contact patterns* (pairs of compatible fragments, two from each structure) are formed (which, in effect, is the product graph). The contact patterns are used for finding *seeds* for starting the clustering, which is done by a simulated annealing (Monte Carlo) search. An elastic score for the similarity of corresponding distances is used. Each step in the search is either adding or removing a pair of tetrapeptides (one hexapeptide generating three overlapping tetrapeptides). Several searches are done in parallel so that one can obtain different solutions at the end. A number of techniques are used for speeding up the computations.

Kleywegt and Jones (1997) describe a program, DEJAVU, which searches for similarities in structures where the elements are SSEs. The elements are described by type, number of residues, and the C_α coordinates of the first and last residue. Two elements are compatible if they have the same type. The user can specify which type of distances are to be used (center-to-center, etc.), as well as whether directional and topological constraints are to be included. Weights to a scoring scheme can also be given. The clustering algorithm is a depth-first search. In a final step, elaborate iterative refinements to optimize the superpositioning of the two structures is performed.

Escalier *et al.* (1998) use fragments, and compatibility is decided by structure similarity. A pattern-driven (simultaneous) searching procedure is used for finding compatible fragments from two structures, based on the 3D coordinates. One node in the search tree represent a subfragment (by 3D coordinates) from one of the structures (query), and every occurrences of this in the other structures (scene) is saved in a list. Two search nodes, of equal length (k atoms) and with $k - 1$ common atoms, are then used to generate a search node of $k + 1$ atoms. The relations used for clustering are the positions of corresponding atoms (residues) and the distances between the atoms in two fragments: the pairs (A_i, B_k) and (A_j, B_l) are consistent if (1) every atom a appearing in both A_i and A_j should correspond to the same atom in B_k

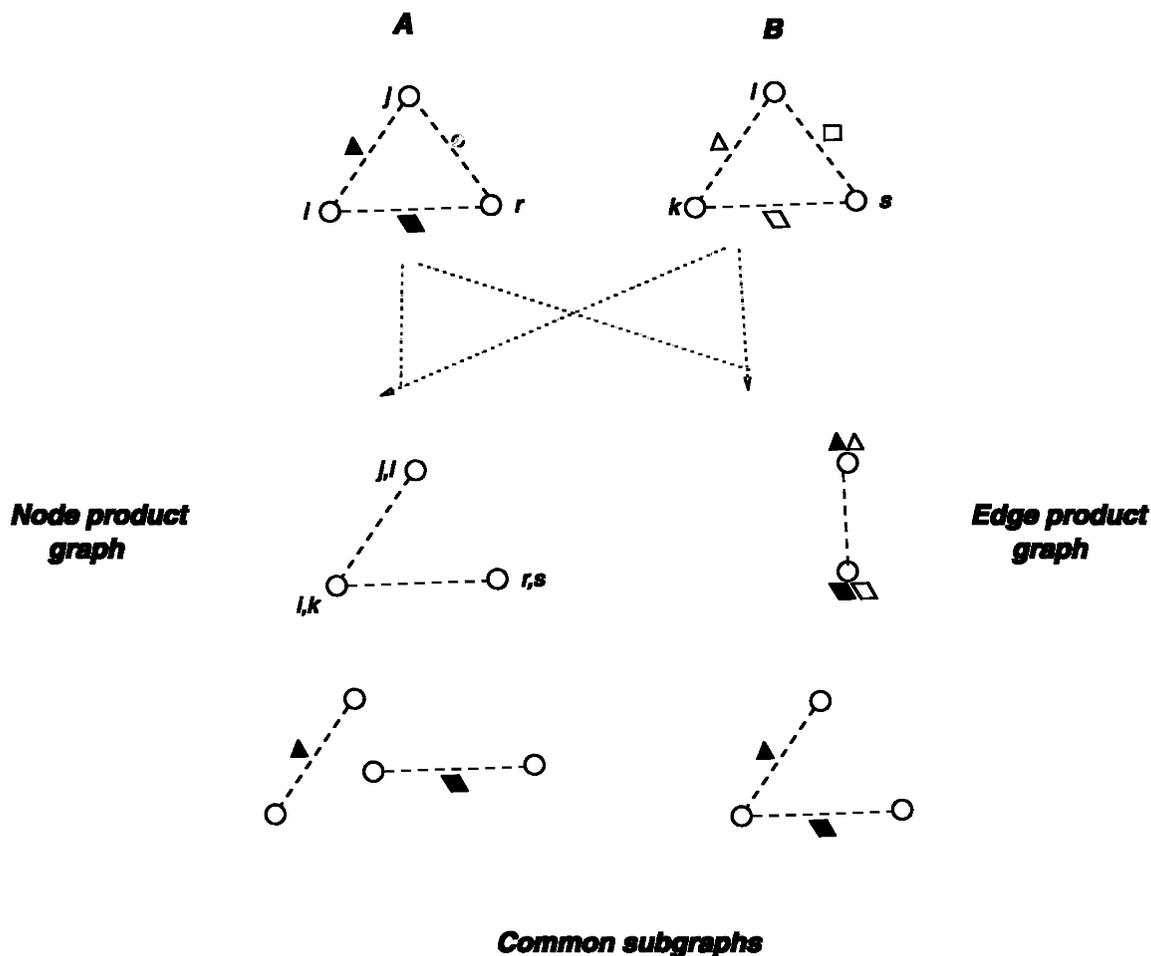


FIG. 9. The difference between the node product graph and edge product graph is illustrated by an example with three nodes in each substructure; $(A_i, B_k), (A_j, B_l), (A_r, B_s)$ are compatible, and the relations $\rho_{ij} \sim \rho_{kl}$ and $\rho_{ir} \sim \rho_{ks}$ but not $\rho_{jr} \sim \rho_{ls}$. The node product graph contains two cliques and hence results in two common substructures, each with two nodes and an edge product graph in the common substructures $(\rho_{ij}, A_i, \rho_{ir}), (\rho_{kl}, B_k, \rho_{ks})$. The edge product graph contains only one clique corresponding to one shared subgraph containing two edges.

and B_l , and (2) the distance between atoms in A_i, A_j should not deviate much from the distance between corresponding atoms in B_k, B_l . The score of the clustering problem is the number of corresponding atoms. The clustering is performed by a branch-and-bound search procedure.

Russell (1998) describes a method for finding similarities at residue level which are possibly related to functional sites (active and binding sites). The method finds sets of residue pairs (one from each structure) with amino acids from the same group and with similar orientation of the side chains. Sequence order need not be preserved. Only conserved residues are considered (requiring a multiple *sequence* alignment); amino acids with only carbon and hydrogen in their side chains and residues involved in disulphide bindings are excluded. Compatibility between the remaining residues is decided by amino acids from the same group (Russell defines nine groups). Each residue is represented by three atoms, and for consistency it is required that the distances between corresponding atoms in two residues (in the same structure) fall below a limit and that corresponding interatomic distances in the two structures are sufficiently close. The clustering is incremental, implemented as a recursive depth-first search. A new pair is tested for consistency against all pairs in the running cluster. A weighted RMSD measure is calculated, and the statistical significance of the similarity is assigned from analysis of randomly generated side-chain patterns.

4.8. Statistical approach

These methods are used when the proteins are characterized by some distributions of biological attributes and are compared by using a (protein-length independent) function that computes a distance between the proteins. The three methods classified as Content methods in Brown *et al.* (1996) are all statistical.

The methods using space-based description are all statistical (Section 3.2). A common pattern can be found comparing sets of corresponding cells. This is used by Bagley and Altman (1995), where the cells are shells around a common center. The substructures are described by the set of atoms within the cells, along with their 3D coordinates. In addition, they include user-defined properties, such as the types of atoms, chemical groups, amino acids, secondary structures, charge, polarity, mobility and solvent accessibility.

The statistical approach is also used by Kastenmüller (1998) to compare structures as a whole. Histograms are used for representing the distribution of the number of atoms in each cell and the applicable properties.

5. MULTIPLE STRUCTURE COMPARISON—PATTERN DISCOVERY

The problem of discovering a pattern from a set of biological objects (structures or sequences) can be looked upon as a machine learning problem (e.g., Brazma *et al.* (1998) and Conklin (1995)).

Brazma *et al.* (1998) classify the methods for discovering deterministic sequence pattern as *sequence driven* and *pattern driven*. Another term for sequence driven would be *comparison based*:

1. *Pattern driven*—methods that search a solution space for patterns with a high fitness value with respect to the input (training) objects.
2. *Comparison based*—methods that find common patterns through (pairwise) comparisons of the input objects.

A survey treating both pattern-driven and comparison-driven methods for sequences is presented by Brazma *et al.* (1998). For structures, the majority of the methods are comparison based.

Extension of a pairwise method to the multiple (comparison based) case ($n > 2$ structures) can be done mainly by three approaches (see also Escalier *et al.* (1998)):

1. Pivot (PI) uses one object as the pivot and compares it successively to all the other objects. The results are then compared to find the part of the pivot that has similarities to all the other objects, or in at least $k - 1$ of them.
2. Linear progressive (LP) starts with one object and successively compares the other objects to the results. The results are descriptions of the similarities of the objects involved.
3. Tree progressive (TP) compares (results) using a (possibly implicit) tree, where the leafs are objects.

All pairwise methods can be extended to the multiple case by the PI approach. Note, however, that a description of the similarities is not necessarily given; only the part of the pivot which occurs in $(k - 1)$ other objects, and where, are given. It may be possible to describe the similarities as a pattern (for different pattern descriptions, see Section 5.1). Ideas from the pivot approach can be combined with pattern-driven, search-based methods to obtain a pattern common to the pivot and at least $k - 1$ other objects simultaneously (e.g., Jonassen (1997)).

If a pairwise comparison method can express each of the (one or several) best local similarities as a pattern, and if these patterns can be compared to other objects, then the method can be extended to the multiple case by using the LP approach. Also, if patterns can be compared to find common generalizations, the TP approach can be used.

As noted in Brazma *et al.* (1998), methods for multiple sequence comparison based on pairwise methods have the weakness that, in each pairwise comparison, only information about the two sequences/structures is available and so alignments/patterns optimal for the whole set can be missed if they are not also optimal for every pair. Therefore, pattern driven methods may first be used to find elements common to all structures and these can then guide the pairwise comparisons in a more-detailed structure-based method. For example, output of the local motif discovery tool SPPratt (Jonassen *et al.*, 1999) has been used to guide the multiple alignment program MSAP (Taylor *et al.*, 1994).

An alternative to PI, LP, and TP is the following:

- All-by-all (AA) calls for all $n(n-1)/2$ pairwise comparisons to be performed to find common pairwise similarities. Comparisons (e.g., intersections) are done on these results to find similarities occurring in, for example, k of the objects (a clique of similarity relations).

Using AA, the result of all pairwise similarities can be used to identify similarities between all objects (or some minimum proportion of the objects). In a sequence comparison method by Vingron and Argos (1991), the pairwise similarities are represented by dot matrices (reflecting residue similarities) and they are combined using matrix multiplication to filter out similarities not shared by all objects. We are not aware of methods using this approach for comparison of structures.

In principle, PI and AA are able to give the same results if the similarity relation used is transitive. Transitivity exists if similarity between A and B and between B and C implies that A is similar to C . However, the result of the PI and AA approaches will depend on the exact implementations (for example, how similarities are scored and how many are included in the analysis).

The structure of the rest of this section is as follows. First is a discussion on how patterns can be described. Then follows an overview of methods for multiple structure analysis: first superpositioning methods, then alignment- and clustering-based methods, and finally search methods which take into account information from all structures simultaneously.

5.1. Pattern representation

In this work, a pattern is an object to which a structure description can be matched. If the pattern is deterministic, the structure description either matches or does not match a pattern; if the pattern is probabilistic, the result of the matching is a number (e.g., a probability). In a pattern discovery method, one needs to define the class of patterns to be considered. For a structure comparison method, it is possible to define pattern classes which allow one to describe the outcome of the comparison (e.g., a consensus representation of the aligned substructures).

5.1.1. Sequence patterns. If only primary structure is considered, the patterns will be sequence patterns. A common form of deterministic sequence patterns is regular expression type patterns, such as those used in the PROSITE database. A simple example taken from this database is the pattern used for the classical zinc finger C2H2 motif, 'C-x(2,4)-C-x(3)-[LIVMFYWC]-x(8)-H-x(3,5)-H' that matches all sequences containing a C followed by between two and four arbitrary residues followed by a C and three arbitrary residues, each one being L, I, V, M, F, Y, W, or C, etc. (For a fuller treatment of regular expression type sequence patterns, see Brazma *et al.* (1998).) Probabilistic sequence patterns include profiles, hidden Markov models, neural networks, and stochastic context-free grammars. The Probabilistic pattern formalisms give greater expressive power, but on the other hand have more parameters to be estimated and are less easily interpreted by humans.

Some protein families are difficult to describe adequately using sequence patterns; that is (1) the pattern does not perfectly discriminate between family members and nonmembers and (2) it does not give a complete description of the features critical to the function/structure class associated with the family. When structural information is available, it is possible to define a structure pattern.

5.1.2. Structure patterns describing sites. Kasuya and Thornton (1999) studied the three-dimensional structure of protein fragments matching PROSITE patterns. They matched each PROSITE pattern against all structures in the PDB databank and studied cases where PROSITE patterns give false positives (matches in unrelated proteins). By comparing the structure of each matching fragment to that of a true matching fragment and requiring similarity to accept the match, they achieved increased specificity. Jonassen *et al.* (2000), in an independent but closely related study, have arrived at equivalent results. Additionally, they found that sensitivity can sometimes be improved by relaxing the PROSITE pattern by allowing for approximate matching but requiring structural similarity to a known positive.

The functionality of protein active sites arises from their three-dimensional structures. The residues forming a site need not be close in sequence, and the sequences of the proteins sharing an active site are sometimes extremely divergent. For example, the bacterial and eukaryotic serine proteases share a

very similar active site while having very different structures. Especially when the residues involved in the site are nonlocal, sequence patterns (deterministic and probabilistic) cannot recognize such structures. However, one can define structure motifs describing the spatial configuration of the active site residues.

Wallace *et al.* (1996, 1997) have constructed the PROCAT database which contains manually defined patterns, called "3D coordinate templates," which describe active sites (<http://www.biochem.ucl.ac.uk/bsm/PROCAT/PROCAT.html>). Each template specifies the spatial position of a set of atoms in each residue making up the site and the amino acids allowed for each residue. The sequence positions of the residues are not constrained. A similar description (called "local packing patterns") is proposed by Jonassen *et al.* (1999), where for each residue the allowed amino acids are given together with one coordinate set (mean side chain atom). The order (but not distance) of residues along the sequence is given. Similar patterns discarding the sequence order of the aligned residues can be discovered by the methods using geometric hashing (see Nussinov and Wolfson (1991), Fischer *et al.* (1994), and Fischer *et al.* (1995), Section 4.6).

Yet another similar description formalism, named "Fuzzy Functional Forms" (FFFs), is proposed by Fetrow and Skolnick (1998). FFFs specify, for a set of residues, which amino acids are allowed, the C_{α} - C_{α} distance (average and variance) for each pair of residues, and constraints on the sequence distance between the residues. For example, they specify a FFF for the disulfide oxidoreductase activity of the glutareoxin/thioredoxin protein family: two cystein residues separated by two residues and an α -carbon distance of 5.5\AA ($\pm 0.5\text{\AA}$) close to a proline residues (distance requirements included).

5.1.3. Structure patterns describing folds or architectures. We say that two proteins (or domains) share the same architecture if they have secondary structure elements (SSEs) with similar spatial relative positions. If, additionally, the SSEs are connected in the same way along the sequence, the two proteins have the same fold. The shared architecture or fold of a set of proteins can be described using structure patterns. This can be done at residue, fragment, or SSE level.

Gerstein and Altman (1995) define the core of a family of aligned protein structures as a collection of residues for each of which position (mean and variation) information is given. The descriptions are made from multiple structure alignments and therefore colinearity of aligned residues is maintained. This is equivalent to the internal representation in the extension of the SSAP comparison method to multiple structures (Taylor *et al.*, 1994).

A number of methods for structure comparison perform their analysis at SSE level. Typically, the structures are described using graphs (Section 4.7) and the common substructures are found as subgraphs common to the graphs representing the individual structures. When finding common subgraphs, one can choose whether the solutions should be restricted to obey the sequence order of the SSEs. Different methods have been reported using similar structure representations and matching algorithms (see, for example, Koch *et al.* (1996), Alexandrov and Fischer (1996), White *et al.* (1994), and Gilbert *et al.* (1999)).

5.2. Superposition

Several methods for multiple superposition based on a correspondance (alignment) exist. To avoid biasing the superposition towards a specific (pivot) structure, a simultaneous superposition should be used, resulting also in an average structure.

MNYFIT is a program by Sutcliffe *et al.* (1987) that uses pairwise superposition to iterate to a common average structure. Weights can be used, giving more strength to topologically equivalent positions. Kearsley (1990) building on the work by Gerber and Müller (1987), optimizes simultaneously the superposition of all possible pairs by minimizing the sum $\sum_i^{n-1} \sum_{j=i+1}^n \sum_k^{ns} (Q^i x_k^i - Q^j x_k^j)^2$, where Q^i is the rotation of structure i (weights can be allowed for). One structure is selected as static, and all the centroids are shifted to the origin of coordinates. Minimizing is done by Rational Function Optimization, using iteration for convergence. Shapiro *et al.* (1992) and Diamond (1992) describe other methods for optimally superposing all pairs of structures simultaneously. The method of Diamond has time complexity $O(n)$, and the optimization is done without referring to the original coordinates of the structures.

5.3. Multiple structure alignment

5.3.1. Progressive alignment. The multiple alignment of Sali and Blundell (1990), Russell and Barton (1992), Ding *et al.* (1994), and May and Johnson (1995) all follow the same algorithmic scheme:

- Perform all $n(n - 1)/2$ pairwise alignments, and give a similarity score to each of them.
- Use the scores to generate a dendrogram.
- Perform alignments following the dendrogram from the leaves to the root. When aligning one structure with a subalignment, or two subalignments, the same procedures as for the pairwise alignments are used. The score when aligning two subalignments is *in principle* done in the same manner for all the methods:

$$W_{ij} = \frac{1}{n_1 n_2} \sum_{l=1}^{n_1} \sum_{k=1}^{n_2} \begin{cases} W_{i'l'j'k}^{lk} & \text{if two residues are compared} \\ \text{some value} & \text{if a residue and a gap are compared} \\ \text{some value} & \text{if two gaps are compared} \end{cases}$$

where n_1 and n_2 are the number of structures in the subalignments, i, j are the positions (columns) compared, and i', j' are the true indices for residues from structures l, k , occupying the columns i, j .

Taylor *et al.* (1994) have also adopted the progressive approach in a fusion of the pairwise method SSAP (see Section 4.4) and the multiple sequence alignment program MULTAL (Taylor, 1988). First, each pair of structures are compared using SSAP to assess all pairwise similarities. The most similar structure pairs are aligned independently into a consensus structure and then the pairwise similarities between all single and consensus structures are calculated so as to progressively bring together the most similar of these at each stage. Note that this approach differs from the TP approach in that the progression of the pairwise alignments is not determined by a precomputed dendrogram; instead, at each stage the similarities are recomputed and the most similar objects are aligned first.

The main problem when going from sequence to structure is to define the structural equivalent of a consensus sequence (or profile) for use when aligning alignments. In SSAP, interatomic vectors are used, and the consensus used is an average vector and its variance (error measure). The alignment process results in multiple sets of interatomic vectors (one vector set for each residue) that are not necessarily mutually consistent (due to, say, relative domain movements). Procedures for making them consistent and then calculating the 3D coordinates of the consensus are incorporated.

Gerstein and Levitt (1998) use a similar but simpler method. All pairs of structures are aligned, and the structure that is on average closest to all other structures, the “median structure,” is identified. A multiple alignment is formed by the alignment of all the other structures onto this median structure by consistently combining the alignments. They argue that this simple approach is adequate when the number of structures is small.

5.3.2. Search for global multiple alignment. All types of searching can, in principle, be used to search for multiple alignment of structures, though they do not guarantee finding the best. Each point in the search space represents an alignment and is evaluated by use of a scoring scheme. Godzik and Skolnick (1994) use a multidimensional lattice chain and a Monte Carlo (simulated annealing) approach for making a multiple alignment using different scoring schemes: Dayhoff substitution matrix, RMSD, C_α - C_α distance difference, and contact map overlap. Hence, the algorithm can be used for both sequence and structure comparison. An alignment is represented as a lattice chain, which for d structures consists of a list of elementary vectors of length d : $[1, 0, \dots, 0]$, $[0, 1, \dots, 0]$, \dots , $[0, 0, \dots, 1]$. There is a one-to-one correspondance between an alignment and such a chain. In a single step of the search algorithm, a chain fragment is chosen at random and replaced by another.

5.3.3. Finding patterns from multiple alignment. When a multiple alignment is done (establishing correspondance between residues), superpositions can be performed to find those residue positions with small spatial variations, possibly being a *core*. This can be done as by Gerstein and Altman (1995): the structural variation of each atom (residue) is computed, and a cutoff value is used to remove atoms with high structural variation. Then they iterate by making a new multiple superposition on the atoms not removed. New atoms can be added if they now satisfy the cutoff value, and this removal/adding can be done until convergence.

Rinaldis *et al.* (1998) develop *profiles* representing the most-conserved residues on the surfaces. A multiple superposition is done, and the structures are represented in a 3D grid, whose cells have 2\AA sides, a residue represented by a pseudoatom averaging the side chain. Only exposed residues are retained. For

each cell a column is constructed in a profile, as described by Gribskov *et al.* (1987), with the 20 values for each amino acid and the 3D coordinates of an average pseudo atom.

5.4. Clustering

The result of a pairwise comparison using clustering of transformations is a list of corresponding residues (an equivalence). For the multiple case, a pivot can be chosen and equivalences between the pivot and each of the other structures can be combined. Extension by a progressive method requires calculating an average structure at each iteration, representing the result of the comparison. We are not aware of any methods using clustering and progressive extension to multiple structure comparison.

A method using the pivot approach is presented by Leibowitz *et al.* (1999). The pivot is compared to each of the other structures using geometric hashing to find sets of k residues in the pivot for which similar sets of k residues are found in each of the other structures. Next, the pivot is compared to each of the other structures. A superposition is performed for each of the identified k residue equivalences and the equivalence is extended, if possible. The resulting equivalences are clustered by the number of shared residue pairs. Now equivalences between the pivot and the other structures can be found by taking the intersection of equivalences between the pivot and each of the structures. In principle, all combinations of pairwise equivalences must be tried, but techniques are devised for limiting the work. In conclusion, information from all structures is used in a pivot step before the pairwise comparisons, which are performed between a pivot and the remaining structures. The results of the pairwise comparisons can be combined to find sets of residues present in all the structures.

Two of the methods using relations (Section 4.7.2) are extended to the multiple case, both using the pivot-based extension. Koch *et al.* (1996) first compute all sets $E(G_1, G_i)$, $i = 2, \dots, n$, where $E(G_1, G_i)$ is the set of all connected maximal common substructures in G_1 and G_i represented with the edges. Then they successively intersect each edge set from $E(G_1, G_2)$ with each edge set from $E(G_1, G_i)$, $i = 3, \dots, n$ to find subgraphs of G_1 that are also contained in one or more G_2, \dots, G_n . These subgraphs may be disconnected. Other types of analyses can be done.

Escalier *et al.* (1998) extend to the multiple case in a way analogous to that of Koch *et al.* The same procedure as for the pairwise case is used, but $n - 1$ lists of occurrences are kept in the step for finding compatible elements (one list for each structure compared to the pivot).

5.5. Simultaneous search

A few methods using simultaneous search are published. The Spratt method of Jonassen *et al.* (1999) is a pattern-driven method; the elements are residues. The method finds patterns which occur in at least k of the structures. Each residue's 3D coordinate is represented by a pseudo atom calculated from the side chain atoms. For each residue a string of amino acids is generated from the residues falling inside a sphere around the pseudo atom, satisfying the sequential order. A sequence pattern discovery method called Pratt (Jonassen *et al.*, 1995; Jonassen, 1997) is then used to discover common patterns in these strings, using a simultaneously heuristic depth-first search method. The occurrences are further constrained to have all pairwise RMSD values below a given threshold. The patterns are described in a PROSITE-like manner and give constraints on a set of residues. The constraints are on the coordinates, the allowed amino acid types, and their relative ordering along the backbone.

Wako and Yamato (1998) divide the space around a structure into nonoverlapping tetrahedrons (using Delaunay tessellation). The tetrahedrons are constructed using C_α atoms as vertices. On the bases of each tetrahedron and its up-to-four neighboring tetrahedrons, constructed from 20 atoms, a code of up to 32 digits is found (each digit being between 1 and 8). Thus, the codes represent nonsequential overlapping substructures. For each structure a list of such codes is found. A search procedure then searches for codes occurring in all or several of the structures. Each such code constitutes a potential motif. The codes do not represent the 3D structure explicitly, and superposition of the substructures is necessary to assess the degree of resemblance.

Su and coworkers (Su *et al.*, 1999; Cook *et al.*, 1996), represent the structures as linear graphs (hence, sequential) with one version on residue level and one on SSE-level. The SSEs constitute the nodes of the graph labeled with type, length, and direction (right- or left-handed helix, and parallel or antiparallel sheets). The linear graphs representing the set of structures are concatenated, and the edges inside a structure are

represented by a label. The procedure for finding patterns (substructures) is a beam search of width k , and the evaluation of each node in the search tree is done by use of minimum description length: each occurrence of a pattern is replaced by a link, and the length of this compressed representation is found. The search starts with k search nodes representing one element each. At each iteration, each of the k nodes are extended by all possible neighbors (one from each occurrence), and the k best are kept for the next iteration. The matches may be inexact, and by combining the search with hill climbing, they obtain a polynomial run time.

Gilbert *et al.* (1999) have developed a program for discovering supersecondary motifs in TOPS-cartoons (Flores *et al.*, 1994) by a pattern-driven approach. The TOPS-cartoons are two-dimensional representations of protein folds showing the secondary structure elements and illustrating the relative spatial position and direction of these elements.

5.6. Learning the description of known sites

Fetrow *et al.* (1998) represent active sites by three-dimensional descriptors, termed “fuzzy functional forms” (FFF), based on the geometry, residue identity, and conformation of protein-active sites. An FFF contains the mean and allowed variations of the internal distances between the functionally important residues. They start with multiple sequence alignment and available information about residues important for functions. Then they superpose these and iterate to learn a specific and unique form.

Bagley and Altman (1995) use a machine learning system, FEATURE, with sites and nonsites, to learn description of sites (the distribution of the properties over the cells). The method then learns the description of known sites. Wei *et al.* (1997) present an application to calcium binding sites.

Generally, finding motifs in statistically described structures can be done by characterizing a set of structures by some “middle” distribution, dividing the structure space into subspaces, and comparing the subspaces from the different structures. Since the number of subspaces is indefinite (if real values), a finite number must be chosen, and the goodness of such a method would depend on how the subspaces are chosen.

6. ASSESSMENT

When comparing structures, it is widely known from practical experience (e.g., Taylor and Orengo (1989)) and from more systematic investigations (Godzik, 1996; May, 1996) that, beyond close similarity, there is no uniquely correct structural alignment of two proteins. Different alignments are achieved depending on which biological properties and relations are emphasized in the comparison. This adds a complicating element to the assessment of the result of a comparison.

Most assessments are done by comparing the result to manually determined alignments. This is systematically done by Gerstein and Levitt (1998), where an assessment of the ALIGN program (Cohen, 1997) is made using the SCOP classification (Murzin *et al.*, 1995) as a “gold standard.” The program was tested for 2107 pairs of structures, where each structure in a pair was from the same superfamily (hence, similar structures). They made a plot of the results, plotting a normalized RMSD against the number N of matched residues. The RMSD was normalized by $225\text{RMSD}/(N + 135)$, and a horizontal demarcation line was drawn at 4Å. Only 32 of the 2107 pairs got a normalized RMSD value above the demarcation line. They also assessed their multiple method against 9 SCOP superfamilies.

6.1. Statistical assessment

Several methods are assessed statistically. The scoring found for an alignment of native structures is then compared against what is expected by chance, most often implemented as what is expected by aligning random structures or using fragments of nonrelated proteins. RMSD is used alone or as part of most scoring schemes; hence, knowing the behavior of RMSD of random structures is important.

6.1.1. Structural fragments models. Alexandrov and Go (1994) present an analysis for finding the significance of similar SARFs. For a fixed length L , they pick up all fragment pairs of this length in two unrelated structures and find the value R_L from the condition that only 1% of pairs have smaller RMSD.

Assuming normal distribution, this corresponds to the deviation of 4.3ρ . They then plot R_L against L and show that the result could be approximated by the curve $R_L = 1.37 + (1.16L - 15.1)^{1/2}$. Russell (1998) did an analysis using distance RMSD, which is related to his method for detecting side-chain patterns. Random pairs of structures with different folds were chosen, and groups of two to six residues within interaction distance were picked randomly, requiring conservation of amino acid type. For $P(x)$ being the proportion of the results with $\text{RMSD} < x$ (for a fixed group size), it was found that $\log P$ was linearly related to $1/x$ in the upper part of $P(x)$.

When comparing two sequences, the P-value, E-value, and Z-value are used for measuring the significance of a score S . The P-value gives the probability of achieving a similar score by chance, the E-value is the expected number of comparisons achieving the same score by chance in a data base search, and the Z-value is the number of standard deviations S is from the mean value in an appropriate probability distribution.

Alexandrov and Fischer (1996) use the Z-value for the statistical significance. Comparing a structure A with all others in a structural database, they get a distribution of the scores. From this they can, for each structure B, use the Z-value (number of standard deviations) $Z(A, B)$. As $Z(A, B)$ might be unequal $Z(B, A)$, they define $Z_{AB} = Z(A, B) + Z(B, A)$, resulting in a symmetrical significance value. Holm and Sander (1993b, 1996) also express the significance of an alignment in terms of the Z-value, using the distribution from an all-versus-all comparison. Gibrat *et al.* (1996) in their VAST program compute a P-value for an alignment based on how many secondary structure elements are aligned as compared with the chance of aligning so many elements randomly.

6.1.2. Random structural models. The choice of the best random model against which native/native comparison scores should be compared is not simple and depends on the degree to which the inherent nonrandom features of protein structure in general should be considered significant. Some structural models are discussed by Taylor (1997). The best random models would be those generated with secondary structures. Ideally, these models should be calculated for each comparison to match the length of the native comparison and the secondary structure composition. However, these models are complex to generate and cannot be "tailor made" for each individual comparison without excessive computation. Taylor describes several random models: constrained random walk, random models from distance geometry (DRAGON), combinatorial models, combinatorial reconnection, chain reversal, and chain reflection.

The latter models involving symmetry operations on the protein (reversal and reflection) can only be used in situations where the comparison method restricts its calculation to the α -carbon atoms of the protein, as the arrangement of the other main-chain atoms is directional. Considering just α -carbons, the conformations of local structural features (such as secondary structure and their chirality of connection) in the reversed chain is virtually indistinguishable from a forward running "native" chain. This principle of reversal applies equally at the level of the sequence and has been used previously to provide a random model for sequence pattern matching (Taylor, 1986, 1977). In both sequence and structural data, the reversed model preserves the length and composition of the protein, including directionally symmetric correlations associated with secondary structure, while additionally in the reversed structural model, the bulk properties of packing density and inertial axes are also preserved. The latter are difficult to maintain in randomly generated structures.

The reflected chain is clearly not an ideal model for proteins as they contain both large and small scale chiral features which will change hand under reflection. However, Maiorov and Crippen (1994) used greatly simplified lattice models to avoid this problem and based on this analysis, they proposed a self-referential nonstatistical definition of the significance of RMSD. They take two conformers to be intrinsically similar if their RMSD is smaller than it is when one of them is mirror inverted.

6.1.3. Randomized alignment models. In general, the closer the random model is to preserving the properties of the native proteins, the more difficult it becomes to generate plausible alternatives. This problem is particularly acute for the reversed-chain random model discussed above since, for any given protein, there is only one reversal. This problem can be partially circumvented, however, at the stage of calculating the alignment. At this point, the alignment with each random model can be expanded into a population of variants by introducing "noise" into the score matrix and repeating the calculation of the alignment path from each noisy matrix. This generates a family of near-optimal subalignments and the spread of scores for this population can provide a measure of the stability or uniqueness of the answer.

An advantage of this approach is that it can be applied not only to structures belonging to the set of randomized models but also to the native structure itself and the two resulting score distributions can be tested statistically to see if they are distinct.

6.2. Comparing scoring schemes

Levitt and Gerstein (1998) made a comparison of the scoring of their iterated dynamic programming/superposition program to RMSD. The P-value of a scoring S for fixed N (number of matched residues) can be found by fitting a function $Z(S)$ to an extreme-value distribution. Unrelated structures from the SCOP database are used. The same statistics can be developed for use of RMSD. These two can then be compared by a method of Brenner *et al.* (1998). The E-values of each structure pair were sorted, and the E-value giving 1% false-positives was recognized, as was the number of true positives, calculated for the found E-value. The scores used by Levitt and Gerstein both show a much better E-value ($\log(E)$ equal -1.58 compared to -32.8) and larger number of true positives (627 compared to 202). This confirms that the scoring scheme used is much better than RMSD.

The introduction of “noise” into the alignment calculation can lead to the generation of a population of alignments (as described above). If there is sufficient noise and the population is large enough, then almost all reasonable alignments for a pair of proteins can be sampled. Plotting these solutions by their number of aligned positions against RMSD revealed a “cloud” of points which was diffuse at high RMSD but had a sharp boundary on its lower edge (Taylor, 1999). This edge represents the limit, for a given number of aligned positions, below which a smaller RMSD cannot be found. As judged by the “hard” edge to the distribution, this limit is not restricted by the method of comparing the proteins and so provides an absolute standard against which other methods can be compared. For a few protein pairs, the results of other methods (gathered by Godzik (1996)) were plotted and compared to these lines. Most of these results were found to lie above the line, indicating that the optimal solution in terms of minimum RMSD had not been attained. Only a few results lay on the line and these mostly involved fewer equivalent pairs of positions.

It should be noted that assessing methods by the use of the RMSD value is sometimes unfair since the aim of many of the methods is not to minimize the RMSD value, since this is not a good measure when all equivalent parts of the protein cannot be simultaneously superposed. An advantage of the use of large randomized alignment populations is that the alignments can be scored by any evaluation function, and so long as the evaluation function does not radically contradict what would be considered a reasonable alignment (that is: one with a low, if not minimal RMSD), there is then a good chance that some members of the population will correspond closely with the minimum of the evaluation function. This can also be used to check if a method attains the minimum of its own scoring function.

6.3. Scoring and biological significance

When a structure is compared to every other structure (or to a representative selection), then scores will result ranging from the clear relationships of homologous proteins to a large number of poor scores for obviously unrelated pairs. Between these extremes lies a twilight zone within which it is very difficult to assess the significance of the score. This problem is exacerbated because many proteins contain similar substructures, such as secondary and supersecondary structures, and the problem is to decide when a similarity is just a consequence of being proteinlike and when it indicates a more specific relationship between the two proteins.

Because of its common currency, most considerations of this problem have focused on the significance of the RMSD measure based on comparison of proteins or protein fragments of equal length (see above). Others, such as the DALI method, have adopted a similar approach based on the scores achieved over matches of protein fragments (Holm and Sander, 1993a, 1997). Both these approaches require that the selected fragments be unrelated to the proteins being assessed; however, this raises the problem of what criterion can be used to make this distinction and, in principle, it should not be a weaker method than that used for the current comparison. It is not acceptable, either, to consider completely unrelated proteins since, to take an extreme example, if the two proteins being compared contained only α -helices and the clearly unrelated control set contained only β -structure, then the two α proteins would appear more related than they should do.

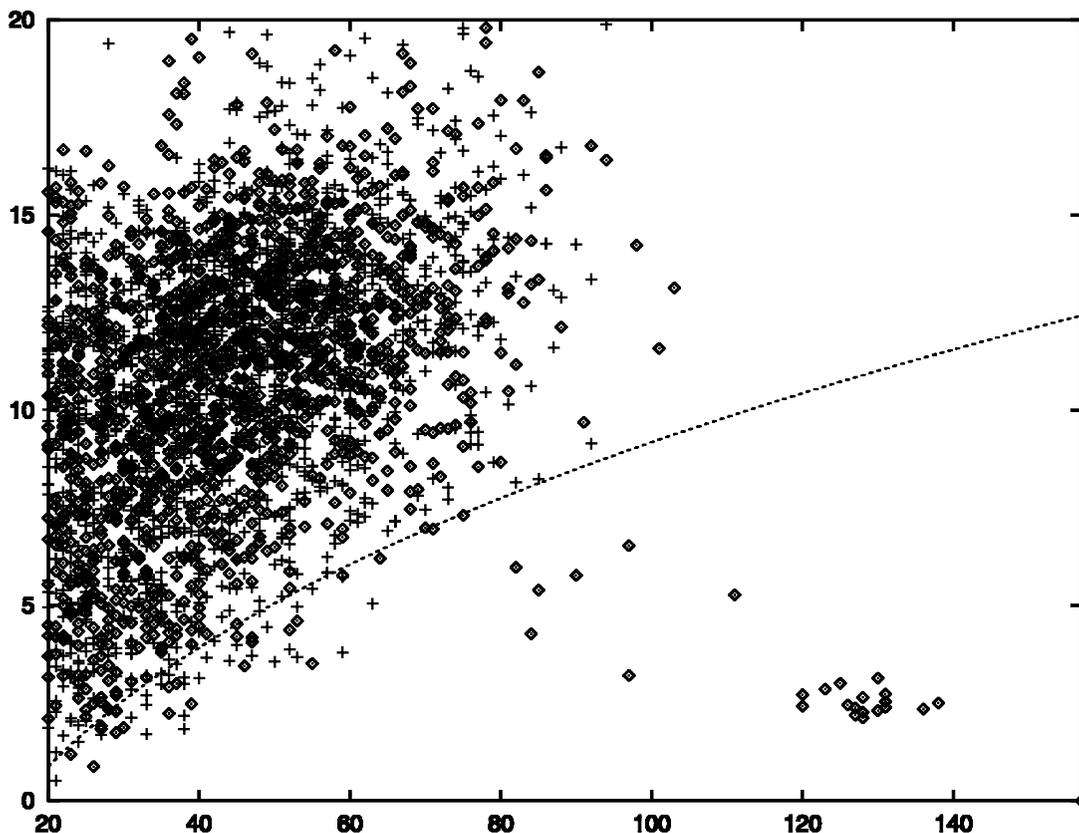


FIG. 10. Match of a myoglobin probe against a nonredundant PDB. The RMSd of each comparison (Y) is plotted against the number of residues aligned (X) for comparisons using the native probe structure (\diamond) and for comparisons using the reversed probe as a control (+). The dashed line follows a function of Meirov and Crippen (1994) that has been scaled to include 99% of the reversed controls. The matches lying below this line include the probe matching itself (150 residues with 0 RMSd), other globins (120–140, respectively) and the phycocyanin family (80–115, respectively) which have a similar fold.

An alternate approach to this problem is to use the reversed structure (as described above). When this is matched against the structure databank, a similar range of scores should result—since the reversed structure has exactly the same length, overall shape, and secondary structure content as the native probe. What will be lost is any specific overall similarity to proteins that are homologous to the native probe. In addition, if the probe structure is a particularly simple fold (such as four α -helices), then the reversed structure will also embody this property so a specific match will need to capture more than a few matched helices to gain significantly over the background of scores derived from the reversed structure (see Figure 10).

7. CONCLUSION

From the wide variety of methods surveyed above, it should be clear that there is no single best method for protein structure comparison. The choice of method depends on the level of data representation, which in turn depends on the nature of the original question. Possibilities range from considering just fragments to an overall match, atomic detail to rough fold, and even considering or neglecting chain direction. For all these aspects there is a method—and often there is more than one.

From the algorithmic viewpoint, the problem of structure comparison remains of interest perhaps because it has no exact solution (neglecting exhaustive enumeration). This has allowed a variety of standard (and some less standard) optimization methods to be applied. Often these have a stochastic component which can bring useful information on the uniqueness of the resulting solution.

Moving to multiple structure comparison presents yet more complexities which have again been overcome using a variety of heuristic methods—often following the equivalent multiple sequence alignment methods. Some innovative developments are being made in this direction (also following sequence-based methods) that allow a simultaneous approach to the multiple structure data rather than a conventional progressive approach.

ACKNOWLEDGMENTS

I.J. was supported by grants from the Norwegian Research Council. The authors wish to thank Dr. David R. Gilbert for helpful discussions.

REFERENCES

- Alesker, V., Nussinov, R., and Wolfson, H.J. 1996. Detection of non-topological motifs in protein structures. *Protein Eng.* 9, 1103–1119.
- Alexandrov, N., and Fischer, D. 1996. Analysis of topological and nontopological structural similarities in the PDB: new examples with old structures. *Proteins: Struct., Func. Gen.* 25, 354–365.
- Alexandrov, N., and Go, N. 1994. Biological meaning, statistical significance, and classification of local spatial similarities in nonhomologous proteins. *Prot. Sci.* 3, 866–875.
- Alexandrov, N., Takahashi, K.J., and Go, N. 1992. Common spatial arrangements of backbone fragments in homologous and non-homologous proteins. *J. Mol. Biol.* 225, 5–9.
- Altschul, S., Gish, W., Miller, W., Myers, E., and Lipman, D.J. 1990. Basic local alignment search tool. *J. Mol. Biol.* 215, 403–410.
- Artymiuk, P., Porrette, A., Grindley, H., Rice, D., and Willett, P. 1994. A graph-theoretic approach to the identification of three-dimensional patterns of amino acid side-chains in protein structures. *J. Mol. Biol.* 243, 327–344.
- Bagley, S., and Altman, R. 1995. Characterizing the microenvironment surrounding protein sites. *Prot. Sci.* 4, 622–635.
- Baxevanis, A.D. 2000. The Molecular Biology Database Collection: An online compilation of relevant database resources. *Nucl. Acids Res.* 28, 1–7.
- Bernstein, F.C., Koetzle, T.F., Williams, G.J., Jr., E.E.M., Brice, M.D., Rodgers, J.R., Kennard, O., Shimanouchi, T., and Tasumi, M. 1977. The protein data bank: A computer-based archival file for macromolecular structures. *J. Mol. Biol.* 112, 535.
- Brazma, A., Jonassen, I., Eidhammer, I., and Gilbert, G. 1998. Approaches to the automatic discovery of patterns in biosequences. *J. Comp. Biol.* 5 (2), 279–305.
- Brenner, S., Chothia, C., and Hubbard, T. 1998. Assessing sequence comparison methods with reliable structurally identified distant evolutionary relationships. *Proc. Natl. Acad. Sci.* 95 (11), 6073–6078.
- Brown, C., and Kerbosch, J. 1973. Algorithm 457—finding all cliques of an undirected graph. *Comm ACM* 16, 575–577.
- Brown, N., Orengo, C., and Taylor, W. 1996. A protein structure comparison methodology. *Comp. Chem.* 20, 359–380.
- Chew, L., Huttenlocher, D., Kedem, K., and Kleinberg, J. 1999. Fast detection of common geometric substructure in proteins. *Proceedings of the Third International Conference on Computational Molecular Biology, RECOMB'99.*
- Cohen, G.H. 1997. ALIGN: A program to superimpose protein coordinates, accounting for insertions and deletions. *J. Appl. Cryst.* 30, 1160–1161.
- Conklin, D. 1995. Machine discovery of protein motifs. *Machine Learning* 21, 125–150.
- Cook, D., Holder, L., and Djoko, S. 1996. Scalable discovery of informative structural concepts using domain knowledge. *IEEE EXPERT* 11 (5), 59–68.
- Dayhoff, M.O. 1978. *Atlas of Protein Sequence and Structure*, vol. 5. National Biomedical Research Foundation.
- de Rinaldis, M., Ausiello, G., Cesareni, G., and Helmer-Citterich, M. 1998. Three-dimensional profiles: A new tool to identify protein surface similarities. *J. Mol. Biol.* 284, 1211–1221.
- Diamond, R. 1992. On the multiple simultaneous superposition of molecular structures by rigid body transformations. *Prot. Sci.* 1, 1279–1287.
- Diederichs, K. 1995. Structural superposition of proteins with unknown alignment and detection of topologically similarity using a six-dimensional search algorithm. *Proteins: Struct., Func., Gen.* 23, 187–195.
- Ding, D., Qian, J., and Feng, Z. 1994. A differential geometric treatment of protein structure comparison. *Bull. Math. Biol.* 56, 923–943.
- Escalier, V., Pothier, J., Soldano, H., and Viari, A. 1998. Pairwise and multiple identification of three-dimensional common substructures in proteins. *J. Comp. Biol.* 5, 41–56.

- Falicov, A., and Cohen, F. 1996. A surface of minimum area metric for the structural comparison of proteins. *J. Mol. Biol.* 258, 871–892.
- Fetrow, J.S., and Skolnick, J. 1998. Method for prediction of protein function from sequence using the sequence-to-structure-to-function paradigm with application to glutaredoxins/thioredoxins and t_1 ribonucleases. *J. Mol. Biol.* 281, 949–968.
- Fischer, D., Tsai, C., Nussinov, R., and Wolfson, H. 1995. A 3D sequence-independent representation of the protein data bank. *Protein Eng.* 8, 981–997.
- Fischer, D., Wolfson, H., Lin, S., and Nussinov, R. 1994. Three-dimensional, sequence order-independent structural comparison of a serine protease against the crystallographic database reveals active site similarities: Potential implications to evolution and to protein folding. *Prot. Sci.* 3, 769–778.
- Flores, T., Moss, D., and Thornton, J. 1994. An algorithm for automatically generating protein topology cartoons. *Protein Eng.* 7 (1), 31–37.
- Gerber, P., and Müller, K. 1987. *Act. Cryst.* A43, 426.
- Gerstein, M., and Altman, R. 1995. Using a measure of structural variation to define a core for the globins. *CABIOS* 11 (6), 633–644.
- Gerstein, M., and Levitt, M. 1998. Comprehensive assessment of automatic structural alignment against a manual standard, the scop classification of proteins. *Prot. Sci.* 7, 445–456.
- Gibrat, J., Madej, T., and Bryant, S. 1996. Surprising similarities in structure comparison. *Curr. Op. Struct. Biol.* 6, 377–385.
- Gilbert, D., Westhead, D., Nagano, N., and Thornton, J. 1999. Motif-based searching in tops protein topology databases. *Bioinformatics* 15, 317–326.
- Gilbert, D., Westhead, D., Thornton, J., and Viksna, J. 1999. Tops cartoons: formalisation, searching and comparison. Poster at RECOMB 99.
- Godzik, A. 1996. The structural alignment between two proteins: is there a unique answer? *Prot. Sci.* 5, 1325–1338.
- Godzik, A., and Skolnick, J. 1994. Flexible algorithm for direct multiple alignment of protein structures and sequences. *CABIOS* 10 (6), 587–596.
- Gribskov, M., McLachlan, A., and Eisenberg, A. 1987. Profile analysis: Detection of distantly related proteins. *Proc. Natl. Acad. Sci. USA* 84, 4355–4358.
- Grindley, H., Artymiuk, P., Rice, D., and Willett, P. 1993. Identification of tertiary structure resemblance in proteins using a maximal common subgraph isomorphism algorithm. *J. Mol. Biol.* 229, 707–721.
- Henikoff, S., and Henikoff, J.G. 1992. Amino acid substitution matrices from protein blocks. *Proc. Natl. Acad. Sci. USA* 89, 100915–100919.
- Holm, L., and Sander, C. 1993a. Protein-structure comparison by alignment of distance matrices. *J. Mol. Biol.* 233, 123–138.
- Holm, L., and Sander, C. 1993b. Protein structure comparison by alignment of distance matrices. *J. Mol. Biol.* 233, 123–138.
- Holm, L., and Sander, C. 1995. 3-d lookup: Fast protein structure database searches at 90% reliability. *Proceedings of the Third International Conference on Intelligent Systems for Molecular Biology*, 179–187.
- Holm, L., and Sander, C. 1996. Mapping the protein universe. *Science* 273, 595–602.
- Holm, L., and Sander, C. 1997. Dali/FSSP classification of three-dimensional protein folds. *Nucleic Acids Research* 25, 231–234.
- Jonassen, I. 1997. Efficient discovery of conserved patterns using a pattern graph. *CABIOS* 13, 509–522.
- Jonassen, I., Collins, J.F., and Higgins, D.G. 1995. Finding flexible patterns in unaligned protein sequences. *Prot. Sci.* 4 (8), 1587–1595.
- Jonassen, I., Eidhammer, I., Taylor, W., and Grindhaug, S. 2000. Searching the protein structure databank with weak sequence patterns and structural constraints. *J. Mol. Biol.* 304(4), 597–617.
- Jonassen, I., Eidhammer, I., and Taylor, W.R. 1999. Discovery of local packing motifs in protein structures. *Proteins: Struct., Func., Gen.* 34, 206–219.
- Kabsch, W. 1978. A solution for the best rotation to relate two sets of vectors. *Acta Cryst.* A32, 922–923.
- Kastenmüller, G., Kriegl, H.P., and Seidl, T. 1998. Similarity search in 3d protein databases. *German Conference on Bioinformatics, GCB* 98.
- Kasuya, A., and Thornton, J. 1999. Three-dimensional structure analysis of prosite patterns. *J. Mol. Biol.* 286, 1673–1691.
- Kearsley, S. 1990. An algorithm for the simultaneous superposition of a structural series. *J. Comp. Chem.* 11 (10), 1187–1192.
- Kleywegt, G., and Jones, T. 1997. Detecting folding motifs and similarities in protein structures. In: *Methods in Enzymology*, 525–545. Academic Press.
- Koch, I., Lengauer, T., and Wanke, E. 1996. An algorithm for finding maximal common subtopologies in a set of protein structures. *J. Comp. Biol.* 3 (2), 287–306.

- Lathrop, R. 1994. The protein threading problem with sequence amino acid interaction preferences is NP-complete. *Protein Eng.* 7, 1059–1068.
- Lathrop, R., and Smith, T. 1996. Global optimum protein threading with gapped alignment and empirical pair score functions. *J. Mol. Biol.* 255 (4), 641–665.
- Leibowitz, N., Fligerman, Z.Y., Nussinov, R., and Wolfson, H.J. 1999. Multiple structural alignment and core detection by geometric hashing. *Proc. of the Seventh International Conference on Intelligent Systems for Molecular Biology, ISMB'99*, 169–177.
- Levitt, M., and Gerstein, M. 1998. A unified statistical framework for sequence comparison and structure comparison. *Proc. Natl. Acad. Sci USA* 95, 5913–5920.
- Lipman, D., and Pearson, W. 1985. Rapid and sensitive protein similarity searches. *Science* 227, 1435–1441.
- Madej, T., Gibrat, J., and Bryant, S. 1995. Threading a database of protein core. *Proteins: Struct., Func., Gen.* 23, 356–369.
- Maiorov, V., and Crippen, G. 1994. Significance of root-mean-square deviation in comparing three-dimensional structures of globular proteins. *J. Mol. Biol.* 235, 625–634.
- Matsuda, H., Taniguchi, F., and Hashimoto, A. 1997. An approach to detection of protein structural motifs using an encoding scheme of backbone conformation. *Pacific Symposium on Biocomputing '97*.
- May, A. 1996. Pairwise iterative superposition of distantly related proteins and assessment of the significance of 3-D structural similarity. *Protein Eng.* 9, 1093–1101.
- May, A.C.W., and Johnson, M.S. 1995. Improved genetic algorithm-based protein structure comparisons: Pairwise and multiple superpositions. *Protein Eng.* 8, 873–882.
- Mizuguchi, K., and Go, N. 1995. Comparison of spatial arrangements of secondary structural elements in proteins. *Protein Eng.* 8, 353–362.
- Murzin, A., Brenner, S., 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.
- Needleman, S., and Wunsch, C. 1970. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *J. Mol. Biol.* 48, 443–454.
- Nussinov, R., and Wolfson, H. 1991. Efficient detection of three-dimensional structural motifs in biological macromolecules by computer vision techniques. *Proc. Natl. Acad. Sci USA* 88, 10495–10499.
- Orengo, C., Brown, N., and Taylor, W. 1992. *Protein: Struct., Funct., Gen.* 14, 139.
- Orengo, C.A., and Taylor, W.R. 1990. A rapid method for protein structure alignment. *J. Theor. Biol.* 147, 517–551.
- Pennec, X., and Ayache, N. 1998. A geometric algorithm to find small but highly similar 3D substructures in proteins. *Bioinformatics* 14, 516–522.
- Petitjean, M. 1998. Interactive maximal common 3d substructure searching with the combined sdm/rms algorithm. *Comp. Chem.* 22 (6), 463–465.
- Rao, S.T., and Rossmann, M.G. 1973. Comparison of super-secondary structures in proteins. *J. Mol. Biol.* 76, 241–256.
- Rossmann, M.G., and Argos, P. 1975. A comparison of the heme binding pocket in globins and cytochrome b5*. *J. Biol. Chem.* 250, 7523–7532.
- Rossmann, M.G., and Argos, P. 1976. Exploring structural homology of proteins. *J. Mol. Biol.* 105, 75–96.
- Rufino, S., and Blundell, T. 1994. Structure-based identification and clustering of protein families and superfamilies. *J. Computer-Aided Molecular Design* 8, 5–27.
- Russell, R. 1998. Detection of protein three-dimensional side-chain patterns: New examples of convergent evolution. *J. Mol. Biol.* 279, 1211–1227.
- Russell, R., and Barton, G. 1992. Multiple protein sequence alignment from tertiary structure comparison: Assignment of global and residue confidence levels. *Proteins: Struct., Func., Gen.* 14, 309–323.
- Sali, A., and Blundell, T. 1990. Definition of general topological equivalence in protein structures: A procedure involving comparison of properties and relationships through simulated annealing and dynamic programming. *J. Mol. Biol.* 212, 403–428.
- Sankoff, D., and Kruskal, J.B. 1983. *Time Warps: String Edits, and Macromolecules: The Theory and Practice of Sequence Comparison*. Addison-Wesley.
- Satow, Y., Cohen, G.H., Padlan, E.A., and Davies, D.R. 1986. *J. Mol. Biol.* 190, 593–604.
- Schulz, G.E., and Schirmer, R.H. 1979. *Principles of Protein Structure*. Springer-Verlag, Germany.
- Shapiro, A., Botha, J.D., Pastore, A., and Lesk, A. 1992. A method for multiple superposition of structures. *Acta Cryst.* A48, 11–14.
- Shinohara, T., and Arikawa, S. 1995. Pattern inference. In: *Algorithmic Learning for Knowledge-Based Systems, GOSLER Final Report*, 259–291. Springer-Verlag.
- Su, S., Cook, D., and Holder, L. 1999. Applications of knowledge discovery to molecular biology: Identifying structural regularities in proteins. In: *Pacific Symposium on Biocomputing '99*.
- Sutcliffe, M., Haneef, I., Carney, D., and Blundell, T.L., 1987. Knowledge-based modeling of homologous proteins 1: 3-dimensional frameworks derived from the simultaneous superposition of multiple structures. *Prot. Eng.* 1, 377.

- Taylor, W. 1988. A flexible method to align large numbers of biological sequences. *J. Mol. Evol.* 28, 161–169.
- Taylor, W. 1997. Random structural models for double dynamic programming score evaluation. *J. Mol. Evol.* 44, 174–180.
- Taylor, W. 1999. Protein structure comparison using iterated double dynamic programming. *Prot. Sci.* 8, 654–665.
- Taylor, W., Flores, T., and Orengo, C. 1994. Multiple protein structure alignment. *Prot. Sci.* 3, 1858–1870.
- Taylor, W., and Orengo, C. 1989. Protein structure alignment. *J. Mol. Biol.* 208, 1–22.
- Taylor, W.R. 1986. Identification of protein sequence homology by consensus template alignment. *J. Molec. Biol.* 188, 233–258.
- Verbitsky, G., Nussinov, R., and Wolfson, H. 1999. Structural comparison allowing hinge bending, swiveling motions. *Proteins: Struct., Func., Gen.* 34 (2), 232–254.
- Vingron, M., and Argos, P. 1991. Motif recognition and alignment for many sequences by comparison of dot-matrices. *J. Mol. Biol.* 218, 33–43.
- Vriend, G., and Sander, C. 1991. Detection of common three-dimensional substructures in proteins. *Proteins: Struct., Func., Gen.* 11, 52–58.
- Wako, H., and Yamato, T. 1998. Novel method to detect a motif of local structures in different protein conformations. *Protein Eng.* 11 (11), 981–990.
- Wallace, A., Borkakoti, N., and Thornton, J. 1997. TESS: A geometric hashing algorithm for deriving 3D coordinate templates for searching structural databases: Applications to enzyme active sites. *Prot. Sci.* 6, 2308–2323.
- Wallace, A., Laskowski, R., and Thornton, J. 1996. Derivation of 3D coordinate templates for searching structural databases: Applications to ser-his-asp catalytic triads in the serine proteinases and lipases. *Prot. Sci.* 5, 1001–1013.
- Wei, L., Altman, R.B., and Chang, J.T. 1997. Using the radial distribution of physical features to compare amino acid environments and align amino acid sequences. *Pacific Symposium on Biocomputing'97*.
- White J., Muchnik, I., and Smith, T. 1994. Modeling protein cores with markov random fields. *Math. Biosc.* 124, 149–179.
- Zu-Kang, F., and Sippl, M.J. 1996. Optimum superimposition of protein structures: Ambiguities and implications. *Folding and Design* 1, 123–132.

Address correspondence to:

Ingvar Eidhammer
Department of Informatics
University of Bergen
Høyteknologisentret
N-5020 Bergen, Norway

E-mail: Ingvar.Edihammer@ii.uib.no