

# [42] Phylogenetic Relationships from Three-Dimensional Protein Structures

By MARK S. JOHNSON, ANDREJ ŠALI, and TOM L. BLUNDELL

## Introduction

As evolution progresses, accumulated changes in DNA and RNA lead to differences in amino acid sequences and corresponding alterations in the tertiary structure of proteins. Although these changes take place most often on the protein surface exposed to the solvent, mutations can be accepted within the inaccessible hydrophobic interior formed from packed secondary structure elements. Thus, relative translations and rotations of  $\alpha$  helices and  $\beta$  strands do occur while their general spatial relationships remain highly conserved.<sup>1,2</sup> As a result, when comparisons are made among more distantly related structures, fewer topologically equivalent positions are found and lead to greater root mean square (RMS) deviations.<sup>3,4</sup>

Consequently, one should be able to chart the evolution of proteins from a comparison of their structures. Indeed, as protein structures are generally more conserved in evolution than are amino acid sequences,<sup>5</sup> they can be used to infer relationships among proteins where an alignment of their sequences is not statistically significant.<sup>6</sup> The first phylogenetic tree derived from structural information was based on the number of topologically equivalent positions in several dinucleotide- and mononucleotide-binding proteins<sup>7</sup> (for a review, see Matthews and Rossmann<sup>8</sup>). Following this work, Johnson *et al.*<sup>9</sup> compared six families of homologous structures and sequences (immunoglobulins,  $c$ -type cytochromes, globins, serine proteinases, eyeless  $\gamma$ -crystallins, and nucleotide-binding domains) and showed that trees based on sequence and structure are generally congruent.

Pattern matching and comparison methods for proteins can be employed to derive both the equivalences between residues and the overall

similarity scores which can then be used to construct trees. It is important to recognize what information about the proteins is included in the comparison. This information generally focuses on only a few aspects, and there are two major groups of comparison methods: those using sequences of the amino acid residues and those employing three-dimensional structures of proteins.

The most familiar sequence comparison methods are the dynamic programming procedures based on the algorithm of Needleman and Wunsch,<sup>9</sup> which has recently been exploited to obtain multiple sequence alignments.<sup>10-13</sup> These methods usually consider the mutation rates of amino acid residues<sup>14,15</sup> to derive optimal comparison scores and corresponding alignments. The utility of many amino acid properties for the alignment of amino acid sequences was systematically explored by Argos<sup>16</sup> using his own optimization algorithm. Physical properties of amino acid residues are also considered in the pattern matching technique of Taylor,<sup>17</sup> which can align several protein sequences simultaneously. The hierarchical aspects of sequence and secondary structure or organization were implemented in the program ARIADNE<sup>18</sup> for the matching of a given sequence pattern to a protein sequence and in the approach developed by Rawlings *et al.*<sup>19</sup> for reasoning about protein topology.

Alternatively, methods which compare tertiary structures of proteins are dominated by the rigid-body least-squares superposition of the  $\alpha$ -carbon ( $C_\alpha$ ) positions (see Matthews and Rossmann<sup>8</sup> for a review). However, Rao and Rossmann<sup>20</sup> also included the main chain direction in their pairwise comparison procedure, which enabled the alignment of more divergent protein structures. The rigid-body approach was recently extended by Sutcliffe *et al.*<sup>21</sup> for the simultaneous comparison of several

- <sup>1</sup> A. M. Lesk and C. Chothia, *J. Mol. Biol.* **136**, 225 (1980).
- <sup>2</sup> A. M. Lesk and C. Chothia, *J. Mol. Biol.* **160**, 325 (1982).
- <sup>3</sup> C. Chothia and A. M. Lesk, *EMBO J.* **5**, 823 (1986).
- <sup>4</sup> T. J. P. Hubbard and T. L. Blundell, *Protein Eng.* **1**, 159 (1987).
- <sup>5</sup> M. Bajaj and T. L. Blundell, *Annu. Rev. Biophys. Bioeng.* **13**, 453 (1984).
- <sup>6</sup> M. S. Johnson, M. J. Sutcliffe, and T. L. Blundell, *J. Mol. Evol.*, in press (1990).
- <sup>7</sup> W. Evenoff and M. G. Rossmann, *CRC Crit. Rev. Biochem.* **3**, 111 (1975).
- <sup>8</sup> B. W. Matthews and M. G. Rossmann, this series, Vol. 115, p. 397.

- <sup>9</sup> S. B. Needleman and C. D. Wunsch, *J. Mol. Biol.* **48**, 443 (1970).
- <sup>10</sup> D.-F. Feng and R. F. Doolittle, *J. Mol. Evol.* **25**, 351 (1987).
- <sup>11</sup> D.-F. Feng and R. F. Doolittle, this volume, [23].
- <sup>12</sup> G. Barton and M. J. E. Sternberg, *J. Mol. Biol.* **198**, 327 (1987).
- <sup>13</sup> G. J. Barton, this volume, [25].
- <sup>14</sup> R. M. Schwartz and M. O. Dayhoff, in "Atlas of Protein Sequence and Structure" (M. O. Dayhoff, ed.), Vol. 5, Suppl. 3, p. 353. National Biomedical Research Foundation, Washington, D.C., 1978.
- <sup>15</sup> D.-F. Feng, M. S. Johnson, and R. F. Doolittle, *J. Mol. Evol.* **21**, 112 (1985).
- <sup>16</sup> P. Argos, *J. Mol. Biol.* **193**, 385 (1987).
- <sup>17</sup> W. R. Taylor, *J. Mol. Biol.* **188**, 233 (1986).
- <sup>18</sup> R. H. Lathrop, T. A. Webster, and T. F. Smith, *Commun. ACM* **30**, 909 (1988).
- <sup>19</sup> C. J. Rawlings, W. R. Taylor, J. Nyakairu, J. Fox, and M. J. E. Sternberg, *J. Mol. Graph.* **3**, 151 (1985).
- <sup>20</sup> S. T. Rao and M. G. Rossmann, *J. Mol. Biol.* **76**, 241 (1973).
- <sup>21</sup> M. J. Sutcliffe, I. Hanef, D. Carney, and T. L. Blundell, *Protein Engineer.* **1**, 377 (1987).

protein structures. Sippl<sup>22</sup> has used information from intramolecular  $C_\alpha$  distance matrices to compare protein structures, while Murthy<sup>23</sup> applied dynamic programming to compare the secondary structure organization of proteins, taking into account the absolute angles and distances between the idealized secondary structure segments in approximately superposed proteins. Similarly, Richards and Kundrot<sup>24</sup> took into account the internal relationships between the elements of secondary structure to search for a given pattern in a protein structure database. Sheridan *et al.*<sup>25</sup> considered the residue secondary structure type to define the weights for each residue pair, which were later used in a dynamic programming procedure to find the alignment of two proteins. Recently, Barton and Sternberg<sup>26</sup> applied a dynamic programming procedure and an intermolecular distance matrix for roughly superposed loops to obtain the alignment of hypervariable regions.

In this chapter, we first describe a method for the multiple rigid-body superposition of structures and show the usefulness of the pairwise rigid-body superposition in determining relationships among homologous protein structures.<sup>6</sup> We continue with a description of a more flexible alignment procedure that compares a number of structural properties and relationships through simulated annealing and dynamic programming algorithms.<sup>27</sup> This latter technique escapes the limitations imposed by rigid-body alignments: it can taken into account deformations and translocations of secondary structure elements such as those illustrated in Fig. 1 between the two sets of aspartic proteinase domains and in Fig. 2 for the cytochromes *c*. From both approaches, the rigid-body and multifit methods for the comparison of structures, phylogenetic trees are derived for the proteins listed in Table I. In general, these trees are isomorphous to those that may be obtained from the alignment of the corresponding amino acid sequences.<sup>28</sup>

<sup>22</sup> M. J. Sippl, *J. Mol. Biol.* **156**, 359 (1982).

<sup>23</sup> M. R. N. Murthy, *FEBS Lett.* **168**, 97 (1984).

<sup>24</sup> F. M. Richards and C. E. Kundrot, *Proteins* **3**, 71 (1988).

<sup>25</sup> R. P. Sheridan, J. S. Dixon, and R. Venkatarang, *van, Int. J. Peptide Protein Res.* **25**, 132 (1986).

<sup>26</sup> G. Barton and M. J. E. Sternberg, *J. Mol. Graph.* **6**, 190 (1988).

<sup>27</sup> A. Sali and T. L. Blundell, *J. Mol. Biol.*, in press (1990).

<sup>28</sup> The programs described here were written in the C and FORTRAN programming languages on a MICROVAX II computer (Digital Equipment Corporation) running the VMS operating system, version 4.5 or higher. These programs will be available from the authors in the near future.

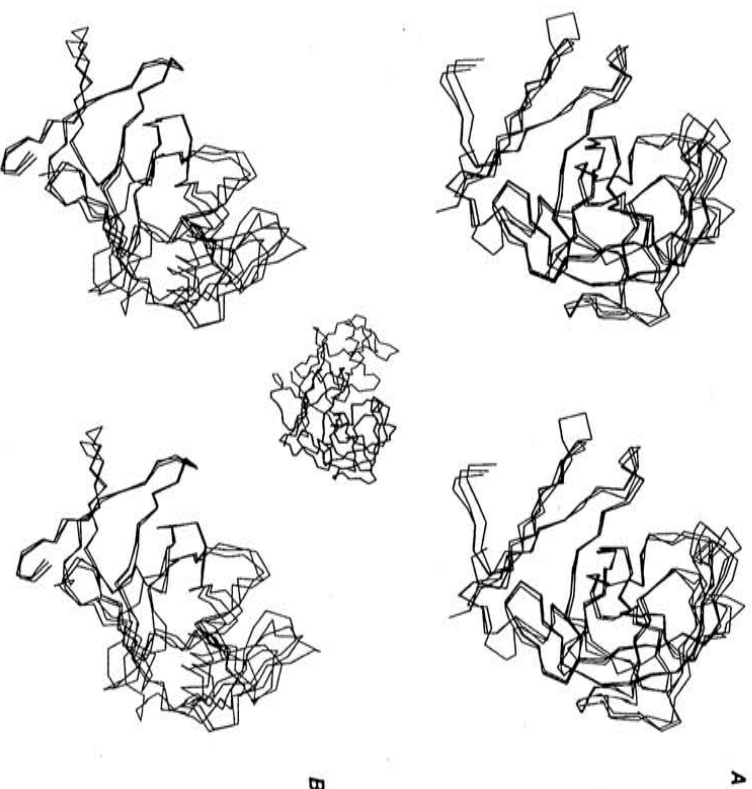


FIG. 1. Stereo views for the rigid-body superposed aspartic proteinase domains: (A) amino-terminal halves and (B) carboxy-terminal halves of endothiasepsin, penicillopepsin, and rhizopuspepsin (Table I). The amino- and carboxy-terminal domains are themselves a result of a duplication event and pack with  $C_2$  symmetry. To show this, the  $C_\alpha$  backbone for the entire two-domain structure of endothiasepsin (amino-terminal domain: thick line; carboxy-terminal: thin line) is inserted in the center; the  $C_2$  rotation axis is located at the domain interface, approximately vertical and in the plane of the paper. Endothiasepsin, penicillopepsin, and rhizopuspepsin were split into domains at residues 174-175, 174-175, and 178-179 (crystallographic numbering), respectively.

## Protein Structure Comparison by Rigid-Body Superposition

### Topological Equivalence

The optimal superposition of two sets of coordinates is a common problem where the goal is to obtain the "best" fit of an object *A* to an object *B* over some set of coordinates said to be topologically equivalent for

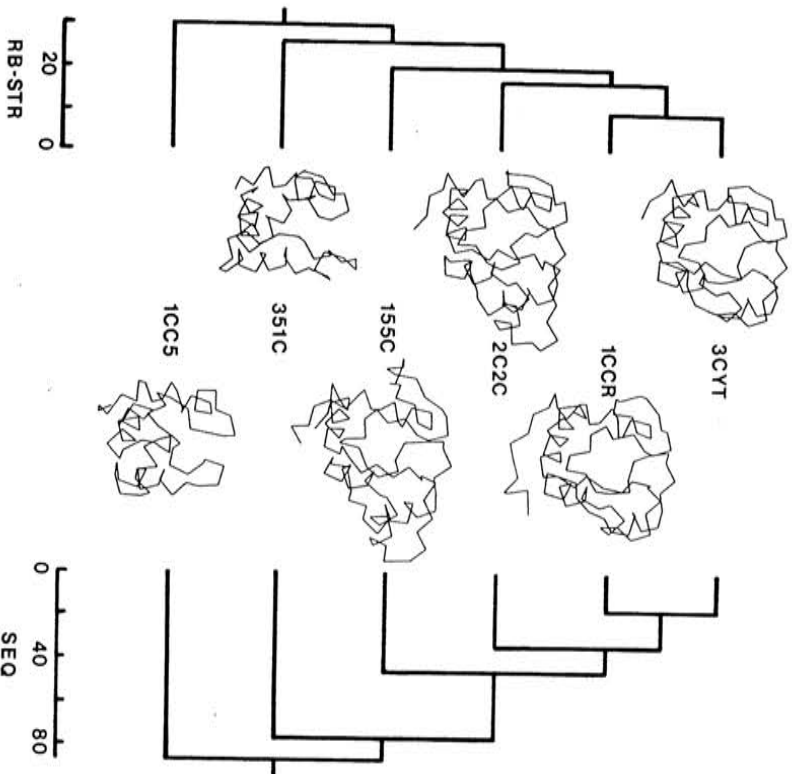


Fig. 2. Cladograms derived for the  $c$ -type cytochromes from both pairwise rigid-body comparisons (RB-STR) and multiple alignment of sequences (SEO).  $C_{\alpha}$  backbones are depicted for each of the structures after they have been fit with the multiple rigid-body superposition procedure; the structures have been translated only within the plane of the page. See Table I for identification of structures and sequences. The  $C_{\alpha}$  coordinates belonging to the "unknown" (UNK) residues, as listed in the Brookhaven file for 155C, were excluded from the comparisons.

the two objects.<sup>7,8,20,29,30</sup> For proteins, these topologically equivalent positions can be defined as those  $C_{\alpha}$  atoms in the two superposed structures that lie within a specified distance of each other, provided that these equivalences are colinear and hence obey the "no-knot" constraint.

<sup>29</sup> D. R. Ferro and J. Hermans, *Acta Crystallogr.* A33, 345 (1977).

<sup>30</sup> A. D. McLachlan, *Acta Crystallogr.* A34, 871, 882.

TABLE I  
PROTEIN STRUCTURES ALIGNED IN THIS STUDY

Brookhaven code <sup>a</sup>	Description	Resolution (Å)
Aspartic proteinases		
4APE	Endothiapepsin	2.1
2APR	Penicillopepsin	1.8
2APR	Rhizopuspepsin	1.8
Cytochromes $c$		
3CYT	Albacore tuna heart ferricytochrome $c$ (oxidized)	1.8
1CCR	Rice embryo ferricytochrome $c$	1.5
2C2C	<i>Rhodospirillum rubrum</i> ferricytochrome $c_2$	2.0
155C	<i>Paracoccus denitrificans</i> cytochrome $c$ -550	2.5
351C	<i>Pseudomonas aeruginosa</i> ferricytochrome $c$ -551	1.6
1CC5	<i>Azotobacter vinelandii</i> ferricytochrome $c_3$	2.5
Globins		
2HHB	Human deoxyhemoglobin $\alpha$ and $\beta$ chains	1.7
2HCO	Human carbonmonoxyhemoglobin $\alpha$ and $\beta$ chains	2.7
1HHO	Human oxyhemoglobin $\alpha$ and $\beta$ chains	2.1
1HBS	Human sickle cell hemoglobin	3.0
1FDH	Human deoxyhemoglobin $\alpha$ and $\gamma$ fetal chains	2.5
2DHB	Horse deoxyhemoglobin $\alpha$ and $\beta$ chains	2.8
1HDS	Deer sickle cell hemoglobin	2.0
2LHB	Sea lamprey hemoglobin V (cyano/met)	2.0
2MBN	Sperm whale metmyoglobin	2.0
3MBN	Sperm whale deoxymyoglobin	2.0
1ECD	<i>Chironomus thummi thummi</i> erythrocytochrome	1.4
1LH1	<i>Lupinus luteus</i> leghemoglobin	2.0

<sup>a</sup> F. C. Bernstein, T. F. Koetzle, G. J. B. Williams, E. F. Meyer, M. D. Brice, J. R. Rodgers, O. Kennard, T. Shimanouchi, and M. Tasumi, *J. Mol. Biol.* 112, 535 (1977).

Typically for proteins, one does not have *a priori* knowledge as to the extent of the topological equivalence, but only an idea of positions that are likely to be equivalent and which may be useful as a starting point for a comparison. This initial set of equivalences can be determined from an examination of atomic coordinates on a graphics device, from residues highly conserved in a sequence alignment, or from positions that are known to be crucial to the structural integrity, catalysis, or ligand binding of a protein. For the initial fit of two structures, three or more  $C_{\alpha}$  positions must be specified to orient the structures. After the initial superposition of the structures, topological equivalences are redetermined, and the structures are again superposed based on this newly determined set of equivalent

lences. The procedure iterates until stability in both the topological equivalences and the RMS distance over the equivalent positions is obtained.

### The Residual

The superposition of the structures hinges on the minimization of a function,  $\delta$  (the residual), of the general form

$$\delta = \sum_{i=1}^a w_i (\bar{X}_i - \mathcal{R} \bar{Y}_i)^2 \quad (1)$$

where  $w_i$  is the weight for the  $i$ th pair of equivalent positions  $\bar{X}_i$  and  $\bar{Y}_i$  of the two structures;  $\mathcal{R}$  is the 3 by 3 rotation matrix that superposes structure  $B$  onto structure  $A$ . The use of this function presupposes that the coordinates for the two structures have been translated so that the centers of gravity for the two sets of equivalent positions are located at the origin of the coordinate system. For the rigid-body comparison of structures in this chapter we have used a modification of a program MNVFT<sup>21</sup> that was designed for the multiple alignment of structures by superposition.

### Iterative Weighted Superposition

In the iterative weighted superposition procedure, one of the structures is chosen initially as the first approximation to the average of all structures (the framework), and each of the other structures is then fitted to it pairwise with the rapid procedure of McLachlan.<sup>20</sup> A new framework can then be calculated for  $b$  molecules from

$$\bar{F}_i^k = \frac{\sum_{j=1}^b w_{ij}^{k-1} \bar{Z}_{ij}}{\sum_{j=1}^b w_{ij}^{k-1}} \quad (2)$$

where  $\bar{F}_i^k$  are the coordinates of point  $i$  on the framework at iteration  $k$  and  $\bar{Z}_{ij} = \mathcal{R}_{ij} \bar{Y}_{ij}$  are the coordinates of atom  $i$  from molecule  $j$  fitted to the previous framework  $F^{k-1}$ . The weight factors

$$w_{ij}^{k-1} = \frac{1}{\alpha_j^2 + \sigma_i d_{ij}} \quad (3)$$

include  $\alpha_j$  as an estimate of the error in the coordinates of molecule  $j$ ,  $\sigma_i$ , the standard deviation of the distance from the framework for the  $i$ th set of topologically equivalent positions  $a_i$  and  $d_{ij}$ , the distance between atom  $i$  of molecule  $j$  to point  $i$  on  $F^{k-1}$ .

The residual is then calculated as

$$\delta^k = \frac{\sum_{i=1}^a \sum_{j=1}^b w_{ij}^{k-1} (\bar{F}_i^k - \bar{Z}_{ij})^2}{\sum_{j=1}^b \sum_{i=1}^a w_{ij}^{k-1}} \quad (4)$$

If the RMS distance between  $F^k$  and  $F^{k-1}$  is less than  $10^{-5}$  Å and the difference between  $\delta^k$  and  $\delta^{k-1}$  is less than  $10^{-5}$  Å then a minimum has been obtained and the equivalences can be updated.

These topologically equivalent atoms are determined from the optimal path through a matrix of Euclidean distances between all main-chain  $C_\alpha$  positions from the proteins. To trace this path, a dynamic programming technique is used.<sup>9,31</sup> If the equivalences have not changed from the previous iteration, the superposition is complete. Otherwise, the molecules are fitted to  $F^k$  pairwise, a new framework determined, the residual calculated, and this process repeated until convergence is attained.

### Distance Metric from Rigid-Body Superposition

From the pairwise rigid-body comparisons, two pieces of information can be obtained: the topological equivalences and the RMS distance over these equivalent  $C_\alpha$  positions (Table II). The number of topologically equivalent positions is converted to a pairwise fractional topological equivalence (*PETE*) by dividing by the length of the smaller structure. The RMS is a distance measure and is converted to a similarity score, the *SRMS*, calculated as  $1 - \text{RMS}(\text{\AA})/3.5(\text{\AA})$ ; a 3.5-Å cutoff is used in the definition of topological equivalence. The distance metric,  $D$ , employed in this study involves a weighted contribution [Eq. (5)] of these two parameters.<sup>6</sup>

$$D = -100 \ln(w_1 \text{PETE} + w_2 \text{SRMS}) \quad (5)$$

where the weights,  $w_1$  and  $w_2$ , are calculated from

$$w_1 = [(1 - \text{PETE}) + (1 - \text{SRMS})]/2 \quad (6)$$

$$w_2 = (\text{PETE} + \text{SRMS})/2 \quad (7)$$

with

$$w_1 + w_2 = 1 \quad (8)$$

For more closely related sequences, our experience has shown that the fraction of topologically equivalent positions may not differentiate be-

<sup>31</sup> M. L. Fredman, *Bull. Math. Biol.* **46**, 553 (1984).



TABLE II  
FEATURES USED IN COMPARISON OF PROTEIN STRUCTURES<sup>a</sup>

Multifeature structural comparisons:	
Residues	Segments
Rigid-body structural comparisons: Number of topologically equivalent positions RMS distance over the topological equivalences	
Multifeature structural comparisons:	
Residues	Segments
Properties	
Identity	Secondary structure type
Physical properties	Amphipathicity
Local conformation	Improper dihedral angle
Distance from gravity center	Distance from gravity center
Side-chain orientation	Orientation relative to gravity center
Main-chain orientation	Side-chain accessibility
Side-chain accessibility	Main-chain accessibility
Main-chain accessibility	Position in space
Position in space	Global orientation
Global direction in space	
Main-chain dihedral angles	
Relations	
Hydrogen bond	Distances to one or more nearest neighbors
Disulfide bond	Relative orientation of two or more segments
Ionic bond	
Hydrophobic cluster	

<sup>a</sup> Structural features that are considered by the rigid-body and the multifeature approach to the comparison of protein structures and the determination of phylogenetic relationships are given. For the multifeature approach, various features are represented by rows and different levels of protein organization by columns. Only residue and secondary structure levels are shown here. The term property is used for all protein features that imply comparison of only one element from each protein. Conversely, the term relationship is used for a feature that implies comparison of at least two elements from each protein.

tween the structures.<sup>4,6</sup> Conversely, the RMS distance does provide a good measure of the difference between structures where the relationship is close. As a result, the weights,  $w_1$  and  $w_2$ , are used to inversely modulate the contribution of the *PETE* and the *SRMS* to the distance score [Eqs. (6)–(8)].<sup>6</sup> A matrix containing all pairwise distances can then be used directly by clustering or tree-generating techniques to display the relationships derived from the rigid-body structural comparisons.

### Multifeature Comparison of Proteins

Although proteins within homologous families have the same tertiary folding, the elements of secondary structure undergo deformations, relative translations, and rotations to optimize packing of side chains and to adapt to evolutionary pressure. Thus, for two proteins with 30% sequence identity, the topologically equivalent residues defined by rigid-body superposition have a root mean square difference of approximately 1.5 Å and may comprise as few as one-third of the total number of residues.<sup>4,6</sup> This may not provide a sufficient basis for structural comparisons and emphasizes the requirement for a more flexible procedure for defining topological equivalence, one that can take into account relative movements and distortions of the secondary structure elements.

### Structural Aspects of Proteins: Properties, Relationships, and Hierarchy

To achieve this flexibility, we include in the comparison method a number of protein features from several levels of the protein structure hierarchy (Table II). The protein is treated as a sequence of elements where each element is associated with a series of properties and may be engaged in a number of relationships with other elements. Additionally, these elements may exist at any level of the hierarchy of protein structure: residue, secondary structure, supersecondary structure, motif, domain, or globular protomer. For example, at the residue level, properties such as the local conformation and relationships like hydrogen bonds can be included. At the level of secondary structure, properties like segment solvent accessibility and relationships such as the relative spatial orientation of two segments can be incorporated.

### Alignments from Dynamic Programming and Simulated Annealing

The comparison method is based on the dynamic programming technique generally used for sequence alignments.<sup>9</sup> In this method, one starts with the calculation of an  $N$  by  $M$  weight matrix  $W$  where  $N$  and  $M$  are the numbers of residues in the two compared proteins. This matrix is calculated in such a way that every element  $W_{ij}$  is proportional to the sum of the differences between various features of the residues  $i$  and  $j$ :

$$W_{ij} = \sum_p \left( \sum_l \rho^l w_{ij}^l + \sum_r \rho^r w_{ij}^r \right) \quad (9)$$

The contributions  $w_{ij}$  are the differences between individual features of the residues  $i$  and  $j$ , and factors  $\rho$  determine their relative weights. Superscript  $l$  runs over all levels of protein structure, superscript  $p$  stands for properties, and superscript  $r$  for relationships. When features at the second-

any structure level are considered, the residues inherit the weights from the secondary structure segments involved, the corresponding  $\alpha$  helices or  $\beta$  strands. It is trivial to define  $w_{ij}^{\beta}$  for properties  $p_i$  for example, the seventh property at the first level of structure,  $w_{ij}^{\beta 7}$ , describes the difference in the residue main-chain solvent accessibilities and is simply an absolute difference in the fractional main-chain accessibilities for the residues  $i$  and  $j$  from the first and second protein, respectively. In addition to properties, specific relationships such as hydrogen bonding interactions, which tend to be conserved in protein folds, can also be used in our comparison method. However, a relationship by its very nature affects more than one element in a sequence, and this precludes the simple procedure for the inclusion of this information into the residue-by-residue weight matrix  $\mathcal{W}$ .

To incorporate the information about relationships into the derivation of the final equivalences, we first use simulated annealing optimization<sup>32</sup> to obtain pairwise alignments based on relationships alone.<sup>27</sup> The underlying goal in the implementation of simulated annealing optimization is to maximize the number of equivalent relationships and minimize violations of a "no-knot" constraint. Since simulated annealing does not necessarily produce a global optimum, the optimization for every pair of structures is repeated several times. The fractional numbers of matching of residues  $i$  and  $j$  from proteins  $A$  and  $B$ , which can be obtained from several relationship alignments of proteins  $A$  and  $B$  in a straightforward way, are used to define the relationship weights  $w_{ij}^r$ . These weights can be introduced directly into the residue by residue weight matrix  $\mathcal{W}$  [Eq. (9)].

The dynamic programming algorithm then uses the matrix  $\mathcal{W}$  to derive the most parsimonious alignment of the two structures. The overall distance score, which reflects the dissimilarity in selected features of the two proteins, is also obtained. A detailed description of these algorithms and their implementation in the program COMPARER may be found elsewhere.<sup>27</sup>

### Multiple Structural Alignments

In the above description of the multifeature alignment method, we have assumed that the three-dimensional structures would be compared in a pairwise manner. However, such pairwise comparisons of several proteins may not be self-consistent, in the same way as pairwise sequence-based alignments may not be self-consistent. For this reason, we proceed by simultaneously aligning all structures. In COMPARER, we have adopted a strategy that employs a combination of the approaches by Feng and Doolittle<sup>10,11</sup> and Barton and Sternberg.<sup>12,13</sup> The procedure is divided

into two parts. The first part is the construction of a dendrogram relating the homologous proteins, either *ad hoc* or from distance scores from pairwise comparisons.<sup>10,11</sup> The second part involves the gradual addition of new proteins, as imposed by the tree topology, into a growing multiple alignment. The most similar proteins and groups of proteins are structurally aligned first, and the gaps that are introduced do not change in later stages.<sup>10,11</sup> The weight matrix for the dynamic alignment of the two groups of previously aligned proteins is defined on the basis of the pairwise weight matrices relating the proteins from the two groups.

### Distance Metric from Multifeature Comparison

With COMPARER, a pairwise distance score for each protein pair is obtained from the corresponding pairwise alignment implied by the multiple alignment. First, a sum of the weights  $\mathcal{W}_{ij}$  that relate the residues equivalent in the pairwise comparison is found. This sum is then normalized via division by the number of equivalent residues in the pairwise comparison to give the intermediate score  $e$ . The final pairwise distance,  $E$ , that is used in the clustering procedure is then defined as

$$E = -100 \ln(1 - e/D_c) \quad (10)$$

$D_c$  is a constant equal to the random value of the distance score  $e$  and is obtained for each protein pair by increasing the average of the weight matrix elements  $\mathcal{W}_{ij}$  by three standard deviations of these elements divided by the square root of the number of equivalent residues. In addition to the distance score  $E$ , which does not incorporate information about gaps in the alignment, we considered a similar score that does incorporate gap penalties. Tree topologies for the two distance measures were the same in all cases.

Trees reflecting the evolution of different aspects of the proteins can be obtained by calculating the pairwise score  $E$  from the weights  $\mathcal{W}_{ij}$  that were derived from different combinations of protein features. Thus, evolutionarily variable sequence features such as residue identity can be used for classification of similar proteins, and more conserved structural features, like hydrogen bonding, can be used for more divergent structures. Conversely, the clustering can also be used to infer the variability of a given protein feature in evolution.

### Methods

The alignment of sequence data was produced by the "historical" multiple alignment procedure of Feng and Doolittle,<sup>10,11</sup> trees were derived from the distance metric of Feng *et al.*<sup>15</sup> Structures and sequences

<sup>32</sup> S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi, *Science* **220**, 671, (1983).

were obtained from the structure files of the Brookhaven Protein Data Bank<sup>33</sup> with the following exceptions: the sequence of cytochrome *c*-550 from *Paracoccus denitrificans* is as revised by Ambler *et al.*,<sup>34</sup> the globin sequences from *Parasponia*<sup>35</sup> and *V. vesicula*,<sup>36</sup> for which there are no structures, were obtained from the NI-VAT<sup>37</sup> sequence data bank.

Tree topologies and branch lengths were determined from the distance matrices using the program KITSCH from the phylogeny inference package (PHYLIIP) of Felsenstein.<sup>38</sup> This procedure is a modification of the original Fitch-Margoliash<sup>39</sup> method and accounts for unequal rates of change among the proteins by adjusting distances so that the total branch lengths from the root of the tree to the tips of each of the leaves are equivalent. In addition, numerous topologies are explored by swapping branches locally. The "best" tree is defined as the one that minimizes the sum of the squared differences between the equivalent distances from the tree and the input matrix, where each squared difference is also normalized by the corresponding squared distance from the input matrix.

### Phylogenetic Trees from Structural Comparisons

In this chapter, we concentrate on the phylogenies that can be inferred from distance scores obtained from the pairwise rigid-body superposition and the alignment of structures based on many features. The results stemming from these analyses of three-dimensional protein structures are then compared with those obtained from an alignment of the amino acid sequences. Three homologous families serve as examples; these include the amino- and carboxy-terminal domains of the aspartic proteinases, eukaryotic and microbial *c*-type cytochromes, and globins.

### Aspartic Proteinases

Structures of three fungal aspartic proteinases (Table I) have been solved to high resolution: endothiapsin (4APE), penicillopepsin (2APP),

<sup>33</sup> F. C. Bernstein, T. F. Koetzle, G. J. B. Williams, E. F. Meyer, M. D. Brice, J. R. Rodgers, O. Kennard, T. Shimanouchi, and M. Tasumi, *J. Mol. Biol.* **112**, 535 (1977).

<sup>34</sup> R. P. Ambler, T. E. Meyer, M. D. Kamen, S. A. Schenman, and L. Sawyer, *J. Mol. Biol.* **147**, 351 (1981).

<sup>35</sup> A. A. Kort, J. E. Burns, M. J. Trinick, and C. A. Appleby, *FEBS Lett.* **180**, 55 (1985).

<sup>36</sup> S. Wakabayashi, H. Matsubara, and D. A. Webster, *Nature (London)* **322**, 481 (1986).

<sup>37</sup> R. F. Doolittle, *Science* **214**, 149 (1981).

<sup>38</sup> J. Felsenstein, *Evolution* **39**, 783 (1985).

<sup>39</sup> W. M. Fitch and E. Margoliash, *Science* **15**, 279 (1967).

and rhizopuspepsin (2APR). The sequence identity between any pair of these three proteins is roughly 40%. Additionally, Tang *et al.*<sup>40</sup> have shown that amino-terminal and carboxy-terminal domains, each comprising about one-half of the molecule, are related by *C*<sub>2</sub> symmetry.

The multifeature alignment of amino- and carboxy-terminal domains (Fig. 3a) gives equivalences that are generally identical to those obtained from careful inspection of the structures on a graphics terminal (Fig. 1). This is so even though fragments of secondary structure have undergone translations, rotations, distortions, and numerous changes in sequence. This alignment stands in sharp contrast to that obtained from either the rigid-body multiple structure superposition or from the multiple sequence alignment (Fig. 3). The multiple structure rigid-body technique locates only 43 topologically equivalent positions among the six domains (asterisked positions in Fig. 3a), each domain consisting of approximately 150 residues. These positions are in complete agreement with the COMPARE alignment (Fig. 3a). Pairwise rigid-body superposition of structures was used to derive the distances for tree construction and led to between 67 and 77 topologically equivalent positions between the two sets of domains. A comparison of the sequences of the amino-terminal domains with the carboxy-terminal domains of the aspartic proteinases aligns the conserved active-site amino acid triad Asp-Thr-Gly that is present in both domains. Outside of the neighborhood surrounding the catalytic region, only one other section is aligned similarly to that obtained from the multifeature comparison method (Fig. 3).

The trees (Fig. 4) derived from the structural comparisons, either the rigid-body or the multifeature procedure, are congruent: the amino-terminal domains clearly branch apart from the carboxy-terminal domains. For the cluster of either domain, the shorter distance is between endothiapsin and penicillopepsin (Fig. 4) and is consistent with an alignment of the three full-length sequences. By sequence, the branch order within each domain's cluster is well determined. However, the sequence similarity between the domains is not statistically significant at the level of 3σ; this illustrates the power of the structural comparison method where unequivocal relationships are found. The numerous structural features that are common to both sets of domains, most notably, the hydrogen bonding patterns, are consistent with the notion that the *c* domains result from gene duplication.<sup>40</sup>

<sup>40</sup> J. Tang, M. N. G. James, I. N. Hsu, J. A. Jenkins, and T. L. Blundell, *Nature (London)* **271**, 618 (1978).





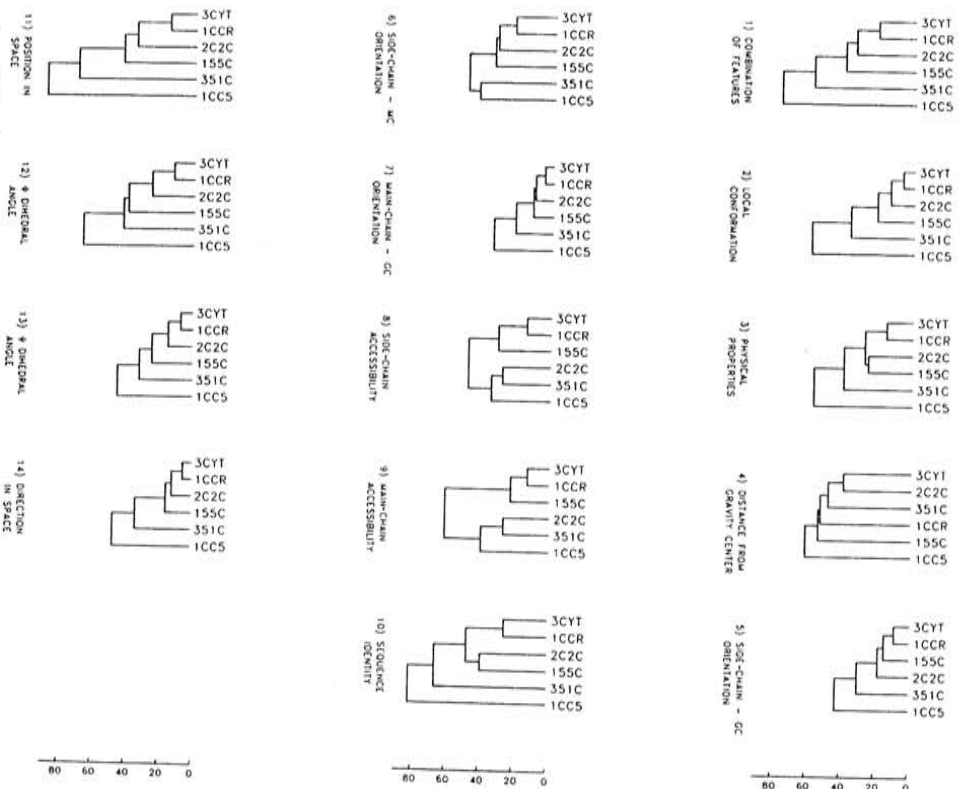


Fig. 5. Cladograms constructed from a multifeature alignment of the c-type cytochromes. Trees 2–14 were derived from the pairwise distances  $E$  [Eq. (10)] that were calculated using the individual features only. On the other hand, tree 1 was constructed from the same weight matrix elements  $W_{ij}$  that were used to derive the multiple alignment (data not shown). Features included were physical properties of amino acid residues (0.20), distance of the  $C_\alpha$  from the molecular gravity center (0.20), residue identities (0.10), absolute distance in space (0.30), and absolute main-chain directions in space (0.20). GC, Molecular gravity center; MC, main chain.

common among the structures than the small number of equivalences would suggest. The multifeature comparison method provides a technique for establishing a complete alignment of the structures (results not shown).

Trees reflecting various aspects of c-type cytochromes can be derived from their multiple alignment (Fig. 5). Trees 2 through 14 were obtained from the pairwise distances  $E$  that were calculated by considering only one feature in the derivation of weight matrix elements  $W_{ij}$  [Eq. (9)]. Tree 1 was constructed from the distances  $E$  obtained from the combination of features used to align the structures (Fig. 5). Of 14 trees, 7 have the same topology as the trees based on the multiple sequence and pairwise rigid-body comparison (Fig. 2); these include the trees derived from the local conformation of the main chain, orientation of the main chain relative to the molecular gravity center, absolute position of  $\alpha$  atoms, main-chain direction in approximately superposed structures,  $\phi$  and  $\psi$  dihedral angles, and the tree constructed from a combination of distances calculated from the features used to obtain the multiple alignment.

Trees that have a topology different from the most frequent one include the two trees derived from sequence information: the first of these two trees is based on the five physical characteristics of amino acid residues (such as hydrophobicity) that were found useful in construction of sequence alignments by Argos,<sup>16</sup> and the second tree is derived from a consideration of residue identities only. It may be noted that the clustering in these two trees is the same and corresponds to the subjective impression obtained from a consideration of the shape of the cytochrome structure in Fig. 2. The two trees reflecting similarities in the main-chain and side-chain accessibilities are also congruent with each other, but they are different from the two topologies mentioned above. In contrast, the unique and self-inconsistent topologies of the two trees that involve the orientation of side chains relative to the main chain and relative to the molecular gravity center imply that the orientation of the side chains is not a useful indicator for establishing relationships between divergent protein structures.

The three tree topologies (Fig. 5) based on the combination of features, sequence criteria, and solvent accessibility, demonstrate that evolutionary pressure does not act on all aspects of protein structure in the same way; thus, different criteria may be better for different purposes. For example, trees constructed from rigid-body superpositions are suitable in the selection of structures for determination of a framework<sup>6</sup> in homology-based protein modeling,<sup>21</sup> sequence-based trees are convenient for the description of evolutionary relationships among relatively similar proteins, while trees based both on pairwise rigid-body and multifeature comparison may be better for the analysis of more divergent structures.



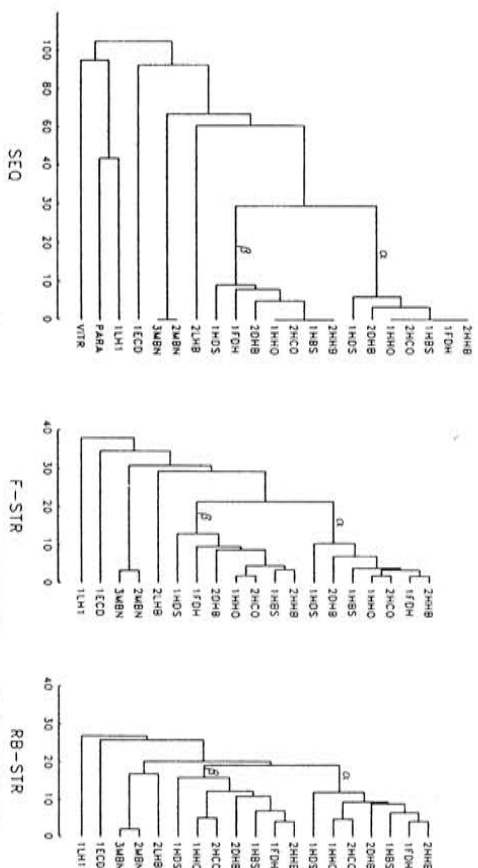


Fig. 7. Trees constructed for 19 globin structures (Table I) from pairwise rigid-body (RB-STR), multifactor (F-STR), and multiple sequence (SEC) alignments. The tree for the multifactor approach was calculated from the same combination of properties that was used to derive the alignment in Fig. 6.  $\alpha$ ,  $\alpha$  chains;  $\beta$ ,  $\beta$  and  $\gamma$  chains.

(2) Trees that could not be obtained from the comparison of sequences alone may be established from structural comparisons. Evolutionary trees for divergent proteins, where the sequence relationships are not statistically significant, can be derived from structures because tertiary structure is more conserved in evolution than sequence. For example, even though the comparison of the amino- and carboxy-terminal domains of the aspartic proteinases reveals many common structural features, the sequence similarity is not statistically significant. Structural procedures can distinguish between different crystal structures for the same sequence. For example, perturbations of hemoglobin structures induced by ligand binding are seen in the tree constructed from rigid-body superpositions, where all deoxy chains cluster apart from those with bound ligands.

## Acknowledgments

We thank our colleague John P. Overington for valuable discussions. We also thank Michael J. Sutcliffe. We thank the American Cancer Society (M.S.J.) for a fellowship and the ORS Awards Scheme and the Research Council of Slovenia for support (A.S.).

Numbers in parentheses are footnote reference numbers and indicate that an author's work is referred to although the name is not cited in the text.

## Author Index

- |   |   |
|---|---|
| <p><b>A</b></p> <p>Aaronson, S. A., 102</p> <p>Abdulaev, N. G., 104, 108(15)</p> <p>Abel, Y., 433, 636</p> <p>Abola, E., 56</p> <p>Ahmed, S., 445</p> <p>Akaike, H., 558</p> <p>Akira, M., 270</p> <p>Albert, P., 105</p> <p>Alberts, B. M., 5</p> <p>Altschul, S. F., 134, 135, 136(9), 137(9), 144, 347, 489</p> <p>Ambler, R. P., 682</p> <p>Amemura, A., 454</p> <p>Amemura, M., 123</p> <p>Anderson, W. F., 411, 418(23), 438, 440, 444(6), 445</p> <p>Andrews, P., 559</p> <p>Antonides, H. G., 102</p> <p>Avila, S., 112, 168</p> <p>Applebury, M. L., 105, 108(16)</p> <p>Appleby, C. A., 682</p> <p>Aquadro, C. F., 535</p> <p>Argast, M., 123</p> <p>Argos, P., 165, 352, 353, 355, 360, 362, 363, 364(15), 404, 421(4), 439, 671, 687(16)</p> <p>Arnberg, A. C., 241</p> <p>Arratia, R., 137, 222, 224(3), 487</p> <p>Artamanov, I. D., 104, 108(15)</p> <p>Atencio, E. J., 18</p> <p>Ausubel, F. M., 132, 223</p> <p>Axel, R., 105</p> <p>Ayala, F. J., 650</p> | <p><b>B</b></p> <p>Baas, F., 241</p> <p>Baba, M. L., 606</p> <p>Bachman, B. J., 20</p> <p>Backer, K. D., 260</p> <p>Bacon, D. J., 411, 418(23), 440, 444(6), 445(6)</p> <p>Becht, W., 105, 108(16)</p> <p>Bains, W., 411</p> <p>Baird, S., 452</p> <p>Baird, A., 238</p> <p>Bajaj, M., 670</p> <p>Baldwin, A. S., 399</p> <p>Baldwin, T. O., 124</p> <p>Baltimore, D., 107</p> <p>Banaszak, L. J., 154</p> <p>Bandelt, H. J., 530</p> <p>Barclay, A. N., 104</p> <p>Bardwell, J. C. A., 399</p> <p>Barker, W. C., 21, 32, 37, 56, 58, 144, 282, 341, 342, 344, 345(16), 346, 348, 352, 440, 457, 464(1), 461, 465(6), 466(6), 469(6), 472</p> <p>Barnabas, J., 601, 604</p> <p>Barnes, D., 21, 56</p> <p>Barton, G. J., 34, 35(4), 358, 370, 371(9), 406, 408, 414(1), 419, 421, 671, 680(13)</p> <p>Barton, G., 671, 672, 680(12)</p> <p>Barton, J. G., 463</p> <p>Bastford, D., 423</p> <p>Baut, J., 128, 130(26), 132(26)</p> <p>Baudin, F., 284</p> <p>Beattie, W. G., 74, 128(25), 400</p> <p>Beaud, G., 223</p> <p>Beck, C. F., 12</p> <p>Belfort, M., 29, 293, 295(26)</p> <p>Bell, G. I., 5, 20</p> <p>Bennett, C. D., 85, 105</p> <p>Benoist, C., 238, 253</p> <p>Benoiv, J. L., 85, 105</p> <p>Benstein, S. I., 260</p> <p>Benton, D., 8, 56, 282</p> <p>Berg, J. M., 107</p> <p>Berg, O. G., 219</p> |
|---|---|