

Sequence Order Independent Comparison of Protein Global Backbone Structures and Local Binding Surfaces for Evolutionary and Functional Inference

Joe Dundas, Bhaskar DasGupta, and Jie Liang

Abstract Alignment of protein structures can help to infer protein functions and can reveal ancient evolutionary relationship. We discuss computational methods we developed for structural alignment of both global backbones and local surfaces of proteins that do not depend on the ordering of residues in the primary sequences. The algorithm for global structural alignment is based on fragment assembly, and takes advantage of an approximation algorithm for solving the maximum weight independent set problem. We show how this algorithm can be applied to discover proteins related by complex topological rearrangement, including circularly permuted proteins as well as proteins related by complex higher order permutations. The algorithm for local surface alignment is based on solving the bi-partite graph matching problem through comparison of surface pockets and voids, such as those computed from the underlying alpha complex of the protein structure. We also describe how multiple matched surfaces can be used to automatically generate signature pockets and basis set that represents the ensemble of conformations of protein binding surfaces with a specific biological function of binding activity. This is followed by illustrative examples of signature pockets and basis set computed for NAD binding proteins, along with a discussion on how they can be used for discriminating NAD-binding enzymes from other enzymes.

Introduction

To understand the molecular basis of cellular processes, it is important to gain a comprehensive understanding of the biological functions of protein molecules. Although an increasing number of sequences and structures of proteins become available, there are many proteins whose biological functions are not known, or knowledge of their biological roles is incomplete. This is evidenced by the existence

J. Dundas (✉)

Bioinformatics Program, Department of Bioengineering, University of Illinois Chicago, Wolcott, Chicago, IL 60612-7340, USA

e-mail: jdunda1@uic.edu

46 of a large number of partially annotated proteins, as well as the accumulation of a
47 large number of protein structures from structural genomics whose biological func-
48 tions are not well characterized [1, 2]. Researchers have turned to *in silico* methods
49 to gain biological insight into the functional roles of these uncharacterized proteins,
50 and there has been a number of studies addressing the problem of computationally
51 predicting the biological function of proteins [3–8].

52 A relatively straightforward method for inferring protein function is to transfer
53 annotation based on homology analysis of shared characteristics between proteins.
54 If a protein shares a high level of sequence similarity to a well characterized family
55 of proteins, frequently the biological functions of the family can be accurately trans-
56 ferred onto that protein [9–11]. At lower levels of sequence similarity, probabilistic
57 models such as profiles can be constructed using local regions of high sequence sim-
58 ilarity [11–13]. The large amount of information of protein such as those deposited
59 in the SWISS-PROT database [14] provides rich information for constructing such
60 probabilistic models.

61 However, limitations to sequence-based homology transfer for function predic-
62 tion arise when sequence identity between a pair of proteins is less than 60% [16].
63 An alternative to sequence analysis is to infer protein based on structural similar-
64 ity. It is now well known that protein structures are much more conserved than
65 protein sequences, as proteins with little sequence identity often fold into similar
66 three-dimensional structures [17].

67 Protein structure and protein function are strongly correlated [18]. Conceptually,
68 knowledge of three-dimensional structures of proteins should enable inference of
69 protein function. Computational tools and databases for structural analysis are indis-
70 pensable for establishing the relationship between protein function and structure.
71 Among databases of protein structures, the SCOP [19] and CATH [20] databases
72 organizes protein structures hierarchically into different classes and folds based on
73 their overall similarity in topology and fold. Such classification of protein structures
74 based on structure generally depends on a reliable structural comparison method.
75 Although there are several widely used methods, including Dali [21] and CE
76 [22], current structural alignment methods cannot guarantee to give optimal results
77 and structural alignment methods do not have the reliability and interpretability
78 comparable to that of sequence alignment methods.

79 Comparing protein structures is challenging. First, it is difficult to obtain a quan-
80 titative measure of structural similarity that is generally applicable to different types
81 of problems. Similar to sequence alignment methods, one can search for global
82 structural similarity between overall folds or focus on local similarity between
83 surface regions of interest. Defining a quantitative measure of similarity is not
84 straightforward as illustrated by the variety of proposed structural alignment scoring
85 methods [23]. Unlike sequence alignment, in which the scoring systems are largely
86 based on evolutionary models of how protein sequence evolve [24, 25], scoring
87 systems of structural alignment must take into account both the three-dimensional
88 positional deviations between the aligned residues or atoms, as well as other charac-
89 teristics that are biologically important. Second, many alignment methods assume
90 the ordering of the residues follows that of the primary sequence when seeking

Comparison of Protein Global Backbone Structures and Local Binding Surfaces

to optimize structure similarity [22, 26]. This assumption can be problematic, as similar three-dimensional placement of residues may arise from residues with different sequential ordering. This problem is frequently encountered when comparing local regions on proteins structures. When comparing global structures of proteins, the existence of circular and higher ordered permutations [27, 28] also poses significant problems. Third, proteins may undergo minor residue side chain structural fluctuations as well as large backbone conformational changes in vivo. These structural fluctuations are not represented in a static snapshot of a crystallized structures in the Protein Data Bank (PDB) [29]. Many structural alignment methods assume rigid bodies and cannot account for structural changes that may occur.

In this chapter, we will first discuss several overall issues important for protein structural alignment. We then discuss a method we have developed for sequence order independent structural alignment at both the global and local level of protein structure. This is followed by discussion on how this method can be used to detect protein pairs that appear to be related by simple and complex backbone permutations. We will then describe the use of local structural alignment in automatic construction of *signature pockets* of binding surfaces, which can be used to construct *basis set* for a specific biological function. These constructs can detect structurally conserved surface regions and can be used to improve the accuracy of protein function prediction.

Structural Alignment

Protein structural alignment is an important problem [23]. It is particularly useful when comparing two proteins with low sequence identity between them. A widely used measure of protein structural similarity is the root mean squared distance (RMSD) between the equivalent atoms or residues of the two proteins. When the equivalence relationship between structural elements are known, a superposition described by a rotation matrix R and a translation vector T that minimizes the root mean squared distances (RMSD) between the two proteins can be found by solving the minimization problem:

$$\min \sum_{i=1}^{N_B} \sum_{j=1}^{N_A} |T + RB_i - A_j|^2, \quad (1)$$

where N_A is the number of points in structure A and N_B is the number of points in structure B and it is assumed that $N_A = N_B$. The least-squares estimation of the transformation parameters R and T in Eq. (1) can be found using the technique of singular value decomposition [30].

However, it is often the case that the equivalences between the structural elements are not known a priori. For example, when two proteins have diverged significantly. In this case, one must use heuristics to determine the equivalence relationship, and the problem of protein structural alignment becomes a multi-objective problem. That is, we are interested in finding the maximum number of equivalent elements as

136 well as in minimizing the RMSD upon superposition of the equivalent elements of
137 the two proteins.

138 A number of methods that are heuristic in nature have been developed for align-
139 ing protein structures [31–38]. These methods can be divided into two categories.
140 *Global* structural alignment methods are suited for detecting similarities between
141 the overall backbones of two proteins, while *local* structural alignment methods are
142 suited for detecting similarities between local regions or sub-structures within the
143 two proteins. As discussed earlier, many structural alignment algorithms are con-
144 strained to find only structural similarities where the order of the structural elements
145 follows their order in the primary sequence. Sequence order independent methods
146 ignore the sequential ordering of the structural elements and are better suited to find
147 more complex global structural similarities. They are also very effective for all atom
148 comparison of protein sub-structures, as in the case of binding surface alignment.
149 Below we discuss methods for both global and local structural alignment.

151 **Global Sequence Order Independent Structural Alignment**

152
153 Global sequence order independent structural alignment is a powerful tool that can
154 be used to detect similarities between two proteins that have complex topological
155 rearrangements, including permuted structures. Permuted proteins can be described
156 as two proteins with similar three-dimensional spatial arrangement of secondary
157 structures, but with a different backbone connection topology. An example of per-
158 muted proteins are proteins with circular permutations. It can be thought of as
159 ligation of the N- and C-termini of a protein, and cleavage somewhere else on the
160 protein. Circular permutations are interesting not only because they tend to have
161 similar three-dimensional structure but also because they often maintain the same
162 biological function [27]. Circularly permuted proteins may provide a generic mech-
163 anism for introducing protein diversity that is widely used in evolution. Detecting
164 circular permutations is also important for homology modeling, for studying protein
165 folding, and for designing protein.

167 168 ***A Fragment Assembly Based Approach to Sequence Order*** 169 ***Independent Structural Alignment***

170
171 We have developed a sequence order independent structural alignment method
172 that is well-suited for detecting circular permutation and more complex topolog-
173 ical rearrangement relationship among proteins [28]. Our algorithm is capable of
174 aligning two protein backbone structures independent of the secondary structure
175 element connectivity. Briefly, the two proteins to be aligned are first separately
176 and exhaustively fragmented. Each fragment $\lambda_{i,k}^A$ from protein structure S_A is then
177 pair-wise superimposed onto each fragment $\lambda_{j,k}^B$ from protein structure S_B , form-
178 ing a set of fragment pairs $\chi_{i,j,k}$, where $i \in S_A$ and $j \in S_B$ are the indices in
179 the primary sequence of the first residue of the two fragment, respectively. Here
180

Comparison of Protein Global Backbone Structures and Local Binding Surfaces

181 $k \in \{5, 6, 7\}$ is the length of the fragment. For each fragment, we assign a similarity
182 score,

$$183 \sigma(\chi_{i,j,k}) = \alpha \left[C - s(\chi_{i,j,k}) \cdot \frac{cRMSD}{k^2} \right] + SCS, \quad (2)$$

184 where $cRMSD$ is the measured RMSD value after optimal superposition, α and C
185 are two constants, $s(\chi_{i,j,k})$ is a scaling factor to the measured RMSD values that
186 depends on the secondary structure of this fragment, and SCS is a BLOSSUM-like
187 measure of similarity in sequence of the matched fragments [25]. Details of the
188 similarity score and the parameters α and C can be found in [28].

189 The goal of structural alignment for the moment seeks to find a consistent set of
190 fragment pairs $\Delta = \{\chi_{i_1,j_1,k_1}, \chi_{i_2,j_2,k_2}, \dots, \chi_{i_t,j_t,k_t}\}$ that minimize the global RMSD.
191 Finding the optimal combination of fragment pairs is a special case of the well
192 known maximum weight independent set problem in graph theory. This problem
193 is MAX-SNP-hard. We employ an approximation algorithm that was originally
194 described for scheduling split-interval graphs [39] and is itself based on a fractional
195 version of the local-ratio approach.

196 Our method begins by creating a conflict graph $G = (V, E)$, where a vertex is
197 defined for each aligned fragment pair. Two vertices are connected by an edge if any
198 of the fragments $(\lambda_{i,k}^A, \lambda_{i',k'}^A)$ or $(\lambda_{j,k}^B, \lambda_{j',k'}^B)$ from the aligned pair is not disjoint,
199 that is, if both fragments from the same protein share one or more residues. For
200 each vertex representing aligned fragment pair, we assign three indicator variables
201 $x_\chi, y_{\chi_{\lambda_A}},$ and $y_{\chi_{\lambda_B}} \in \{0, 1\}$ and a closed neighborhood $Nbr[\chi]$. x_χ indicates whether
202 the fragment pair should be used ($x_\chi = 1$) or not ($x_\chi = 0$) in the final alignment.
203 $y_{\chi_{\lambda_A}}$ and $y_{\chi_{\lambda_B}}$ are artificial indicator values for λ_A and λ_B , which allow us to encode
204 consistency in the selected fragments. The closed neighborhood of a vertex χ of G
205 is $\{\chi' \mid \{\chi, \chi'\} \in E\} \cup \{\chi\}$, which is simply χ and all vertices that are connected to χ
206 by and edge.

207 Our algorithm for sequence order independent structural alignment can now be
208 described as follows. To begin, we initialize the structural alignment Δ equal to the
209 entire set of aligned fragment pairs. We then:

- 210 1. Solve a linear programming (LP) formulation of the problem:

211 *maximize*

$$212 \sum_{\chi \in \Delta} \sigma(\chi) \cdot x_\chi \quad (3)$$

213 *subject to*

$$214 \sum_{a_t \in \lambda^A} y_{\chi_{\lambda_A}} \leq 1 \quad \forall a_t \in S_A \quad (4)$$

$$215 \sum_{b_t \in \lambda^B} y_{\chi_{\lambda_B}} \leq 1 \quad \forall b_t \in S_B \quad (5)$$

$$y_{\chi\lambda_A} - x_\chi \geq 0 \quad \forall \chi \in \Delta \quad (6)$$

$$y_{\chi\lambda_B} - x_\chi \geq 0 \quad \forall \chi \in \Delta \quad (7)$$

$$x_\chi, y_{\chi\lambda_A}, y_{\chi\lambda_B} \geq 0 \quad \forall \chi \in \Delta \quad (8)$$

2. For every vertex $\chi \in V_\Delta$ of G_Δ , compute its *local conflict number* $\alpha_\chi = \sum_{\chi' \in \text{Nbr}_\Delta[\chi]} x_{\chi'}$. Let χ_{\min} be the vertex with the *minimum* local conflict number. Define a new similarity function σ_{new} from σ as follows:

$$\sigma_{\text{new}}(\chi) = \begin{cases} \sigma(\chi), & \text{if } \chi \notin \text{Nbr}_\Delta[\chi_{\min}] \\ \sigma(\chi) - \sigma(\chi_{\min}), & \text{otherwise} \end{cases}$$

3. Create $\Delta_{\text{new}} \subseteq \Delta$ by removing from Δ every substructure pair χ such that $\sigma_{\text{new}}(\chi) \leq 0$. Push each removed substructure on to a stack in arbitrary order.
4. If $\Delta_{\text{new}} \neq \emptyset$ then repeat from step 1, setting $\Delta = \Delta_{\text{new}}$ and $\sigma = \sigma_{\text{new}}$. Otherwise, continue to step 5.
5. Repeatedly pop the stack, adding the substructure pair to the alignment as long as the following conditions are met:
 - a. The substructure pair is consistent with all other substructure pairs that already exist in the selection.
 - b. The *cRMSD* of the alignment does not change beyond a threshold. This condition bridges the gap between optimizing a local similarity between substructures and optimizing the tertiary similarity of the alignment. It guarantees that each substructure from a substructure pair is in the same spatial arrangement in the global alignment.

Detecting Permuted Proteins

This algorithm is used in a large scale study, where a subset with 3,336 protein structures taken from the PDBSELECT 90 data set % [40] are structurally aligned in a pair-wise fashion. Our goal is to determine if we could detect structural similarities with complex topological rearrangements such as circular permutations. From this subset of 3,336 proteins, we aligned two proteins if they met the following conditions: the difference in their lengths was no more than 75 residues, and they had approximately the same secondary structure content (see [28] for details). Within the approximately 200,000 alignments, we found many known circular permutations, and three novel circular permutations previously unknown, as well as a pair of non-cyclic complex permuted proteins. Below we describe in some details the circular permutations we found between a nucleoplasmin-core and an auxin binding protein, as well as details of the more complex non-cyclic permutation.

Comparison of Protein Global Backbone Structures and Local Binding Surfaces

Nucleoplasmin-Core and Auxin Binding Protein

A novel circular permutation was detected between the nucleoplasmin-core protein in *Xenopus laevis* (PDB ID 1k5j, chain E) [41] and the auxin binding protein in maize (PDB ID 1l1rh, chain A, residues 37 through 127) [42]. The structural alignment between 1k5jE (Fig. 1a, top) and 1l1rhA (Fig. 1a, bottom) consisted of 68 equivalent residues superimposed with an RMSD of 1.36 Å. This alignment is statistically significant with a p -value of 2.7×10^{-5} after Bonferroni correction. Details of p -value calculation can be found in reference [28]. The short loop connecting two antiparallel strands in nucleoplasmin-core protein (in circle, top of Fig. 1b) becomes disconnected in auxin binding protein 1 (in circle, bottom of Fig. 1b), and the N- and C- termini of the nucleoplasmin-core protein (in square, top of Fig. 1b) are connected in auxin binding protein 1 (square, bottom of Fig. 1b). For details of other circular permutations we discovered, including permutations between aspartate racemase and type II 3-dehydrogenase and between microphage migration inhibition factor and the C-terminal domain of arginine repressor, please see [28].

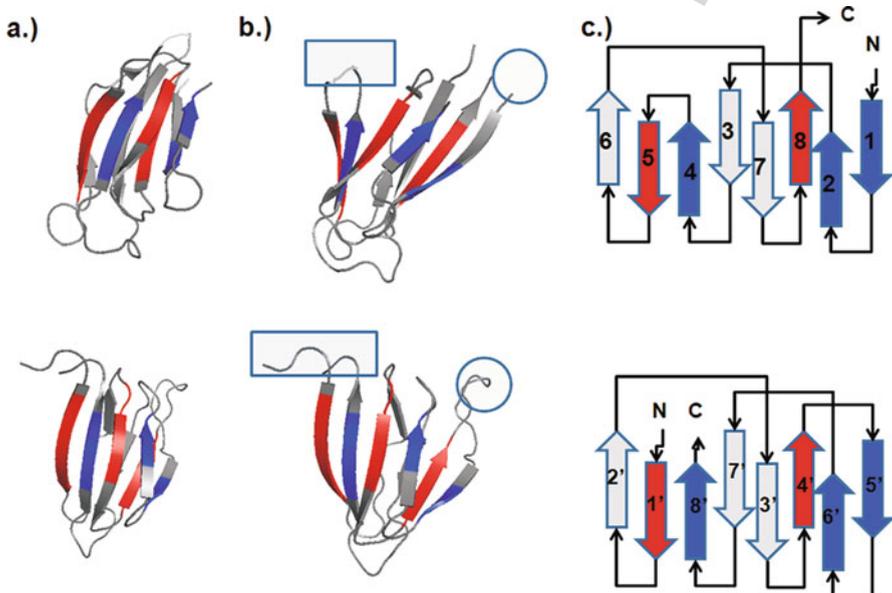


Fig. 1 A newly discovered circular permutation between nucleoplasmin-core (1k5j, chain E, *top panel*), and a fragment of auxin binding protein 1 (residues 37–127) (1l1rh, chain A, *bottom panel*). **a** These two proteins align well with a RMSD value of 1.36 Å over 68 residues, with a significant p -value of 2.7×10^{-5} after Bonferroni correction. **b** The loop connecting strand 4 and strand 5 of nucleoplasmin-core (in *rectangle, top*) becomes disconnected in auxin binding protein 1. The N- and C- termini of nucleoplasmin-core (in *rectangle, top*) become connected in auxin binding protein 1 (in *rectangle, bottom*). To aide in visualization of the circular permutation, residues in the N-to-C direction before the cut in the nucleoplasmin-core protein are colored *red*, and residues after the cut are colored *blue*. **c** The topology diagram of these two proteins. In the original structure of nucleoplasmin-core, the electron density of the loop connecting strand 4 and strand 5 is missing in the PDB structure file. This figure is modified from [28]

Beyond Circular Permutation

Because of its relevance in understanding the functional and folding mechanism of proteins, circular permutations have received much attention [27, 43]. A more challenging class of permuted proteins is that of the non-cyclic permutation with more complex topological changes. Very little is known about this class of permuted proteins, and the detection of non-cyclic permutations is a challenging task [44–47].

Non-cyclic permutations of the Arc repressor were created artificially and were found to be thermodynamically stable. It can refold on the sub-millisecond time scale, and can bind operator DNA with nanomolar affinity [48], indicating that naturally occurring non-cyclic permutations may be as rich as the cyclic permutations. Our database search uncovered a naturally occurring non-cyclic permutation between chain F of AML1/Core Binding Factor (AML1/CBF, PDB ID 1e50, Fig. 2a, top) and chain A of riboflavin synthase (PDB ID 1pkv, Fig. 2a, bottom) [49, 50]. The

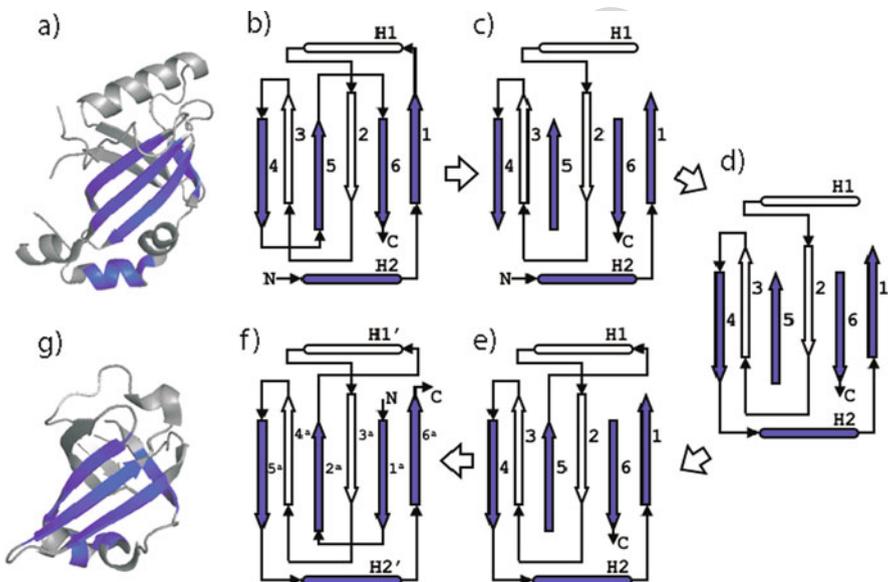


Fig. 2 A non-cyclic permutation discovered between AML1/Core Binding Factor (AML1/CBF, PDB ID 1e50, Chain F, *top*) and riboflavin synthase (PDB ID 1pkv, chain A, *bottom*) **a** These two proteins structurally align with an RMSD of 1.23 Å over 42 residues, and has a significant p -value of 2.8×10^{-4} after Bonferroni correction. The residues that were assigned equivalences from the structural alignment are colored blue. **b** These proteins are related by a complex permutation. The steps to transform the topology of AML1/CBF (*top*) to riboflavin (*bottom*) are as follows: **c** Remove the loops connecting strand 1 to helix 2, strand 4 to strand 5, and strand 5 to helix 6; **d** Connect the C-terminal end of strand 4 to the original N-termini; **e** Connect the C-terminal end of strand 5 to the N-terminal end of helix 2; **f** Connect the original C-termini to the N-terminal end of strand 5. The N-terminal end of strand 6 becomes the new N-termini and the C-terminal end of strand 1 becomes the new C-termini. We now have the topology diagram of riboflavin synthase. This figure was modified from [28]

Comparison of Protein Global Backbone Structures and Local Binding Surfaces

361 two structures align well with an RMSD of 1.23 Å, at an alignment length of 42
362 residues, with a significant p -value of 2.8×10^{-4} after Bonferroni correction.

363 The topology diagram of AML1/CBF (Fig. 2b) can be transformed into that of
364 riboflavin synthase (Fig. 2f) by the following steps: Remove the loops connecting
365 strand 1 to helix 2, strand 4 to strand 5, and strand 5 to strand 6 (Fig. 2c). Connect
366 the C-terminal end of strand 4 to the original N-termini (Fig. 2d). Connect the C-
367 terminal end of strand 5 to the N-terminal end of helix 2 (Fig. 2e). Connect the
368 original C-termini to the N-terminal end of strand 5. The N-terminal end of strand
369 6 becomes the new N-termini and the C-terminal end of strand 1 becomes the new
370 C-termini (Fig. 2f).

371 372 Local Sequence Order Independent Structural Alignment 373

374 The comparison of overall structural folds regardless of topological reconnections
375 can lead to insight into distant evolutionary relationship. However, similarity in
376 overall fold is not a reliable indicator of similar function [51–53]. Several studies
377 suggest that structural similarities between local surface regions where biological
378 function occurs, such as substrate binding sites, are a better predictor of shared
379 biological function [8, 54–58].

380 Substrate binding usually occurs at concave surface regions, commonly referred
381 to as *surface pockets* [56, 59–61]. A typical protein has many surface pockets, but
382 only a few of them present a specific three-dimensional arrangement of chemical
383 properties conducive to the binding of a substrate. This protein must maintain this
384 physiochemical environment throughout evolution in order to maintain its biological
385 function. For this reason, shared structural similarities between *functional surfaces*
386 among proteins may be a strong indicator of shared biological function. This has
387 lead to a number of promising studies, in which protein functions can be inferred by
388 similarity comparison of local binding surfaces [56, 62–65].

389 A challenging problem with the structural comparison of protein pockets lies in
390 the inherent flexibility of the protein structure. A protein is not a static structure
391 represented by a Protein Data Bank entry. The whole protein as well as the local
392 functional surface may undergo large structural fluctuations. The use of a single
393 surface pocket structure as a representative template for a specific protein function
394 will often result in many false negatives. This is due to the inability of a single
395 representative to capture the full functional characteristics across all conformations
396 of the protein.

397 To address this problem, we have developed a method that can automatically
398 identify the structurally preserved atoms across a family of protein structures that
399 are functionally related. Based on sequence-order independent surface alignments
400 across the functional pockets of a family of protein structure, our method creates
401 *signature pockets* with structurally conserved atoms identified and their fluctuation
402 measured. As more than one signature pocket may result for a single functional
403 class, the signature pockets can be organized into a *basis set* of pockets for that
404 functional family. These signature pockets of the binding surfaces then can be used
405 for scanning a protein structure database for function inference.

Bi-partite Graph Matching Approach to Structural Alignment

Our method for surface alignment is sequence order independent. It is based on a maximum weight bi-partite graph matching formulation of [66] with further modifications. This alignment method is a two step iterative process. First, an optimal set of equivalent atoms under the current superposition are found using a bi-partite graph representation. Second, a new superposition of the two proteins is determined using the new equivalent atoms from the previous step. The two steps are repeated until a stopping condition has been met.

To establish the equivalence relationship, two protein functional pocket surfaces S_A and S_B are represented as a graph, in which a node on the graph represent an atom from one of the two functional pockets. The graph is bi-partite if edges only connect nodes from protein S_A to nodes from protein S_B . In our implementation, directed edges are only drawn from nodes of S_A to nodes of S_B if a similarity threshold is met. The similarity threshold used in our implementation is a function of spatial distances and chemical differences between the corresponding atoms (see [67] for details). Each edge $e_{i,j}$ connecting node i to node j is assigned a weight $w(i,j)$ equal to the similarity score between the two corresponding atoms. A set of equivalence relations between atoms of S_A and atoms of S_B can be found by selecting a subset of the edges connecting nodes of S_A to S_B , with maximized total edge weight, where at most one edge can be selected for each atom [68]. A solution to the maximum weight bi-partite graph matching problem can be found using the Hungarian algorithm [69].

The Hungarian method works as follows. To begin, an overall score $F_{\text{all}} = 0$ is initialized, and an artificial source node s and an artificial destination node d are added to the bi-partite graph. Directed edges with 0-weight from the source node s to each node of S_A and from each node of S_B to the destination node d are also added. The algorithm then proceeds as follows:

1. Find the shortest distance $F(i)$ from the source node s to every other node i using the Bellman-Ford [71] algorithm.
2. Assign a new weight $w'(i,j)$ to each edge that does not originate from the source node s as follows,

$$w'(i,j) = w(i,j) + [F(i) - F(j)]. \quad (9)$$

3. Update F_{all} as $F_{\text{all}}' = F_{\text{all}} - F(d)$
4. Reverse the direction of the edges along the shortest path from s to d .
5. If $F_{\text{all}} > F(d)$ and a path exists between s and d then start again at step 1.

The Hungarian algorithm terminates when either there is no path from s to d or when the shortest distance from the source node to the destination node $F(d)$ is greater than the current overall score F_{all} . The bi-partite graph will now consist of directed edges that have been reversed (point from nodes of S_B to nodes of S_A). These flipped edges represent the current equivalence relationships between atoms of S_A and atoms of S_B .

Comparison of Protein Global Backbone Structures and Local Binding Surfaces

451 The equivalence relations can then be used to superimpose the two proteins. After
452 superposition, a new bi-partite graph is created and the maximum weight bi-partite
453 matching algorithm is called again. This process is repeated iteratively until the
454 change in RMSD upon superposition falls below a threshold.

455
456
457 ***Signature Pockets and Basis Set of Binding Surface***
458 ***for a Functional Family of Proteins***
459

460 Based on the pocket surface alignment algorithm, we have developed a method that
461 automatically generate structural templates of local surfaces, called *signature pock-*
462 *ets*, which can be used to represent an enzyme function or a binding activity. These
463 signature pockets contain broad structural information as well as discriminating
464 ability.

465 A signature pocket is derived from an optimal alignment of precomputed surface
466 pockets in a sequence-order-independent fashion, in which atoms and residues are
467 aligned based on their spatial correspondence when maximal similarity is obtained,
468 regardless how they are ordered in the underlying primary sequences. Our method
469 does not require the atoms of the signature pocket to be present in all member
470 structures. Instead, signature pockets can be created at varying degrees of partial
471 structural similarity, and can be organized hierarchically at different level of binding
472 surface similarity.

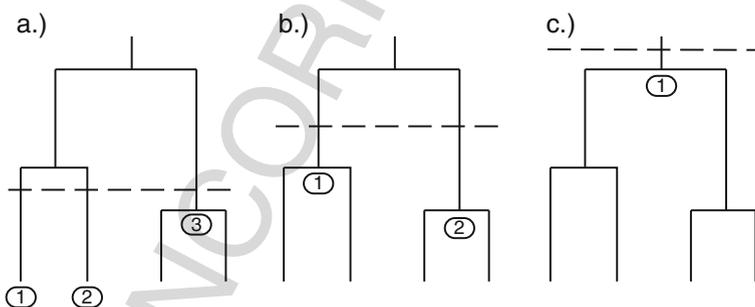
473 The input to the signature pocket algorithm is a set of functional pockets from a
474 pre-calculated database of surface pockets and voids on proteins, such as those con-
475 tained in the CASTp database [61]. The algorithm begins by performing all vs all
476 pair-wise sequence order independent structural alignment on the input functional
477 surface pockets. A distance score, which is a function of the RMSD and the chem-
478 istry of the paired atoms from the structural alignment, is recorded for each aligned
479 pair of functional pockets (see [67] for details). The resulting distance matrix is
480 then used by an agglomerative clustering method, which generates a hierarchical
481 tree. The signature of the functional pockets can then be computed using a recursive
482 process following the hierarchical tree.

483 The process begins by finding the two closest siblings (pockets S_A and S_B), and
484 combining them into a single surface pocket structure S_{AB} . Because of the recursive
485 nature of this algorithm, either of the two structures being combined may themselves
486 already be a combination of several structures. When combining the two structures,
487 we follow the criteria listed below:

- 488
489 1. If two atoms were considered equivalent in a structural alignment, a single
490 coordinate is created in the new structure to represent both atoms. The new coord-
491 inate is calculated by averaging the coordinates of all underlying atoms that are
492 currently represented by the two coordinates to be averaged.
- 493 2. If no equivalence was found for an atom during the structural alignment,
494 the coordinates of that atom are transferred directly into the new pocket
495 structure.

496 During each step in combining two surface pockets, a count of the number of
 497 times that an atom at the position i was present in the underlying set of pockets
 498 is recorded, which is then divided by the number of the constituent pockets. This
 499 is the *preservation ratio* $\rho(i)$. In addition, the mean distance of the coordinates of
 500 the aligned atoms to their geometric center is recorded as the *location variation* v .
 501 At the end of each step, the new structure S_{AB} replaces the two structures S_A and
 502 S_B in the hierarchical tree, and the process is repeated on the updated hierarchical
 503 tree. At a specific height of the hierarchical tree, different signature pockets can be
 504 created with different extents of structural preservation by selecting a ρ threshold
 505 value.

506 The signature pocket algorithm can be terminated at any point during its traversal
 507 of the hierarchical tree. Figure 3 illustrates this point by showing three different
 508 stopping thresholds (horizontal dashed lines). Depending on the choice of the
 509 threshold, one or multiple signature pockets may result. Figure 3a shows a low
 510 threshold which results in a set of 3 signature pockets. Raising the threshold can
 511 produce fewer signature pockets (Fig. 3b). A single signature pocket that repre-
 512 sents all surface pockets in the data set can be generated by raising the threshold
 513 even further (Fig. 3c). Since clusters from the hierarchical tree represent a set of
 514 surface pockets that are similar within certain threshold, if a stopping threshold is
 515 chosen such that there exist multiple clusters in the hierarchical tree, a signature
 516 pocket will be created for each cluster. The set of signature pockets from differ-
 517 ent clusters collectively form a *basis set* of signature pockets, which represent the
 518 ensemble of differently sampled conformations for a functional family of proteins.
 519 As a basis set of signatures can represent many possible variations in shapes and
 520 chemical textures, it can represent structural features of an enzyme function with
 521 complex binding activities, and can also be used to accurately predict enzymes
 522 function.



537 **Fig. 3** Different basis sets of signature pockets can be produced at different levels of structural
 538 similarity by raising or lowering the similarity threshold (*vertical dashed line*). **a** A low threshold
 539 will produce more signature pockets. **b** As the threshold is raised, fewer signature pockets will be
 540 created. **c** A single signature pocket can in principle be created to represent the full surface pocket
 data set by raising the threshold

Signature Pockets of NAD Binding Proteins

To illustrate how signature pockets and basis set help to identify key structural elements important for binding and how they can facilitate function inference, we discuss a study of the nicotinamide adenine dinucleotide (NAD) binding proteins. NAD consists of two nucleotides, nicotinamide and adenine, which are joined by two phosphate groups. NAD plays essential roles in metabolism where it acts as a coenzyme in redox reactions, including glycolysis and the citric acid cycle.

Using a set of 457 NAD binding proteins of diverse fold structures and diverse evolutionary origin, we first extracted the NAD binding surfaces from precomputed CASTp database of protein pockets and voids [61]. Based on similarity values from a comprehensive all-against-all sequence order independent surface alignment, we obtain a hierarchical tree of NAD binding surfaces. The resulting 9 signature pockets of the NAD binding pocket form a basis set, which are shown in Fig. 4.

These signature pockets contain rich biological information. Among the NAD-binding oxio-reductase, three signature pockets (Fig. 4e, h, and i) are for clusters of oxio-reductases that act on the CH-OH group of donors (alcohol oxio-reductases), one signature pocket (Fig. 4j) is for a cluster that act on the aldehyde group of donors, and the remaining two signature pockets (Fig. 4f and g) are for oxio-reductases that act on the CH-CH group of donors. For NAD-binding lyase, one of the two signature pockets (Fig. 4d) represent lyase that cleave both C-O and P-O bonds. The other signature pocket (Fig. 4b) represent lyases that cleave both C-O and C-C bonds. These two signatures come from two clusters of lyase conformations, each with a very different class of conformations of the bound NAD cofactor.

We found that the structural fold and the conformation of the bound NAD cofactor are the two major determinants of the formation of the clusters of the NAD binding pockets (Fig. 4a). It can be seen in Fig. 4b-j that there are two general conformations of the NAD coenzyme. The NAD coenzymes labeled C (Fig. 4b, c, f, g, h, and j) have a closed conformation, while the coenzymes labeled X (Fig. 4d, e, and i) have an extended conformation. This indicates that the binding pocket may take multiple conformations yet bind the same substrate in the same general structure. For example, the two structurally distinct signature pockets shown in Fig. 4f, g are derived from proteins that have the same biological function and SCOP fold. All of these proteins bind to the same NAD conformation.

We have further evaluated the effectiveness of the NAD binding site basis set by determining its accuracy in correctly classifying enzymes as either NAD-binding or non-NAD-binding. We constructed a test data set of 576 surface pockets from the CASTp database [61] independent of the training set of 457 NAD binding proteins. These 576 surface pockets were selected by taking the top 3 largest pockets in volume from 142 randomly chosen proteins and 50 proteins that have NAD bound in the PDB structure, with the further constraint that they were not in our training data set. We then structurally aligned all 576 pockets in our test data set against each of the nine NAD signature pockets in the resulting basis set. The testing pocket was assigned to be an NAD binding pocket if it structurally aligned to one of the nine NAD signature pockets, with its distance under a predefined threshold. Otherwise it

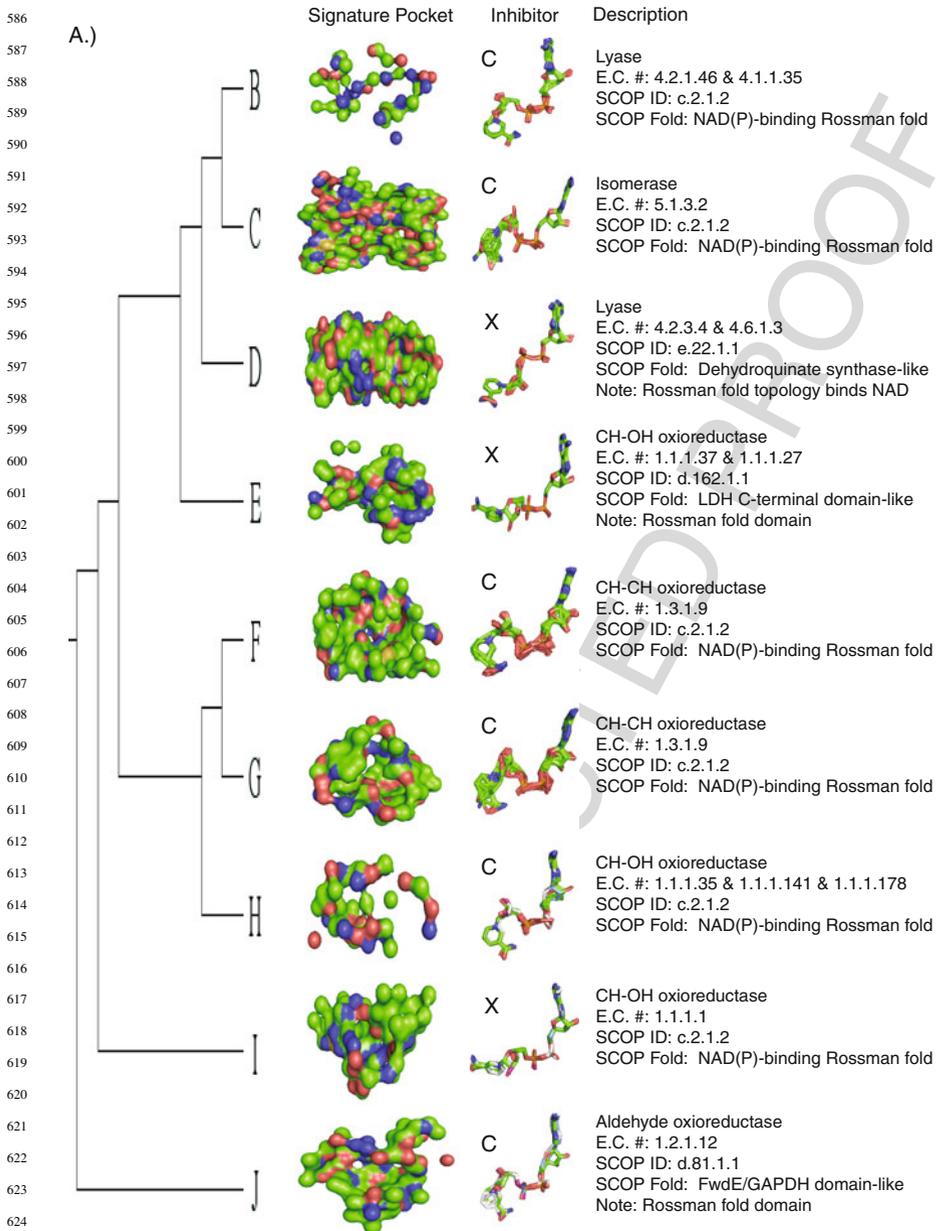


Fig. 4 The topology of the hierarchical tree and signature pockets of the NAD binding pockets. **a** The resulting hierarchical tree topology. **b–j** The resulting signature pockets of the NAD binding proteins, along with the superimposed NAD molecules that were bound in the pockets of the member proteins of the respective clusters. The NAD coenzymes have two distinct conformations. Those in an extended conformation are marked with an X and those in a compact conformation are marked with a C

Comparison of Protein Global Backbone Structures and Local Binding Surfaces

631 was classified as non-NAD binding. The results show that the basis set of 9 signature
632 pockets can classify the correct NAD binding pocket with sensitivity and specificity
633 of 0.91 and 0.89, respectively. We performed further testing to determine whether a
634 single representative NAD binding pocket, as opposed to a basis set, is sufficient for
635 identifying NAD-binding enzymes. We chose a pocket representative pocket from
636 one of the 9 clusters that were used to construct the 9 signature pockets. Here, a
637 testing pockets was classified as NAD-binding if its structural similarity to the single
638 representative pocket was above the same pre-defined threshold used in the basis
639 set study. We repeat this exercise nine times, each time using a different representa-
640 tive from a different cluster. We found that the results deteriorated significantly, with
641 an average sensitivity and specificity of only 0.36 and 0.23, respectively. This study
642 strongly indicates that the construction of a basis set of signatures as a structural
643 template provides significant improvement for a set of proteins binding the same
644 co-factor but with diverse evolutionary origin. Further details of the NAD-binding
645 protein study can be found in [67], along with an in-depth study of the metalloen-
646 dopeptidase, including the construction of its signatures and basis set, as well as
647 their utility in function prediction.

648

649

650 Conclusion

651

652 In this chapter, we have discussed methods that provide solutions to the problem of
653 aligning protein global structures as well as aligning protein local surface pockets.
654 Both methods disregard the ordering of residues in the protein primary sequences.
655 For global alignment of protein structures, such a method can be used to address
656 the challenging problem of identifying proteins that are topologically permuted but
657 are spatially similar. The approach of fragment assembly based on the formulation
658 of a relaxed integer programming problem and an algorithm based on scheduling
659 split-interval graphs works well, and is characterized by a guaranteed approximation
660 ratio. In a scaled up study, we showed that this method works well in discovery
661 of circularly permuted proteins, including several previously unrecognized protein
662 pairs. It also uncovered a case of two proteins related by higher order permutations.

663 We also described a method for order-independent alignment of local spatial sur-
664 faces that is based on bi-partite graph matching. By assessing surface similarity
665 for a group of protein structures of the same function, this method can be used to
666 automatically construct signatures and basis set of binding surfaces characteristic
667 of a specific biological function. We showed that such signatures can reveal use-
668 ful mechanistic insight on enzyme function, and can correlate well with substrate
669 binding specificity.

670 In this chapter, we neglect an important issue in our discussion of comparing
671 protein local surfaces for inferring biochemical functions, namely, how to detect
672 evolutionary signals and how to employ such information for protein function pre-
673 diction. Instead of going into details, we first point readers to the general approach of
674 constructing continuous time Markovian models to study protein evolution [72, 73].
675 In addition, a Bayesian Monte Carlo method that can separate selection pressure due

676 to biological function from selection pressure due to the constraints of protein fold-
677 ing stability and folding dynamics can be found in [58] and in [74]. The Bayesian
678 Monte Carlo approach can be used to construct customized scoring matrices that are
679 specific to a particular class of proteins of the same function. Details of how such
680 method works and how it can be used to accurately predict enzyme functions from
681 structure with good sensitivity and specificity for 100 enzyme families can be found
682 in a recent review [74] and original publications [8, 58]. The task of computing
683 surface pockets and voids using alpha shape is discussed in a recent review [75].
684

685 **Acknowledgements** This work was supported by NIH grants GM079804, GM081682,
686 GM086145, NSF grants DBI-0646035 and DMS-0800257, ONR grant N00014-09-1-0028.
687

688 References

- 691
- AQ2
- 692 1. Binkowski, A., Joachimiak, A., Liang, J. Protein surface analysis for function annotation in
693 high-throughput structural genomics pipeline. *Protein Sci.* **14**: 2972–2981 (2005).
 - 694 2. Pazos, F., Sternberg, M.J.E. Automated prediction of protein function and detection of
695 functional sites from structure. *PNAS* **101**:14, 14754–14759 (2004).
 - 696 3. Andrade, M.A., Brown, N.P., Leroy, C., Hoersch, S., de Daruvar, A., Reich, C., Franchini, A.,
697 Tamames, J., Valencia, A., Ouzounis, C., Sander, C. Automated genome sequence analysis and
698 annotation. *Bioinformatics* **15**: 391–412 (1999).
 - 699 4. Jensen, L.J., Gupta, R., Blom, N., Devos, D., Tamames, J., Kesmir, C., Nielsen, H., Staerfeldt,
700 H.H., Rapacki, K., Workman, C., Andersen, C.A.F., Knudsen, S., Krogh, A., Valencia, A.,
701 Brunak, S. Prediction of human protein function from post-translational modifications and
702 localization features. *J. Mol. Biol.* **319**: 1257–1265 (2002).
 - 703 5. Pal, D., Eisenberg, D. Inference of protein function from protein structure. *Structure* **13**:
704 121–130 (2005).
 - 705 6. Laskowski, R.A., Watson, J.D., Thornton, J.M. ProFunc: a server for predicting protein
706 function from 3D structure. *Nucleic Acids Res.* **33**: W89–93 (2005).
 - 707 7. Deng, M., Zhang, K., Mehta, S., Chen, T., Sun, F. Prediction of protein function using protein-
708 protein interaction data. *J. Comput. Biol.* **10**(6): 947–960 (2003).
 - 709 8. Tseng, Y.Y., Dundas, J., Liang, J. Predicting protein function and binding profile via
710 matching of local evolutionary and geometric surface patterns. *J. Mol. Biol.* **387**(2): 451–
711 464 (2009).
 - 712 9. Shah, I., Hunterm, L. Predicting enzyme function from sequence: a systematic appraisal.
713 *ISMB* **5**: 276–283 (1997).
 - 714 10. Altschul, S.F., Warren, G., Miller, W., Myers, E.W., Lipman, D.J. Basic local alignment search
715 tool. *J. Mol. Biol.* **215**: 403–410 (1990).
 - 716 11. Altschul, S.F., Madden, T.L., Schaffer, A.A., Zhang, J., Zhang, Z., Miller, W., Lipman, D.J.
717 Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.
718 *Nucleic Acids Res.* **25**(17): 3389–3402 (1997).
 - 719 12. Karplus, K., Barret, C., Hughey, R. Hidden Markov Models for detecting remote protein
720 homologues. *Bioinformatics* **14**: 846–856 (1998).
 13. Hulo, N., Sigrist, C.J.A., Le Saux, V. Recent improvements to the PROSITE database. *Nucleic
Acids Res.* **32**: D134–D137 (2004).
 14. Boeckmann, B., Bairoch, A., Apweiler, R., Blatter, M.C., Estreicher, A., Gasteiger, E., Martin,
M.J., Michoud, K., O'Donovan, C., Phan, I., Pilbout, S., Schneider, M. The SWISS-PROT
protein knowledgebase and its supplement TrEMBL in 2003. *Nucleic Acids Res.* **31**: 365–370
(2003).

Comparison of Protein Global Backbone Structures and Local Binding Surfaces

- AQ3** 721 15. Chung, J.L., Wang, W., Bourne, P.E. High-throughput identification of interacting protein-
722 protein binding sites. *BMC Bioinformatics* **8**: 223 (2007).
- 723 16. Weidong, T., Skolnick, J. How well is enzyme function conserved as a function of pairwise
724 sequence identity. *J. Mol. Biol.* **333**: 863–882 (2003).
- 725 17. Rost, B. Twilight zone of protein sequence alignments. *Protein Eng.* **12**: 85–94 (1999).
- 726 18. Hegyi, H., Gerstein, M. The relationship between protein structure and function: a com-
727 prehensive survey with application to the yeast genome. *J. Mol. Biol.* **288**: 147–164
(1999).
- 728 19. Murzin, A.G., Brenner, S.E., Hubbard, T., Chothia, C. SCOP: a structural classification of
729 proteins database for the investigation of sequences and structures. *J. Mol. Biol.* **247**: 536–540
(1995).
- 730 20. Orengo, C.A., Michie, A.D., Jones, D.T., Swindells, M.B., Thornton, J.M. CATH: a hierar-
731 chical classification of protein domain structures. *Structure* **5**: 1093–1108 (1997).
- 732 21. Holm, L., Sander, C. Protein structure comparison by alignment of distance matrices. *J. Mol.*
733 *Biol.* **233**: 123–138 (1993).
- 734 22. Shindyalov, I.N., Bourne, P.E. Protein structure alignment by incremental combinatorial
735 extension (CE) of the optimal path. *Protein Eng.* **11**(9): 739–747 (1998).
- 736 23. Hasegawa, H., Holm, L. Advances and pitfalls of protein structural alignment. *Curr. Opin.*
Struct. Biol. **19**: 341–348 (2009).
- AQ4** 737 24. Dayhoff, M.O., Schwartz, R.M., Orcutt, B.C. Atlas Protein Seq. *Struct.* **5**(3): 345–352 (1978).
- 738 25. Henikoff, S., Henikoff, J.G. Amino acid substitution matrices from protein blocks. *PNAS*
739 **89**(22): 10915–10919 (1992).
- 740 26. Teichert, F., Bastolla, U., Porto, M. SABERTOOTH: protein structure comparison based on
741 vectorial structure representation. *BMC Bioinformatics* **8**: 425 (2007).
- 742 27. Lindqvist, Y., Schneider, G. Circular permutations of natural protein sequences: structural
743 evidence. *Curr. Opin. Struct. Biol.* **7**: 422–427 (1997).
- 744 28. Dundas, J., Binkowski, T.A., DasGupta, B., Liang, J. Topology independent protein structural
745 alignment. *BMC Bioinformatics* **8**(388) doi:10.1186/1471-2105-8-388 (2007).
- 746 29. Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H.,
747 Shindyalov, I.N., Bourne, P.E. The protein data bank. *Nucleic Acids Res.* **28**:
748 235–242 (2000).
- 749 30. Umeyama, S. Least-squares estimation of transformation parameters between two point
750 patterns. *IEEE Trans. Pattern Anal. Mach. Intell.* **13**(4): 376–380 (1991).
- 751 31. Veeramalai, M., Gilbert, D. A novel method for comparing topological models of protein
752 structures enhanced with ligand information. *Bioinformatics* **24**(23): 2698–2705 (2008).
- 753 32. Aghili, S.A., Agrawal, D., El Abbadi, A. PADS: protein structure alignment using directional
754 shape signatures. In DASFFA (2004).
- 755 33. Szustakowski, J.D., Weng, Z. Protein structure alignment using a genetic algorithm. *Proteins:*
756 *Struct. Funct. Genet.* **38**: 428–440 (2000).
- 757 34. Standley, D.M., Toh, H., Nakamura, H. Detecting local structural similarity in proteins by
758 maximizing number of equivalent residues. *Proteins: Struct. Funct. Genet.* **57**: 381–391
759 (2004).
- 760 35. Roach, J., Sharma, S., Kapustina, M., Cater Jr., C.W. Structure alignment via delaunay
761 tetrahedralization. *Proteins: Struct. Funct. Genet.* **60**: 66–81 (2005).
- 762 36. Teyra, J., Paszkowski-Rogacz, M., Anders, G., Pisabarro, M.T. SCOWLP classifica-
763 tion: structural comparison and analysis of protein binding regions. *BMC Bioinformatics*
764 doi:10.1186/1471-2105-9-9 (2008).
- 765 37. Gold, N.D., Jackson, R.M. Fold independent structural comparisons of protein-ligand binding
766 sites for exploring functional relationships. *J. Mol. Biol.* **355**: 1112–1124 (2006).
38. Zhu, J., Weng, Z. A novel protein structure alignment algorithm. *Proteins: Struct. Funct.*
Bioinform. **58**: 618–627 (2005).
- AQ5** 39. Bar-Yehuda, R., Halldorsson, M.M., Naor, J., Shacknai, H., Shapira, I. Scheduling split
767 intervals. 14th ACM-SIAM Symposium on Discrete Algorithms, pp. 732–741 (2002).

- 766 40. Hobohm, U., Sander, C. Enlarged representative set of protein structures. *Protein Sci.* **33**: 522
767 (1994).
- 768 41. Dutta, S., Akey, I.V., Dingwall, C., Hartman, K.L., Laue, T., Nolte, R.T., Head, J.F.,
769 Akey, C.W. The crystal structure of nucleoplasmin-core implication for histone binding and
770 nucleosome assembly. *Mol. Cell* **8**: 841–853 (2001).
- 771 42. Woo, E.J., Marshall, J., Bauly, J., Chen, J.G., Venis, M., Napier, R.M., Pickersgill, R.W.
772 Crystal structure of the auxin-binding protein 1 in complex with auxin. *EMBO J.* **21**:
773 2877–2885 (2002).
- 774 43. Uliel, S., Fliess, A., Amir, A., Unger, R. A simple algorithm for detecting circular permuta-
775 tions in proteins. *Bioinformatics* **15**(11): 930–936 (1999).
- 776 44. Alexandrov, N.N., Fischer, D. Analysis of topological and nontopological structural similar-
777 ities in the PDB: new examples with old structures. *Proteins* **25**: 354–365 (1996).
- 778 45. Dror, O., Benyamini, H., Nussinov, R., Wolfson, H.J. MASS: multiple structural alignment
779 by secondary structures. *Bioinformatics* **19**: i95–i104 (2003).
- 780 46. Shih, E.S., Hwang, M.J. Alternative alignments from comparison of protein structures.
781 *Proteins* **56**: 519–527 (2004).
- 782 47. Ilyin, V.A., Abyzov, A., Leslin, C.M. Structural alignment of proteins by a novel TOPOFIT
783 method, as a superimposition of common volumes at a topomax point. *Protein Sci.* **13**: 1865–
784 1874 (2004).
- 785 48. Tabtiang, R.K., Cezairliyan, B.O., Grant, R.A., Cochran, J.C., Sauer, R.T. Consolidating critical
786 binding determinants by noncyclic rearrangement of protein secondary structure. *PNAS*
787 **7**: 2305–2309 (2004).
- 788 49. Warren, A.J., Bravo, J., Williams, R.L., Rabbitts, T.H. Structural basis for the heterodimeric
789 interaction between the acute leukemia-associated transcription factors AML1 and CBFbeta.
790 *EMBO J.* **19**: 3004–3015 (2000).
- 791 50. Meining, W., Eberhardt, S., Bacher, A., Ladenstein, R. The structure of the N-terminal domain
792 of riboflavin synthase in complex with riboflavin at 2.6Å resolution. *J. Mol. Biol.* **331**: 1053–
793 1063 (2003).
- 794 51. Lichtarge, O., Bourne, H.R., Cohen, F.E. An evolutionary trace method defines binding
795 surfaces common to protein families. *J. Mol. Biol.* **7**: 39–46 (1994).
- 796 52. Norel, R., Fischer, H., Wolfson, H., Nussinov, R. Molecular surface recognition by computer
797 vision-based technique. *Protein Eng.* **7**(1): 39–46 (1994).
- 798 53. Fischer, D., Norel, R., Wolfson, H., Nussinov, R. Surface motifs by a computer vision-
799 technique: searches, detection, and implications for protein-ligand recognition. *Proteins* **16**:
800 278–292 (1993).
- 801 54. Meng, E., Polacco, B., Babbitt, P. Superfamily active site templates. *Proteins* **55**: 962–967
802 (2004).
- 803 55. Orengo, C., Todd, A., Thornton, J. From protein structure to function. *Curr. Opin. Struct. Biol.*
804 **9**: 374–382 (1999).
- 805 56. Binkowski, A., Adamian, L., Liang, J. Inferring functional relationship of proteins from local
806 sequence and spatial surface patterns. *J. Mol. Biol.* **332**: 505–526 (2003).
- 807 57. Jeffery, C. Molecular mechanisms for multi-tasking: recent crystal structures of moon-lighting
808 proteins. *Curr. Opin. Struct. Biol.* **14**: 663–668 (2004).
- 809 58. Tseng, Y.Y., Liang, J. Estimation of amino acid residue substitution rates at local spatial
810 regions and application in protein function inference: a Bayesian Monte Carlo approach. *Mol.*
Biol. Evol. **23**: 421–436 (2006).
59. Liang, J., Edelsbrunner, H., Woodward, C. Anatomy of protein pockets and cavities: measure-
ment of binding site geometry and implications for ligand design. *Protein Sci.* **7**: 1884–1897
(1998).
60. Edelsbrunner, H., Facello, M., Liang, J. On the definition and the construction of pockets in
macromolecules. *Disc. Appl. Math.* **88**: 83–102 (1998).
61. Dundas, J., Ouyang, Z., Tseng, J., Binkowski, A., Turpaz, Y., Liang, J. CASTp: computed atlas
of surface topography of proteins with structural and topographical mapping of functionally
annotated residues. *Nucleic Acids Res.* **34**: W116–W118.

Comparison of Protein Global Backbone Structures and Local Binding Surfaces

- 811 62. Lee, S., Li, B., La, D., Fang, Y., Ramani, K., Rustamov, R., Kihara, D. Fast protein tertiary
812 structure retrieval based on global surface shape similarity. *Proteins* **72**: 1259–1273 (2008).
- 813 63. Binkowski, T.A., Joachimiak, A. Protein functional surfaces: global shape matching and local
814 spatial alignments of ligand binding sites. *BMC Struct. Biol.* **8**: 45 (2008).
- 815 64. Bandyopadhyay, D., Huan, J., Liu, J., Prins, J., Snoeyink, J., Wang, W., Tropsha, A. Functional
816 Neighbors: Inferring relationships between non-homologous protein families using family-
817 specific packing motifs. *Proc. IEEE Int. Conf. Bioinform. Biomed.* (2008).
- AQ7 818 65. Moll, M., Kavvaki, L.E. A flexible and extensible method for matching structural motifs. *Nat.*
819 *Proc.* (2008).
- AQ8 820 66. Chen, L., Wu, L.Y., Wang, R., Wang, Y., Zhang, S., Zhang, X.S. Comparison of protein
821 structures by multi-objective optimization. *Genome Inform.* **16**(2): 114–124 (2005).
- AQ9 822 67. Dundas, J., Adamian, L., Liang, J. Signatures and basis sets of enzyme binding surfaces by
823 sequence order independent surface alignment. Manuscript submitted for publication (2010).
- 824 68. Corment, T.H., Leiserson, C.E., Rivest, R.L., Stein, C. *Introduction to algorithms*, 2nd edn.
825 Cambridge, MA: MIT Press (2001).
- 826 69. Kuhn, H.W. The hungarian method for the assignment problem. *Nav. Res. Logist. Q.* **2**: 83–97
827 (1995).
- 828 70. Wang, L., Wang, J.G., Chen, L. Align protein surface structures to identify evolutionarily
829 and structurally conserved residues. *The Third International Symposium on Optimization and*
830 *Systems Biology*, pp. 296–303 (2009).
- AQ10 831 71. Bellman, R. On a routing problem. *Q. Apply Math.* **16**(1): 87–90 (1958).
- 832 72. Yang, Z., Nielsen, R., Hasegawa, M. Models of amino acid substitution and applications to
833 mitochondrial protein structures. *Mol. Biol. Evol.* **15**: 1600–1611 (1998).
- 834 73. Huelsenbeck, J.B., Ronquist, R., Nielsen, R., Bollback, J. Bayesian inference of phylogeny
835 and its impact on evolutionary biology. *Science* **294**: 2310–2314 (2001).
- 836 74. Liang, J., Tseng, Y.Y., Dundas J., Binkowski, A., Joachimiak, A., Ouyang, Z., Adamian, L.
837 [Chapter 4](#): predicting and characterizing protein functions through matching geometric and
838 evolutionary patterns of binding surfaces. *Adv. Protein Chem.* **75**: 107–141 (2008).
- 839 75. Liang, J., Kachalo, S., Li, X., Ouyang, Z., Tseng, Y.Y., Zhang, J. Geometric structures of pro-
840 teins for understanding folding, discriminating natives and predicting biochemical functions.
841 *The World is a Jigsaw*. van de Weygaert R. (ed.). Springer (2009).
- 842
- 843
- 844
- 845
- 846
- 847
- 848
- 849
- 850
- 851
- 852
- 853
- 854
- 855