

ADVANCES IN PROTEIN CHEMISTRY AND STRUCTURAL BIOLOGY

EDITED BY

FREDERIC M. RICHARDS

Department of Molecular Biophysics
and Biochemistry
Yale, University
New Haven, Connecticut

DAVID S. EISENBERG

Department of Chemistry and Biochemistry
Center for Genomics and Proteomics
University of California, Los Angeles
Los Angeles, California

JOHN KURIYAN

Department of Molecular and Cellular Biology
University of California, Berkeley
Berkeley, California

VOLUME 75

Structural Genomics, Part A

EDITED BY

ANDRZEJ JOACHIMIAK

Structural Biology Center and Midwest Center
for Structural Genomics, Biosciences Division,
Argonne National Laboratory, Illinois, USA



ELSEVIER

AMSTERDAM • BOSTON • HEIDELBERG • LONDON
NEW YORK • OXFORD • PARIS • SAN DIEGO
SAN FRANCISCO • SINGAPORE • SYDNEY • TOKYO

Academic Press is an imprint of Elsevier



1 **PREDICTING AND CHARACTERIZING PROTEIN FUNCTIONS** 1
2 **THROUGH MATCHING GEOMETRIC AND EVOLUTIONARY** 2
3 **PATTERNS OF PROTEIN-BINDING SURFACES** 3
4 4

5 **By JIE LIANG,* YAN-YUAN TSENG,† JOSEPH DUNDAS,*** 5
6 **ANDREW BINKOWSKI,‡ ANDRZEJ JOACHIMIAK,‡** 6
7 **ZHENG OUYANG,* AND LARISA ADAMIAN*** 7
8 8

9 ***Program of Bioinformatics, Department of Bioengineering,** 9
10 **University of Illinois at Chicago, Chicago, Illinois 60607;** 10
11 **†Department of Ecology and Evolution, University of Chicago,** 11
12 **Chicago, Illinois 60637;** 12
13 **‡Structural Biology Center and Midwest Center for Structural Genomics,** 13
14 **Biosciences Division, Argonne National Laboratory,** 14
15 **Argonne, Illinois 60439** 15

16	I. Introduction	102	16
17	II. Voids and Pockets in Protein Structures and Their Origins.....	103	17
18	III. Identifying Functional Surfaces of Proteins.....	106	18
19	IV. Matching Local Binding Surfaces	108	19
20	A. Comparison of Sequence Patterns of Surface Pockets and Voids	109	20
21	B. Comparison of Shapes of Surface Pockets and Voids	112	21
22	C. Statistical Significance	114	22
23	V. Uncovering Evolutionary Patterns of Local Binding Surfaces.....	115	23
24	A. Evolution Model	116	24
25	B. Estimating Model Parameters Q and Bayesian Monte Carlo.....	118	25
26	C. Deriving Scoring Matrices from Rate Matrix	119	26
27	D. Validity of the Evolutionary Model.....	120	27
28	E. Evolutionary Rates of Binding Surfaces and Other Surfaces		28
29	are Different	120	29
30	VI. Predicting Protein Function by Detecting Similar Biochemical		30
31	Binding Surfaces.....	120	31
32	VII. Adaptive Patterns of Spectral Tuning of Proteorhodopsin from		32
33	Metagenomics Projects	125	33
34	VIII. Generating Binding Site Negative Images for Drug Discovery	127	34
35	IX. Summary and Conclusion	130	35
36	References.....	131	36

37 **ABSTRACT** 37

38 Predicting protein functions from structures is an important and chal- 38
39 lenging task. Although proteins are often thought to be packed as tightly 39
40 as solids, closer examination based on geometric computation reveals that 40

1 they contain numerous voids and pockets. Most of them are of random 1
2 nature, but some are binding sites providing surfaces to interact with other 2
3 molecules. A promising approach for function inference is to infer func- 3
4 tions through discovery of similarity in local binding pockets, as proteins 4
5 binding to similar substrates/ligands and carrying out similar functions have 5
6 similar physical constraints for binding and reactions. In this chapter, we 6
7 describe computational methods to distinguish those surface pockets that 7
8 are likely to be involved in important biological functions, and methods to 8
9 identify key residues in these pockets. We further describe how to predict 9
10 protein functions at large scale (millions) from structures by detecting 10
11 binding surfaces similar in residue make-ups, shape, and orientation. We 11
12 also describe a Bayesian Monte Carlo method that can separate selection 12
13 pressure due to biological function from pressure due to protein folding. 13
14 We show how this method can be used to reconstruct the evolutionary 14
15 history of binding surfaces for detecting similar binding surfaces. In addi- 15
16 tion, we briefly discuss how the negative image of a binding pocket can be 16
17 casted, and how such information can be used to facilitate drug discovery. 17
18

19
20

21 I. INTRODUCTION 21

22 The structural genomics projects have made significant contributions to 22
23 our current body of knowledge of protein structures (Chandonia and 23
24 Brenner, 2006). They have further facilitated the establishment of a 24
25 comprehensive view of the global universe of protein structures, and 25
26 have provided a foundation with a wealth of information for developing 26
27 model and computational tools that can be used to understand the 27
28 molecular mechanism how individual proteins carry out their biological 28
29 roles and how protein functions evolve. 29
30

31 Functional characterization of proteins with unassigned functions is an 31
32 important task. By design, a large number of newly determined protein 32
33 structures from structural genomics are not related to other known pro- 33
34 teins, and bioinformatics tools based on sequence alignment often cannot 34
35 provide accurate information about the functional roles of these proteins. 35
36 Several early studies showed that reliable functional assignment will re- 36
37 quire sequence identity of 60–70% between the protein of unknown 37
38 function and a well-studied protein (Rost, 2002; Tian and Skolnick, 2003). 38
39

40 Recently, the approach of inferring protein functions by detecting local 39
40 spatial regions on protein structures with similar patterns has been shown 40

1 to be very effective (Binkowski *et al.*, 2003a; Glaser *et al.*, 2003; Gold and 1
2 Jackson, 2006; Laskowski *et al.*, 2005; Najmanovich *et al.*, 2005; Pazos and 2
3 Sternberg, 2004; Russell, 1998; Torrance *et al.*, 2005; Tseng and Liang, 3
4 2006). The rationale behind this approach is intuitive and appealing. 4
5 For proteins binding to similar substrates or ligands and carrying out 5
6 similar functions, they are constrained by the requirement of providing 6
7 the necessary microenvironment for similar binding and biochemical 7
8 reactions to occur. These physical constraints are reflected by similarity 8
9 in the shape of local binding surfaces and in the physicochemical texture 9
10 of the binding surfaces. In order for similar functions to occur, the 10
11 evolution of residues involved in binding and reaction will be constrained 11
12 and this results in similarly allowed and forbidden residue substitution on 12
13 binding surfaces (Tseng and Liang, 2006). 13
14

15 In this chapter, we discuss our approach to predict and characterize 15
16 protein functions from protein structures by comparing local surfaces. We 16
17 first discuss the existence of voids and pockets, and their distribution in 17
18 proteins (Liang and Dill, 2001). We then describe how to identify those 18
19 that are likely to be functionally important, as well as the key residues on 19
20 them (Tseng and Liang, 2007). This is followed by a discussion on how to 20
21 match local surfaces and how to assess their similarity in both sequence 21
22 order-dependent and -independent fashion (Binkowski *et al.*, 2003a). 22
23 Next we discuss how to extract evolution patterns of small local regions 23
24 directly related to protein function and unaffected by folding requirement 24
25 using a Bayesian Monte Carlo method, and how this approach improves 25
26 protein function prediction (Tseng and Liang, 2006). We then describe 26
27 three examples of protein function prediction and characterizations using 27
28 proteins generated from the Midwest Center for Structural Genomics 28
29 (Binkowski *et al.*, 2005). This is followed by a brief discussion on how 29
30 further information from computed protein local binding pockets can be 30
31 extracted in the form of negative image to guide for selecting inhibitors 31
32 from a collection of candidate compounds (Ebalunode *et al.*, 2008). 32
33

34 35 II. VOIDS AND POCKETS IN PROTEIN STRUCTURES AND THEIR ORIGINS 35 36

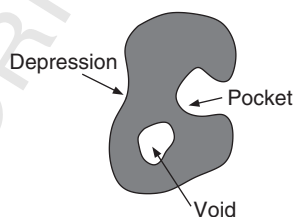
37 Protein structure is known to be packed tightly. The packing density of 37
38 protein interior is comparable to that of solid, with low compressibility 38
39 (Gavish *et al.*, 1983). Protein packing has been described as a jig-saw puzzle 39
40 (Richards and Lim, 1994). However, detailed study using the technique of 40

1 alpha shape (Edelsbrunner and Mücke, 1994; Edelsbrunner *et al.*, 1998; 1
2 Liang *et al.*, 1998a,b) revealed that there are numerous voids and pockets 2
3 in protein structures (Fig. 1) (Liang and Dill, 2001). 3

4 Here, voids are enclosed empty space that is inaccessible to a water 4
5 molecule modeled as a probe of 1.4 Å radius, and pocket is an empty space 5
6 in the protein that has a constricted opening to the bulk exterior and is 6
7 accessible to a water molecule (Fig. 1). The size of the void or pocket in 7
8 this study is required to be large enough to contain at least one water 8
9 molecule. In fact, there is a scaling relationship between the number of 9
10 voids and pocket and the chain length of the protein (Fig. 2A). On 10
11 average, there is an increase of 15 voids or pockets for every 100 amino 11
12 acid residues (Liang and Dill, 2001). For example, the binding sites of 12
13 HIV-1 protease and phosphatidylinositol transfer protein (PITP) both 13
14 correspond to well-defined surface pockets (Fig. 3). 14
15

16 Various scaling relationships suggest that protein packing is of random 16
17 nature (Liang and Dill, 2001). For example, if we use a simple solid ball 17
18 packing as a model of protein, we would expect that the volume 18
19 $V = 4\pi r^3/3$ and the area $A = 4\pi r^2$ should have a scaling relationship of 19
20 $V = A^{3/2}$. In reality, this scaling relationship is linear (Fig. 2B). This linear 20
21 relationship is reminiscent of the scaling relationship of clustered random 21
22 spheres in off-lattice and on-lattice models (Lorenz *et al.*, 1993; Stauffer, 22
23 1985). 23

24 To further investigate the nature of protein packing and the origin of 24
25 voids and pockets, we have studied the packing behavior of random chain 25
26 polymer in off-lattice three-dimensional space (Zhang *et al.*, 2003a). Other 26
27 than the requirement that these polymer chains are compact and self- 27
28



31
32
33
34
35
36
37
38
39
40
FIG. 1. Pockets and voids in proteins. There are three types of unfilled space on protein surfaces. *voids* are fully enclosed and have no outlet, *pockets* are accessible from the outside but with constriction at mouths, and shallow *depressions* have wide openings. We use the general term *surface pockets* to include both pockets and voids. Adapted from (Liang and Dill 2001).

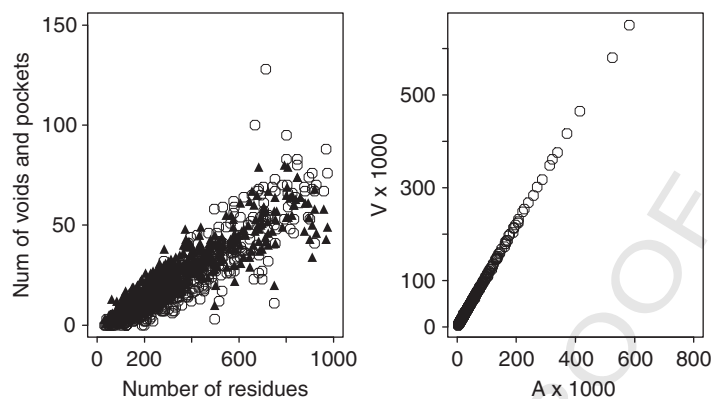


FIG. 2. Voids and pockets in protein structures. (A) Number of voids and pockets scale roughly linearly with protein length for a representative set of 636 proteins. Here, circles and solid triangles represent the numbers of voids and pockets, respectively. (B) The volume of protein as calculated using van der Waals model scales linearly with the van der Waals area of protein. Adapted from Liang and Dill (2001).

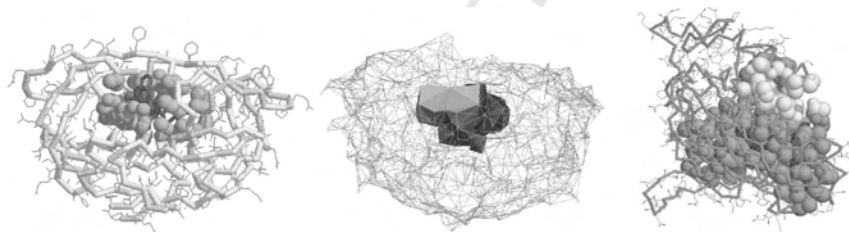


FIG. 3. The binding pockets on HIV-1 protease and phosphatidylinositol transfer protein (PITP). Left: binding pocket (yellow) on HIV-1 shown in van der Waals space filling model. Ligand is colored red. Middle: the alpha shape of the HIV-1-binding site. Its mouth opening is colored gold. Right: Binding pocket (green) on PITP for phospholipid (red) and a regulatory site on a different region (yellow) of the same protein.

avoiding, there is no relationship between these studied chains and real protein. The task of assessing the ensemble properties of packing of these chain polymers in a statistically accurate manner is technically very challenging, as one needs to generate adequate samples that are independent and properly weighted. This relates to the well-known attrition problem: the success rate of generating self-avoiding chain polymers is rapidly diminishing with the increase of chain length, as it becomes exponentially

1 difficult to maintain the self-avoiding requirement. For example, even for 1
2 a short chain of length 48, the success rate of using simple growth method 2
3 would be only 0.79% (Liu, 2001). 3

4 Using the sequential Monte Carlo method (Doucet *et al.*, 2001; Liu and 4
5 Chen, 1998), we have overcome this technical difficulty, and succeeded in 5
6 generating properly weighted ensemble of thousands of self-avoiding 6
7 chains up to length 2000 (Zhang *et al.*, 2003a). We have carried out 7
8 the same geometric analysis on these chain polymer structures, just as we 8
9 did with protein structures. The results indicate that both the scaling 9
10 relationship of the coordination number, and the packing density with 10
11 the chain length show characteristically the same scaling relationship as 11
12 that of proteins (Zhang *et al.*, 2003a). Altogether, these findings provide 12
13 strong evidence that proteins are not optimized by evolution to eliminate 13
14 voids and pockets. Rather, the majority of the voids and pockets simply 14
15 emerge from the requirement of packing self-avoiding chains in a compact 15
16 space. 16
17

18 19 20 III. IDENTIFYING FUNCTIONAL SURFACES OF PROTEINS 20

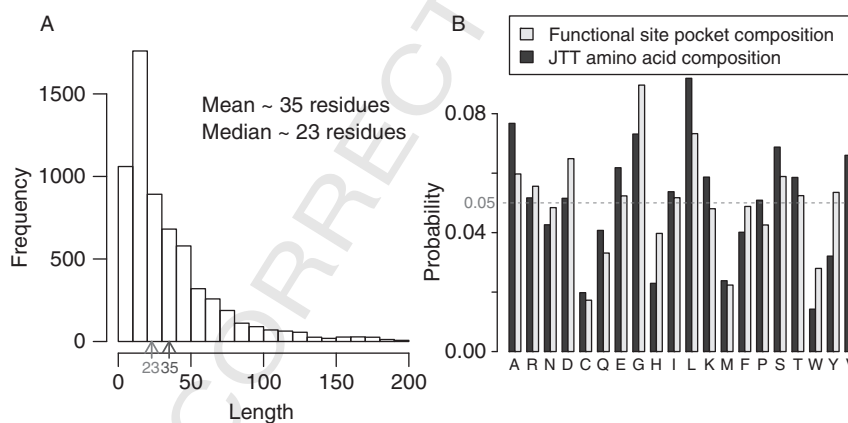
21 The existence of numerous voids and pockets poses two challenging 21
22 problems. First, how do we identify the void(s) and pocket(s) that are 22
23 biologically important, for example, how to distinguish those involved in 23
24 binding and biochemical reactions from those formed by random chance. 24
25 Second, for a given pocket or voids found on a protein structure, how do 25
26 we know if it is important for some biological functions known or yet to be 26
27 discovered? 27

28 We have developed a method to address these problems for enzymes. 28
29 In this method, we do not directly compare the structure or function of a 29
30 well-characterized protein with the protein in question. Rather, we seek to 30
31 recognize pocket or void that might be involved in enzyme function based 31
32 on general characteristics. We discuss in later sections the comparative 32
33 approach when the unknown query protein is compared with a database of 33
34 protein structures. 34
35

36 Typically, about 10–30% of all residues in an enzyme participate in the 36
37 formation of the binding pocket (Tseng and Liang, 2007). Compared to 37
38 the full length primary sequences, the usage of residues in forming pocket 38
39 is biased. Often His, Asp, Glu, Ser, and Cys account for the most important 39
40 active site residues (Bartlett *et al.*, 2002; Binkowski *et al.*, 2003a; Laskowski 40

1 *et al.*, 2005; Tseng and Liang, 2007). These are residues known to be 1
2 important for catalytic functions. On the other hand, nonpolar residues 2
3 such as Val, Leu, Pro are far less frequent in enzyme-binding pocket 3
4 (Tseng and Liang, 2007). Although these hydrophobic residues are fre- 4
5 quently conserved for maintaining protein structures and for protein 5
6 folding, they are often not directly involved in molecular functions of 6
7 enzymes. In fact, the composition of residue on binding surfaces of 7
8 enzyme is very different from that of the overall sequences (Fig. 4). 8
9

10 In our method for identifying functional region from enzyme structures 10
11 (Tseng and Liang, 2007), we examine the occurrence of the *atomic pattern* 11
12 of a residue with exposed surface in the binding pocket. That is, we record 12
13 the residue type and all of the exposed atoms from this residue, along with 13
14 the secondary structure environment this residue belongs to. A probability 14
15 function for each atom pattern, residue type, and secondary structure is 15
16 then constructed based on statistical analysis of a database of annotated 16
17 key residues of enzymes. After evaluating this probability function for each 17
18 residue in a candidate pocket, we can sum up the probability values for all 18
19 residues in the identified pocket, and if it is above a threshold value, a 19
20 functional binding pocket is predicted, and the few residues with the 20
21



22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
FIG. 4. The length distribution and residue composition of functional surfaces for 3275 enzyme proteins containing known functional key residues. (A) Functional surfaces usually consist of 8–200 residues, with the mean at 35 residues. (B) The amino acid residue composition of functional surfaces is different from the composition of sequences used to construct the Jones–Taylor–Thornton (JTT) model. Adapted from Tseng and Liang (2007).



18 FIG. 5. The binding surface (green) and key residues predicted from a structure of 18
19 alpha amylase. Here, the predicted four key residues are colored yellow (D176), cyan 19
20 (H180), pink (N208), and blue (D269). They contain several high propensity atomic 20
21 patterns from our library of 1031 functional atomic patterns. Their classes of secondary 21
22 structural environment (sheet *s*, helix *h*, and coil *c*) are also listed. The substrate 22
23 molecule is colored red. Adapted from Tseng and Liang (2007). 23

24 highest probability values are further predicted to be functionally impor- 24
25 tant key residues. 25

26 This method has been shown to work well in a 10-fold crossvalidation 26
27 test of 3503 protein surfaces from 70 proteins, with a sensitivity of 92.9% 27
28 and specificity of 99.88% (Tseng and Liang, 2007). We have also shown 28
29 that for four enzyme families (2,3-dihydroxybiphenyl dioxygenase, E.C. 29
30 1.13.11.39; adenosine deaminase, E.C. 3.5.4.4; 2-haloacid dehalogenase, 30
31 E.C. 3.8.1.2; and phosphopyruvate hydratase, E.C. 4.2.1.11), the key resi- 31
32 dues predicted are also consistent with annotated information contained 32
33 in the Structure–Function Linkage Database (SFLD) (Pegg *et al.*, 2006). 33
34 Figure 5 illustrates the example of predicted binding surface and key 34
35 residue on a structure of alpha amylase. 35
36

37 IV. MATCHING LOCAL BINDING SURFACES 38

39 A different approach that can potentially yield rich information is to 39
40 compare the local surface of a binding pocket to a database of local 40
surfaces, some of which have known biological characterization. Figure 6

1 illustrates an example. The cAMP-dependent protein kinase (1cdk) and
2 Tyr protein kinase c-src (pdb 2src) share only 13% sequence identity. 2
3 However, the ATP-binding pockets have similar shape and chemical texture. 3
4 Once these ATP-binding pockets are identified and computed from 4
5 their structures, we can select the residues located on the wall of the 5
6 binding pocket, and remove residues on the loops connecting these wall 6
7 residues. It is clear that the remaining sequence fragments have much 7
8 higher sequence identity (51%). In both cases, the residues forming the 8
9 pocket wall come from diverse regions in the primary sequences. 9

10 The simple example shown in Fig. 6 suggests an effective strategy that 10
11 can rapidly decide if two pocket surface are similar. We can derive surface 11
12 patterns from the residues forming the walls of pockets (called pvSOAR 12
13 patterns for pocket and void surface patterns of amino acid residues), and 13
14 rapidly compare these patterns. Once a pair of protein surfaces are found 14
15 to be similar, we can further examine their shape and chemical texture in 15
16 detail, and determine the statistical significance of their overall similarity. 16
17 This approach is generally applicable to any two surface patterns of 17
18 pockets and voids (Binkowski *et al.*, 2003a). 18

19 There are several technical problems to be solved for this approach to 19
20 be generally useful. We need to identify and generate local surfaces 20
21 automatically and accurately. This can be achieved by applying void and 21
22 pocket algorithm for exhaustive identification and measurement of voids 22
23 and pockets from protein structures (Edelsbrunner *et al.*, 1998; Liang 23
24 *et al.*, 1998a,b). We also need to rapidly and accurately assess surface 24
25 similarity. Once a pair of similar local surfaces are found, we need to 25
26 evaluate whether the similarity is statistical significant. 26
27

28 29 30 A. Comparison of Sequence Patterns of Surface Pockets and Voids 30

31 *Sequence order-dependent method.* By concatenating wall residues of a pocket 31
32 or void on a peptide chain, we have compiled a database of pvSOAR sequence 32
33 patterns for all protein structures in the protein data bank (PDB). This 33
34 database is part of the CASTp database (Binkowski *et al.*, 2003b; Dundas 34
35 *et al.*, 2006). It currently (August, 2008) contains 46,071 protein structures, 35
36 with 1,582,472 voids and 1,555,994 pockets. We can rapidly query a protein 36
37 surface pocket against CASTp database through alignment of sequence 37
38 fragments using standard dynamic programming technique, allowing gap 38
39 insertion (Binkowski *et al.*, 2003a). In this approach, we assume that the 39
40



FIG. 6. Functional surfaces on the catalytic domains of cAMP-dependent protein kinase (1cdk) and tyrosine protein kinase (2src). (A) In both cases, the active sites are computed as surface pockets. (B) Residues defining the pockets are well dispersed throughout the primary sequences (full sequence identity = 16%). (C) The identity of their surface sequence patterns is much higher (51%).

residues in the sequence pattern are positioned following their order of the primary sequence.

Sequence order-independent comparison. The alignment of pvSOAR sequence fragments through dynamic programming can discover many similar binding pockets. However, there are many cases where two proteins with similar placement of amino acids in their tertiary structures have different relative positioning of these amino acids in their primary structures (see Fig. 7 for stromelysin). When comparing two local surface pockets, we also need to detect similar residue patterns while ignoring

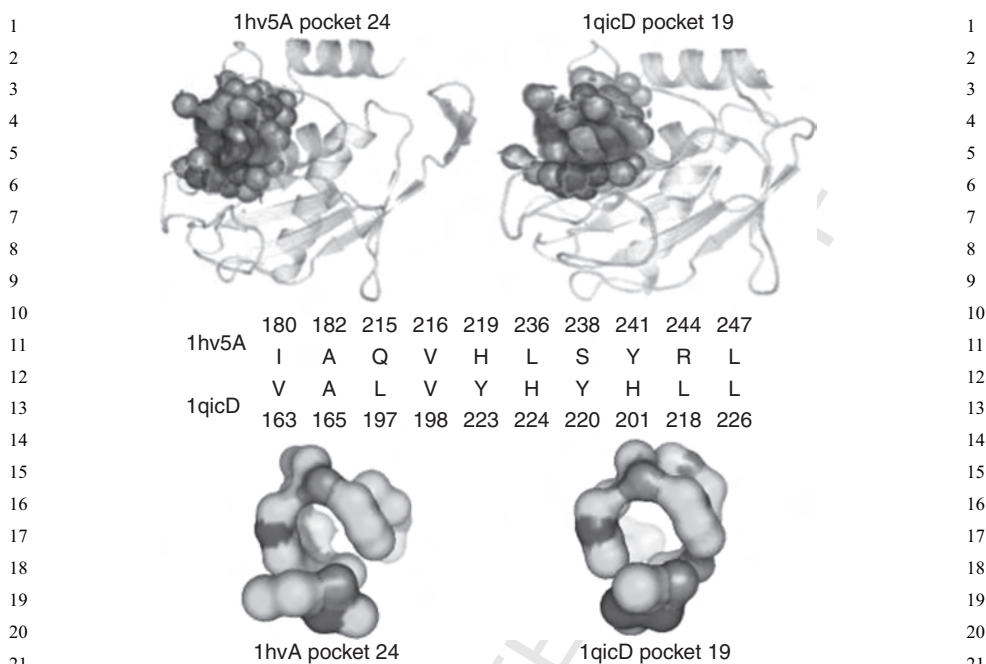


FIG. 7. The binding pockets from two different stromelysin catalytic domains (pocket 29 from pdb 1hv5.A and pocket 19 from 1qic.D). They are aligned in a sequence order-independent fashion with a cRMSD of 0.76 Å for 29 atoms from 10 residues. Top: the binding pockets on the two protein structures, with pocket atoms shown in space filling form. The aligned atoms are colored in red. Middle: the alignment of residues of these two surface pockets. Atomic details of the alignment are not shown. Sequence numbers are listed above and below the residue names for 1hv5 and 1qic, respectively. Residues in 1hv5 are arranged in order, but it is clear that the aligned residues in 1qic are not in sequence order. This residue alignment is derived from detailed alignment of atoms from surface pockets. Bottom: aligned atoms from these two surface pockets, with N atoms in blue, O in red, and C in green.

their strict positioning in the primary structures. This is the problem of finding which amino acid on the query protein surface pocket is equivalent to which amino acid on the target protein surface pocket.

Sequence order-independent matching of pockets can be formulated as a maximum weight bipartite matching problem, where graph nodes represent amino acids (e.g., using C_{α} atoms) from the two protein pockets. Directed edges are used to connect nodes from the query protein to nodes of the target protein, if the two nodes share some similarity (e.g., by a

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

1 scoring function based on shape and chemistry). Each edge is given a 1
2 weight that is based on the similarity measure. The problem is to find a set 2
3 of edges connecting nodes in query pocket to nodes in target pocket, with 3
4 maximized total edge weight, while insisting only at most one edge is 4
5 selected for each residue (Cormen *et al.*, 2001). 5

6 One way to solve this problem is by using the Hungarian algorithm 6
7 (Kuhn, 1955) as described in (Chen *et al.*, 2005) with modifications. This is 7
8 an iterative method that uses the Bellman–Ford algorithm (Bellman, 8
9 1958). First, we add a fictitious source node s that connects to every 9
10 query node with 0-weight. We then add a fictitious destination node d 10
11 that connects to every target node with 0-weight. The Bellman–Ford 11
12 algorithm computes the distance $F(i)$ of the shortest path(s) from the 12
13 source node to each of the remaining node i . The weight for each edge 13
14 that does not contain the source node is then updated. The new weight 14
15 $w'(i, j)$ for edge $e(i, j)$ starting from node i to node j is 15
16

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

17
18
19 An overall score F_{all} , initialized to 0, is now updated as $F'_{\text{all}} = F_{\text{all}} - F(d)$. 19
20 Next, we flip the directions of all edges in the shortest path from the 20
21 source s to the destination d . 21

22 We then apply the Bellman–Ford algorithm on this new graph, and this 22
23 is repeated until either there is no directed path from s to d as edges have 23
24 been flipped, or the shortest distance $F(d)$ to the destination is greater 24
25 than the current overall score F_{all} . The output of the Hungarian algorithm 25
26 includes a set of directed edges starting from target nodes to query nodes, 26
27 and these provide the equivalence relationship, namely, which residue in 27
28 the target pocket should be aligned to which residue in the query pocket. 28
29 Based on this equivalence relationship, we can then compute the shape 29
30 similarity between these two surface pockets at atomic details, as described 30
31 below. When we use atoms as nodes instead of residues, the results will be 31
32 atomic alignment of pocket surfaces. 32
33

34 35 B. Comparison of Shapes of Surface Pockets and Voids 35 36

37 Once two voids or pockets are found to have significant sequence 37
38 similarity, we then follow up with more detailed shape analysis using two 38
39 methods. First, we compute the coordinate root mean square distance 39
40 (cRMSD) between the subset of equivalent residues or atoms. This 40

1 equivalence relationship is established by the local alignment of pocket 1
2 sequence fragments. The cRMSD distance is measured when the subset of 2
3 residues are optimally aligned with rigid motion and has the least RMSD 3
4 value. This alignment and the cRMSD value can be computed from the 4
5 singular value decomposition of the correlation matrix of the coordinates 5
6 of the point sets (Umeyama, 1991). 6

7 cRMSD is not a perfect measure of shape similarity. It works well when 7
8 two structures are similar, but is sensitive to outliers. If a protein experi- 8
9 ences conformational change, its binding pocket may expand or shrink 9
10 and its residues may retain the relative orientational relationship, but with 10
11 significantly altered Euclidean distances. To address this deficiency, we 11
12 can use the orientational RMSD (oRMSD) measure (Binkowski *et al.*, 12
13 2003a). We first place a unit sphere at the geometric center of the pocket. 13
14 The location of each residue is then projected onto the unit sphere along 14
15 the direction of the vector from the geometric center. The projected 15
16 pocket is therefore represented by a set of unit vectors on the unit sphere, 16
17 which preserves the original orientational relationship. The RMSD of the 17
18 two sets of unit vectors for the two pockets in comparison can then be 18
19 measured, which gives the oRMSD value (Binkowski *et al.*, 2003a). 19
20

21 For sequence order-independent comparison of two surface pockets, we 21
22 start from a crude initial equivalence relationship that represents the 22
23 initial correspondence between residues from query and target pockets. 23
24 We then apply the optimal rotation matrix and translation vector com- 24
25 puted using (Umeyama, 1991) to this initial alignment. The Euclidean 25
26 distances between residues (or atoms) in the query pocket and target 26
27 pocket are then computed after the optimal superposition. Those that 27
28 are below a threshold are updated with new weights computed using a 28
29 similarity scoring function. The Bellman–Ford algorithm and the SVD- 29
30 based optimal alignment and update of Euclidean distances are then 30
31 repeated iteratively. One can stop this iterative process if the improvement 31
32 is less than a threshold. As the overall alignment shape score may deterio- 32
33 rate temporarily when a new equivalence relationship is found and new 33
34 superposition applied, simulated annealing allowing a probability that 34
35 structural alignment may temporarily deteriorate can also be applied 35
36 here (Chen *et al.*, 2005). 36
37

38 As an illustration, the sequence order-independent alignment of surface 38
39 pockets in two structures of stromelysin shown in Fig. 7. It has an overall 39
40 cRMSD of 0.76 Å for 29 atoms from 10 residues. The C_α atoms from these 10 40

1 residues align with a cRMSD of 1.05 Å. The alignment obtained in a 1
2 sequence order-dependent fashion contains 16 residues. If we select the 2
3 subset of 10 residues from these 16 residues that overlap most with that of 3
4 the sequence order-independent alignment, the alignment of their C_α 4
5 atoms has a cRMSD value of 3.71 Å. This example illustrates that this 5
6 method of sequence order-independent comparison of two surface pockets 6
7 works well, and often can identify excellent surface matches that are chal- 7
8 lenging for other methods (J. Dundas and J. Liang, unpublished results). 8
9

10 11 12 13 *C. Statistical Significance* 13

14
15 After the similarity of two surface pockets is calculated, we need to assess 15
16 its statistical significance to aid in biological interpretation. pvSOAR 16
17 sequence patterns are typically short, and are of different composition 17
18 from the full chain sequences. In addition, frequently the two pocket 18
19 sequence patterns in comparison have different number of residues. 19
20 Although the theoretical model of extreme value distribution (EVD) pro- 20
21 vides accurate description of gapless local alignment of random sequences 21
22 (Karlin and Altschul, 1990), no exact theoretical models are known in 22
23 general for local sequence alignment of very short sequences with gaps. 23

24 We have developed a heuristic approach to assess the statistical signifi- 24
25 cance of two pocket pvSOAR sequences aligned in sequence order. By 25
26 removing the largest peak in the low-score region of the distribution of 26
27 alignment scores of random short sequences which often contain just one 27
28 or two matched residues, we found that the remaining distribution can be 28
29 described by an EVD well (Binkowski *et al.*, 2003a). Specifically, the Smith- 29
30 Waterman scores of the search results of a query sequence pvSOAR 30
31 pattern to a database of randomly shuffled pocket sequences are collected. 31
32 They are then fitted to an EVD distribution, and the goodness of fit is then 32
33 evaluated using the Kolmogorov–Smirnov test (Pearson, 1991). If the 33
34 observed Kolmogorov–Smirnov statistic does not indicate that the random 34
35 scores are inconsistent with an EVD distribution, we further estimate the 35
36 statistical significance *p*-value using the calculated *z*-score $z = (S - \mu)/\sigma$, 36
37 where *S* is the similarity score, μ is the mean of random scores, and σ is the 37
38 standard deviation. The *p*-value can be estimated from the *z*-score as 38
39 (Binkowski *et al.*, 2003a) 39
40

$$p(Z > z) = 1 - \exp(-e^{-1.282z-0.5772}).$$

The expected number E of random pocket sequences with the same or better score can be calculated as

$$E = p \times N_r,$$

where N_r is the number of randomly shuffled sequence fragments. The p -value or E -value can be used to exclude matched pairs of pocket pvSOAR sequences that are unlikely to be biologically relevant.

Once the cRMSD or oRMSD value is calculated for two surface pockets, we also need to evaluate the statistical significance of shape comparison. As illustrated above, a common practice for determining statistical significance is to assume the similarity score are drawn randomly from a specific underlying distribution. The parameters of the assumed distribution are then estimated by curve-fitting the distribution of scores from the random comparison of protein pockets. The derived parameters can then be used to find the Z -score or p -value of a given similarity score (Jia *et al.*, 2004; Levitt and Gerstein, 1998; Ye and Godzik, 2004; Zhu and Weng, 2005). We found that the distribution of both cRMSD and oRMSD for random surfaces on protein structures do not follow known parametric model such as the EVD (Binkowski *et al.*, 2003a). We empirically estimate the probability p of obtaining a specific cRMSD or oRMSD value for n number of matched positions from a set of randomly generated surface pockets and voids. By collecting cRMSD and oRMSD values of millions of randomly matched pockets with different number of selected matched residues, we can estimate the p -value of a specific cRMSD or oRMSD with a specific number of matched residues. This can be found by finding the closest value of the rank order statistic in the randomly collected cRMSD or oRMSD data of the same number of residues (Binkowski *et al.*, 2003a; Russell, 1998).

V. UNCOVERING EVOLUTIONARY PATTERNS OF LOCAL BINDING SURFACES

Fast comparison of pvSOAR sequence fragments is a key step when querying a specific surface pocket/void against a database of precomputed pocket/voids, as the database can contain hundreds of thousands or millions of entries. This is possible by applying fast dynamic programming method to align the sequence fragments representing the two pockets/

1 voids. This step is carried out before promising hits are identified and 1
2 further detailed shape comparison is carried out. 2

3 The specific scoring matrix used to assess the similarity of two aligned 3
4 pocket/void sequence fragments is critical for detecting functionally 4
5 related binding pockets/voids. A convenient choice is to adopt widely 5
6 used PAM matrices or BLOSUM matrices (Dayhoff *et al.*, 1978; Henikoff 6
7 and Henikoff, 1992). A disadvantage of this approach is that these are 7
8 precomputed matrices and have implicit parameters with values prede- 8
9 termined from the analysis of large quantities of sequences, which 9
10 contain little information of the protein of interest. Another approach 10
11 is to use position-specific scoring matrix (PSSM) such as those gener- 11
12 ated by the PSI-BLAST program (Altschul *et al.*, 1997). The drawback of 12
13 this latter approach is that it often leads to serious bias as the PSSM is 13
14 derived from all sequences aligned to the query sequence satisfying 14
15 certain statistical significance requirement. Bias comes from the fact 15
16 that all aligned sequences contribute equally to the derivation of 16
17 PSSM, regardless how closely or distantly they are related. This is 17
18 particularly problematic if the query result from the database is domi- 18
19 nated by closely related proteins. 19
20

21 22 23 *A. Evolution Model* 23

24 To resolve these issues, we have adopted an approach that models the 24
25 evolutionary process using a continuous time Markov process and an 25
26 explicit phylogenetic tree (Tseng and Liang, 2006). Markovian evolution- 26
27 ary models are parametric models and do not have prespecified parameter 27
28 values. These values are instead estimated from specific sequence data 28
29 relevant to the protein of interests (Whelan *et al.*, 2001). This approach 29
30 has been shown to be more effective in deriving informative rate matrices 30
31 with significant advantage over matrices obtained from other methods 31
32 (Whelan *et al.*, 2001). 32
33

34 We assume that a reasonably accurate phylogenetic tree T , the branch 34
35 lengths of the tree representing divergence time, and an accurate multiple 35
36 sequence alignment are known. These can be computed using maximum 36
37 likelihood method or Bayesian method (Adachi and Hasegawa, 1996; 37
38 Huelsenbeck *et al.*, 2001; Yang, 1997). The subset of columns in the 38
39 multiple sequence alignment corresponding to the residues in the bind- 39
40 ing pocket are then identified based on pocket calculation (Binkowski 40

1 *et al.*, 2003a; Liang *et al.*, 1998c; Tseng and Liang, 2006). Our model 1
 2 assumes that the evolution of the residues in the binding pocket can 2
 3 be modeled by a Markovian process characterized by a 20×20 matrix 3
 4 $\mathbf{Q} = \{q_{ij}\}$ of instantaneous substitution rates. The divergence time t is 4
 5 measured in the unit of the expected number of residue changes per 5
 6 100 sites between the sequences. 6

7 Once the instantaneous substitution rate matrix $\mathbf{Q} = \{q_{ij}\}$ is known, the 7
 8 matrix of probabilities of substitution of residue i by residue j in the time 8
 9 interval t can be computed as 9
 10

$$11 \quad P(t) = \{p_{ij}(t)\} = \exp(\mathbf{Q} \cdot t). \quad 11$$

12
 13 For symmetric \mathbf{Q} , the matrix exponential can be conveniently computed 13
 14 as 14

$$15 \quad \exp(\mathbf{Q} \cdot t) = \mathbf{U} \exp(\Lambda t) \mathbf{U}^{-1}, \quad 15$$

16
 17 where \mathbf{U} is the matrix of right eigenvectors of \mathbf{Q} , and \mathbf{U}^{-1} is that of the left 17
 18 eigenvectors. A technique to construct a more general nonsymmetric 18
 19 instantaneous rate matrix \mathbf{Q} that can be symmetrized can be found in 19
 20 Tseng and Liang (2006) and Whelan and Goldman (2001). 20
 21

22 For a column in the multiple sequence, we follow the phylogenetic tree 22
 23 \mathbf{T} and compute the transition probability $p_{x_i, x_j}(t_{ij})$ for each of the edge in 23
 24 the tree, whose length denotes the time interval $t_{i,j}$. Here, x_i and x_j are the 24
 25 residues at the positions corresponding to the nodes connected by 25
 26 the edge. If we knew all the ancestral sequences (corresponding to the 26
 27 internal nodes in the phylogenetic tree) of the extant sequences 27
 28 (corresponding to the leaf nodes), the likelihood given the tree \mathbf{T} and 28
 29 the instantaneous rates \mathbf{Q} for this column h can be obtained by combining 29
 30 probabilities along all edges: 30

$$31 \quad p(\mathbf{x}_h | \mathbf{T}, \mathbf{Q}) = \pi_{x_k} \prod p_{x_i, x_j}(t_{ij}). \quad 31$$

32
 33 Here, the π_{x_k} is the prior probability of an arbitrarily chosen node k as 34
 35 the starting node taking its residue as type x_k at column h . π_{x_k} typically can 35
 36 be computed as the composition of the aligned sequences. The product 36
 37 sign \prod is over all edges in the phylogenetic tree. Since in reality we do not 37
 38 know the identities of the residues in ancestral sequences, we sum over all 38
 39 possible values the ancestral sequence might take in this column, and the 39
 40

1 probability $p(x_h|T, Q)$ of observing this particular column h in the multiple 1
 2 sequence alignment is 2

$$3 \quad p(x_h|T, Q) = \pi_{x_h} \sum \prod p_{x_i x_j}(t_{ij}). \quad 3$$

4
 5
 6 Here, the summation sign Σ is overall all possible residues in this 6
 7 column for each of the ancestral sequences. 7

8 Treating each column independently, the probability $P(S|T, Q)$ of ob- 8
 9 serving all residues in the selected columns for the functional region S is 9

$$10 \quad P(S|T, Q) = P(x_1, \dots, x_s|T, Q) = \prod p(x_h|T, Q). \quad 10$$

11
 12 Here, the product Π sign is over all columns. 12
 13

14 *B. Estimating Model Parameters Q and Bayesian Monte Carlo* 14

15
 16 We adopt a Bayesian framework, and each model parameter is described 16
 17 with a distribution instead of a single value. The *posterior probability* 17
 18 $\pi(Q|S, T)$ of the rate matrix for a given aligned pocket region S and the 18
 19 phylogenetic tree T integrates our prior information (represented by the 19
 20 prior distribution $\pi(Q)$) on the model parameters, and the likelihood 20
 21 function-related probability $P(S|T, Q)$ derived from the observed data: 21
 22

$$23 \quad \pi(Q|S, T) \propto \int P(S|T, Q) \cdot \pi(Q) dQ. \quad 23$$

24
 25 Once this posterior distribution is known, we can calculate the posterior 24
 26 mean of the parameters: 25
 27

$$28 \quad E_\pi(Q) = \int Q \cdot \pi(Q|S, T) dQ. \quad 28$$

29
 30 In practice, we generate correlated samples from the posterior distribu- 29
 31 tion, and the posterior means of the model parameters are estimated from 30
 32 these samples: 31
 33

$$34 \quad E_\pi(Q) \approx \sum Q_i \cdot \pi(Q_i|S, T). \quad 34$$

35
 36 Samples drawn from the desired posterior distribution $\pi(Q|S, T)$ are 35
 37 generated by running a Markov chain. Briefly, we start with an initial set of 36
 38 37
 39 40

parameter values for Q . The new parameter set Q_{t+1} at time $t + 1$ is generated from a proposal transition function $T(Q_t, Q_{t+1})$. It will be either accepted or rejected by following the acceptance rule denoted as $r(Q_t, Q_{t+1})$. The criterion in designing the acceptance rule is to ensure that the detailed balance

$$\pi(Q_t|S, T) \cdot A(Q_t, Q_{t+1}) = \pi(Q_{t+1}|S, T) \cdot A(Q_{t+1}, Q_t)$$

is observed. This is necessary for the samples generated by the Markov chain to follow the desired posterior probability distribution $\pi(Q|S, T)$. The move set behind the proposal transition function that generates new trial parameter set is very important for efficient computation. Its design is discussed in Tseng and Liang (2006).

The Metropolis–Hastings acceptance rule

$$r(Q_t, Q_{t+1}) = \min \left\{ 1, \frac{\pi(Q_{t+1}|S, T) \cdot T(Q_{t+1}, Q_t)}{\pi(Q_t|S, T) \cdot T(Q_t, Q_{t+1})} \right\}$$

is a rule that ensures detailed balance. It either accepts or rejects the proposed new parameter set Q_{t+1} by evaluating whether a random number u generated from the uniform distribution between 0 and 1 is no greater than $r(Q_t, Q_{t+1})$.

C. Deriving Scoring Matrices from Rate Matrix

Once the expected values for the rate matrix Q are obtained, we follow the framework by Karlin and Altschul and derived scoring matrix used for assessing the similarity between residues at different time interval (Altschul *et al.*, 1997). For residue i and residue j at time interval t , the similarity score $b_{ij}(t)$ can be computed as

$$b_{ij}(t) = \frac{1}{\lambda} \log \frac{p_{ij}(t)}{\pi_j} = \frac{1}{\lambda} \log \frac{m_{ij}(t)}{\pi_i \pi_j},$$

where $m_{ij}(t)$ is the joint probability of observing both residue type i and j at the two nodes separated by time t , and λ is a scalar (Altschul *et al.*, 1997).

1 *D. Validity of the Evolutionary Model* 1

2 The validity of this approach is confirmed by extensive simulation test. 2
3 In Tseng and Liang (2006), an explicit phylogenetic tree and 16 artificially 3
4 evolved sequences of carboxypeptidase A2 are used to test if the underlying 4
5 model of substitution rate parameters of Jones, Taylor, and Thornton 5
6 (JTT) (Jones *et al.*, 1992) used to generate the artificial sequences can be 6
7 recovered. In 50 independent simulations, the recovered rates and the 7
8 true JTT parameters all have the weighted mean error (as defined in 8
9 Mayrose *et al.*, 2004) less than 0.0045. In addition, the parameters can 9
10 be recovered with acceptable accuracy when only about 20 residues in total 10
11 size are used (Tseng and Liang, 2006). 11
12 12
13 13
14 14
15 15

16 *E. Evolutionary Rates of Binding Surfaces and Other Surfaces are Different* 16

17 We have calculated the substitution rate matrix for both the binding 17
18 surface region and the remaining surface region of alpha amylase. The 18
19 distinct selection pressure for functional surface is also clearly evident in 19
20 the different patterns of the inferred substitution rates for binding region 20
21 and for the rest of the protein surface region (Fig. 8) (Tseng and Liang, 21
22 2006). In addition, both substitution patterns are also very different from 22
23 the precomputed JTT model (Jones *et al.*, 1992). This example illustrates 23
24 the need of extracting evolution pattern specific to the functional surfaces 24
25 of a particular protein for constructing sensitive and specific scoring 25
26 matrix for detecting functionally related protein surfaces. It also indicates 26
27 that selection pressure specific for protein function can be extracted 27
28 without being altered by selection pressure due to folding. 28
29 29
30 30
31 31

32 VI. PREDICTING PROTEIN FUNCTION BY DETECTING SIMILAR 32
33 BIOCHEMICAL BINDING SURFACES 33
34 34

35 *Amylase and other enzymes.* Alpha amylase (Enzyme Classification number 35
36 3.3.1.1) is an enzyme that breaks down starch, glycogen, and other related 36
37 polysaccharides and oligosaccharides. An objective test for protein func- 37
38 tion prediction is to take a known amylase structure and ask if it is used as a 38
39 template, whether we can find all other amylase structures in the PDB and 39
40 nothing else. This is a challenging task, as amylase exist in diverse species, 40

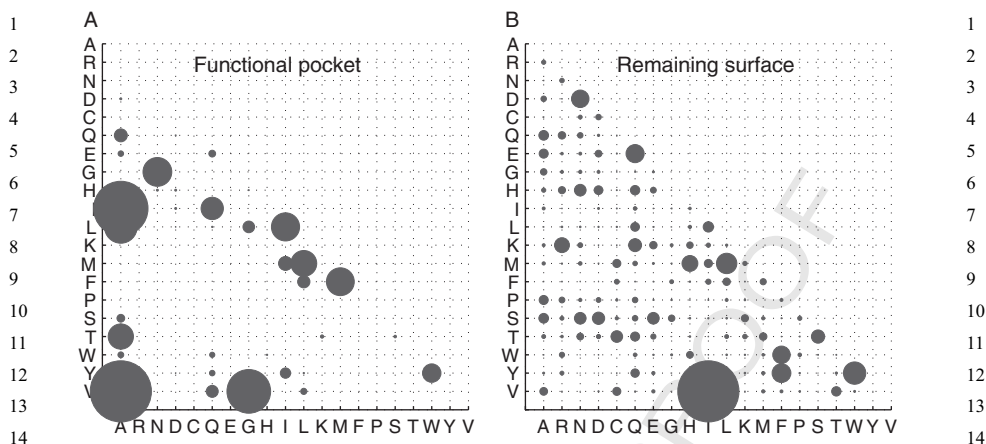


FIG. 8. Substitution rates of residues in the functional binding surface and the remaining surface of alpha-amylase (pdb 1bag). (A) Substitution rates of residues on functional binding surface (values represented by bubble sizes). (B) Substitution rates of residues on the remaining surface on 1bag. The values and overall pattern of substitutions that appear in both surface regions are very different. Adapted from Tseng and Liang (2006).

and some of them have very low sequence identity (<25%), which is challenging for function inference.

Using the template structure 1bag from *B. subtilis*, we are able to identify one of the computed pocket-containing 18 residues as the binding pocket (Fig. 9). With multiple sequence alignment of 14 sequences homologous to the template 1bag, all with <90% sequence identity to the template or to each other, we have constructed a phylogenetic tree using the Molphy package (Fig. 9A) (Adachi and Hasegawa, 1996). The rate matrix Q for the binding region (which corresponds to the positions of the 18 residues) is then estimated using the Bayesian Monte Carlo method we developed (Tseng and Liang, 2006). Scoring matrices of different divergence time are then generated from this rate matrix Q . These scoring matrices are then used to evaluate the similarity for each of the >2 million precomputed pocket/void sequence fragment contained in the pvSOAR database (Binkowski *et al.*, 2004) with the query sequence fragment. This comparison is carried out using the Smith–Waterman method as implemented in the FASTA package (Pearson, 1991). Promising hits with E -value <0.1 are then selected for further shape analysis. Those with cRMSD or oRMSD

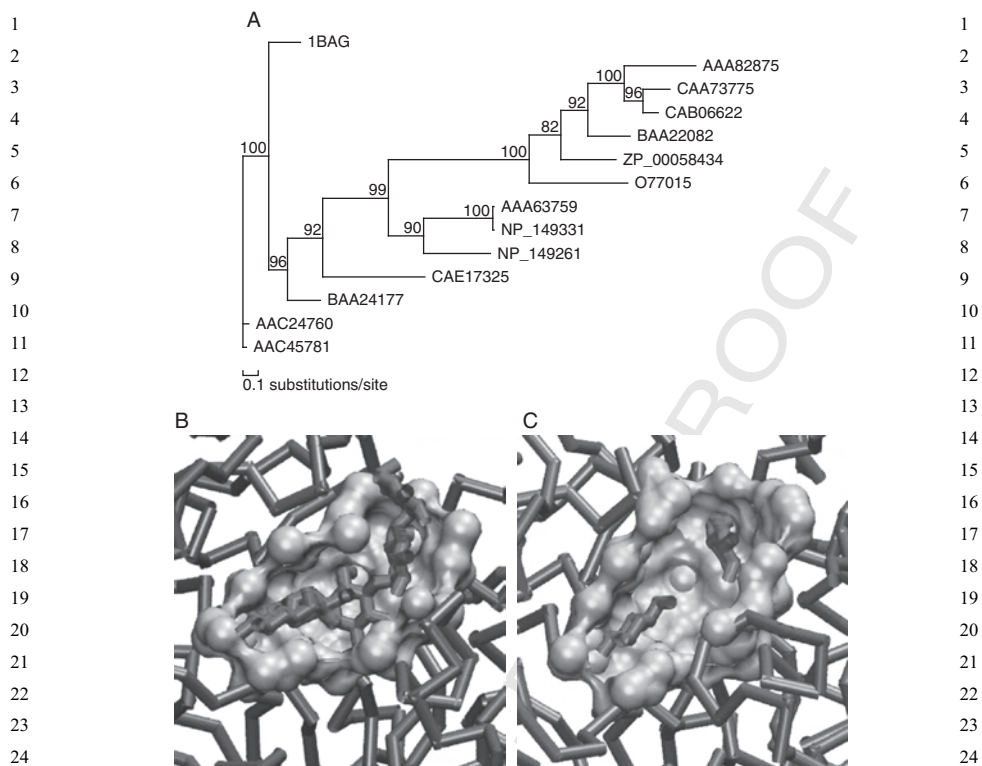


FIG. 9. Function prediction of alpha amylases. (A) The phylogenetic tree for PDB structure 1bag from *B. subtilis*. (B) The functional binding pocket of alpha amylase on 1bag. (C) A matched binding surface on a different protein structure (1b2y from human, full sequence identity 22%) obtained by querying with the binding surface of 1bag. Adapted from Tseng and Liang (2006).

values with the template surface pocket at a statistical significance of $p < 0.01$ (Binkowski *et al.*, 2003a) are then chosen as predicted hits, namely, proteins that are predicted as alpha amylase.

Using this template, we are able to predict 58 other PDB structures as alpha amylase. Indeed, all of them are found to have the same EC number as that of 1bag. When following the same procedure but using a different PDB template 1bg9 from the plant barley, we can predict 48 other PDB structures to be alpha amylase, again in this case all are of the same E.C. number as that of 1bg9 and 1bag (Tseng and Liang, 2006). Combining the

1 hits using these two templates together, we are able to identify 69 PDB 1
2 structures of alpha amylases among the 75 known alpha amylase structures. 2
3 This method using specific matrix estimated by Bayesian Monte Carlo 3
4 compares more favorably than using the general JTT matrix, and than 4
5 using the iterative dynamic programming sequence alignment method 5
6 PSI-BLAST. Details can be found in Tseng and Liang (2006). 6

7 This method has been tested for other enzymes. The results for 7
8 2,3-dihydroxybiphenyl dioxygenase (E.C. 1.13.11.39), adenosine deami- 8
9 nase (E.C. 3.5.4.4), 2-haloacid dehalogenase (E.C. 3.8.1.2), and phospho- 9
10 pyruvate hydratase (E.C. 4.2.1.11) are described in (Tseng and Liang, 2006), 10
11 where all other protein structures of the same E.C. numbers are correctly 11
12 predicted. In a recent study, we have selected a set of 100 enzyme families 12
13 with about 6000 structures and 770,000 precomputed binding surface pock- 13
14 ets/voids for testing. By taking the structure with the best resolution and *R*- 14
15 factor as template, we test if our method can identify other members of the 15
16 same protein family and nothing else. After calculating the overall sensitivity 16
17 and specificity of predictions of all 100 protein families, the accuracy of 17
18 predictions for the functions of all 6000+ structures from the 100 protein 18
19 family is 92%, and the best Mathews coefficient is 86.6% (Y. Y. Tseng and J. 19
20 Liang, unpublished results). 20
21

22 *Identifying metal cofactor of YecM from E. coli.* The problem of predicting 22
23 ion specificity of YecM protein structure is studied in (Binkowski *et al.*, 23
24 2005). YecM protein (pdb 1k4n) from *E. coli* was chosen as a structural 24
25 genomics target, as it does not have recognizable similarity to other 25
26 proteins of known structures. Structural analysis indicates that YecM shares 26
27 some similarity to an isomerase and several oxidoreductases (Zhang *et al.*, 27
28 2003b). As these proteins all contain a divalent metal cation, it was pre- 28
29 dicted that YecM is a metal-binding protein, but the preferred metal ions 29
30 were not known. 30

31 To predict the metal cofactor more accurately, the putative metal-binding 31
32 pocket on the YecM structure was compared against all known metal- 32
33 binding surfaces in the PDB database using pvSOAR (Binkowski *et al.*, 33
34 2004, 2005). The results of surface alignment indicate that several zinc- 34
35 binding surfaces from diverse species (*Rattus norvegicus*, *Bacillus thermopro-* 35
36 *teolyticus*, and *Bacillus anthracis*) share strong similarity to that of YecM, all 36
37 with significant *p*-values (Binkowski *et al.*, 2005). In fact, the top 30% of a 37
38 rank ordered list of all significant hits are zinc-binding surfaces. In contrast, 38
39 binding surfaces for other metal ions (i.e., Co, Mn, Fe, and Mg) have less 39
40

1 significant similarity to that of YecM. This result suggests that YecM is likely 1
2 to have zinc as its preferred metal cofactor. 2

3 *Locating the active site of ribose-5-phosphate isomerase.* pvSOAR analysis 3
4 helped to identify the active site of another protein from structural 4
5 genomics project (Binkowski *et al.*, 2005). RpiB protein from *E. coli* (pdb 5
6 Inn4) is known to have ribose-5-phosphate isomerase activity. However, 6
7 the active site on this protein is unknown (Zhang *et al.*, 2003c). Although 7
8 RpiA and RpiB have similar function, these two proteins belong to two 8
9 different structural folds (Binkowski *et al.*, 2005). The active site of RpiA as 9
10 identified by mutagenesis and cocrystal structure with inhibitor is absent 10
11 on RpiB structure (Zhang *et al.*, 2003c). A ligand docking study suggested 11
12 that the active site of RpiB from *M. tuberculosis* is located at the dimer 12
13 interface (Binkowski *et al.*, 2005). 13
14

15 Pairwise comparisons of the active sites using pvSOAR show that the active 15
16 sites of RpiA and RpiB from *E. coli* and *M. tuberculosis* have similar area and 16
17 volume, and the active sites on RpiB from *E. coli* and *M. tuberculosis* have 17
18 almost identical geometry measured in both cRMSD and oRMSD, with 18
19 strongly conserved phosphate-binding residues. Detailed analysis further 19
20 reveals that the most notable difference between RpiA and RpiB is in the 20
21 composition of basic residues, where His/Arg in RpiB are replaced by Lys in 21
22 RpiA. The surface patches of positively charged residues, and the orienta- 22
23 tion of acidic and basic residues important for catalysis are all conserved for 23
24 these proteins to carrying out similar functions. 24

25 Although biochemical assays clearly indicate that all three proteins have 25
26 the same substrate, and they are likely to have very similar binding 26
27 surfaces, the location and identities of the binding surfaces cannot be 27
28 detected without surface comparison, as RpiA and RpiB have no detect- 28
29 able similarity in overall sequence and structural fold. This study indicates 29
30 that pvSOAR analysis can help to understand how two seemingly different 30
31 binding surfaces performed the same function. 31

32 *Putative adenine nucleotide-binding site on CBS domain.* CBS domains are 32
33 present in many species and have unknown specific functions, but are 33
34 thought to be part of an energy status sensor complex (Scott *et al.*, 2004). 34
35 They appear in AMP-activated protein kinase, IMP dehydrogenase-2, and 35
36 chloride channel CLC2-binding adenosyl moieties (such as AMP, ATP, or 36
37 S-adenosyl methionine), and are often found in tandem pairs (Bateman, 37
38 1997; Scott *et al.*, 2004). Their biochemical roles and the locations of the 38
39 active sites are uncharacterized. 39
40

1 In the study of Binkowski *et al.* (2005), three structures of different 1
2 proteins from different species of archaea and bacteria-containing CBS 2
3 domains are analyzed (Fig. 10). These domains have about 20% sequence 3
4 identities, which is insufficient for functional inference. Surface patches 4
5 from the structures of these domains are identified and searched against a 5
6 library of AMP- and ATP-binding surfaces for potential matches. Among 6
7 these, well-defined interface pockets are identified by CASTp computa- 7
8 tion, and strong hits of diverse AMP- and ATP-binding surfaces are found 8
9 that are similar to these interface surfaces (Binkowski *et al.*, 2005). The 9
10 results suggest that both tandem CBS domains from protein mt1622 (pdb 10
11 1pbj from *M. thermoautotrophicum*) and inosine-5'-monophosphate dehy- 11
12 drogenase (IMPDH from *S. pyogenes*, pdb 1zjf) can bind to AMP and ATP, 12
13 consistent with experimental studies (Scott *et al.*, 2004). 13
14

15 An unexpected finding for hypothetical protein Ta549 CBS from 15
16 *T. acidophilum* is that an alternative binding surface is found to have formed 16
17 by a C-terminal additional insert of the singleton CBS domain, and a CBS 17
18 domain tandem pair on a different chain. This binding surface has only weak 18
19 similarity to the above-mentioned binding surface of the tandem CBS pairs, 19
20 but showed strong similarity to ATP-binding surface on saicar synthase from 20
21 *S. cerevisiae*. This finding suggests the existence of multiple-binding sites in a 21
22 CBS-binding domain, stabilized by a third CBS domain. 22
23

24 25 VII. ADAPTIVE PATTERNS OF SPECTRAL TUNING OF PROTEORHODOPSIN 26 FROM METAGENOMICS PROJECTS 26

27 Our method can also be applied to protein sequences with only limited 27
28 structural information to gain biological insight (Adamian *et al.*, 2006). 28
29 Proteorhodopsins (PR) are a class of newly discovered retinal-containing 29
30 rhodopsins with structural and functional similarities to archaeal bacter- 30
31 iorhodopsins (Beja *et al.*, 2000, 2001). They are found in numerous marine 31
32 bacteria and archaea through metagenomics studies of the communities 32
33 of marine organisms. A number of homologous proteorhodopsins were 33
34 functionally expressed in *E. coli* and found to form active, light-driven 34
35 proton pumps in the presence of retinal (Beja *et al.*, 2000; Friedrich *et al.*, 35
36 2002; Kim *et al.*, 2008; Sabehi *et al.*, 2005). 36
37

38 The absorption maxima of light wavelength of several subfamilies of 38
39 proteorhodopsins span the spectral range from blue (490 nm) to green 39
40 (525 nm) (Man *et al.*, 2003). The absorption maxima correlate with the 40

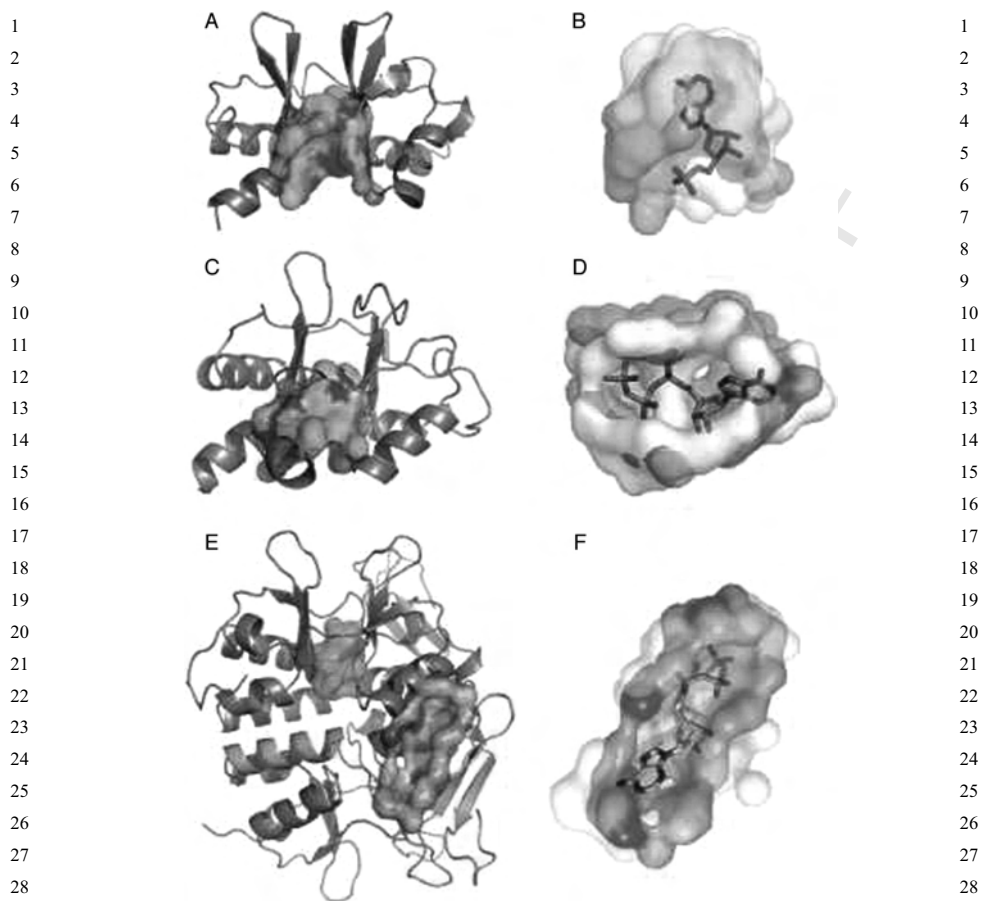


FIG. 10. Structures containing the CBS domain: (A) CBS domain protein mt1622 from *M. thermoautotrophicum* (PDB ID = 1pbj), (C) inosine-5'-monophosphate dehydrogenase (IMPDH) from *S. pyogenes* (PDB ID = 1zlj), and (E) conserved hypothetical protein Ta549 from *T. acidophilum* (PDB ID = 1pvm). The proposed nucleotide-binding surface of mt1622 (CASTp ID = 9, cyan, A) is shown superpositioned to a flavoprotein (PDB ID = 1efp, white) with bound AMP molecule (B). The IMPDH-binding surface (CASTp ID = 31, yellow) is show superpositioned with ATP bound cyclin-dependent kinase 2 (PDB ID = 1b38, white) (D). Ta549 contains an additional C-terminus CBS domain (C, orange) opposite the tandem domain interface surface (CASTp ID = 27, C, green). The domain insert creates a novel surface (CASTp ID = 30, orange) that shares similarity to an ATP-binding surface from saicar synthase (PDB ID = 1obd, white) (F).

1 depth at which the samples were collected, for example, green-absorbing 1
2 pigments (GPR) are found at the surface, and blue-absorbing pigments 2
3 (BPR) are found at the deeper waters (Beja *et al.*, 2001). Spectroscopic and 3
4 mutagenesis analyses indicate that a single residue difference at the 4
5 position 105 (Leu in GPR and Gln in BPR) functions as a spectral tuning 5
6 switch and accounts for most of the spectral differences (Man *et al.*, 2003). 6
7 Residues A, E, M, and V also appear at the position 105 in the family of 7
8 green-absorbing pigments, each with a specific absorption maximum 8
9 (Gomez-Consarnau *et al.*, 2007; Man *et al.*, 2003). 9

10 Based on sequence similarity to the archaeal bacteriorhodopsin with 10
11 known structures, we have mapped out 13 nonredundant putative 11
12 retinal-binding pocket sequence fragments from 99 sequences of proteor- 12
13 hodopsins (Adamian *et al.*, 2006). The substitution rates for the amino acid 13
14 residues forming the putative retinal-binding pocket are then calculated 14
15 using the Bayesian Markov chain Monte Carlo method (Tseng and Liang, 15
16 2006). Figure 11 shows the putative proteorhodopsin retinal-binding pock- 16
17 et sequences, along with the phylogenetic tree and the bubble plot of amino 17
18 acid substitution rates. The amino acid substitution rates indicate very fast 18
19 exchange rate between the pairs of amino acid residues at position 105 19
20 (Fig. 11C), such as A/E, A/L, A/V, E/Q, L/Q, E/L, and E/V, indicating 20
21 that this position of the retinal-binding pocket is the important location of 21
22 the functional adaptation of the proteorhodopsin. Results from this analy- 22
23 sis support the model that proteorhodopsins experience fast adaptation to 23
24 the environmental conditions (ocean depth) of their habitat by mutating at 24
25 position 105, rather than acquiring a new function (such as signal trans- 25
26 duction). As light is at a premium at ocean depth, spectral tuning is very 26
27 important, as a well-tuned pigment would be more effective at capturing 27
28 light (Beja *et al.*, 2001; Man *et al.*, 2003; Sabehe *et al.*, 2003). 28
29 29

30 31 32 33 VIII. GENERATING BINDING SITE NEGATIVE IMAGES 34 FOR DRUG DISCOVERY 35

36 We can also construct the negative image of a binding pocket, and use it 36
37 as a shape template for understanding substrate/ligand and protein bind- 37
38 ing. With additional chemical texture mapped on the template, negative 38
39 images of binding pockets can be used for rapid screening of compounds 39
40 to identify those that might bind to the proteins (Ebalunode *et al.*, 2008). 40

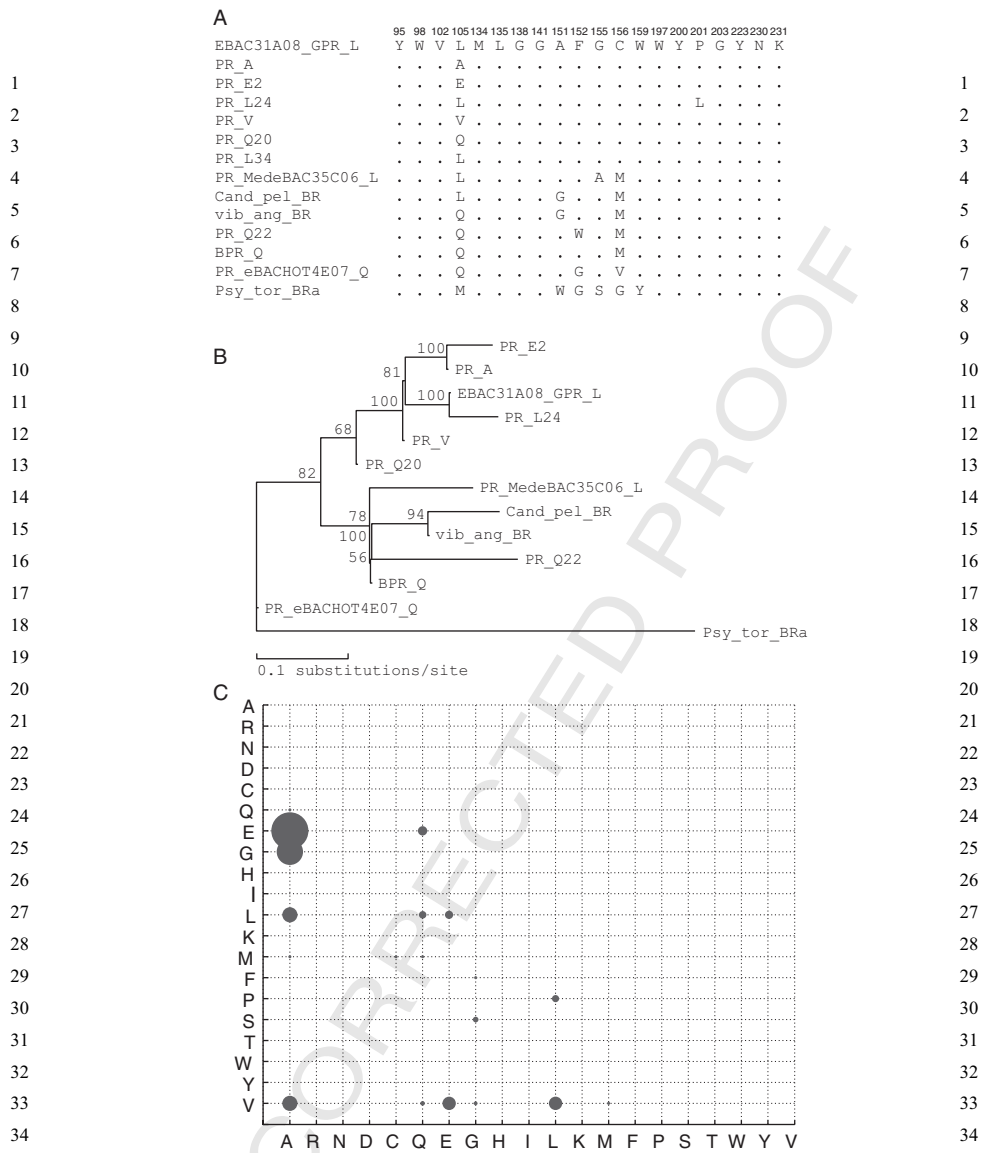
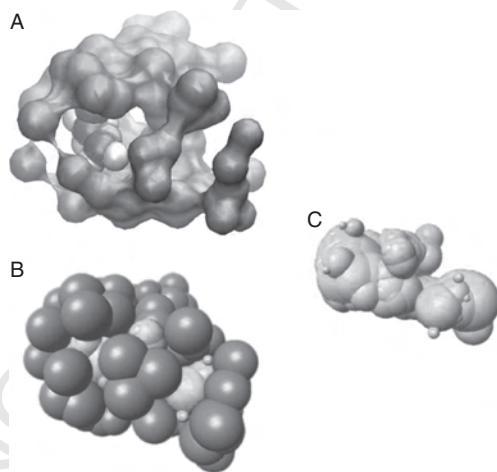


FIG. 11. Amino acid substitution rates in the putative retinal-binding pockets of proteorhodopsins. (A) Alignment of putative pocket sequences. The 20 pocket residue positions are mapped from retinal-binding pocket in bacteriorhodopsin structure 1KGB. Residues that are identical with the residues in the first sequence are substituted with "dots." (B) Phylogenetic tree of the full-length proteorhodopsin sequences. (C) The plot of amino acid substitution rates for residues in the putative retinal-binding pocket. The area of the circles is proportional to the substitution rate. The exchange pairs with the fastest rates are found at positions 93 and 137 in PR (following BR numbering). These are A/L, A/V, A/E, E/Q, E/L, L/Q, L/V, and M/T. Adapted from Adamian *et al.* (2006).

1 The negative image of a binding pocket can be constructed using a set
2 of circumscribing spheres for the discrete set of Delaunay tetrahedra and
3 triangles that defines the binding pocket (Ebalunode *et al.*, 2008;
4 Edelsbrunner *et al.*, 1998). First, the orthogonal centers of each Delaunay
5 tetrahedron contained in the binding pocket are calculated. Circum-
6 scribed spheres are then generated with the orthogonal centers taken as
7 their spherical centers. The radii of the circumscribed spheres are then
8 further optimized so the resulting collection of spheres most faithfully
9 represents the negative shape of the binding pocket (Ebalunode *et al.*,
10 2008). Figure 12 gives an example of the negative image computed for the
11 isoflurane-binding pocket in apoferritin, which provides the only soluble
12 protein model known to contain the structural motif thought to be
13 important for strong anesthetic binding (Liu *et al.*, 2005).
14

15 When combined with pharmacophore information, the negative images
16 of protein-binding pockets are found to be very effective in enriching
17 inhibitors when examining and ranking a long list of chemical compounds
18 for potential binding activities (Ebalunode *et al.*, 2008). Results for HIV-1
19 protease, phosphodiesterase 4B, estrogen receptor alpha, HIV-1 reverse
20 transcriptase, and thymidine kinase show that the enriched compounds
21



22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38 FIG. 12. The generation of a negative image of a binding pocket. (A) The surface
39 pocket in apoferritin that binds isoflurane. (B) The atoms forming the binding pocket
40 and its computed negative image. (C) Negative image of the binding pocket.

1 are of generally diverse chemical nature (Ebalunode *et al.*, 2008). This 1
2 offers an advantage for further development of drug-like compounds 2
3 based on these leads. 3
4 4
5 5
6 6

7 IX. SUMMARY AND CONCLUSION 7

8 Structural genomics projects have significantly advanced our under- 8
9 standing of the structural basis of the protein universe. It provides a wealth 9
10 of information for tackling the challenging problem of understanding 10
11 protein functions. By providing a large amount and standardized data, the 11
12 success of structural genomics enables development of new and well- 12
13 tailored computational methodology to interrogate a variety of problems 13
14 in functional understanding of the biological roles of protein molecules. 14
15 15

16 In this chapter, we have discussed our approach of studying protein 16
17 local surfaces for function inference and function characterization. The 17
18 approach described in this chapter combines computational geometric 18
19 characterization of protein structure, sequence and shape matching, and 19
20 uncovers evolutionary signal of protein function. Our results suggest that 20
21 this approach is effective in detecting enzyme functional surfaces, in 21
22 inferring and characterizing protein functions, and in gaining biological 22
23 insight of the relevant cellular processes. An important advantage of 23
24 this integrated approach is that it gives clear location information about 24
25 the region of protein surfaces where biological function occurs. Another 25
26 important advantage is that by generating well-defined surface pockets 26
27 and interior voids, by identifying those surfaces related to binding, and by 27
28 applying the Bayesian Monte Carlo method as developed in (Tseng and 28
29 Liang, 2006), we are now able to achieve the important task of separating 29
30 selection pressure due to protein function from that due to protein 30
31 stability and folding. This is evidence by the improved ability in predicting 31
32 protein functions when using customized scoring matrices computed 32
33 using our approach versus using precomputed scoring matrices. 33
34 34

35 It is envisioned that this approach of local surface analysis and comparison 35
36 can be generalized to study the challenging problem of physical protein- 36
37 protein interactions. Additional development in surface partition, shape 37
38 matching, and evolutionary signal detection will likely to yield new insight. 38
39 39
40 40

ACKNOWLEDGMENTS

We thank Drs. Rong Chen, Ken Dill, Rod Eckenhoff, Herbert Edelsbrunner, Jinfeng Zhang, Weifan Zheng, and Clare Woodward for fruitful collaborations. This work is supported by grants from National Science Foundation (CAREER DBI0133856, DBI0078270, DBI0646035, and DMS-0800257), National Institute of Health (GM68958, GM079804, and GM081682), Office of Naval Research (N000140310329), and Whitaker Foundation (TF-04-0023).

REFERENCES

- Adachi, J., and Hasegawa, M. (1996). MOLPHY, version 2.3. Programs for molecular phylogenetics based on maximum likelihood. *Comput. Sci. Monogr. Inst. Stat. Math. Tokyo* **28**, 1–150.
- Adamian, L., Ouyang, Z., Tseng, Y. Y., and Liang, J. (2006). Evolutionary patterns of retinal-binding pockets of type I rhodopsins and their functions. *Photochem. Photobiol.* **82**(6), 1426–1435.
- Altschul, S. F., Madden, T. L., Schaffer, A. A., Zhang, J., Zhang, Z., Miller, W., and Lipman, D. J. (1997). Gapped BLAST and PSI-BLAST: A new generation of protein database search programs. *Nucleic Acids Res.* **25**, 3389–3402.
- Bartlett, G. J., Porter, C. T., Borkakoti, N., and Thornton, J. M. (2002). Analysis of catalytic residues in enzyme active sites. *J. Mol. Biol.* **324**, 105–121.
- Bateman, A. (1997). The structure of a domain common to archaeobacteria and the homocystinuria disease protein. *Trends Biochem. Sci.* **22**(1), 12–13.
- Beja, O., Aravind, L., Koonin, E. V., Suzuki, M. T., Hadd, A., Nguyen, L. P., Jovanovich, S. B., Gates, C. M., Feldman, R. A., Spudich, J. L., Spudich, E. N., and DeLong, E. F. (2000). Bacterial rhodopsin: Evidence for a new type of phototrophy in the sea. *Science* **289**, 1902–1906.
- Beja, O., Spudich, E. N., Spudich, J. L., Leclerc, M., and DeLong, E. F. (2001). Proteorhodopsin phototrophy in the ocean. *Nature* **411**, 786–789.
- Bellman, R. (1958). On a routing problem. *Q. Appl. Math.* **16**(1), 87–90.
- Binkowski, T. A., Adamian, L., and Liang, J. (2003a). Inferring functional relationships of proteins from local sequence and spatial surface patterns. *J. Mol. Biol.* **332**, 505–526.
- Binkowski, T. A., Naghibzadeh, S., and Liang, J. (2003b). CASTp: Computed atlas of surface topography of proteins. *Nucleic Acids Res.* **31**, 3352–3355.
- Binkowski, T. A., Freeman, P., and Liang, J. (2004). pvSOAR: Detecting similar surface patterns of pocket and void surfaces of amino acid residues on proteins. *Nucleic Acids Res.* **32**, W555–W558.
- Binkowski, T. A., Joachimia, A., and Liang, J. (2005). Protein surface analysis for function annotation in high-throughput structural genomics pipeline. *Protein Sci.* **14**, 2972–2981.
- Chandonia, J. M., and Brenner, S. E. (2006). The impact of structural genomics: Expectations and outcomes. *Science* **311**(5759), 347–351.

- 1 Chen, L., Wu, L. Y., Wang, R., Wang, Y., Zhang, S., and Zhang, X. S. (2005). Compari- 1
2 son of protein structures by multi-objective optimization. *Genome Informatics* **16**(2), 2
3 114–124. 3
- 4 Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C. (2001). “Introduction to 4
5 Algorithms,” 2nd edn. MIT Press, Cambridge, MA. 5
- 6 Dayhoff, M. O., Schwartz, R. M., and Orcutt, B. C. (1978). “Atlas of Protein Sequence 6
7 and Structure.” National Biomedical Research Foundation, Washington, DC. 6
- 8 Doucet, A., De Freitas, N., Gordon, N., and Smith, A. (2001). “Sequential Monte Carlo 7
9 Methods in Practice.” Springer-Verlag, Berlin. 8
- 10 Dundas, J., Ouyang, Z., Tseng, J., Binkowski, A., Turpaz, Y., and Liang, J. (2006). 9
11 CASTp: Computed atlas of surface topography of proteins with structural and 10
12 topographical mapping of functionally annotated residues. *Nucleic Acids Res* **34**, 11
13 (Web Server issue), W116–W118. 12
- 14 Ebalunode, J. O., Ouyang, Z., Liang, J., and Zheng, W. (2008). Novel approach to 13
15 structure-based pharmacophore search using computational geometry and shape 14
16 matching techniques. *J. Chem. Inf. Model.* **48**(4), 889–901. 14
- 17 Edelsbrunner, H., and Mücke, E. P. (1994). Three-dimensional alpha shapes. *ACM* 15
18 *Trans. Graphics* **13**, 43–72. 16
- 19 Edelsbrunner, H., Facello, M., and Liang, J. (1998). On the definition and the con- 17
20 struction of pockets in macromolecules. *Discrete Appl. Math.* **88**, 83–102. 18
- 21 Friedrich, T., Geibel, S., Kalmbach, R., Chizhov, I., Ataka, K., Heberle, J., 19
22 Engelhard, M., and Bamberg, E. (2002). Proteorhodopsin is a light-driven proton 20
23 pump with variable vectoriality. *J. Mol. Biol.* **321**(5), 821–838. 20
- 24 Gavish, B., Gratton, E., and Hardy, C. J. (1983). Adiabatic compressibility of globular 21
25 proteins. *Proc. Natl Acad. Sci. USA* **80**, 750–754. 22
- 26 Glaser, F., Pupko, T., Paz, I., Bell, R. E., Shental, D., Martz, E., and Tal, N. (2003). 23
27 Consurf: Identification of functional regions in proteins by surface-mapping of 24
28 phylogenetic information. *Bioinformatics* **19**, 163–164. 25
- 29 Gold, N. D., and Jackson, R. M. (2006). Fold independent structural comparisons of 26
30 protein–ligand binding sites for exploring functional relationships. *J. Mol. Biol.* 27
31 **355**, 1112–1124. 27
- 32 Gomez-Consarnau, L., Gonzalez, J. M., Coll-Llado, M., Gourdon, P., Pascher, T., 28
33 Neutze, R., Pedros-Alio, C., and Pinhassi, J. (2007). Light stimulates growth of 29
34 proteorhodopsin-containing marine flavobacteria. *Nature* **445**(7124), 210–213. 30
- 35 Henikoff, S., and Henikoff, J. G. (1992). Amino acid substitution matrices from protein 31
36 blocks. *Proc. Natl Acad. Sci. USA* **89**, 10915–10919. 31
- 37 Huelsenbeck, J. P., Ronquist, F., Nielsen, R., and Bollback, J. P. (2001). Bayesian 32
38 inference of phylogeny and its impact on evolutionary biology. *Science* **294**, 33
39 2310–2314. 34
- 40 Jia, Y., Dewey, G., Shindyalov, I. N., and Bourne, P. E. (2004). A new scoring function 35
and associated statistical significance for structure alignment by CE. *J. Comput. Biol.* 36
11(5), 787–799. 37
- Jones, D. T., Taylor, W. R., and Thornton, J. M. (1992). The rapid generation of 38
mutation data matrices from protein sequences. *CABIOS* **8**, 275–282. 39
40

- 1 Karlin, S., and Altschul, S. F. (1990). Methods for assessing the statistical significance of 1
2 molecular sequence features by using general scoring schemes. *Proc. Natl Acad. Sci.* 2
3 *USA* **87**, 2264–2268. 3
- 4 Kim, S. Y., Waschuk, S. A., Brown, L. S., and Jung, K. H. (2008). Screening and 4
5 characterization of proteorhodopsin color-tuning mutations in *Escherichia coli* 5
6 with endogenous retinal synthesis. *Biochim. Biophys. Acta* **1777**(6), 504–513. 6
- 7 Kuhn, H. W. (1955). The Hungarian method for the assignment problem. *Nav. Res.* 7
8 *Logist. Q.* **2**, 83–97. 8
- 9 Laskowski, R. A., Watson, J. D., and Thornton, J. M. (2005). Protein function prediction 9
10 using local 3D templates. *J. Mol. Biol.* **351**, 614–626. 10
- 11 Levitt, M., and Gerstein, M. (1998). A unified statistical framework for sequence 11
12 comparison and structure comparison. *Proc. Natl Acad. Sci. USA* **95**, 5913–5920. 12
- 13 Liang, J., and Dill, K. A. (2001). Are proteins well-packed? *Biophys. J.* **81**(2), 751–766. 13
- 14 Liang, J., Edelsbrunner, H., Fu, P., Sudhakar, P. V., and Subramaniam, S. (1998a). 14
15 Analytical shape computing of macromolecules I: Molecular area and volume 15
16 through alpha-shape. *Proteins* **33**, 1–17. 16
- 17 Liang, J., Edelsbrunner, H., Fu, P., Sudhakar, P. V., and Subramaniam, S. (1998b). 17
18 Analytical shape computing of macromolecules II: Identification and computation 18
19 of inaccessible cavities inside proteins. *Proteins* **33**, 18–29. 19
- 20 Liang, J., Edelsbrunner, H., and Woodward, C. (1998c). Anatomy of protein pockets 20
21 and cavities: Measurement of binding site geometry and implications for ligand 21
22 design. *Protein Sci.* **7**, 1884–1897. 22
- 23 Liu, J. S. (2001). “Monte Carlo Strategies in Scientific Computing.” Springer-Verlag, 23
24 New York. 24
- 25 Liu, J. S., and Chen, R. (1998). Sequential Monte Carlo methods for dynamic systems. 25
26 *J. Am. Stat. Assoc.* **93**, 1032–1044. 26
- 27 Liu, R., Loll, P. J., and Eckenhoff, R. G. (2005). Structural basis for high-affinity 27
28 volatile anesthetic binding in a natural 4-helix bundle protein. *FASEB J.* **19**(6), 28
29 567–576. 29
- 30 Lorenz, B., Orgzall, I., and Heuer, H. O. (1993). Universality and cluster structures in 30
31 continuum models of percolation with two different radius distributions. *J. Phys.* 31
32 *A: Math. Gen.* **26**, 4711–4722. 32
- 33 Man, D. L., Wang, W. W., Sabehi, G., Aravind, L., Post, A. F., Massana, R., 33
34 Spudich, E. N., Spudich, J. L., and Beja, O. (2003). Diversification and spectral 34
35 tuning in marine proteorhodopsins. *EMBO J.* **22**(3), 1725–1731. 35
- 36 Mayrose, I., Graur, D., Tal, N., and Pupko, T. (2004). Comparison of site-specific rate- 36
37 inference methods for protein sequences: Empirical Bayesian methods are superior. 37
38 *Mol. Biol. Evol.* **21**, 1781–1791. 38
- 39 Najmanovich, R. J., Torrance, J. W., and Thornton, J. M. (2005). Prediction of protein 39
40 function from structure: Insights from methods for the detection of local structural 40
41 similarities. *BioTechniques* **38**(6), 847–851. 41
- 42 Pazos, F., and Sternberg, M. J. (2004). Automated prediction of protein function and 42
43 detection of functional sites from structure. *Proc. Natl Acad. Sci. USA* **101**(41), 43
44 14754–14759. 44

- 1 Pearson, W. R. (1991). Searching protein sequence libraries: Comparison of the sensi- 1
2 tivity and selectivity of the Smith–Waterman and FASTA algorithms. *Genomics* **11**, 2
3 635–650. 3
- 4 Pegg, S. C., Brown, S. D., Ojha, S., Seffernick, J., Meng, E. C., Morris, J. H., Chang, P. J., 4
5 Huang, C. C., Ferrin, T. E., and Babbitt, P. C. (2006). Leveraging enzyme structure– 5
6 function relationships for functional inference and experimental design: The 6
7 structure–function linkage database. *Biochemistry* **45**, 2545–2555. 7
- 8 Richards, F. M., and Lim, W. A. (1994). An analysis of packing in the protein folding 8
9 problem. *Q. Rev. Biophys.* **26**, 423–498. 9
- 10 Rost, B. (2002). Enzyme function less conserved than anticipated. *J. Mol. Biol.* **318**, 10
11 595–608. 11
- 12 Russell, R. B. (1998). Detection of protein three-dimensional side-chain patterns: New 12
13 examples of convergent evolution. *J. Mol. Biol.* **279**, 1211–1227. 13
- 14 Sabehi, G., Massana, R., Bielawski, J. P., Rosenberg, M., Delong, E. F., and Bj, O. 14
15 (2003). Novel proteorhodopsin variants from the Mediterranean and Red Seas. 15
16 *Environ. Microbiol.* **5**(10), 842–849. 16
- 17 Sabehi, G., Loy, A., Jung, K. H., Partha, R., Spudich, J. L., Isaacson, T., Hirschberg, J., 17
18 Wagner, M., and Beja, O. (2005). New insights into metabolic properties of marine 18
19 bacteria encoding proteorhodopsins. *PLoS Biol.* **3**(8), e273. 19
- 20 Scott, J. W., Hawley, S. A., Green, K. A., Anis, M., Stewart, G., Scullion, G. A., 20
21 Norman, D. G., and Hardie, D. G. (2004). CBS domains from energy-sensing 21
22 modules whose binding of adenosine ligands is disrupted by disease mutations. 22
23 *J. Clin. Invest.* **113**(2), 274–284. 23
- 24 Stauffer, D. (1985). “Introduction to Percolation Theory.” Taylor & Francis, London. 24
25 Tian, W., and Skolnick, J. (2003). How well is enzyme function conserved as a function 25
26 of pairwise sequence identity? *J. Mol. Biol.* **333**, 863–882. 26
- 27 Torrance, J. W., Bartlett, G. J., Porter, C. T., and Thornton, J. M. (2005). Using a library 27
28 of structural templates to recognise catalytic sites and explore their evolution in 28
29 homologous families. *J. Mol. Biol.* **347**, 565–581. 29
- 30 Tseng, Y. Y., and Liang, J. (2006). Estimation of amino acid residue substitution rates at 30
31 local spatial regions and application in protein function inference: A Bayesian 31
32 Monte Carlo approach. *Mol. Biol. Evol.* **23**(2), 421–436. 32
- 33 Tseng, Y. Y., and Liang, J. (2007). Predicting enzyme functional surfaces and locating 33
34 key residues automatically from structures. *Ann. Biomed. Eng.* **35**(6), 1037–1042. 34
- 35 Umeyama, S. (1991). Least-square estimation of transformation parameters between 35
36 two point patterns. *IEEE Trans. Pattern Anal. Mach. Intell.* **13**, 376–380. 36
- 37 Whelan, S., and Goldman, N. (2001). A general empirical model of protein evolution 37
38 derived from multiple protein families using a maximum-likelihood approach. 38
39 *Mol. Biol. Evol.* **18**, 691–699. 39
- 40 Whelan, S., Liò, P., and Goldman, N. (2001). Molecular phylogenetics: State-of-the-art 40
41 methods for looking into the past. *Trends Genet.* **17**, 262–272. 41
- 42 Yang, Z. (1997). PAML: A program package for phylogenetic analysis by maximum 42
43 likelihood. *Comput. Appl. Biosci.* **13**, 555–556. 43
- 44 Ye, Y., and Godzik, A. (2004). Database searching by flexible protein structure align- 44
45 ment. *Protein Sci.* **13**, 1841–1850. 45

- 1 Zhang, J., Chen, R., Tang, C., and Liang, J. (2003a). Origin of scaling behavior of 1
2 protein packing density: A sequential Monte Carlo study of compact long chain 2
3 polymers. *J. Chem. Phys.* **118**, 6102–6109. 3
4 Zhang, R. G., Duke, N., Laskowski, R., Evdokimova, E., Skarina, T., Edwards, A., 4
5 Joachimiak, A., and Savchenko, A. (2003b). Conserved protein YecM from *Escher-* 5
6 *ichia coli* shows structural homology to metal-binding isomerases and oxygenases. 6
7 *Proteins* **51**(2), 311–314. 7
8 Zhang, R. G., Andersson, C. E., Skarina, T., Evdokimova, E., Edwards, A. M., 8
9 Joachimiak, A., Savchenko, A., and Mowbray, S. L. (2003c). The 2.2 Å resolution 9
10 structure of RpiB/AlsB from *Escherichia coli* illustrates a new approach to the ribose- 10
11 5-phosphate isomerase reaction. *J. Mol. Biol.* **332**(5), 1083–1094. 11
12 Zhu, J., and Weng, Z. (2005). A novel protein structure alignment algorithm. *Proteins:* 12
13 *Struct. Funct. Bioinf.* **14**, 417–423. 13
14 14
15 15
16 16
17 17
18 18
19 19
20 20
21 21
22 22
23 23
24 24
25 25
26 26
27 27
28 28
29 29
30 30
31 31
32 32
33 33
34 34
35 35
36 36
37 37
38 38
39 39
40 40