

# Protein–Protein Binding-Sites Prediction by Protein Surface Structure Conservation<sup>†</sup>

Janez Konc and Dušanka Janežič\*

National Institute of Chemistry, Hajdrihova 19, SI-1000 Ljubljana, Slovenia

Received November 20, 2006

A new algorithm to predict protein–protein binding sites using conservation of both protein surface structure and physical–chemical properties in structurally similar proteins is developed. Binding-site residues in proteins are known to be more conserved than the rest of the surface, and finding local surface similarities by comparing a protein to its structural neighbors can potentially reveal the location of binding sites on this protein. This approach, which has previously been used to predict binding sites for small ligands, is now extended to predict protein–protein binding sites. Examples of binding-site predictions for a set of proteins, which have previously been studied for sequence conservation in protein–protein interfaces, are given. The predicted binding sites and the actual binding sites are in good agreement. Our algorithm for finding conserved surface structures in a set of similar proteins is a useful tool for the prediction of protein–protein binding sites.

## 1. INTRODUCTION

The number of proteins, that are known to interact, grows much faster<sup>1,2</sup> than the number of structures of protein complexes in the Protein Data Bank (PDB).<sup>3</sup> This creates an opportunity for the development of computational approaches to predict binding sites on these proteins. The common strategy for predicting protein–protein binding sites is to analyze interfaces of a set of existing protein complexes and to determine parameters which differentiate a binding site from the rest of the protein surface. It is known that certain residues appear more often in interfaces than in the rest of the protein and that certain conserved residues, namely, hotspots, contribute the most to the binding free energy.<sup>4–8</sup> A prediction model is usually constructed on a set of known protein–protein interfaces and is then used to predict binding sites for proteins excluded from this set. Approaches for predicting protein–protein binding sites can be divided into sequence and structural approaches. Among the latter, a combination of surface parameters can be used to calculate the probability of a surface patch forming protein–protein interactions.<sup>9</sup> Support vector machines (SVMs) have been used to predict interaction sites, where a SVM is trained to distinguish between interacting and noninteracting surface patches using various surface properties.<sup>10</sup> A similar structural method, which samples the surface and assigns a probability to be a part of an interface to each residue, is ProMate.<sup>11</sup> Approaches where only sequences are used to predict protein–protein binding sites have also been reported.<sup>12</sup> In both structural and sequence approaches, a set of interacting proteins is partitioned into the training set and test set. The methods are trained on the first set of proteins, and predictions are made for the test set. Docking methods have also been used to predict protein–protein interactions. For a review, see the literature.<sup>13</sup>

Previously, it has been shown that binding-site residues of proteins are more conserved among proteins than the rest of the surface residues. Although it has been argued that

sequence conservation is rarely sufficient for complete and accurate prediction of a protein–protein interface,<sup>14</sup> other methods successfully predict protein–protein interactions on the basis of the conservation of sequence and structure.<sup>15</sup>

In this paper, we develop an algorithm which predicts protein–protein binding sites using conserved protein surface structure together with physical–chemical properties in structurally similar proteins. It is based on the idea that the most conserved part of the protein surface in terms of the physical–chemical properties must be related either to the binding of small endogenous ligands or to that of other proteins. To find the conserved part of the protein surface, the algorithm takes this query protein and compares it to one or more of its structural neighbors. We represent functional groups of the surface residues with labeled points (called also pseudocenters), which was previously used in a method to detect related functions among proteins.<sup>16</sup> The main idea in that method is to detect functional relationships among proteins independent of a given sequence or fold homology, assuming that proteins with similar functions must have conserved recognition features, that is, common physical–chemical properties inside binding cavities for endogenous ligands. Being more flat, protein–protein interfaces are more difficult to detect than binding sites for small ligands. As opposed to the approach of Schmitt et al.,<sup>16</sup> who considered only protein cavities, our algorithm compares whole surfaces of proteins.

The algorithm is tested in predicting protein–protein binding sites on a set of protein complexes from the literature.<sup>14</sup> Each protein complex in the set is split into its constituent chains, and one chain from each complex is compared with one or more of its structural neighbors. The surface that we find to be conserved in both the chain and its neighbor structure(s) is then predicted to be the binding site for the second chain in the complex. To verify the predicted binding site, we compare it against the actual binding site.

## 2. METHODS

Our algorithm for predicting protein–protein binding sites by protein structure conservation is schematically depicted

<sup>†</sup> Dedicated to Professor Nenad Trinajstić on the occasion of his 70th birthday.

\* Corresponding author email: dusa@cmm.ki.si.



same as in the previous step. Higher values of this parameter correspond to a greater similarity.

**D. Generate a Product Graph around Each Pair of Similar Points.** We generate one product graph for each pair of similar points ( $p$ ,  $r$ ), where  $p$  is from the first and  $r$  is from the second protein, for which we calculated in the previous step that similarity  $> 2.8$ . A product graph consists of a set of vertices, where each vertex is a pair of similar points, and a set of edges connecting these vertices. Besides vertex ( $p$ ,  $r$ ), other vertices in this product graph are all pairs of points ( $p_i$ ,  $r_i$ ); again,  $p_i$  is from the first and  $r_i$  is from the second protein, which have similarity  $> 1.9$  and both distance( $p$ ,  $p_i$ ) and distance( $r$ ,  $r_i$ ) within 12 Å. We then connect each two vertices of a product graph ( $p_i$ ,  $r_i$ ) and ( $p_j$ ,  $r_j$ ) by an edge if  $|\text{distance}(p_i, p_j) - \text{distance}(r_i, r_j)| < 0.5$  Å.

**E. Find a Maximum Clique in Each Product Graph.** Product graphs are taken one by one, and each is searched for a maximum clique. A clique is a subgraph of a graph in which each vertex is connected to all other vertices. A maximum clique in a product graph generated with the above rules corresponds to a surface patch that is common to both proteins.<sup>18,19</sup>

**F. Calculate the Rotation and Translation of Protein 1 to Protein 2 for Each Maximum Clique Found.** We rotate and translate points of the common surface patch on the first protein, so that they are aligned with the points of the common surface patch on the second protein.<sup>20</sup> We repeat this for each maximum clique found.

**G. Join Cliques that Give Similar Transformations into Clusters.** We join maximum cliques with similar rotation matrices and translation vectors into clusters. Each cluster represents a larger similarity in shape and physical–chemical properties of the two protein surfaces.

**H. Sort Clusters by Their Size (Number of Aligned Points) and Their Root-Mean-Squared Deviation (RMSD).** Finally, we sort the list of clusters with respect to their size and the RMSD of the two sets of aligned points. We output the clusters, for each cluster the two sets of aligned points, their corresponding two sets of residues, and rotation matrix/translation vector. In this study, we use only the residues from the first and largest cluster in the list, which gives the best prediction of a binding site.

### 3. RESULTS

Our algorithm for finding conserved regions on protein surfaces was tested for predicting protein–protein binding sites on a set of protein complexes partially adapted from the literature<sup>14</sup> and augmented with our own test proteins. This set was chosen because it has been used by other authors to study the conservation of protein–protein interface sequences. Instead of the sequences, we studied conservation of the interface structures, which we also extended to predictions of binding sites. The algorithm, shown in Figure 1, is used. From start to finish, this algorithm took  $< 10$  s on a 1.6 GHz AMD Opteron processor.

Each complex in the set was split into its constituent chains. We then compared one of the chains to one or more of this chain's structural neighbors. The surface that was conserved (the first cluster of residues or its part) in both the chain and its neighbor structures was then predicted to be the binding site for a second chain. The predicted binding-

site residues were then compared with the actual binding-site residues, which we extracted from each protein complex. We used the definition of an interface to be the region between two polypeptide chains that are not covalently linked.

Two residues were defined as being a part of the protein–protein interface if the distance between any two atoms of the two residues from different chains was less than the sum of their van der Waals radii plus 3.0 Å. This is more than 0.5 Å as used by other authors,<sup>21</sup> which we found to be too limiting. As the conservation of residues is not limited only to direct contacting residues, we also wanted to sample nearby residues, which provide a structural scaffold to the interface and may also be predicted with our algorithm.

The agreement of predicted binding sites with the actual binding sites was measured with *specificity* and *sensitivity*, which are standard measures used in this field.<sup>10</sup> *Specificity* indicates the proportion of the predicted residues that are also interface residues and is defined as specificity = the number of predicted residues in the actual interface/number of predicted residues. *Sensitivity* tells us the proportion of the interface that was predicted. It is defined as sensitivity = number of predicted residues in the actual interface/number of interface residues.

The polypeptide chain for which we predicted the binding site was compared to one or more of its structural neighbors. A list of structural neighbors was provided by the VAST (vector alignment search tool) Web page, which offers structure–structure alignments of publicly available protein structures.<sup>22</sup> We used a medium redundancy list to avoid too similar structures which would give biased predictions. All of the protein structures in the list share some sequence identity (%ID = 0–100) with the polypeptide chain for which the prediction is performed. Our program settings are tuned to find similarities in proteins with fair sequence identity, so we used proteins from the list with %ID = 20–50. We also took care that the alignments stretched over whole sequences. The results of the protein–protein binding-site predictions are shown in Table 2.

Chain A of the phycobiliprotein allophycocyanin (1al1A) forms an interface with chain B. To predict the binding site on chain A, we compared this chain with each of its top structural neighbors in the list. All share %ID = 30–38 sequence identity with chain A. Their PDB codes with chain identifiers, together with the specificity and the sensitivity of the predicted binding sites, are 1kn1B (24%, 29%), 1jboA (33%, 63%), 2c71A (74%, 33%), 2bv8A (69%, 42%), and 1b8dA (26%, 47%). These results are presented in Table 2. The conserved residues of all five predicted binding sites were then mapped onto the surface of chain A of the protein in question in an effort to improve the specificity and the sensitivity. This time we counted a residue to be a part of the binding site if it was conserved in at least four structural neighbors. This prediction is shown in Table 2. The predicted conserved surface contains a separate binding site for the chromophore, which reduces the specificity of this prediction.

An interesting case is chain A in the coagulation factor X (1hcgA) from the superfamily of trypsin-like serine proteases. It is known that the active site in this family is very well conserved. When we compared chain A with the first protein in the list of structural neighbors (1q3xA), we found two distinct conserved surface patches. The first patch

**Table 2.** Results from Protein-Protein Binding-Site Predictions<sup>a</sup>

PDB code and chain	compared against	interface size (no. of residues)	predicted interface size (no. of residues)	specificity (%)	sensitivity (%)
<b>Heterodimer</b>					
1allA	1kn1B	43	88	24	49
1allA	1jboA	43	83	33	63
1allA	2c71A	43	19	74	33
1allA	2bv8A	43	26	69	42
1allA	2b8dA	43	76	26	47
1allA	1kn1B&1jboA&2c71A&2bv8A&1b8bA	43	49	37	42
		32	21	47	31
1hcgA	1q3xA				
1lucA	1bs1B	64	95	52	77
1tcoB	2ct9B	62	38	47	29
1tcoB	1dguA	62	87	39	55
1tcoB	1uhnA	62	15	67	16
1tcoB	2ct9B&1dguA&1uhnA	62	94	41	63
<b>Homodimer</b>					
1bncA	1ulzA	40	146	16	58
1bncA	1ulzA	40	32	35	28
1daaA	1i2kA	60	48	52	42
<b>Transient</b>					
1azeA	1jegA	15	20	45	60
1azeA	1e6hA	15	25	40	67
1azeA	1ju5C	15	26	46	80
1azeA	1jegA&1e6hA&1ju5C	15	8	63	33
1lw6I	1cseI	18	23	61	78

<sup>a</sup> Each chain is assigned a code that consists of the PDB code and the identifier of the chain that was used in structural comparison.

overlaps with the actual binding site for chain B. The second conserved surface patch overlaps with the protein's active site. The specificity and the sensitivity of the predicted protein-protein binding site are shown in Table 2.

The interface between chains A and B in luciferase (1luc) is big and flat, extending one whole side of each chain. Two distinct surface patches, with 95 and 11 residues, are found in the first cluster. The larger of the two overlaps with the binding site for chain B on chain A. The results are presented in Table 2.

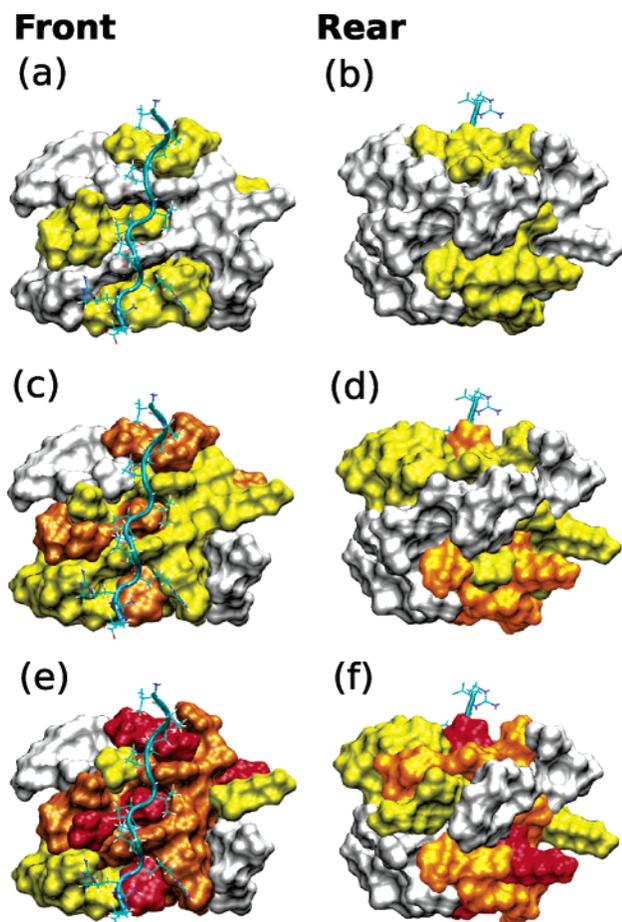
The calmodulin-dependent phosphatase (1tco) is known to interact with several different proteins.<sup>23,24</sup> We compare chain B (1tcoB) with each of the first three proteins in the list of structural neighbors with %ID = 30–45. Their PDB codes and chain identifiers, together with the specificity and the sensitivity of the predicted binding sites, are 2ct9B (47%, 29%), 2dguA (39%, 55%), and 1uhnA (67%, 16%). These results are again presented in Table 2. Chain B also forms a small interface with chain C, which we did not consider. To increase the reliability of predictions, we also mapped the residues of the three predicted binding sites onto the surface of chain A. This time, we obtained the best prediction of the binding site, when we counted a residue to be a part of the binding site if it was conserved in at least two of the structural neighbors, which is shown in Table 2. The interface between chains A and B in this protein is also found to be significantly conserved in a previous study.<sup>14</sup>

For chain A of the homodimer acetyl-CoA carboxylase (1bncA), we predicted the binding site with only 16% specificity. The number of predicted binding-site residues (146) far exceeds the actual binding-site size (40), which inevitably lowers the specificity of this prediction. This failure could be explained by the previous findings that the interface in this protein is less conserved than the rest of the exposed surface.<sup>14</sup> It may also be due to the high sequence

identity (%ID = 50) of the structural neighbor (1ulzA) that was used, which causes the program to find too many conserved residues. To test if this was the case, we modified the parameters of the algorithm to make it less sensitive and repeated the calculation. The specificity improved to 35%, but now the sensitivity was only 28%, which is still a better prediction. We show both predictions in Table 2.

The prediction of the binding site on chain A of homodimer D-amino acid aminotransferase (1daaA), which forms an interface with chain B of this complex, is shown in Table 2. In this case, the protein-protein binding site partially overlaps with the active site, which may also contribute to good conservation of this region.

We used our algorithm to predict binding sites on the two chains involved in transient complexes, the Grb2 SH3 domain (1azeA) and chymotrypsin inhibitor 2a (1lw6I). We compared chain A (1azeA) with three of its structural neighbors and mapped the conserved residues onto the surface of this chain. The PDB codes and the chain identifiers of the first, second, and third structural neighbor, together with the specificity and sensitivity of each prediction, are 1jegA (45%, 60%), 1e6hA (40%, 67%), and 1ju5C (46%, 80%), respectively. This mapping of the surface is shown in Figure 2 in which panels a and b show the surface of 1azeA mapped with the residues which are conserved in both the query protein and its first structural neighbor. In panels c and d, the residues that are conserved in the query protein and the first two structural neighbors have been mapped. Finally, in panels e and f, the residues that are conserved in each of the four proteins, including the third structural neighbor (1ju5C), were again mapped, and residues which were conserved in all four proteins (red) were predicted to be in the binding site. The residues that are conserved in all four proteins are colored red; those conserved in only three of the proteins are colored orange, and the ones conserved



**Figure 2.** In panels a, c, and e, the front and, in b, d, and f, the back side of chain A (lazeA) forming the interface with chain B (line representation) are shown. Chain A was compared with an increasing number of its structural neighbors: (a and b) with one (1jegA); (c and d) with two (1jegA and 1e6hA); (e and f) with three (1jegA, 1e6hA, and 1ju5C). The conserved residues are mapped on the surface of lazeA; residues, which are conserved in lazeA and one, two, or three of its structural neighbors, are colored yellow, orange, and red, respectively.

in only two of the proteins are yellow. From Figure 2, we observe that, by increasing the number of compared structural neighbors, the conserved surface corresponds better to the actual binding site (specificity is 63%) than when taking each of the single comparisons. These results, together with the results for the chymotrypsin inhibitor 2a (1lw6I), are shown in Table 2.

#### 4. CONCLUSIONS

We describe an algorithm which predicts the protein–protein binding site in a protein by finding the most conserved surface between this protein and one or more of its structural neighbors. Our algorithm differs from others in that only the structure of a protein and a couple of its structural neighbors is needed. This approach may give more unbiased predictions of protein–protein binding sites than predictions obtained by other methods, which are trained on a set of existing interfaces. Since it uses a different paradigm, our approach may be best when used in combination with these methods. Our algorithm can also be used to reduce

the search space of docking algorithms and can provide new targets for potential inhibitors of protein–protein interactions.

#### ACKNOWLEDGMENT

The financial support through grants P1-0002 of the Ministry of Higher Education, Science, and Technology of Slovenia is acknowledged.

#### REFERENCES AND NOTES

- (1) Salwinski, L.; Miller, C. S.; Smith, A. J.; Pettit, F. K.; Bowie, J. U.; Eisenberg, D. The Database of Interacting Proteins: 2004 Update. *Nucleic Acids Res.* **2004**, *32*, 449–451.
- (2) Bader, G. D.; Donaldson, I.; Wolting, C.; Ouellette, B. F. F.; Pawson, T.; Hogue, C. W. V. BIND-The Biomolecular Interaction Network Database. *Nucleic Acids Res.* **2001**, *29*, 242–245.
- (3) Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T. N.; Weissig, H.; Shindyalov, I. N.; Bourne, P. E. The Protein Data Bank. *Nucleic Acids Res.* **2000**, *28*, 235–242.
- (4) Cunningham, B. C.; Wells, J. A. Rational Design of Receptor-Specific Variants of Human Growth Hormone. *Proc. Natl. Acad. Sci. U.S.A.* **1991**, *88*, 3407–3411.
- (5) Clackson, T.; Wells, J. A. A Hot Spot of Binding Energy in a Hormone-Receptor Interface. *Science* **1995**, *267*, 383–386.
- (6) Ma, B.; Elkayam, T.; Wolfson, H.; Nussinov, R. Protein–Protein Interactions: Structurally Conserved Residues Distinguish between Binding Sites and Exposed Protein Surfaces. *Proc. Natl. Acad. Sci. U.S.A.* **2003**, *100*, 5772–5777.
- (7) Li, X.; Keskin, O.; Ma, B.; Nussinov, R.; Liang, J. Protein–Protein Interactions: Hot Spots and Structurally Conserved Residues often Locate in Complemented Pockets that Pre-Organized in the Unbound States: Implications for Docking. *J. Mol. Biol.* **2004**, *344*, 781–795.
- (8) Keskin, O.; Ma, B.; Nussinov, R. Hot Regions in Protein–Protein Interactions: The Organization and Contribution of Structurally Conserved Hot Spot Residues. *J. Mol. Biol.* **2005**, *345*, 1281–1294.
- (9) Jones, S.; Thornton, J. M. Prediction of Protein–Protein Interaction Sites Using Patch Analysis. *J. Mol. Biol.* **1997**, *272*, 133–143.
- (10) Bradford, J. R.; Westhead, D. R. Improved Prediction of Protein–Protein Binding Sites Using a Support Vector Machines Approach. *Bioinformatics* **2005**, *21*, 1487–1494.
- (11) Neuvirth, H.; Raz, R.; Schreiber, G. ProMate: A Structure Based Prediction Program To Identify the Location of Protein–Protein Binding Sites. *J. Mol. Biol.* **2004**, *338*, 181–199.
- (12) Gallet, X.; Charlotiaux, B.; Thomas, A.; Brasseur, R. A Fast Method to Predict Protein Interaction Sites from Sequences. *J. Mol. Biol.* **2000**, *302*, 917–926.
- (13) Smith, G. R.; Sternberg, M. J. Prediction of Protein–Protein Interactions by Docking Methods. *Curr. Opin. Struct. Biol.* **2002**, *12*, 28–35.
- (14) Caffrey, D. R.; Somaroo, S.; Hughes, J. D. Are Protein–Protein Interfaces More Conserved in Sequence than the Rest of the Protein Surface? *Protein Sci.* **2004**, *13*, 190–202.
- (15) Aytuna, A. S.; Gursoy, A.; Keskin, O. Prediction of Protein–Protein Interactions by Combining Structure and Sequence Conservation in Protein Interfaces. *Bioinformatics* **2005**, *21*, 2850–2855.
- (16) Schmitt, S.; Kuhn, D.; Klebe, G. A New Method to Detect Related Function Among Proteins Independent of Sequence and Fold Homology. *J. Mol. Biol.* **2002**, *323*, 387–406.
- (17) Konc, J.; Hodošček, M.; Janežič, D. Molecular Surface Walk. *Croat. Chem. Acta* **2006**, *79*, 237–241.
- (18) Konc, J.; Janežič, D. A Maximum Clique Problem Revisited. *Eur. J. Oper. Res.* (submitted).
- (19) Konc, J.; Janežič, D. A Branch and Bound Algorithm for Matching Protein Structures. *Lect. Notes Comput. Sci.* **2007**, *4432*, 399–406.
- (20) Kabsch, W. A Solution for the Best Rotation to Relate Two Sets of Vectors. *Acta Crystallogr., Sect. A* **1976**, *32*, 922–923.
- (21) Tsai, C. J.; Lin, S. L.; Wolfson, H. J.; Nussinov, R. A Dataset of Protein–Protein Interfaces Generated with a Sequence-Order-Independent Comparison Technique. *J. Mol. Biol.* **1996**, *260*, 604–620.
- (22) Gibrat, J. F.; Madej, T.; Bryant, S. H. Surprising Similarities in Structure Comparison. *Curr. Opin. Struct. Biol.* **1996**, *6*, 377–385.
- (23) Griffith, J. P.; Kim, J. L.; Kim, E. E.; Sintchak, M. D.; Thomson, J. A.; Fitzgibbon, M. J.; Fleming, M. A.; Caron, P. R.; Hsiao, K.; Navia, M. A. X-ray Structure of Calcineurin Inhibited by the Immunophilin–Immunosuppressant FKBP12-FK506 Complex. *Cell* **1995**, *82*, 507–522.
- (24) Moroiyanu, J. Nuclear Import and Export Pathways. *J. Cell. Biochem.* **1998**, *33*, 76–83.