

PocketMatch: A new algorithm to compare binding sites in protein structures

Yeturu Kalidas and Nagasuma Chandra*

Bioinformatics Centre and Supercomputer Education and Research Centre, Indian Institute of Science, Bangalore – 560 012, INDIA

Email: Yeturu Kalidas - kalidas@rishi.serc.iisc.ernet.in; Nagasuma Chandra* - nchandra@serc.iisc.ernet.in;

*Corresponding author

Abstract

Background: Recognizing similarities and deriving relationships among protein molecules is a fundamental requirement in present-day biology. Similarities can be present at various levels which can be detected through comparison of protein sequences or their structural folds. In some cases similarities obscure at these levels could be present merely in the substructures at their binding sites. Inferring functional similarities between protein molecules by comparing their binding sites is still largely exploratory and not as yet a routine protocol. One of the main reasons for this is the limitation in the choice of appropriate analytical tools that can compare binding sites with high sensitivity. To benefit from the enormous amount of structural data that is being rapidly accumulated, it is essential to have high throughput tools that enable large scale binding site comparison.

Results: Here we present a new algorithm *PocketMatch* for comparison of binding sites in a frame invariant manner. Each binding site is represented by 90 lists of sorted distances capturing shape and chemical nature of the site. The sorted arrays are then aligned using an incremental alignment method and scored to obtain *PMScores* for pairs of sites. A comprehensive sensitivity analysis and an extensive validation of the algorithm have been carried out. Perturbation studies where the geometry of a given site was retained but the residue types were changed randomly, indicated that chance similarities were virtually non-existent. Our analysis also demonstrates that shape information alone is insufficient to discriminate between diverse binding sites, unless combined with chemical nature of amino acids.

Conclusions: A new algorithm has been developed to compare binding sites in accurate, efficient and high-throughput manner. Though the representation used is conceptually simplistic, we demonstrate that along with the new alignment strategy used, it is sufficient to enable binding comparison with high sensitivity. Novel methodology has also been presented for validating the algorithm for accuracy and sensitivity with respect to geometry and chemical nature of the site. The method is also fast and takes about $1/250^{th}$ second for one comparison on a single processor. A parallel version on BlueGene has also been implemented.

Background

Much of present day biology is dependent on sequence-structure-function relationships in protein molecules, insights obtained for one protein heavily influencing understanding of other proteins in the family. Recognizing similarities and deriving relationships therefore is a fundamental objective in bioinformatics. Some of these similarities are obvious at the sequence level while some are detected at the structure level [1, 2]. It is in fact well established now that the conservation at the structure level of related proteins can be higher and hence much more detectable than at the sequence level [3]. In this context, there are a number of examples in the literature, which illustrate that structures often convey the ‘meaning’, more efficiently than sequences, here ‘meaning’ referring to the ‘function’ of the protein. On the other hand, there are also a number of instances, which illustrate that a particular ‘function’ is achieved by proteins whose sequences and structures are dis-similar. For example, at least three different proteins with different folds and architectures recognize mannose and exhibit mannose-mediated physiology [4]. In other words, structures also fail to convey the ‘meaning’ in many cases. We do not yet know if this failure is because of our inability to recognize any similarities in such seemingly dis-similar proteins or it is simply because no similarities actually exist among them. What ultimately matters for a protein molecule however, is its function and not what means it uses to achieve it [5]. A given function could be conserved simply by having similarities in some elements of the structure, such as the binding site residues [6]. A classic example is the large family of serine proteases which are classified into different sequence and structural families, but all come under the functional class of serine proteases due to the presence of the catalytic triad [7]. Comparison of binding sites differ from comparison of whole structures for two main reasons (a) the binding sites are small containing only a few residues and (b) these residues are often not

contiguous in sequence. Alignment of two sites containing discrete sets of atoms involves evaluation of a huge number of mappings. This makes it important to have efficient algorithms with low time and space complexities that are capable of identifying and ranking different extents of similarities appropriately. With several structural genomics projects as well as advances in computational methods for structure prediction, the structural databases are growing at a rapid pace, providing experimental structures of thousands and confident homology models of millions of protein molecules [8,9] <http://www.ncbi.nlm.nih.gov/Structure/MMDB/mmdb.shtml>. Algorithms for identification of binding pockets with reasonable confidence have also been developed [10–13]. The need for large scale comparisons of binding sites is hence accentuated. Some methods are already available for such a purpose based on ideas well established in the field of image processing. For example SitesBase [14] that uses ‘Geometric hashing’ involves selection of triads of points representing atomic types and positions in each site and comparing the triangles formed by triads; PINTS [15] uses a depth first traversal strategy, adopted to find common set of nodes between a pair of graphs connected by similar pattern of edges; Spherical harmonics based algorithm [16], captures distribution of points representing the site in terms of coefficients of a square integrable function on a unit sphere; CavBase [17] identifies maximal common sub-graphs between pairs of sites; SPASM and RIGOR [18] compares distribution of residues from the centroid of the binding sites. Each of the methods have their own merits and demerits, warranting exploration of newer methods for site comparison. Here we present a new algorithm called *PocketMatch* for representing a binding site in a frame invariant manner and comparison of pairs of sites based on alignment of sorted sequences of distances between pairs of points representing sites.

Results and Discussion

A new efficient and accurate algorithm *PocketMatch* has been developed for comparing binding sites in protein structures. A comparison between a pair of sites of an average size of 50 atoms takes only $1/250^{th}$ of a second on a single processor. The algorithm has been used for large scale database searches and all-vs-all comparisons. A typical database search of a query site against those from a large dataset comprising about 20000 sites takes in the order of 4 minutes to complete on the same machine with MPI-C version using 4 processors.

Algorithm

Comparison of a pair of binding sites involves three aspects, (a) representation of each site as sorted lists of distances between chosen points, (b) alignment of two sets of distance lists and (c) choosing a scoring scheme for arriving at a final score. Our representation scheme for the site is based on capturing the geometry of a 3D object in a 1D representation by a set of all pairs of distances between points. Such a set of distances would become a frame-invariant representation, a highly desirable scheme for any general shape comparison method. Two sets of distances can be compared for similarity by considering a suitable mapping between distances whose dissimilarity is bounded by a small amount. As the number of such mappings of distances from one set to the other would be huge, we represent them as sorted sequences in ascending order. An alignment strategy for comparison of two sequences of distances is also presented. The different steps involved are described.

Representation of binding site

Step-1 A set of residues whose one or more atoms surround a given crystallographic ligand within a specified distance (4Å as default) from each atom of the ligand is taken as the binding site. **We chose to consider atoms of the complete residues corresponding to each of these atoms for shape-representation.** For any 3D object represented by a set of points, the set of all pairs of distances between the points encodes the shape of the object. These points can then be flagged with specific chemical properties. 20 amino acids were considered in 5 groups - Group-0:(A,V,I,L,G,P);Group-1:(K,R,H);Group-2:(D,E,Q,N); Group-3:(Y,F,W);Group-4:(C,S,T). Grouping is implemented as a user defined parameter enabling use of other types of grouping.

Step-2 Represent each complete residue of the site by 3 types of points - C_α , C_β and $C_{centroid}$ corresponding to C-Alpha, C-Beta and centroid of atoms of the side chain of the residue. Centroid is computed as

$$\bar{x} = \sum_{i=1}^n \bar{x}_i$$

where \bar{x}_i indicates coordinates of the i^{th} atom of the side chain and n is the number of atoms in the side chain.

Step-3 Binning of distances and representation format.

- (a) Compute distances between all pairs of points and bin the distances into 90 sets corresponding to group-type-pair and point-type-pairs.

- (i) There are 5 types of residue-groups. Therefore $5 * (5 - 1)/2 + 5 \rightarrow 15$ pairs of groups are possible.
- (ii) There are 3 types of points. Therefore $3 * (3 - 1)/2 + 3 \rightarrow 6$ pairs of point types are possible.
- (iii) Considering both residue-group and point-type information, total number of possible sets of distances would be $15 * 6 \Rightarrow 90$.

(b) Write the above representation in the following format to a file.

$$\begin{array}{c}
 NGP \\
 NTP \\
 \\
 \left\{ \begin{array}{c} ND_1 \\ d_1, \quad d_2, \quad \dots \quad d_j, \quad \dots \end{array} \right\} \\
 \\
 \left\{ \begin{array}{c} ND_2 \\ d_1, \quad d_2, \quad \dots \quad d_j \quad \dots \end{array} \right\} \\
 \\
 \vdots \\
 \left\{ \begin{array}{c} ND_{90} \\ d_1, \quad d_2, \quad \dots \quad d_j \quad \dots \end{array} \right\}
 \end{array}$$

Where, NGP : Number of pairs of group-types, NTP : Number of pairs to point-types, ND_i : Number of distances in the i^{th} bin, d_j : distance between j^{th} pair of points.

Step-4 The distances in each list were sorted in ascending order.

Alignment of a pair of distance-sets

To compute similarity score between two binding sites, each of the 90 sets of one site has to be compared with its corresponding set of the other site. In the next step, similarity between sets was computed using a *greedy strategy* for alignment of sorted distance-sequences.

Step-5: Let the given pair of sorted(ascending) distance lists be S_1 and S_2 where each element is indexed by $S[i]$ for the list, S . Let the threshold for alignment of two distances be τ i.e., $|S_1[i] - S_2[j]| \leq \tau$. Let $m = |S_1|$ and $n = |S_2|$ denote the cardinalities of the sets S_1 and S_2 respectively. Let the variable *counter* hold the number of matched distances between two sequences. Then alignment strategy is as indicated in the (*Sub-routine 1*) below.

Sub-routine 1 Alignment of a pair of sorted distance sequences

```
i=0; j=0; counter=0;
while (i ≤ m) ∧ (j ≤ n) do
  if |S1[i] - S2[j]| ≤ τ then
    i ← i + 1; j ← j + 1
    counter ← counter + 1
  else
    if S1[i] < S2[j] then
      i ← i + 1;
    else
      j ← j + 1;
    end if
  end if
end while
```

Scoring for similarity between sites

Step-6: Scoring of the alignment: An alignment score between a pair of sites is the net average of the number of matching distance elements in the 90 lists as a fraction of total number of distance elements in the bigger set, for a chosen threshold τ , as shown

$$PM\text{Score} = \frac{\sum_{i=1}^{90} \text{Count}_i}{\text{maximum}(|S_1|, |S_2|)}$$

where $|S|$ indicates cardinality of the set.

This measure of similarity, referred to as *PMScore* is used as the default scoring scheme in this study. A variant of this scoring scheme in which the denominator was taken as the $\text{minimum}(|S_1|, |S_2|)$, was also explored, since it emphasises local structural similarity ignoring the relative size mismatch between the two sites being compared. However we observed that the comparison was too insensitive with such a scheme and did not use it further. τ is an important parameter that governs the alignment, which decides when a given pair of distances constitute a match. Different values of τ (please see section on sensitivity analysis) were tried and a default of 0.5\AA was used in this study.

Testing

PocketMatch has been validated for two aspects, (a) first to validate how well known similarities are being reproduced and (b) how sensitive is the algorithm with respect to minor perturbations in geometry and residue types at the site. Comparison with known binding sites that are similar to each other is achieved by three ways (i) by considering proteins belonging to the same SCOP classification up to the family level, since such proteins are likely to have similar binding sites, given high similarity in their overall folds and their inferred homology within a family (ii) by comparing sites for the same ligand of multiple subunits in

the same protein in a dataset of tetrameric proteins and finally (iii) by performing an all-vs-all comparison in a curated set of 51 sites for 4 ligand types corresponding to 27 proteins described in Table 1 and testing if the sites for each of the ligands cluster separately.

Based on SCOP classification:

399171 all-pair-site comparisons in PDBBind(Table 1-i) was carried out. Figure 1(a) illustrates XOR between the *PMScore* and the SCOP matrices as a function of *PMScore* threshold. As expected the XOR values start at a value close to 1 when the threshold is very low (at 0%) and inch towards 0 when the thresholds are increased to 100%. A sharp decline in XOR can be seen at around a threshold of 40%. This means that when a pair of sites match with a score of 40% or higher, they exhibit the same SCOP classification. At lower thresholds, for example at 20%, i.e., when two sites are considered similar to each other even when they match to only 20%, the correspondence of such a metric with the SCOP matrix is poor. This analysis serves to validate two points, (a) that higher *PMScores* are reflective of true-positives in the context of the family level SCOP classification and (b) that lower *PMScores* are reflective of true-negatives by the same logic. The graph also shows that a score of 40% is discriminatory enough to identify a true-positive. Figure 1(b) illustrates superposed matrices, at a threshold of 50% for which the dis-similarity and similarity scores are 0.0284 and 0.9716 respectively.

Similarities in tetramers:

Next, we tested the *PocketMatch* algorithm by comparing multiple binding sites of the same ligand in the same protein. The ligand sites in the different subunits in the tetramer dataset (Table 1-iv) were compared. In all 174,372 comparisons were carried out for the 11301 sites in 1525 tetramers. In most subunits, binding sites match with sites of other subunits of the same protein with a score of 90% to 100%(Figure 2). The two validation tests demonstrate high prediction accuracies, both in terms of assigning high scores to similar sites and low scores to dis-similar sites.

Using a curated dataset of known similarities:

A dataset of 21 proteins containing 57 sites for ligand types CIT, MTX, MK1 and PGA (Table 1-iii) was curated and all vs all comparisons were carried out. The sites were then clustered based on their extent of dissimilarities. As evident from Figure 3 illustrating the cluster tree, sites corresponding to each ligand type cluster separately forming distinct groups among them. This serves to demonstrate that similarities

among sites for a given ligand type are high in different proteins and the similarity scores are also sufficient to discriminate between sites for other ligand types.

Sensitivity analysis

Sensitivity with respect to random perturbation of positions of site points:

The basic idea behind this validation is to understand the sensitivity of our algorithm, with respect to minor perturbations in the positions of binding site residues. This would encompass scenarios where the proteins being analyzed (a) have a minor error in the crystallographically determined coordinates, or (b) where the given site residue is flexible or (c) where the structure in hand is a homology model with minor changes in the site, such as due to mutations or binding of slightly different ligands. In all these cases the nature of the binding sites of the corresponding proteins is essentially the same and should be identified as similar by any site comparison algorithm. We chose two ligands of different sizes : phenylpropanoic acid methyl ester(PP8) consisting of 54 atoms and methotrexate(MTX) consisting of 34 atoms. We perturbed the sites to varying extents within bounds of 0 to 5Å about 1000 times for each site individually. We then computed *PMScores* for the original and the perturbed pair in each instance. A plot of *PMScore* vs RMSD between the sites is shown in Figures 4 (a), (b) for the two ligands. Four different scoring schemes were used which differed in the τ parameter. While testing for sensitivity, this analysis also tests for an appropriate definition of similarity.

For a small perturbation of say $< 2.0\text{\AA}$ the *PMScores* between the original and the perturbed site was seen to be very high with all scoring schemes for both the ligands. This indicates that the algorithm is robust enough to recognize two sites as similar even when minor perturbations to the extent of 2\AA are present. When the perturbation was increased to 5\AA the *PMScores* showed lower values as expected. With scoring scheme-1(green, $\tau = 1.0$) which was most liberal in terms of its similarity definition, the *PMScores* varied very little even with high perturbation where as when scheme-5(yellow, $\tau = 0.01$) was used, *PMScores* moved to near zero even for small perturbations of about 0.1\AA . The 2nd, 3rd and 4th scoring schemes corresponding to τ of 0.5, 0.25 and 0.125 respectively, showed in-between trends with a reasonable balance between sensitivity and robustness. The figure shows that the trend followed by the various schemes is consistent for different sizes of ligands. Given that the average coordinate error in crystal structures is in the order 0.5\AA we have used the scoring scheme-2(red) as the default scheme for large scale analysis. τ is however implemented as a user defined parameter. This analysis indicates that the scores obtained are (a) reflective of the extent of similarities, (b) resistant to minor perturbations in the site, (c) scoring schemes

are self consistent validating the basic logic used in the algorithm and (d) perturbed sites where some atoms have moved even up to 5Å are recognizable as similar to their original sites, albeit with lower *PMScores*, because of retaining the overall nature of the site and high similarities with respect to remaining atoms in the site.

Sensitivity with respect to random perturbation of residue types of site points:

Given that spatial arrangement of specific amino acid residues at a given binding site dictates its recognition properties, we felt it was important to test the sensitivity of our algorithm to perturbations in the nature of the residues in the binding site without disturbing their spatial arrangement. We carried out this analysis on a pair of sites which were known to be similar to each other and another pair of sites which were known to be dissimilar to each other. To minimize site-specific biases that may arise during comparison, we chose one protein to be common between the two sites and all three to be nucleotide binding sites. One of the sites was kept constant while the other was perturbed. Of a possible 5^N perturbations for a site a site of N residues, 1000 random perturbations were carried out and computed the *PMScore* of each of the perturbed sites with respect to the original site. The *PMScores* for the unperturbed sites of the two chosen pairs (1H8H-ATP and 1W0K-ADP) and (1H8H-ATP and 1H8H-ADP) were 80.9% and 25.8% respectively whose superpositions are shown in Figures 5 (a) and 5 (b). The distribution of new *PMScores* over the set of perturbed sites with respect to the original site is shown in Figure 6(a) for the first pair and Figure 6(b) for the second pair. Figure 6(a) shows that when perturbed, the similarity between two sites disappears and the scores get poorer. Figure 6(b) shows a similar trend indicating no high similarities were seen during random perturbations. Both these suggest that the arrangement of the site is specific and has been derived for a purpose and not just by chance. Our algorithm is sufficiently sensitive to detect changes in the nature of residues at the binding site. Our current implementation considers 20 amino acids classified into 5 groups and any change within the same group will obviously not be detected here. However the classification type is also implemented as a user defined parameter so as to consider each amino acid as a separate group to overcome this loss of sensitivity where more stringent analysis is required.

All vs all comparison - Site geometry vs residue types

In order to estimate the sensitivity of our algorithm in terms of finding similarities for a given site in a large dataset, we performed an all vs all comparison in the PDBBind dataset Table 1-ii. The distribution of

PM*Scores* across all pairs is shown in Figure 7(a). As evident from the histogram, majority of the pairwise scores are in the range of 0 to 40%. This means that any randomly chosen pair of two different sites will only have a score less than 40%. The histogram shows that only a small percentage of 0.8% of all possible pairs have scores higher than 50% which are indicative of true positives. Examination of 657231 such pairs indeed shows that they are true positives and in fact belong to the same SCOP class indicating that they are similar not only in their binding site architectures but also in their overall folds. This also shows that there are no false positives. It must be noted that this analysis was carried out with the scoring scheme, which we recommend as default that uses amino acid group information apart from the sorted distances. The same analysis was then repeated with another scoring scheme that differed from the previous one by not differentiating among amino acid types. Such a scoring scheme would identify similarities purely based on geometric features of the binding sites without considering their chemical nature explicitly. The distribution of PM*Scores* by this scheme is shown in Figure 7(b). Surprisingly, the same dataset showed a very different distribution of scores as compared to the previous analysis in Figure 7(a). Here a majority of pairs exhibited high scores indicating that many of the sites appear similar, obviously leading to a large number of false positives. However, when their chemical group information is added, differences between types of sites emerge. This analysis shows that shapes by themselves without considering chemical information is not sufficient to discriminate between different types of sites.

Implementation

The software *PocketMatch* was developed on *gcc (GCC) 3.4.3 20041212version 4.1.2* on Linux 2.6.9-5.ELsmp machine. A parallelized version of the software using MPI-C libraries was also developed and implemented both on a standard Quad-core machine (Intel (R) Core (TM) 2 Quad CPU @ 2.4GHz; Address space : 32 bits physical, 48 bits virtual; 4GB RAM) and also on a IBM BlueGene cluster and tested with 512 processors. *Matlab* version 7.1 from Mathworks <http://www.mathworks.com/> was used for generating plots and histograms. Cluster tree was constructed using neighbour-joining method of *phylip-3.67* <http://evolution.genetics.washington.edu/phylip.html> and viewed using *PhyloDraw* [19]. *PyMol* www.pymol.org version 0.99rc1 was used for visualizing various structures.

Methods

Datasets

PocketMatch has been validated on a variety of datasets - PDBBind, a set of tetrameric proteins and a curated dataset containing known similarities of four ligand types (Table 1). To eliminate noise in the datasets, sites corresponding to *small ligands* with less than 6 non-hydrogen atoms and *covalently bound* ligands were not considered here. The PDBBind dataset [20] contains a comprehensive curated set of 1091 protein-ligand complexes determined crystallographically. Using this, two sub-datasets were derived - Table 1-i: A dataset which has only one ligand site for one ligand type in each protein, amounting to 786 proteins and 893 sites and Table 1-ii: A dataset in which ligands suggested to contribute to noise by Jackson and co-workers [14] were removed, but all sites for all ligand types were considered, amounting to 456 proteins and 1146 sites corresponding to 289 ligand types. Tetramers (Table 1-iv), obtained from PQS server, containing 3768 proteins has been curated to yield a dataset of 1525 proteins having 11301 sites has been chosen for studying the sensitivity of *PocketMatch* with respect to recognizing known highly similar sites. Another dataset representing multiple sites for four known ligands was curated (Table 1-iii). 51 sites from 27 different proteins in PDBBind for Citrate (CIT), Methotrexate (MTX), Indinavir (MK1) and phosphoglycerate (PGA) were chosen for the dataset. A distance metric measuring the dis-similarity between these sites was computed using *PocketMatch* and their clustering was studied.

Scheme for Validation

We have validated the algorithm based on (a) SCOP equivalences between pairs of proteins at 4th level of SCOP [1], (b) similarities among multiple sites for a given ligand in the tetramer dataset Table 1-iv and (c) clusters formed by sites whose similarities have been previously identified by independent analysis. To compare with SCOP, we constructed a matrix, M_1 where each of the rows and columns correspond to list proteins and each element $M_{1_{ij}}$ corresponds to a score assigned to that pair of proteins, i and j. A score of 1.0 was assigned to those pairs that had the same SCOP class. All other pairs were assigned a score of 0. Then we constructed another matrix, M_2 of the same dimension where a score of 1.0 was assigned to those pairs whose top scoring site-pair had a *PMScore* greater than a threshold and score of 0 for all others. Given two matrices M_1 and M_2 - one for SCOP and one for *PMScores* respectively, the dissimilarity between them has been calculated as *XOR*

$$XOR = \frac{\sum_{\forall(i,j)} M_{1_{ij}}! = M_{2_{ij}}}{N^2}$$

where N is the dimension of the matrix.

Schemes for sensitivity analysis

Sensitivities of *PocketMatch* (a) with respect to geometry and (b) with respect to specificity of residue types at the site, were studied. For (a), we applied tiny perturbations to coordinates to each of the points representing sites, leading to generation of altered structures within bounded a RMSD and computed *PMScores* between pairs of perturbed and reference sites. These plots are obtained for various sizes of binding sites. For (b), we randomly assigned residue group information to each of site points and computed *PMScores* .

Method for random perturbation to positions of site-points:

Our method for random perturbation applies random displacement from a uniform distribution to each of the points C_α , C_β , and $C_{centroid}$ corresponding to each residue as described in (*Sub-routine 2*).

Sub-routine 2 Performing random perturbation to positions of points of the site

- a Let the ‘net required RMSD’ be ρ and number of intermediate randomly perturbed point-sets be K .
 - b In order to generate, K perturbed sets of points,
for EACH $\delta = 0$ till ρ in steps of $\frac{\rho}{K}$; **do**
 - Perturb the actual set of points to generate a new point-set at an RMSD of δ .
 - Let each point be represented by $\bar{x} = (x_1, x_2, x_3)$ and let the random vector be $\bar{r} = (r_1, r_2, r_3)$.
 - Initialize each r_i to $\frac{\cos(2.0*\pi*random())}{1+RAND_MAX}$; where *random()* function generates a random number between 0 and *RAND_MAX* from ‘uniform distribution’.
 - Normalise \bar{r} by reinitializing each $r_i = \frac{r_i}{\sqrt{\sum_{i=1}^3 r_i^2}}$.
 - Apply the present displacement, $\bar{r} * \delta$ to each atom which set each $\bar{x}^{new} = \bar{x}^{old} + \bar{r} * \delta$
 - These steps generate a perturbed set of points with a net RMSD of δ .**end for**
 - c For each of the perturbed point-set generated, compute RMSD with respect to the unperturbed version and measure *PMScores* .
-

Method for random perturbation of types of site-residues:

Our method for random perturbation of group numbers assigned to residues of site involves assignment of random integers between 0 and 4 from a ‘uniform’ distribution to the residues represented by three points, C_α , C_β , and $C_{centroid}$ as described in (*Sub-routine 3*).

Conclusions

A new algorithm has been developed to compare binding sites in accurate, efficient and high-throughput manner, where sites are represented as 90 lists of sorted distances flagged with residue type information. This representation therefore captures both the shape and chemical nature of amino acid types at the site.

Sub-routine 3 Random perturbation of types of site residues

- a Get C_α 's alone for a given site represented in C_α , C_β , and $C_{centroid}$ - format.
 - b To each of the C_α 's assign a random number between 0 and 4 using $\frac{5*random()}{(RAND_MAX+1.0)}$
 - c Copy back the modified group information of C_α atoms back to the file representing site by C_α , C_β , and $C_{centroid}$ points.
 - d The $PMScore$ is then computed for the perturbed site with respect to unperturbed version.
-

Extensive validation has been performed using different datasets. Sensitivity analysis has also been performed to analyze the performance of the algorithm with respect to perturbations of two types - (a) in the actual atomic positions of atoms of the site and (b) in the amino acid type at the site. Several scoring schemes have been analyzed by virtue of which a scoring function with a good balance between sensitivity and ability to detect similarities has been identified and recommended as the default scoring scheme. Perturbation studies where the geometry of a given site was retained but the residue types were changed randomly, indicated that chance similarities were virtually non-existent. Our analysis also suggests that shape information alone is insufficient to discriminate between diverse binding sites. However, combining shape information with chemical grouping of amino acids at the site enables discrimination between different types of sites.

Authors contributions

YK, a graduate student developed and implemented the method under the guidance of his advisor NSC. Both authors discussed and wrote the paper.

Availability:

The software can be accessed at <http://proline.physics.iisc.ernet.in/pocketmatch/>

Acknowledgements

Support for the Centre of Excellence in Bioinformatics by Department of biotechnology(DBT), Govt. of India and facilities at the Supercomputer Education and Research Centre of this institute are gratefully acknowledged. Financial support from DBT computational genomics initiative is also acknowledged.

References

1. Murzin AG, Brenner SE, Hubbard T, Chothia C: **SCOP: a structural classification of proteins database for the investigation of sequences and structures.** *Journal of Molecular Biology* 1995, **247**:536–540.
2. Orengo CA, Pearl FM, Bray JE, Todd AE, Martin AC, Lo Conte L, Thornton JM: **The CATH Database provides insights into protein structure/function relationships.** *Nucleic Acids Res* 1999, **27**:275–279.
3. Holm L, Sander C: **Mapping the protein universe.** *Science* 1996, **273**(5275):595–603.
4. Ramachandriah G, Chandra NR: **Sequence and structural determinants of mannose recognition.** *Proteins* 2000, **39**(4):358–364.
5. Prasad T, Subramanian T, Hariharaputran S, Chaitra HS, Chandra N: **Extracting hydrogen-bond signature patterns from protein structure data.** *Appl Bioinformatics* 2004, **3**(2-3):125–135.
6. Russell R, Sasieni P, Sternberg M: **Supersites within superfolds. Binding site similarity in the absence of homology.** *Journal of Molecular Biology* 1998, **282**:903–18.
7. Dodson G, Wlodawer A: **Catalytic triads and their relatives.** *Trends Biochem Sci* 1998, **23**(9):347–352.
8. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov I, Bourne P: **The Protein Data Bank.** *Nucleic Acids Research* 2000, **28**:235–242.
9. Pieper U, Eswar N, Braberg H, Madhusudhan MS, Davis FP, Stuart AC, Mirkovic N, Rossi A, Marti-Renom MA, Fiser A, Webb B, Greenblatt D, Huang CC, Ferrin TE, Sali A: **MODBASE, a database of annotated comparative protein structure models, and associated resources.** *Nucleic Acids Res* 2004, **32**(Database issue):217–222.
10. Glaser F, Morris RJ, Najmanovich RJ, Laskowski RA, Thornton JM: **A method for localizing ligand binding pockets in protein structures.** *Proteins* 2006, **62**(2):479–488.
11. Laurie AT, Jackson RM: **Q-SiteFinder: an energy-based method for the prediction of protein-ligand binding sites.** *Bioinformatics* 2005, **21**(9):1908–1916.

12. Huang B, Schroeder M: **LIGSITEcsc: predicting ligand binding sites using the Connolly surface and degree of conservation.** *BMC Struct Biol* 2006, **6**:19–19.
13. Kalidas Y, Chandra N: **PocketDepth: A new depth based algorithm for identification of ligand binding sites in proteins.** *Journal of Structural Biology* 2007, **1**:31–42.
14. Gold ND, Jackson RM: **Fold Independent Structural Comparisons of Protein-Ligand Binding Sites for Exploring Functional Relationships.** *Journal of Molecular Biology* 2006, **355**:1112–1124.
15. Stark A, Russell RB: **Annotation in three dimensions. PINTS: Patterns in Non-homologous Tertiary Structures.** *Nucleic Acids Research* 2003, **31**:3341–3344.
16. Morris RJ, Najmanovich RJ, Kahraman A, Thornton JM: **Real spherical harmonic expansion coefficients as 3D shape descriptors for protein binding pocket and ligand comparisons.** *Bioinformatics* 2005, **21**:2347–2355.
17. Kuhn D, Weskamp N, Schmitt S, Hüllermeier E, Klebe G: **From the similarity analysis of protein cavities to the functional classification of protein families using cavbase.** *J Mol Biol* 2006, **359**(4):1023–1044.
18. Kleywegt GJ: **Recognition of Spatial Motifs in Protein Structures.** *Journal of Molecular Biology* 1999, **285**:1887–1897.
19. Choi JH, Jung HY, Kim HS, Cho HG: **PhyloDraw: a phylogenetic tree drawing system.** *Bioinformatics* 2000, **16**(11):1056–1058.
20. Wang, R, Fang, X, Lu, Y, Wang, S: **The PDBBind Database: Collection of Binding Affinities for Protein-Ligand Complexes with known Three-dimensional Structure.** *J. Med. Chem* 2004, **47**:2977–2980.

Figures

Figure 1 - Validation with respect to SCOP

Validation of *PocketMatch* against SCOP. (a) Plot showing the XOR values for different *PMScore* matrices varying in threshold from 0 to 100 percent; (b) Superposition of *PMScore* (red) and SCOP(blue) matrices. *PMScore* matrix was computed for 50%-threshold. Corresponding dis-similarity(XOR) and similarity scores are 0.0284 and 0.9716

Figure 2 - Validation with respect to tetramers

Histogram showing *PMScores* in the tetramer dataset

Figure 3 - Validation with respect to selected set of ligands

Cluster-tree showing similarities in an all vs all comparison of 51 sites for four ligands in dataset-iii (Table 1). The nodes are labelled by their corresponding PDB codes followed by the ligand code.

Figure 4 - Validation with respect to random perturbation of positions of site-points

Random perturbations of site points for (a) a 54 atom ligand(PP8) and (b) a 34 atom ligand(MTX). *PMScores* for perturbed sites with respect to its original site for different extents of perturbations(RMSD) are shown at different values of τ (1.0-green,0.5-red,0.25-cyan,0.125-blue,0.01-yellow)

Figure 5 - Superposition of two ATP/ADP sites

Examples illustrating validation of *PocketMatch*: Superposition of sites with (a) High *PMScores* (80.9% for 1H8H-ATP and 1W0K-ADP) and (b) low *PMScores* (25.8% for 1H8H-ATP and 1H8H-ADP)

Figure 6 - Validation with respect to random perturbation of types of site-residues

Perturbation analysis for the examples shown in Figure 5 (a) for a pair of high scoring sites and (b) for a pair of low scoring sites. Distribution of *PMScores* for 1000 randomly perturbed sites with respect to their original site is shown in both cases.

Figure 7 - Importance of residue group information over pure geometry

All vs all comparison of 1146 sites of PDBBind dataset-ii(Table 1) using (a) default scoring scheme that uses both geometry and residue type information and (b) a scoring scheme that uses only site geometry. Histograms indicate distribution of number of pairs showing different *PMScore* values.

Table 1 - Description of the dataset

No.	Dataset	Initial Size	Remarks	Final set used in the study	
				Proteins	Sites
i	PDBBind	1091	One representative ligand of each type is considered in each protein; filtered to remove sites for small and covalently bound ligands; retained proteins common with the SCOP database	786	893
ii	PDBBind	1091	Filtered to remove sites for small , covalently bound ligand and all suggested by Jackson and co-workers (10) to contribute to noise; considered sites for all ligand types for the remaining	456	1146
iii	Curated dataset of CIT,MTX,MK1 and PGA	27	Ligands for varying sizes and types from PDBBind with multiple sites for a ligand in different proteins	27	51
iv	Tetramer	3768	Multiple sites for each ligand in the same protein; filtered to remove small and covalently bound ligands	1525	11301

Table 1. Different datasets used in the study

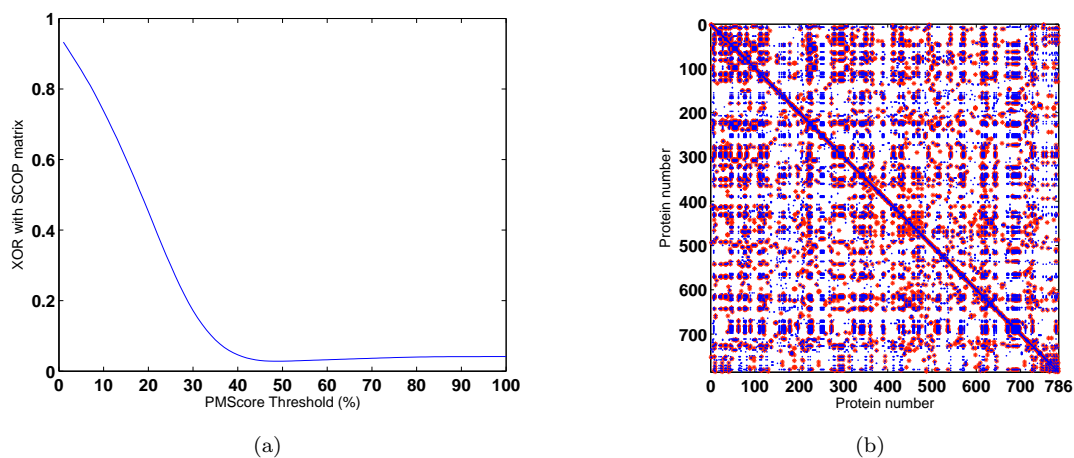


Figure 1:

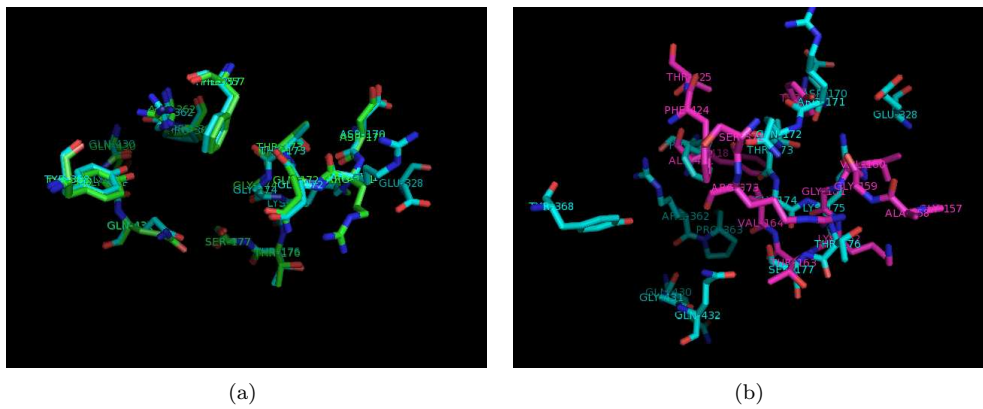


Figure 5:

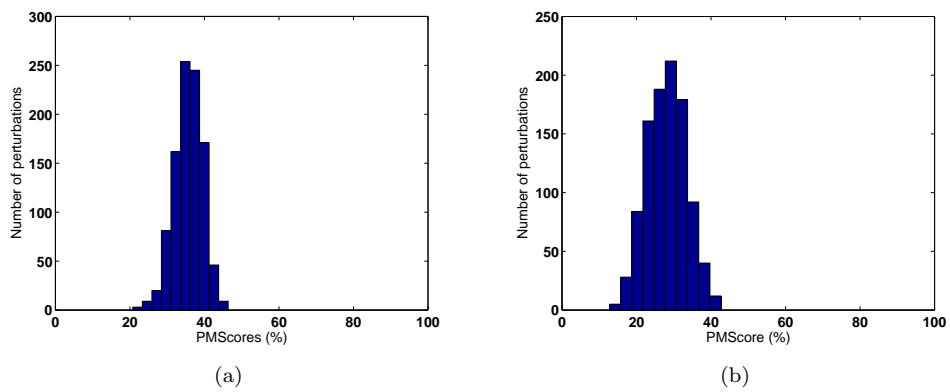


Figure 6:

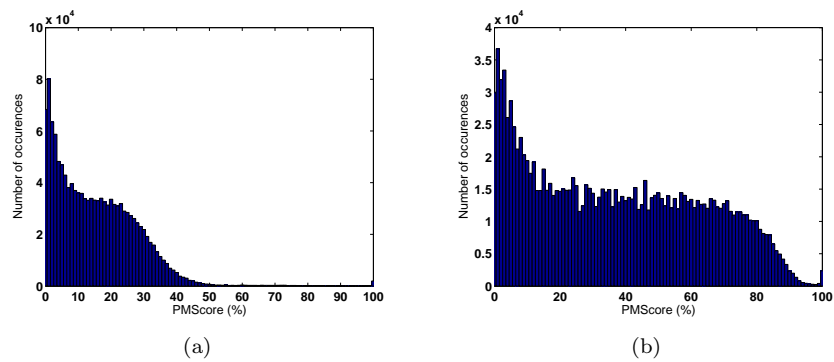


Figure 7: