

A peer-reviewed version of this preprint was published in PeerJ on 3 February 2020.

[View the peer-reviewed version](https://peerj.com/articles/cs-253) (peerj.com/articles/cs-253), which is the preferred citable publication unless you specifically need to cite this preprint.

Sadique N, Ahmed AAN, Islam MT, Pervage MN, Shatabda S. 2020. Image-based effective feature generation for protein structural class and ligand binding prediction. PeerJ Computer Science 6:e253 <https://doi.org/10.7717/peerj-cs.253>

Image based effective feature generation for protein structural class and ligand binding prediction

Nafees Sadique ^{Equal first author, 1}, **Al Amin Neaz Ahmed** ^{Equal first author, 1}, **Md Tajul Islam** ¹, **Md. Nawshad Pervage** ¹, **Swakkhar Shatabda** ^{Corresp. 1}

¹ Department of Computer Science and Engineering, United International University, Dhaka, Dhaka, Bangladesh

Corresponding Author: Swakkhar Shatabda
Email address: swakkhar@cse.uui.ac.bd

Proteins are the building blocks of all cells in both human and all our living creatures of the world. Most of the work in the living organism is performed by Proteins. Proteins are polymers of amino acid monomers which are biomolecules or macromolecules. The tertiary structure of protein represents the three-dimensional shape of a protein. The functions, classification and binding sites are governed by protein's tertiary structure. If two protein structures are alike then the two proteins can be of the same kind implying similar structural class and ligand binding properties. In this paper, we have used protein structure to generate effective features for applications in structural similarity to detect structural class and ligand binding. Firstly, we analyze the effectiveness of a group of image based features to predict the structural class of a protein. These features are derived from the image generated by the distance matrix of the tertiary structure of a given protein. They include local binary pattern histogram, Gabor filtered local binary pattern histogram, separate row multiplication matrix with uniform local binary pattern histogram, neighbour block subtraction matrix with uniform local binary pattern histogram and atom bond. The experiments were done on a standard benchmark dataset. We have demonstrated the effectiveness of these features over a large variety of supervised machine learning algorithms. Experiments suggest Random Forest is the best performing classifier on the selected dataset using the set of features. We believe the excellent performance of Hybrid LBP in terms of accuracy would motivate the researchers and practitioners to use it to identify protein structural class. To facilitate that, a classification model using Hybrid LBP is readily available for use at <http://brl.uui.ac.bd/PL/>.

Protein-Ligand binding is accountable for managing the tasks of biological receptors that helps to cure diseases and many more. So, binding prediction between protein and ligand is important for understanding a protein's activity or to accelerate docking computations in virtual screening-based drug design. Protein-Ligand Binding Prediction requires three-dimensional tertiary structure of the target protein to be searched for ligand binding. In this paper, we've proposed a supervised learning algorithm for predicting Protein-Ligand Binding which is a Similarity-Based Clustering approach using the same set of features. Our algorithm works better than most popular and widely used machine learning algorithms

1 Image Based Effective Feature Generation 2 for Protein Structural Class and Ligand 3 Binding Prediction

4 Nafees Sadique^{1,*}, Al Amin Neaz Ahmed^{1,*}, Md Tajul Islam¹, Md.
5 Nawshad Pervage¹, and Swakkhar Shatabda¹

6 ¹Department of Computer Science and Engineering, , United International University,
7 Bangladesh*

8 Corresponding author:

9 Last Author¹

10 Email address: swakkhar@cse.uui.ac.bd

11 ABSTRACT

12 Proteins are the building blocks of all cells in both human and all our living creatures of the world.
13 Most of the work in the living organism is performed by Proteins. Proteins are polymers of amino acid
14 monomers which are biomolecules or macromolecules. The tertiary structure of protein represents the
15 three-dimensional shape of a protein. The functions, classification and binding sites are governed by
16 protein's tertiary structure. If two protein structures are alike then the two proteins can be of the same
17 kind implying similar structural class and ligand binding properties. In this paper, we have used protein
18 structure to generate effective features for applications in structural similarity to detect structural class
19 and ligand binding. Firstly, we analyze the effectiveness of a group of image based features to predict
20 the structural class of a protein. These features are derived from the image generated by the distance
21 matrix of the tertiary structure of a given protein. They include local binary pattern histogram, Gabor
22 filtered local binary pattern histogram, separate row multiplication matrix with uniform local binary pattern
23 histogram, neighbour block subtraction matrix with uniform local binary pattern histogram and atom bond.
24 The experiments were done on a standard benchmark dataset. We have demonstrated the effectiveness
25 of these features over a large variety of supervised machine learning algorithms. Experiments suggest
26 Random Forest is the best performing classifier on the selected dataset using the set of features. We
27 believe the excellent performance of Hybrid LBP in terms of accuracy would motivate the researchers
28 and practitioners to use it to identify protein structural class. To facilitate that, a classification model using
29 Hybrid LBP is readily available for use at <http://brl.uui.ac.bd/PL/>.

30 Protein-Ligand binding is accountable for managing the tasks of biological receptors that helps to
31 cure diseases and many more. So, binding prediction between protein and ligand is important for
32 understanding a protein's activity or to accelerate docking computations in virtual screening-based drug
33 design. Protein-Ligand Binding Prediction requires three-dimensional tertiary structure of the target
34 protein to be searched for ligand binding. In this paper, we've proposed a supervised learning algorithm
35 for predicting Protein-Ligand Binding which is a Similarity-Based Clustering approach using the same set
36 of features. Our algorithm works better than most popular and widely used machine learning algorithms.

37 INTRODUCTION

38 Protein tertiary structure comparison is very important in many applications of modern structural biology,
39 drug design, drug discovery, in studies of protein-ligand binding, protein-protein interactions and other
40 fields. This is especially significant because the structure of a protein is more protected than the protein
41 sequence (Chothia and Lesk, 1986). Many works have been done to find protein binding (Brady and
42 Stouten, 2000). Comparison of protein structure has been done in many works of literature by alignment
43 of distance matrices (Holm and Sander, 1993), using iterated double dynamic programming (TAYLOR,

*First two authors contributed equally

44 1999), using elastic shape analysis (Srivastava et al., 2016) and many other techniques. The most common
45 way of comparing protein tertiary structure is to treat the protein as a three-dimensional object and
46 superimpose one on another. Different distances are used to calculate the differences between the proteins.

47 The distance matrix of α carbon can be seen extensively used in (Holm and Sander, 1997; Singh and
48 Brutlag, 1997) as a feature which represents the tertiary structure of a protein chain. This feature is used
49 as a feature vector which represents the structure of a protein to measure either similarity or dissimilarity
50 to measure and compare the feature vectors with one another in pattern recognition literature. A mapped
51 two-dimensional feature matrix is created from the 3D coordinate data of protein. The intra-molecular
52 distance is used to make the α carbon distance matrix which mirrors the tertiary structure of a protein and
53 the conserved elements of the secondary structure in it. With an input matrix size of $N \times N$, the distance
54 matrix based exact algorithms run in $O(N)$ time (Karim et al., 2015).

55 An image is basically a matrix of $N \times N$ dimension with corresponding data in each cell. Thus the
56 distance matrix can be used as an image. Basically, three types of features can be generated from an image:
57 pixel based, filter based and computationally generated features. Pixel-based features e.g histograms
58 are simplistic and dependent on the capability of classification algorithms. Filter based methodologies
59 transform the original image to use feature extraction methods. Refined algorithms are used to segment
60 and other various algorithms are used to detect different features. Using ideas from computer vision and
61 utilizing it in protein structure retrieval is not uncommon in the field. ProteinDBS server (Shyu et al.,
62 2004) implement a similar approach in (Chi et al., 2005) by Chi et al. Texture features from the original
63 size images and diagonally partitioned images were extracted by Chi et al. CoMOGrad and PHOG (Karim
64 et al., 2015) also used images to extract their two novel feature whereas we are extracting histograms of
65 local binary pattern images from the original image.

66 Human body uses protein for repairing tissues, making enzymes, hormones, and other biological
67 chemicals. It is an essential building block of bones, muscles, cartilage, skin, and blood. On the
68 other hand, a ligand is a material that has the potentiality to bind to and forms a composite with a
69 biomolecule in order to carry out a biological function. In Protein-Ligand Binding, the ligand is usually
70 a molecule which produces a signal by binding to a locus on a target protein. The binding typically
71 results in a change of conformational isomerism (conformation) of the target protein. The evolution
72 of the protein's responsibility depends on the development of specific sites which are designed to bind
73 ligand molecules. Ligand binding ability is important for the management of biological functions. Ligand
74 binding interactions changes the protein state and function. Protein-Ligand Binding prediction is very
75 important in many applications of modern structural biology, drug design, drug discovery and other fields.

76 Many experimental techniques can be used to investigate various aspects of protein-ligand binding.
77 X-ray crystallography, nuclear magnetic resonance(NMR), Laue X-ray diffraction, small-angle X-ray
78 scattering, and cryo-electron microscopy provide atomic-resolution or near-atomic-resolution structures
79 of the unbound proteins and the protein-ligand complexes, which can be used to study the changes
80 in structure and/or dynamics between the free and bound forms as well as relevant binding events.
81 Although experimental techniques can investigate thermodynamic profiles for a ligand-protein complex,
82 the experimental procedures for determination of binding affinity are laborious, time-consuming, and
83 expensive. Modern rational drug design usually involves the HTS of a large compound library comprising
84 hundreds or thousands of compounds to find the lead molecules, but this is still not realistic using
85 experimental methods alone. Different methods like Isothermal Titration Calorimetry (ITC) (Chaires,
86 2008), Surface Plasmon Resonance (SPR) (Patching, 2014), Fluorescence Polarization(FP) (Rossi and
87 Taylor, 2011), Protein-Ligand Docking (Sousa et al., 2013), Free Energy Calculations (Steinbrecher and
88 Labahn, 2010), etc are being used to predict ligand-binding prediction.

89 In this paper, we propose the combination of local binary pattern histogram, Gabor Filtered Local
90 Binary Pattern Histogram, Separate Row Multiplication Matrix with Uniform Local Binary Pattern
91 Histogram, Neighbour Block Subtraction Matrix with Uniform Local Binary Pattern Histogram and
92 Atom Bond features to be used for protein similarity measurement. We extract the distance matrix of
93 α carbon of a protein from PDB file and use the distance matrix as an image to extract our first four
94 features and Atom Bond is extracted from the PDB files. We have used a large variety of classification
95 algorithms to test the extracted features. We are also going to show the results and comparative study of
96 different implementation methodologies such as wavelet and pyramid histogram based features (Ahmed
97 et al., 2019) and CoMOGrad and PHOG. The method we have proposed is able to produce a better
98 result on some classification algorithm over the previous methods on the same benchmark. In addition

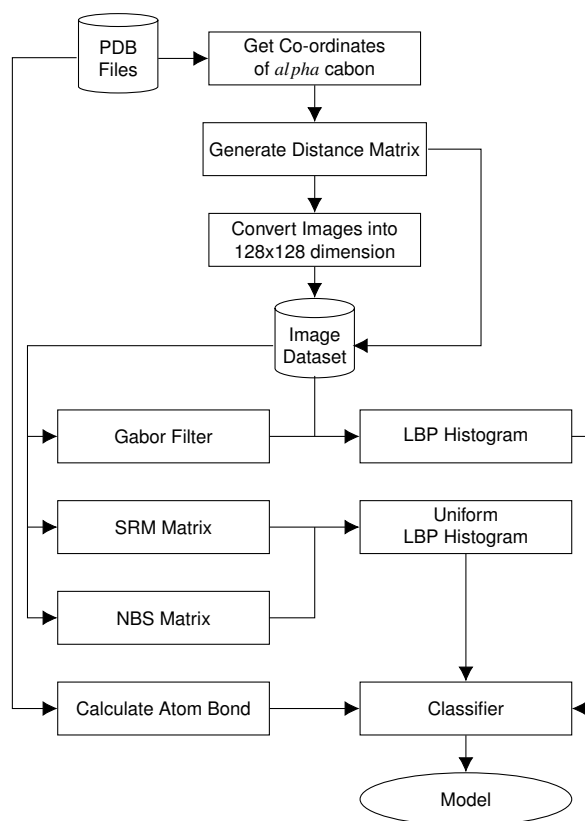


Figure 1. Block diagram of the methodology used in structural class prediction.

99 to that, we've proposed a supervised learning algorithm for predicting Protein-Ligand Binding which
 100 is a Similarity-Based Clustering approach using the same set of features. Our algorithm works better
 101 than most popular and widely used machine learning algorithms. Our proposed method uses the features
 102 proposed in this paper.

103 MATERIALS AND METHODS

104 Our methodology is divided into two parts. Firstly, we have generated image based features using protein
 105 tertiary structures and performed feature analysis based on the prediction power on the structural class
 106 prediction problem. In this section, we present the materials and methods for both of the problems.
 107 For each of the problems the dataset, features, necessary algorithms and performance measurement is
 108 described accordingly.

109 Structural Class Prediction

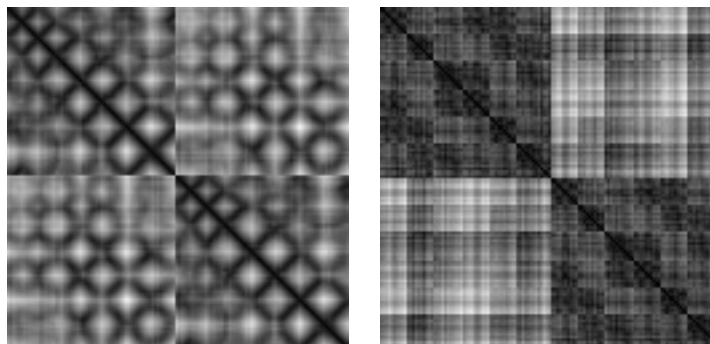
110 In this section, we present the methodology on structural class prediction. Atom bond features are
 111 generated from the protein tertiary structures given as PDB files. Images are created from the distance
 112 matrix calculated using α carbon atom coordinates of the amino acids of the protein structures in the
 113 given dataset. From each image of protein, we have derived five features. There are in total seven different
 114 classes of protein structures. Synthetic minority over-sampling technique (SMOTE) is used to handle
 115 class imbalance problem. K -fold cross-validation with three fold was used to test the capability and
 116 efficiency of the dataset. The block diagram of the methodology is given in Figure 1.

117 Structural Class Prediction Dataset

118 We have used 40 percent ID filtered subset of PDB-style files for SCOPe domains version 2.03 (Fox et al.,
 119 2013) as our dataset. It contains a total of 12119 PDB files. Each PDB files contains SCOP(e) concise
 120 classification string (sccs) which respectively describes class, fold, superfamily, and family. In this paper,
 121 we are going to experiment only with the class of the protein. In the dataset, there are total seven protein

Table 1. Protein Classes and its Corresponding Instances

Class Name	Total Instances
Small Proteins	640
All α Proteins	2195
α and β proteins(a/b)	3305
α and β proteins(a+b)	3006
Membrane and cell surface proteins and peptides	204
All β proteins	1485
Multi-domain proteins(α and β)	219

**Figure 2.** Sample images of protein structures after rescaling.

122 structural classes. For benchmark analysis with CoMOGrad and Phog, the common pdb files were used as
 123 dataset. The common PDB files are total of 11052. The details of the protein structural classes are given
 124 in Table 1. This dataset is widely used as a benchmark in the literature for protein structural similarity
 125 prediction (Karim et al., 2015).

126 **Image Generation**

127 We have generated images of protein structures according to the methodology described in CoMOGrad
 128 and PHOG (Karim et al., 2015). Only α carbons of the amino acids in the protein structure are considered
 129 for image generation. From the three dimensional coordinates of the α carbon atoms a distance matrix
 130 is generated by taking the Euclidean distance among all pairs. Thus only half of the image contains
 131 redundant information due to symmetry.

132 **Scaling of Images**

133 The dimension of protein images is based on the total number of α carbon they have. So, every individual
 134 protein images are different from the other in dimension. Therefore, the images were scaled to the same
 135 dimension. CoMOGrad and PHOG have used Bi-cubic interpolation and wavelet transform to scale all
 136 the protein images into 128 x 128 dimension (Karim et al., 2015). During the Bi-cubic interpolation step,
 137 most of the images were in 128x128 dimension so in the wavelet transform step they scaled all the images
 138 to that dimension. Thus, we have directly scaled the images to 128x128 dimension. We have used both
 139 real and scaled images to examine the differences in their predictive power. Sample rescaled images of
 140 protein structures are given in Figure 2.

141 **Feature Extraction**

142 We have generated five different feature groups. Our first four feature groups are different types of
 143 histograms and the fifth feature group is about the prognosis of the atoms. The histograms were taken
 144 from both scaled and unscaled images.

145 **Local Binary Pattern Histogram** Local binary pattern (LBP) histogram was first proposed by Ojala
 146 et al. (1994) and popularized by the work of Ojala et al. (2002). Local binary pattern computes the
 147 local representation of the texture of an image as a texture descriptor. Comparing each pixel with its

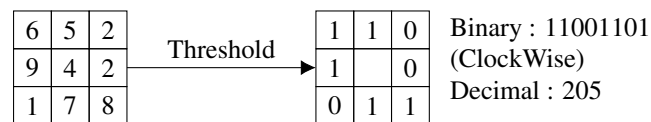


Figure 3. An example of basic LBP

148 neighboring pixels the local representation is created. The image is transformed into a grayscale image. In
 149 a 3×3 neighborhood, the center pixel value is calculated by comparing with its eight neighboring pixels.
 150 Each comparison gives a result of either 0 if the center pixel value is greater then the comparing neighbor
 151 pixel or 1 for the latter. A clockwise direction starting from the top-left one provides a binary number. The
 152 binary number is converted to a decimal number and the value is placed in the center pixel. LBP codes or
 153 Local Binary Patterns are the obtained binary numbers. An example of a basic Local Binary Pattern is
 154 given in Figure 3. After calculating the value for each pixel of the image, a histogram is calculated. A 3×3
 155 neighborhood has $2^8 = 256$ possible patterns, thus the values range from 0 to maximum 255 in each
 156 pixel of the image. The total number of bins of the histogram is thus 256. We would get 256 attributes
 157 from each image. We have used zero padding technique to generate local binary pattern.

158 **Gabor Filtered Local Binary Pattern Histogram (GfLBP-Hist)** Gabor Filter is titled after Dennis
 159 Gabor. It is used for texture segmentation (Jain and Farrokhnia, 1991), optical character recognition (Jain
 160 and Bhattacharjee, 1992), edge detection (Mehrotra et al., 1992) etc. It is a linear filter which examines if
 161 there is any particular frequency content in the image in specific areas in a localized region throughout the
 162 point. The multiplication of a sinusoid and a Gaussian is called the Gabor filter (Eq.1).

$$g(x, y; \lambda, \theta, \phi, \gamma) = \exp\left(-\frac{x'^2 + \gamma^2 y'^2}{2\sigma^2}\right) \cos\left(2\pi\frac{x'}{\lambda} + \phi\right) \quad (1)$$

163 .
 164 Here, λ controls the wavelength of this sinusoid, θ is the angle of the normal to the sinusoid, ϕ is the
 165 phase shift of the sinusoid, γ controls the aspect ratio, The spatial envelope or the standard deviation of
 166 the Gaussian is σ . For our experiments, we have used $\lambda = 10$, $\theta = 0$, $\phi = 0$, $\gamma = 0.02$ and $\sigma = 5$. After
 167 applying the Gabor filter, LBP techniques are applied to the image to get 256 attributes.

168 **Atomic Bond Features** First of all, we've identified unique atoms amidst all the protein PDB files.
 169 From each protein PDB file, we've counted occurrences of each atom. Then we've taken the percentage
 170 as features of each atom among all the atoms that each protein has. Then we've taken first 100 sequential
 171 atoms and used their atomic mass as the feature. Then we've counted the bond that each pair of atoms
 172 has in a particular protein using atomic distance based on a threshold value. Finally, we've taken the
 173 percentage as the feature of the bond of each unique pair of atoms among all the bonds that the protein
 174 has.

175 **Separate Row Multiplication Matrix with Uniform LBP Histogram(SRMMat-ULBP-Hist)** The image
 176 is split into 3×3 matrices. From each matrix, we get 3 rows with the dimension of 1×3 . By multiplying
 177 each row with the same 3×3 matrix, we get three result matrix consisting of 1×3 dimension. Each cell is
 178 divided by 100. The results are then put in the 3×3 matrix in accordance with the row numbers. The color
 179 intensity of an image is between 0 to 255. So, if the value of any cell of the result matrix is greater than
 180 255, then the value is replaced with 255. After applying this technique, the uniform local binary pattern
 181 is applied. From Figure 4, (a) presents a 3×3 section of matrix and the rows, (b) exhibits the result of
 182 multiplication, (c) shows the value after dividing by 100, (d) shows the replacement result of value greater
 183 than 255 and (e) shows a 3×3 matrix section after SRM-Matrix transformation.

184 Another variation of the LBP is called uniform pattern (Ojala et al., 2002). Some binary patterns occur
 185 more generally in texture images. If the binary pattern comprises of at most two 0-1 or 1-0 transitions
 186 when the bit pattern is held circular then the pattern is called uniform. For instance, 01000000 has 2
 187 transitions, 00000111 has 2 transitions which are uniform pattern on the other hand 01010100 has 6
 188 transitions, 11001001 has 4 transitions which are not uniform. A neighborhood with the dimension of
 189 3×3 has $2^8 = 256$ possible patterns with 58 of them being uniform. For estimating the histogram, every

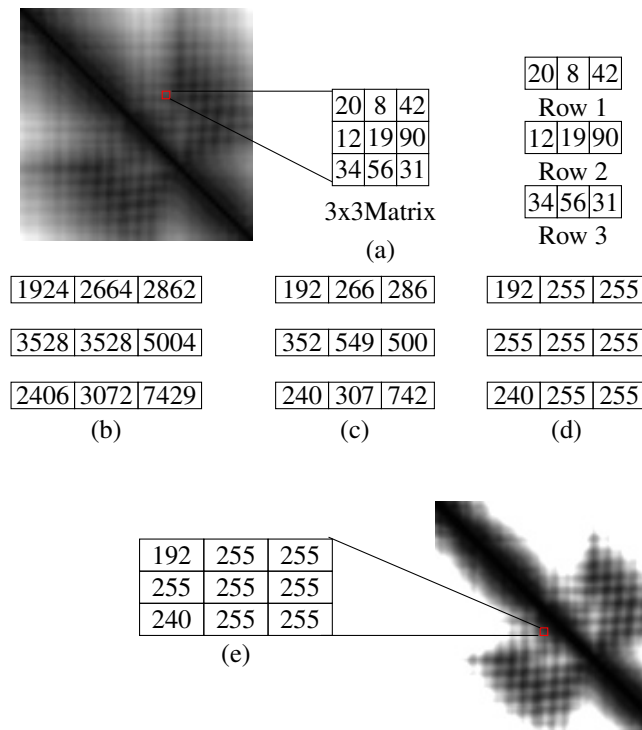


Figure 4. An example of Separate Row Multiplication Matrix with Uniform Local Binary Pattern Histogram.

Table 2. Feature Groups

Identifier	Feature Group Name	Number of Features
A	LBP-Hist	256
B	GfLBP-Hist	256
C	Atom Bond	116
D	SRMMat-ULBP-Hist	59
E	NBSMat-ULBP-Hist	59

190 uniform pattern gets a separate bin while a single bin is allotted for all non-uniform patterns. Therefore,
191 from a uniform binary pattern, we get the histogram of total bin size of 59.

192 **Neighbour Block Subtraction Matrix with Uniform LBP Histogram (NBSMat-ULBP-Hist)** Blocks
193 are of the same dimension, 3x3. Two blocks of matrices are considered neighbors for this method if the
194 center cells are neighboring. Because of this, the value of the last two columns of the first block and first
195 two columns of the second block are same. The two blocks of matrices are subtracted and the result is set
196 in the place of the first block. If any of the cells have any negative number, then 0 is placed instead of
197 the negative value. The replacing of value is made because the histogram bin begins from zero. Uniform
198 local binary pattern is then used to compute the histogram.

199 Summary of all the feature groups used in this paper is given in Table 2.

200 **Handling Imbalance in Data**

201 From Table 1 it can be noted that the classes are imbalanced. To balance the classes, we have used
202 Synthetic Minority Over-sampling Technique (SMOTE) (Chawla et al., 2002). The percentage of SMOTE
203 indicates that how many more instances would be generated. As the highest number of instance a class
204 has is 3305, we have over-sampled our instances close to that number. If x denotes the highest number of
205 instances among all the classes and y denoted by a class which we will SMOTE then the expression for

206 the percentage calculation is $\frac{x-y}{y} * 100$. We have used 5 nearest neighbors to generate the over-sampled
207 instances. After applying SMOTE to all data sets, the total number of instances in the dataset is 23132.

208 **Classifiers Used**

209 We have used five classifiers for the analysis of features applied to solve structural class prediction
210 problem: K-Nearest Neighbor (KNN), Naive Bayesian Classifier, Support Vector Machines (SVM),
211 Adaptive Boosting (AdaBoost) and Random Forest. A concise description of the classifiers is given in
212 this section.

213 **K-Nearest Neighbour (KNN)** K-nearest neighbour algorithm (KNN) (Mohri et al., 2012) is a similarity-
214 based classification technique. It is a lazy classification technique. Distance metrics are used for each
215 instance of the whole dataset for calculating the K nearest neighbors. The labels of the nearest neighbors
216 decide the label of the test instances. It works poorly for high dimensional data. Euclidean distance,
217 Hamming distance, Manhattan distance, Minkowski distance, Tanimoto distance and Jaccard distance are
218 used for similarity measures.

219 **Naive Bayesian Classifier** Naive Bayesian classifier (Mohri et al., 2012) is based on probabilistic
220 inference of samples observed where the decision variable and the features form a very naive structure of
221 Bayesian Network. Naive Bayesian classifiers work best for image recognition and text mining.

222 **Support Vector Machine (SVM)** Support Vector Machine (Mohri et al., 2012) works by creating and
223 separating hyperplane for a given dataset by sampling different classes which are separated by maximum
224 width.

225 **Adaptive Boosting (AdaBoost)** Adaptive Boosting classifier (Mohri et al., 2012) is a meta-classifier
226 which aims to make a strong classifier using a set of weak classifiers. The classifiers whose performance
227 are marginally better than random classifiers are called weak classifiers.

228 **Random Forest** Random Forest (Mohri et al., 2012) is an ensemble classifier. A decision tree is created
229 in each iteration with features taken randomly. It samples selected features using bootstrap aggregating.

230 **Ligand-binding Prediction**

231 Protein Ligand Binding prediction is a binary class classification problem. We've used Image Based
232 Features for each Protein and Ligand dataset. Our methodology learns threshold values from the training
233 data and uses these in test data prediction. We have used the same set of features that were generated and
234 analyzed for the structural class prediction problem to solve the ligand-binding problem. In this section,
235 we present the necessary materials and methods that were used for the ligand binding problem.

236 **Ligand-Binding Dataset**

237 We've used Computer Vision and Pattern Discovery for Bioimages Group @ BII as our dataset. In our
238 dataset, there are 3000 protein-ligand complexes that were determined experimentally with 3D structures
239 available. Each protein and its ligand are of one-to-one correspondence, i.e. they can bind to each other
240 and make Protein-Ligand complex. The dataset has 3000 pairs of protein and ligand where same name/ID
241 of protein and ligand interacts/binds with each other.

242 We've used OpenCV (Bradski and Kaehler, 2008) library to create images from PDB files. For protein,
243 we've considered the coordinates of only the alpha-carbons to generate the distance matrix to create image.
244 Because alpha-carbon can represent the structural information of protein quite well. But the given ligands
245 were small in terms of atom number. So, while creating ligand images, we've considered all the atom's
246 co-ordinates for generating distance matrix.

247 Among the PDB files, 33 ligands have only one atom, which will create 1x1 image having no
248 significance for feature extraction. So, we had to compromise those 33 ligands as well as 33 corresponding
249 proteins from training.

250 **Handling Imbalance**

251 The given dataset has only positive instances (the pairs of protein and ligand where they bind with each
252 other). But there were no negative instances (the pairs of protein and ligand where they do not bind with
253 each other). The missing negative instances have created our dataset highly imbalanced. To overcome this
254 imbalance, we've generated negative instances in two different ways.

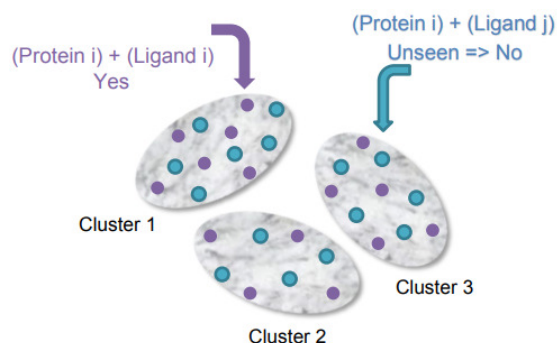


Figure 5. Clustering-Based Undersampling.

- 255 1. **Random Negative Undersampling:** We have 2967 protein PDB and 2967 ligand PDB where
 256 8803089 pairs are possible. Among these, 2967 pairs are given as positive instances and the rest
 257 8800122 pairs are unknown/unseen instances. From the unseen pairs, we've taken 2967 pairs
 258 randomly as negative instances to make our dataset balanced.
- 259 2. **Clustering-Based Undersampling:** Using the positive instances (2967 pairs), we've created 10
 260 clusters. Then we've searched for 2967 unseen pairs randomly as negative instances where they
 261 belong to those 10 clusters. We've made sure that each cluster has exactly same number of positive
 262 and negative instances to make the dataset balanced (See Figure 5).

263 **Similarity Based Classifier**

264 We've developed a similarity-based clustering method to predict the binding class. Distance is used to
 265 measure similarity. Our methodology is given in Figure 6 and the pseudo-code in Algorithm 1.

```

266 Data: A pair  $(p, l)$ , a protein structure and ligand structure in pdb format
Result: Decision, whether they will interact or not
1 for all proteins and ligands do
2   | generate images & extract features
3 end
4 for each of the given pairs of protein-ligand do
5   |  $\mathbb{NP} \leftarrow k\text{-NEARESTPROTEINS}(p)$  of the given protein
6   |  $\mathbb{RL} \leftarrow k\text{-RELATEDLIGANDS}(\mathbb{NP})$ 
7   |  $d_l \leftarrow$  distance between given ligand,  $l$  &  $\mathbb{RL}$ 
8   | if  $d_l < \text{threshold}_l$  then
9     |  $v_l \leftarrow$  vote for positive bind
10  | else
11  |  $v_l \leftarrow$  vote for negative bind
12  | end
13  |  $\mathbb{NL} \leftarrow k\text{-NEARESTLIGANDS}(l)$  of the given ligand
14  |  $\mathbb{RP} \leftarrow k\text{-RELATEDPROTEINS}(\mathbb{NL})$ 
15  |  $d_p \leftarrow$  distance between given protein,  $p$  &  $\mathbb{RP}$ 
16  | if  $d_p < \text{threshold}_p$  then
17  |  $v_p \leftarrow$  vote for positive bind
18  | else
19  |  $v_p \leftarrow$  vote for negative bind
20  | end
21  |  $v \leftarrow$  weighted majority voting between  $(v_l, v_p)$ 
22 end
23 return  $v$ 
  
```

Algorithm 1: Similarity based clustering algorithm.

267 From the PDB dataset of proteins and ligands, firstly we have generated images and converted to
 268 128×128 images for each protein and ligand. From these images we have generated 2 different features.

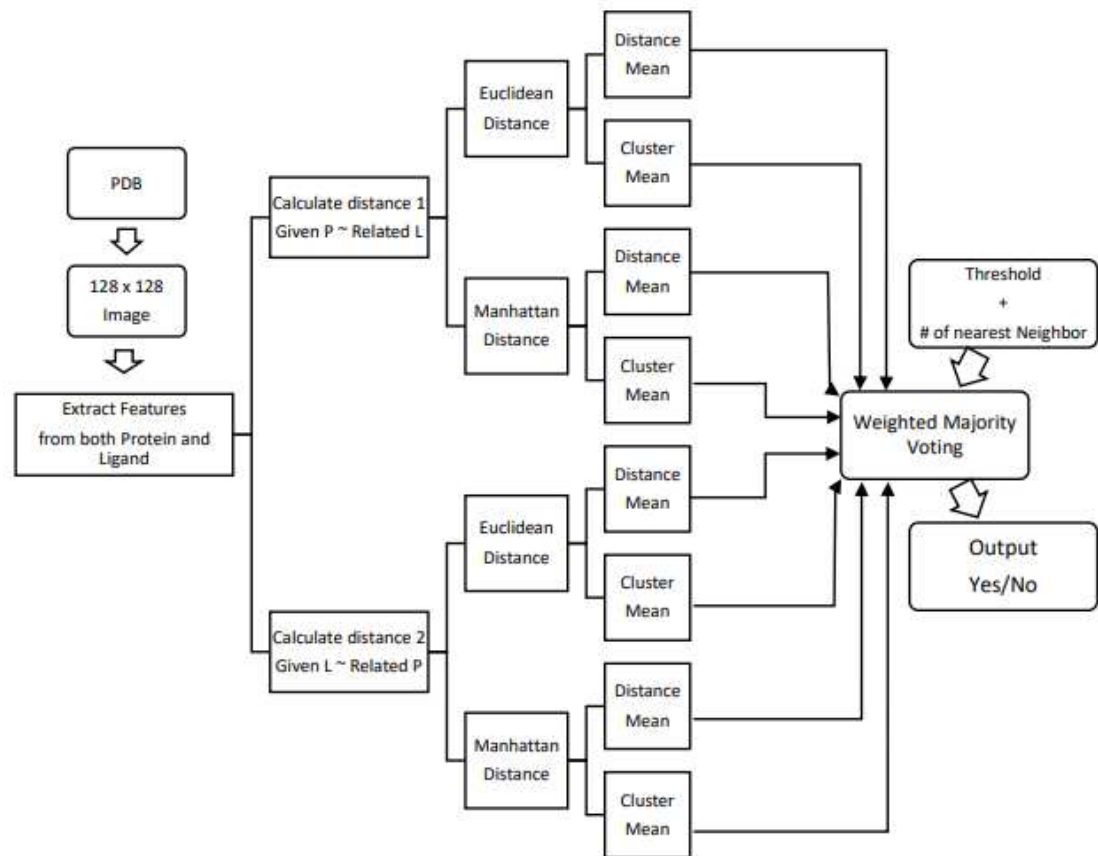


Figure 6. Block Diagram of Similarity Based Clustering.

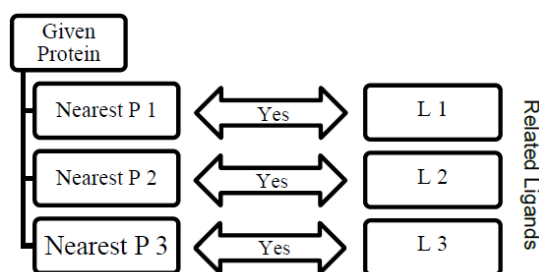


Figure 7. Relation between given protein and related ligands.

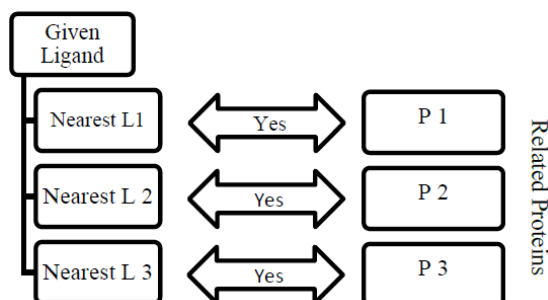


Figure 8. Relation between given ligand and related proteins.

269 1. **CoMOGrad and PHOG:** CoMOGrad stands for Co-occurrence Matrix of the Oriented Gradient
 270 of Distance Matrices and PHOG stands for Pyramid Histogram of Oriented Gradient (Karim et al.,
 271 2015). This methodology also uses the α carbon distance matrix of protein. The dimension of
 272 all distance matrix is converted to 128×128 . In CoMOGrad, the gradient angle and magnitude is
 273 computed from the distance matrix and the values are quantized. Quantization is a compressing
 274 technique which compresses a range of values to a single quantum value. In this methodology, the
 275 values are quantized to 16 bins which produce a co-occurrence matrix which is 16×16 matrix.
 276 The matrix is converted into a vector of size 256. Quadtree from the distance matrix is created
 277 with the desired level in PHOG. Gradient Oriented Histogram of each node is calculated with the
 278 preferred number of bins and bin size. In gradient oriented histogram an image is divided into
 279 small sub-images called cells and histogram of edge orientations are accumulated within the cell.
 280 The combined histogram entries are used as the feature vector describing the object. Total features
 281 which are the multiplication of total nodes and number of bins are incorporated in the vector with
 282 the size of the total number of features. The vector is normalized by dividing it with the sum of its
 283 components.

284 2. **Hybrid Local Binary Pattern (Hybrid LBP):** Local Binary Pattern (Ojala et al., 1994) is a
 285 procedure of local binary pattern histogram. We have used all the five feature groups described in
 286 the last section for structural class prediction problem.

287 Distance can only be calculated between proteins or between ligands. We've used K-nearest neighbor
 288 and Clustering method to calculate these distances.

289 1. **RELATEDLIGANDS(NP):** For a given Protein, find K-nearest proteins. The ligands those binds
 290 with the above nearest proteins, are the Related Ligands for the given protein (See Figure 7).

291 2. **RELATEDPROTEINS(NL):** For a given Ligand, find K-nearest ligands. The proteins those binds
 292 with the above nearest ligands, are the Related Proteins for the given ligand (See Figure 8).

293 To find the distances between pairs of ligands and proteins are calculated using Euclidean and
 294 Manhattan distances. Threshold is the boundary between similarity and dissimilarity in terms of distance.
 295 If distance is less than the threshold, then prediction in positive similarity, else the prediction is negative

296 similarity. Threshold of each category of distances is the average of minimum and maximum distance
297 based on the number of nearest neighbors.

298 For a given pair of Protein and Ligand, we want to predict if they will bind with each other or not. For
299 measuring distance d_l , from the given protein, we searched for k-nearest proteins and found the k related
300 ligands accordingly. Then we've calculated the distance using above mentioned methods. Then we've
301 taken the vote for the binding class by all categories of distances based on their thresholds. Then finally,
302 we've used weighted majority voting mechanism to predict the binding class.

303 **Hyperparameters**

304 There are a number of hyperparameters of our proposed method.

- 305 1. **Number of nearest neighbors:** Our algorithm's prediction accuracy is highly dependent on the
306 number of nearest neighbors for finding both RELATEDLIGANDS(NP) and RELATEDPROTEINS(NL).
307 We've used 5 nearest neighbors in this experiment.
- 308 2. **Threshold:** This is the threshold of distance for determining whether two proteins or two ligands
309 are similar or not. For a higher value of threshold, there is a higher possibility for our algorithm
310 to predict positive binding class for the majority of the Protein-Ligand pairs. And the lower the
311 threshold is, the higher is the possibility of negative binding class prediction. We've taken the
312 average of distances among 5 nearest neighbors as our threshold for each category of the distances.

313 **RESULTS AND DISCUSSION**

314 This section is the description of our experiments performed in this study. Some of the experiments
315 were carried out in a personal desktop computer having Intel Core i3 and 4 GB RAM and others were
316 experimented in a Computing Machine provided by CITS, United International University which was
317 equipped with 8 core processors each having a Dell R 730 Intel Xeon Processor (E5-2630 V3) with
318 2.4 GHz speed and 18.5 GB memory. Java language was used for data preprocessing including feature
319 generation using OpenCV software library, negative data generation and data merging using Eclipse IDE
320 with Java 8 standard edition. Python language was used to implement our algorithm using the Spyder IDE.
321 Weka tool was used to run the traditional classification algorithms for the comparison with our algorithm.
322 We've used Leave-One-Out validation method to get the accuracy of our model.

323 **Analysis of Features**

324 A different set of parameters were used for each classifiers used in this research. A linear searching was
325 used with no distance weighting for KNN. In case of the Naive Bayesian Classifier, SVM, a polynomial
326 kernel was used with $c = 1.0$ and $\epsilon = 1.0w^{-2}$. Data was normalized before supplying to the classifier.
327 J48 decision tree classifier was used in Adaboost classifier as the weak base classifier. Classifier number
328 of iterations was set to 100 for Random Forest.

329 Results in terms of average accuracy in 3-fold cross-validation of protein images are given in Table 3.
330 The highest percentage of correctly classified instances achieved for each of the classifiers are indicated
331 by the boldfaced values of the table.

332 After running the experiments for our five feature groups ABCDE classifies the highest percentage of
333 correct instances in Random Forest, Adaboost and SVM among all other feature groups. Feature scaled B
334 and D individually provides the highest accuracy in Naive Bayesian and KNN. As the whole combination
335 of all feature groups accuracy gives the highest percentage than any other feature group, thus we conclude
336 that the best performing feature group combination is ABCDE and the best classifier is Random Forest
337 classifier.

338 **Effectiveness in structural class prediction**

339 In this section, we compare the performance of our proposed method with CoMOGrad and PHOG (Karim
340 et al., 2015) along with our previous published literature Wavelet and Pyramid Histogram Features for
341 Image Based Leaf Detection (Ahmed et al., 2019). For comparison with our methodology in this literature,
342 we applied CoMOGrad and Phog techniques and Wavelet and Pyramid Histogram techniques in our
343 dataset of 11052 instances and later applied SMOTE for reducing class imbalance problem. We conducted
344 experiments with different classifiers using the same parameters as we did for feature analysis with the
345 feature groups. The results are given in Table 4. From Table 4 it can be comprehended that our feature

Image Type	Feature Type	Classifiers				
		KNN	Naive Baysian	SVM	Ada Boost	Random Forest
Non Scaled	A	77.70	32.48	68.17	84.11	87.53
Scaled	A	78.55	52.55	79.61	84.22	86.16
Non Scaled	B	74.75	35.65	71.23	83.45	86.85
Scaled	B	76.90	60.04	79.26	82.80	84.74
	C	66.96	21.79	44.49	62.26	69.92
Scaled	D	84.11	51.23	71.29	83.17	85.05
Scaled	E	83.76	51.24	71.28	83.23	84.94
Non Scaled	AB	76.70	33.53	78.39	85.87	88.43
Non Scaled	ABC	68.27	34.12	82.50	86.74	88.61
Non Scaled + Scaled	ABCD	73.06	35.46	85.47	87.26	89.21
Non Scaled + Scaled	ABCDE	74.72	37.45	86.12	87.74	89.49

Table 3. Classifier accuracies for different types of feature and groups of features.

Feature Type	Classifiers				
	KNN	Naive Baysian	SVM	Ada Boost	Random Forest
Karim et al. Karim et al. (2015)	87.41	59.50	87.67	84.19	85.49
Ahmed et al. Ahmed et al. (2019)	69.36	36.22	67.30	79.92	84.58
this paper	74.72	37.45	86.12	87.74	89.49

Table 4. Comparison of the proposed features in this paper with Karim et al. (2015) and Ahmed et al. (2019) for structural class prediction.

346 group ABCDE outperforms CoMOGrad and PHOG in Random Forest and in Adaboost. CoMOGrad
 347 and PHOG surpassed our feature groups in KNN, Naive Bayesian and SVM. It can be noted that the
 348 combination of our feature groups are three-fourths of CoMOGrad and PHOG. It also can be discerned
 349 that the accuracy percentage in Random Forest is higher than all the classifier results. Thus, our novel
 350 features can classify more instances than CoMOGrad and PHOG. We have also noticed that our feature
 351 groups outperform the features of our previous literature Ahmed et al. (2019) on all classifiers.

352 We have revealed the precedence of our methodology over CoMOGrad and PHOG (Karim et al.,
 353 2015) and Wavelet and Pyramid Histogram Features for Image Based Leaf Detection (Ahmed et al., 2019).
 354 The same feature groups were used for leaf detection (Ahmed et al., 2019) with the dataset consisting of
 355 RGB images of leaves. Unlike only gray histogram used on this paper, blue, green and red histograms
 356 were used to generate features in each feature group and the accuracy result of each classifier was high.
 357 The distance matrix of α carbons or the protein images were black and white, thus only gray histogram
 358 was used as a feature.

359 We also used Scale-invariant feature transform (SIFT) (Lowe, 2004) methodologies in our experiments.
 360 Each descriptor has a 128-dimensional feature vector. The number of the descriptors of SIFT from every
 361 image is not specific so we cannot use traditional machine learning techniques. Hence to apply traditional
 362 machine learning procedure and specify the feature vector, we split the image into 16 slices and took one

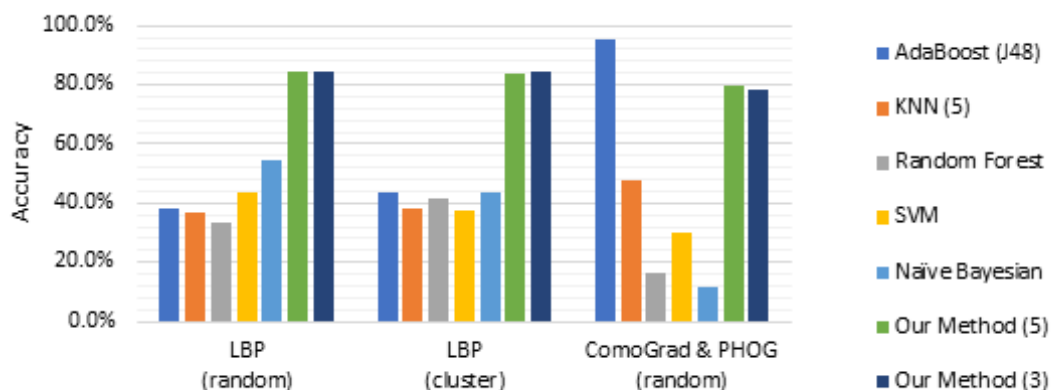


Figure 9. Barplot showing the performance of different algorithms on ligand-binding dataset.

363 descriptor from each of the slice images. Therefore we got 2048 number of attributes(8x16) from each
 364 image. We tested the dataset with the same classifiers mentioned in this paper. The results didn't turn up
 365 to be better or close to our proposed methodology in this literature.

366 Effectiveness in ligand-binding prediction

367 Sensitivity is the true positive rate regarding the positive instances. As we had to generate the negative
 368 data artificially, sensitivity is the actual scale of performance measuring where positive data were the
 369 actual data. Using the thresholds gained using the negative data, sensitivity of our algorithm is very good
 370 comparing to other existing algorithms shown in Table 5 and Figure 9.

Features	AdaBoost	KNN	Random Forest	SVM	Naive Bayesian	Our Method
LBP (random)	40.00%	43.50%	22.00%	36.80%	45.20%	91.33%
LBP (cluster)	51.90%	44.30%	52.20%	49.00%	43.70%	91.60%
CoMOGrad & PHOG (random)	95.20%	47.60%	16.10%	29.70%	11.30%	79.86%

Table 5. Sensitivity Comparison among different methods for ligand-binding prediction.

371 We have generated three different datasets based on three different features. Hybrid LBP gives 736
 372 long feature vectors from protein images and 677 long feature vectors from ligand images. So, for one
 373 protein-ligand pair we've got 1413 (736+677) attributes and one Binding Class value as one instance.
 374 The above mentioned two types of negative data (random and Clustering-Based Undersampling) were
 375 generated using Hybrid LBP for balancing the data. CoMOGrad and PHOG gives 1021 or 1020 long
 376 feature vectors from protein image, but for ligand images, it gives 1020 long feature vectors. We assumed
 377 "0" as the last feature in protein where features were 1020 long, to make it 1021 long feature. So, for
 378 one protein-ligand pair we've got 2041 (1021+1020) attributes and one BINDING Class value as one
 379 instance. Random negative undersampling was used in CoMOGrad and PHOG but Clustering-Based
 380 Undersampling was not possible as some clusters couldn't get any unseen pairs of protein and ligand. Our
 381 method was used based on 5 and 3 nearest neighbors and shown on the above table and chart.

382 We can see that AdaBoost works better than our algorithm in terms of sensitivity in ComoGrad and
 383 PHOG dataset. Because, Ligand data were so small in terms of number of atoms that ComoGrad and
 384 PHOG gave zeros for most of the ligands. But our algorithm's overall performance is better than other
 385 machine learning algorithms in the three different feature datasets.

386 CONCLUSIONS

387 In this paper, we showed how accurately we can detect protein classes using the combination of different
 388 image based feature groups generated from protein images. We also propose a simple similarity-based

389 clustering method to predict Protein-Ligand Binding without using deep-learning or neural-networks.
390 This simple distance-based algorithm is quite effective compared to complex machine learning algorithms.
391 Our main limitation was the missing negative data. If we had the actual negative data, we could've
392 determined the perfect thresholds for each category of distances, and that would give us more accurate
393 prediction. Another problem was dimensions of small Ligands as we're using image-based features.
394 As the advancement of deep learning, neural network, and many other deep learning techniques are
395 being used to classify images, many remarkably interesting applications can be made. For our future
396 advancement, we wish to introduce new features to improve accuracy, use new tools and explore other
397 fields of computer vision such as human emotion detection. In addition, we will try to extract some unique
398 features from the Ligand dataset so that the dimensionality problem doesn't affect our Protein-Ligand
399 binding prediction.

400 ACKNOWLEDGMENTS

401 We thank Rezaul Karim for sharing the SQL dataset files and algorithms to generate Distance Matrix
402 from PDB files for CoMOGrad and PHOG.

403 REFERENCES

- 404 Ahmed, A. A. N., Haque, H. M. F., Rahman, A., Ashraf, M. S., and Shatabda, S. (2019). Wavelet and
405 pyramid histogram features for image-based leaf detection. In Abraham, A., Dutta, P., Mandal, J. K.,
406 Bhattacharya, A., and Dutta, S., editors, *Emerging Technologies in Data Mining and Information*
407 *Security*, pages 269–278, Singapore. Springer Singapore.
- 408 Bradski, G. and Kaehler, A. (2008). *Learning OpenCV: Computer vision with the OpenCV library*. ”
409 O'Reilly Media, Inc.”.
- 410 Brady, G. P. and Stouten, P. F. (2000). Fast prediction and visualization of protein binding pockets with
411 pass. *Journal of computer-aided molecular design*, 14(4):383–401.
- 412 Chaires, J. B. (2008). Calorimetry and thermodynamics in drug design. *Annu. Rev. Biophys.*, 37:135–151.
- 413 Chawla, N. V., Bowyer, K. W., Hall, L. O., and Kegelmeyer, W. P. (2002). Smote: synthetic minority
414 over-sampling technique. *Journal of artificial intelligence research*, 16:321–357.
- 415 Chi, P.-H., Scott, G., and Shyu, C.-R. (2005). A fast protein structure retrieval system using image-based
416 distance matrices and multidimensional index. *International Journal of Software Engineering and*
417 *Knowledge Engineering*, 15(03):527–545.
- 418 Chothia, C. and Lesk, A. M. (1986). The relation between the divergence of sequence and structure in
419 proteins. *The EMBO journal*, 5(4):823–826.
- 420 Fox, N. K., Brenner, S. E., and Chandonia, J.-M. (2013). Scope: Structural classification of pro-
421 teins—extended, integrating scop and astral data and classification of new structures. *Nucleic acids*
422 *research*, 42(D1):D304–D309.
- 423 Holm, L. and Sander, C. (1993). Protein structure comparison by alignment of distance matrices. *Journal*
424 *of molecular biology*, 233(1):123–138.
- 425 Holm, L. and Sander, C. (1997). Dali/fssp classification of three-dimensional protein folds. *Nucleic acids*
426 *research*, 25(1):231–234.
- 427 Jain, A. K. and Bhattacharjee, S. (1992). Text segmentation using gabor filters for automatic document
428 processing. *Machine vision and applications*, 5(3):169–184.
- 429 Jain, A. K. and Farrokhnia, F. (1991). Unsupervised texture segmentation using gabor filters. *Pattern*
430 *recognition*, 24(12):1167–1186.
- 431 Karim, R., Aziz, M. M. A., Shatabda, S., Rahman, M. S., Mia, M. A. K., Zaman, F., and Rakin, S. (2015).
432 Comograd and phog: From computer vision to fast and accurate protein tertiary structure retrieval.
433 *Scientific Reports*, 5:13275 EP –. Article.
- 434 Lowe, D. G. (2004). Distinctive image features from scale-invariant keypoints. *International journal of*
435 *computer vision*, 60(2):91–110.
- 436 Mehrotra, R., Namuduri, K. R., and Ranganathan, N. (1992). Gabor filter-based edge detection. *Pattern*
437 *recognition*, 25(12):1479–1494.
- 438 Mohri, M., Rostamizadeh, A., and Talwalkar, A. (2012). *Foundations of machine learning*. MIT press.
- 439 Ojala, T., Pietikainen, M., and Harwood, D. (1994). Performance evaluation of texture measures with
440 classification based on kullback discrimination of distributions. In *Pattern Recognition, 1994. Vol.*

- 441 *I-Conference A: Computer Vision & Image Processing., Proceedings of the 12th IAPR International*
442 *Conference on*, volume 1, pages 582–585. IEEE.
- 443 Ojala, T., Pietikainen, M., and Maenpaa, T. (2002). Multiresolution gray-scale and rotation invariant
444 texture classification with local binary patterns. *IEEE Transactions on pattern analysis and machine*
445 *intelligence*, 24(7):971–987.
- 446 Patching, S. G. (2014). Surface plasmon resonance spectroscopy for characterisation of membrane
447 protein–ligand interactions and its potential for drug discovery. *Biochimica et Biophysica Acta (BBA)-*
448 *Biomembranes*, 1838(1):43–55.
- 449 Rossi, A. M. and Taylor, C. W. (2011). Analysis of protein-ligand interactions by fluorescence polarization.
450 *Nature protocols*, 6(3):365.
- 451 Shyu, C.-R., Chi, P.-H., Scott, G., and Xu, D. (2004). Proteindbs: a real-time retrieval system for protein
452 structure comparison. *Nucleic Acids Research*, 32(suppl_2):W572–W575.
- 453 Singh, A. P. and Brutlag, D. L. (1997). Hierarchical protein structure superposition using both secondary
454 structure and atomic representations. In *Ismb*, volume 5, pages 284–293.
- 455 Sousa, S. F., Ribeiro, A. J., Coimbra, J., Neves, R., Martins, S., Moorthy, N., Fernandes, P., and Ramos, M.
456 (2013). Protein-ligand docking in the new millennium—a retrospective of 10 years in the field. *Current*
457 *medicinal chemistry*, 20(18):2296–2314.
- 458 Srivastava, S., Lal, S. B., Mishra, D., Angadi, U., Chaturvedi, K., Rai, S. N., and Rai, A. (2016). An
459 efficient algorithm for protein structure comparison using elastic shape analysis. *Algorithms for*
460 *Molecular Biology*, 11(1):27.
- 461 Steinbrecher, T. and Labahn, A. (2010). Towards accurate free energy calculations in ligand protein-
462 binding studies. *Current medicinal chemistry*, 17(8):767–785.
- 463 TAYLOR, W. R. (1999). Protein structure comparison using iterated double dynamic programming.
464 *Protein Science*, 8(3):654–665.