

**A peer-reviewed version of this preprint was published in PeerJ on 17 July 2017.**

[View the peer-reviewed version](https://peerj.com/articles/cs-124) (peerj.com/articles/cs-124), which is the preferred citable publication unless you specifically need to cite this preprint.

Zacharaki EI. 2017. Prediction of protein function using a deep convolutional neural network ensemble. PeerJ Computer Science 3:e124 <https://doi.org/10.7717/peerj-cs.124>

# Prediction of protein function using a deep convolutional neural network ensemble

Evangelia I Zacharaki <sup>Corresp.</sup> <sup>1</sup>

<sup>1</sup> Center for Visual Computing, CentraleSupélec and GALEN Team, INRIA Saclay, Palaiseau, France

Corresponding Author: Evangelia I Zacharaki

Email address: [evangelia.zacharaki@centralesupelec.fr](mailto:evangelia.zacharaki@centralesupelec.fr)

**Background.** The availability of large databases containing high resolution three-dimensional (3D) models of proteins in conjunction with functional annotation allows the exploitation of advanced supervised machine learning techniques for automatic protein function prediction.

**Methods.** In this work, novel shape features are extracted representing protein structure in the form of local (per amino acid) distribution of angles and amino acid distances, respectively. Each of the multi-channel feature maps is introduced into a deep convolutional neural network (CNN) for function prediction and the outputs are fused through Support Vector Machines (SVM) or a correlation-based k-nearest neighbor classifier. Two different architectures are investigated employing either one CNN per multi-channel feature set, or one CNN per image channel.

**Results.** Cross validation experiments on enzymes ( $n = 44,661$ ) from the PDB database achieved 90.1% correct classification demonstrating the effectiveness of the proposed method for automatic function annotation of protein structures.

**Discussion.** The automatic prediction of protein function can provide quick annotations on extensive datasets opening the path for relevant applications, such as pharmacological target identification.

# Prediction of protein function using a deep convolutional neural network ensemble

Evangelia I. Zacharaki<sup>1</sup>

<sup>1</sup>Center for Visual Computing, CentraleSupélec and GALEN Team, INRIA Saclay, France

## ABSTRACT

**Background.** The availability of large databases containing high resolution three-dimensional (3D) models of proteins in conjunction with functional annotation allows the exploitation of advanced supervised machine learning techniques for automatic protein function prediction.

**Methods.** In this work, novel shape features are extracted representing protein structure in the form of local (per amino acid) distribution of angles and amino acid distances, respectively. Each of the multi-channel feature maps is introduced into a deep convolutional neural network (CNN) for function prediction and the outputs are fused through Support Vector Machines (SVM) or a correlation-based k-nearest neighbor classifier. Two different architectures are investigated employing either one CNN per multi-channel feature set, or one CNN per image channel.

**Results.** Cross validation experiments on enzymes ( $n = 44,661$ ) from the PDB database achieved 90.1% correct classification demonstrating the effectiveness of the proposed method for automatic function annotation of protein structures.

**Discussion.** The automatic prediction of protein function can provide quick annotations on extensive datasets opening the path for relevant applications, such as pharmacological target identification.

Keywords:

## 1 INTRODUCTION

Research in metagenomics led to a huge increase of protein databases and discovery of new protein families (Godzik, 2011). While the number of newly discovered, but possibly redundant, protein sequences rapidly increases, experimentally verified functional annotation of whole genomes remains limited. Protein structure, i.e. the 3D configuration of the chain of amino acids, is a very good predictor of protein function, and in fact a more reliable predictor than protein sequence because it is far more conserved in nature (Illergård et al., 2009).

By now, the number of proteins with functional annotation and experimentally predicted structure of their native state (e.g. by NMR spectroscopy or X-ray crystallography) is adequately large to allow learning training models that will be able to perform automatic functional annotation of unannotated proteins. Also, as the number of protein sequences rapidly grows, the overwhelming majority of proteins can only be annotated computationally. In this work enzymatic structures from the Protein Data Bank (PDB) are considered and the enzyme commission (EC) number is used as a fairly complete framework for annotation. The EC number is a numerical classification scheme based on the chemical reactions the enzymes catalyze, proven by experimental evidence (web, 1992).

There have been plenty machine learning approaches in the literature for automatic enzyme annotation. A systematic review on the utility and inference of various computational methods for functional characterization is presented in (Sharma and Garg, 2014), while a comparison of machine learning approaches can be found in (Yadav and Tiwari, 2015). Most methods use features derived from the amino acid sequence and apply Support Vector Machines (SVM) (Cai et al., 2003)(Han et al., 2004)(Dobson and Doig, 2005)(Chen et al., 2006)(Zhou et al., 2007)(Lu et al., 2007)(Lee et al., 2009)(Qiu et al., 2010)(Wang et al., 2010)(Wang et al., 2011)(Amidi et al., 2016), k-Nearest Neighbor (kNN) classifier (Huang et al., 2007)(Shen and Chou, 2007a)(Nasibov and Kandemir-Cavas, 2009a), classification trees/forests (Lee et al., 2009)(Kumar and Choudhary, 2012a)(Nagao et al., 2014)(Yadav and Tiwari, 2015), and neural networks (Volpato et al., 2013). In (Borgwardt et al., 2005) sequential, structural and chemical information was combined into one graph model of proteins which was further classified by SVM. There has been little

47 work in the literature on automatic enzyme annotation based only on structural information. A Bayesian  
48 approach (Borro et al., 2006) for enzyme classification using structure derived properties achieved 45%  
49 accuracy. Amidi et al. (2016) obtained 73.5% classification accuracy on 39,251 proteins from the PDB  
50 database when they used only structural information.

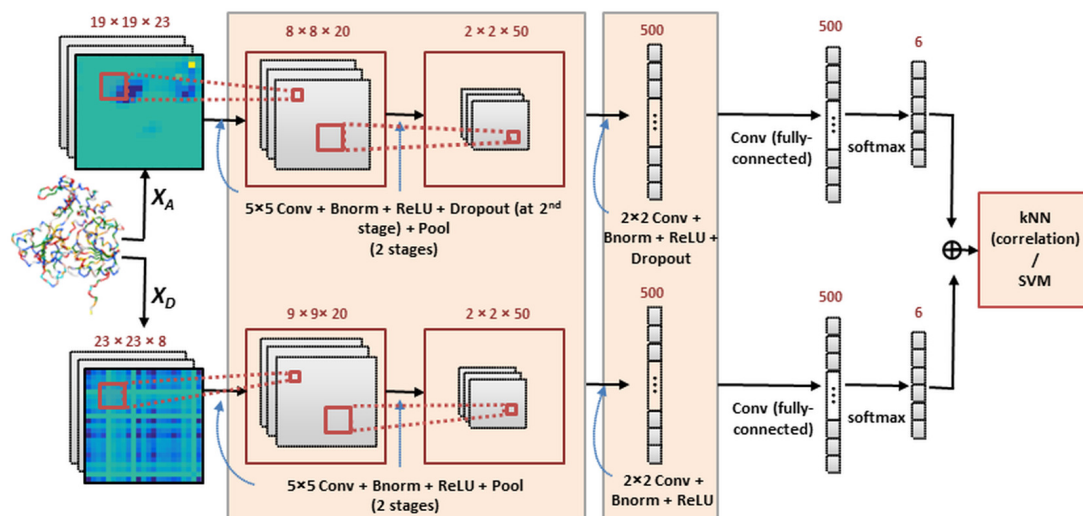
51 In the past few years, deep learning techniques, and particularly convolutional neural networks,  
52 have rapidly become the tool of choice for tackling many challenging computer vision tasks, such as  
53 image classification (Krizhevsky et al., 2012). The main advantage of deep learning techniques is the  
54 automatic exploitation of features and tuning of performance in a seamless fashion, that simplifies the  
55 conventional image analysis pipelines. CNNs have recently been used for protein secondary structure  
56 prediction (Spencer et al., 2015)(Li and Shibuya, 2015). In (Spencer et al., 2015) prediction was based  
57 on the position-specific scoring matrix profile (generated by PSI-BLAST), whereas in (Li and Shibuya,  
58 2015) 1D convolution was applied on features related to the amino acid sequence. Also a deep CNN  
59 architecture was proposed in (Lin et al., 2016) to predict protein properties. This architecture used a  
60 multilayer shift-and-stitch technique to generate fully dense per-position predictions on protein sequences.  
61 To the best of authors's knowledge, deep CNNs have not been used for prediction of protein function so  
62 far.

63 In this work the author exploits experimentally acquired structural information of enzymes and apply  
64 deep learning techniques in order to produce models that predict enzymatic function based on structure.  
65 Novel geometrical descriptors are introduced and the efficacy of the approach is illustrated by classifying  
66 a dataset of 44,661 enzymes from the PDB database into the  $l = 6$  primary categories: oxidoreductases  
67 (EC1), transferases (EC2), hydrolases (EC3), lyases (EC4), isomerases (EC5), ligases (EC6). The novelty  
68 of the proposed method lies first in the representation of the 3D structure as a "bag of atoms (amino acids)"  
69 which are characterized by geometric properties, and secondly in the exploitation of the extracted feature  
70 maps by deep CNNs. Although assessed for enzymatic function prediction, the method is not based  
71 on enzyme-specific properties and therefore can be applied (after re-training) for automatic large-scale  
72 annotation of other 3D molecular structures, thus providing a useful tool for data-driven analysis. In  
73 the following sections more details on the implemented framework are first provided, including the  
74 representation of protein structure, the CNN architecture and the fusion process of the network outputs.  
75 Then the evaluation framework and the obtained results are presented, followed by some discussion and  
76 conclusions.

## 77 2 METHODS

78 Data-driven CNN models tend to be domain agnostic and attempt to learn additional feature bases that  
79 cannot be represented through any handcrafted features. It is hypothesized that by combining "amino acid  
80 specific" descriptors with the recent advances in deep learning we can boost model performance. The  
81 main advantage of the proposed method is that it exploits complementarity in both data representation  
82 phase and learning phase. Regarding the former, the method uses an enriched geometric descriptor that  
83 combines local shape features with features characterizing the interaction of amino acids on this 3D  
84 spatial model. Shape representation is encoded by the local (per amino acid type) distribution of torsion  
85 angles (Bermejo et al., 2012). Amino acid interactions are encoded by the distribution of pairwise amino  
86 acid distances. While the torsion angles and distance maps are usually calculated and plotted for the  
87 whole protein (Bermejo et al., 2012), in the current approach they are extracted for each amino acid  
88 type separately, therefore characterizing local interactions. Thus, the protein structure is represented as  
89 a set of multi-channel images which can be introduced into any machine learning scheme designed for  
90 fusing multiple 2D feature maps. Moreover, it should be noted that the utilized geometric descriptors  
91 are invariant to global translation and rotation of the protein, therefore previous protein alignment is not  
92 required.

93 Our method constructs an ensemble of deep CNN models that are complementary to each other.  
94 The deep network outputs are combined and introduced into a correlation-based k-nearest neighbor  
95 (kNN) classifier for function prediction. For comparison purposes, SVM were also implemented for  
96 final classification. Two system architectures are investigated in which the multiple image channels are  
97 considered jointly or independently, as will be described next. Both architectures use the same CNN  
98 structure (within the highlighted boxes) which is illustrated in Fig.1.



**Figure 1.** The deep CNN ensemble for protein classification. In this framework (*Architecture 1*) each multi-channel feature set is introduced to a CNN and results are combined by kNN or SVM classification. The network includes layers performing convolution (Conv), batch normalization (Bnorm), rectified linear unit (ReLU) activation, dropout (optionally) and max-pooling (Pool). Details are provided in section 2.2.

## 99 2.1 Representation of protein structure

100 The building blocks of proteins are amino acids which are linked together by peptide bonds into a chain.  
 101 The polypeptide folds into a specific conformation depending on the interactions between its amino acid  
 102 side chains which have different chemistries. Many conformations of this chain are possible due to the  
 103 rotation of the chain about each carbon ( $C\alpha$ ) atom. For structure representation, two sets of feature  
 104 maps were used. They express the shape of the protein backbone and the distances between the protein  
 105 building blocks (amino acids). The use of global rotation and translation invariant features is preferred  
 106 over features based on the Cartesian coordinates of atoms, in order to avoid prior protein alignment, which  
 107 is a bottleneck in the case of large datasets with proteins of several classes (unknown reference template  
 108 space). The feature maps were extracted for every amino acid being present in the dataset including the  
 109 20 standard amino acids, as well as asparagine/aspartic (ASX), glutamine/glutamic (GLX), and all amino  
 110 acids with unidentified/unknown residues (UNK), resulting in  $m = 23$  amino acids in total.

111 **Torsion angles density.** The shape of the protein backbone was expressed by the two torsion angles of  
 112 the polypeptide chain which describe the rotations of the polypeptide backbone around the bonds between  
 113 N- $C\alpha$  (angle  $\phi$ ) and  $C\alpha$ -C (angle  $\psi$ ). All amino acids in the protein were grouped according to their type  
 114 and the density of the torsion angles  $\phi$  and  $\psi (\in [-180, 180])$  was estimated for each amino acid type  
 115 based on the 2D sample histogram of the angles (also known as Ramachandran diagram) using equal  
 116 sized bins (number of bins  $h_A = 19$ ). The histograms were not normalized by the number of instances,  
 117 therefore their values indicate the frequency of each amino acid within the polypeptide chain. In the  
 118 obtained feature maps ( $X_A$ ), with dimensionality  $[h_A \times h_A \times m]$ , the number of amino acids ( $m$ ) corresponds  
 119 to the number of channels. Smoothness in the density function was achieved by moving average filtering,  
 120 i.e. by convoluting the density map with a 2D gaussian kernel ( $\sigma = 0.5$ ).

121 **Density of amino acid distances.** For each amino acid  $a_i, i = 1, \dots, m$ , the distances to amino acid  
 122  $a_j, j = 1, \dots, m$ , in the protein are calculated based on the coordinates of the  $C\alpha$  atoms for the residues  
 123 and stored as an array  $d_{ij}$ . Since the size of the proteins varies significantly, the length of the array  $d_{ij}$   
 124 is different across proteins, thus not directly comparable. In order to standardize measurements, the  
 125 sample histogram of  $d_{ij}$  is extracted (using equally sized bins) and smoothed by convolution with a 1D  
 126 gaussian kernel ( $\sigma = 0.5$ ). The processing of all pairs of amino acids resulted to feature maps ( $X_D$ ) of  
 127 dimensionality  $[m \times m \times h_D]$ , where  $h_D = 8$  is the number of histogram bins (considered as number of  
 128 channels in this case).

## 129 2.2 Classification by deep CNNs

130 **Feature extraction stage of each CNN.** The CNN architecture employs three computational blocks of  
131 consecutive convolutional, batch normalization, rectified linear unit (ReLU) activation, dropout (option-  
132 ally) and max-pooling layers, and a fully-connected layer. The convolutional layer computes the output  
133 of neurons that are connected to local regions in the input in order to extract local features. It applies  
134 a 2D convolution between each of the input channels and a set of filters. The 2D activation maps are  
135 calculated by summing the results over all channels and then stacking the output of each filter to produce  
136 the output 3D volume. Batch normalization normalizes each channel of the feature map by averaging over  
137 spatial locations and batch instances. The ReLU layer applies an element-wise activation function, such  
138 as the  $\max(0, x)$  thresholding at zero. The dropout layer is used to randomly drop units from the CNN  
139 during training and reduce overfitting. Dropout was used only for the  $X_A$  feature set. The pooling layer  
140 performs a downsampling operation along the spatial dimensions in order to capture the most relevant  
141 global features with fixed length. The  $\max$  operator was applied within a  $[2 \times 2]$  neighborhood. The last  
142 layer is fully-connected and represents the class scores.

143 **Training and testing stage of each CNN.** The output of each CNN is a vector of probabilities, one for  
144 each of the  $l$  possible enzymatic classes. The CNN performance can be measured by a loss function which  
145 assigns a penalty to classification errors. The CNN parameters are learned to minimize this loss averaged  
146 over the annotated (training) samples. The *softmax* loss function (i.e. the *softmax* operator followed by the  
147 *logistic loss*) is applied to predict the probability distribution over categories. Optimization was based on  
148 an implementation of stochastic gradient descent. At the testing stage, the network outputs after *softmax*  
149 normalization are used as class probabilities.

## 150 2.3 Fusion of CNN outputs using two different architectures

151 Two fusion strategies were implemented. In the first strategy (*Architecture 1*) the two feature sets,  $X_A$   
152 and  $X_D$ , are each introduced into a CNN, which performs convolution at all channels, and then the  $l$  class  
153 probabilities produced for each feature set are combined into a feature vector of length  $l * 2$ . In the second  
154 strategy (*Architecture 2*), each one of the ( $m = 23$  or  $h_D = 8$ ) channels of each feature set is introduced  
155 independently into a CNN and the obtained class probabilities are concatenated into a vector of  $l * m$   
156 features for  $X_A$  and  $l * h_D$  features for  $X_D$ , respectively. These two feature vectors are further combined  
157 into a single vector of length  $l * (m + h_D)$  (=186). For both architectures, kNN classification was applied  
158 for final class prediction using as distance measure between two feature vectors,  $x_1$  and  $x_2$ , the metric  
159  $1 - \text{cor}(x_1, x_2)$ , where *cor* is the sample Spearman's rank correlation. The value  $k = 12$  was selected for  
160 all experiments. For comparison, fusion was also performed with linear SVM classification (Chang and  
161 Lin, 2011). The code was developed in MATLAB environment and the implementation of CNNs was  
162 based on MatConvNet (Vedaldi and Lenc, 2015).

## 163 3 RESULTS

164 The protein structures ( $n = 44,661$ ) were collected from the PDB. Only enzymes that occur in a single  
165 class were processed, whereas enzymes that perform multiple reactions and are hence associated with  
166 multiple enzymatic functions were excluded. Since protein sequence was not examined during feature  
167 extraction, all enzymes were considered without other exclusion criteria, such as small sequence length or  
168 homology bias. The dataset was unbalanced in respect to the different classes. The number of samples per  
169 class is shown in Table 1. The dataset was split into 5 folds. Four folds were used for training and one for  
170 testing. The training samples were used to learn the parameters of the network (such as the weights of the  
171 convolution filters), as well as the parameters of the subsequent classifiers used during fusion (SVM or  
172 kNN model). Once the network was trained, the class probabilities were obtained for the testing samples,  
173 which were introduced into the trained SVM or kNN classifier for final prediction. The SVM model was  
174 linear, thus didn't require any hyper-parameter optimization. Due to lack of hyper-parameters, no extra  
175 validation set was necessary. On the side, the author examined also non-linear SVM with gaussian radial  
176 basis function kernel, but didn't observe any significant improvement, thus the corresponding results are  
177 not reported.

178 A classification result was deemed a true positive if the match with the highest probability was in first  
179 place in a rank-ordered list. The classification accuracy (percentage of correctly classified samples over  
180 all samples) was calculated for each fold and then results were averaged across the 5 folds.

**Table 1.** Cross-validation accuracy (in percentage) in predicting main enzymatic function using the deep CNN ensemble

Class	Samples	<i>Architecture 1</i>		<i>Architecture 2</i>	
		linear-SVM	kNN	linear-SVM	kNN
EC1	8,075	86.4	88.8	91.2	90.6
EC2	12,739	84.0	87.5	88.0	91.7
EC3	17,024	88.7	91.3	89.6	94.0
EC4	3,114	79.4	78.4	84.9	80.7
EC5	1,905	69.5	68.6	79.6	77.0
EC6	1,804	61.0	60.6	73.6	70.4
<b>Total</b>	44,661	84.4	86.7	<b>88.0</b>	<b>90.1</b>

**Table 2.** Confusion matrices for each fusion scheme and classification technique

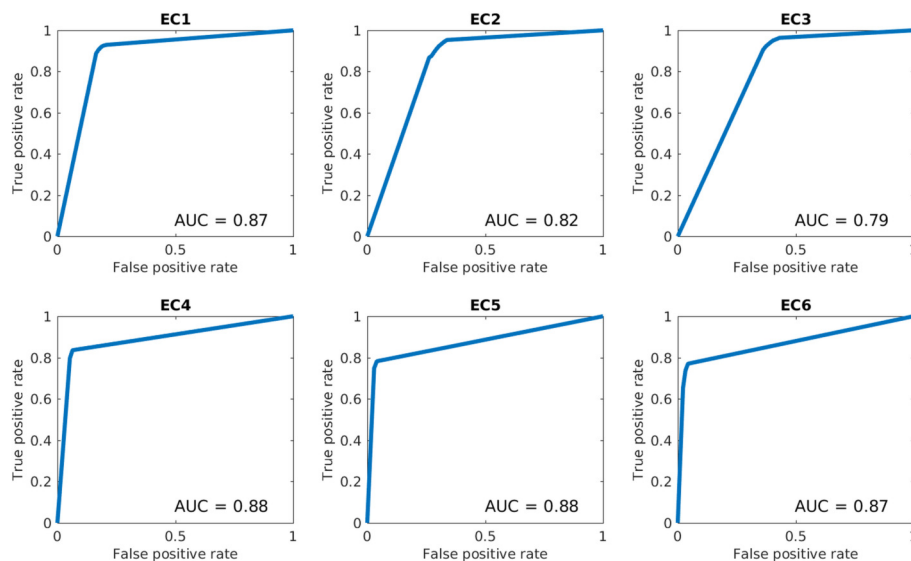
Classifier		prediction by <i>Architecture 1</i>						prediction by <i>Architecture 2</i>					
		1	2	3	4	5	6	1	2	3	4	5	6
linear-SVM	EC1	<b>86.5</b>	4.9	4.8	1.8	1.1	1.0	<b>91.2</b>	2.9	1.9	2.2	1.1	0.7
	EC2	3.4	<b>84.0</b>	7.9	1.9	1.2	1.6	3.6	<b>88.0</b>	3.5	2.2	1.2	1.5
	EC3	2.4	6.1	<b>88.7</b>	1.0	0.8	1.0	2.3	4.1	<b>89.6</b>	1.6	1.2	1.2
	EC4	4.4	7.3	5.7	<b>79.4</b>	1.8	1.3	4.3	4.9	2.7	<b>84.9</b>	1.7	1.4
	EC5	7.0	10.1	9.0	2.9	<b>69.4</b>	1.6	4.5	5.4	4.7	4.4	<b>79.5</b>	1.7
	EC6	5.9	15.5	13.0	2.3	2.3	<b>61.0</b>	5.5	10.3	5.4	3.3	1.9	<b>73.6</b>
kNN	EC1	<b>88.8</b>	5.0	4.5	0.7	0.5	0.5	<b>90.6</b>	4.4	4.6	0.3	0.1	0.0
	EC2	2.5	<b>87.5</b>	7.4	1.0	0.6	1.1	1.7	<b>91.7</b>	5.8	0.3	0.2	0.4
	EC3	1.8	5.4	<b>91.3</b>	0.5	0.4	0.6	1.2	4.4	<b>94.0</b>	0.2	0.1	0.2
	EC4	3.8	9.1	7.2	<b>78.5</b>	1.1	0.4	3.7	8.4	6.9	<b>80.7</b>	0.1	0.1
	EC5	6.1	11.5	10.7	2.3	<b>68.5</b>	1.0	3.5	9.7	8.6	0.9	<b>76.9</b>	0.3
	EC6	4.9	18.8	13.5	1.0	1.3	<b>60.6</b>	4.2	14.1	10.3	0.7	0.3	<b>70.5</b>

### 3.1 Classification performance

Common options for the network were used, except of the size of the filters which was adjusted to the dimensionality of the input data. Specifically, the convolutional layer used neurons with receptive field of size 5 for the first two layers and 2 for the third layer. The stride (specifying the sliding of the filter) was always 1. The number of filters was 20, 50 and 500 for the three layers, respectively, and the learning rate 0.001. The batch size was selected according to information amount (dimensionality) of input. It was assumed (and verified experimentally) that for more complicated the data, a larger number of samples is required for learning. One thousand samples per batch were used for *Architecture 1*, which takes as input all channels, and 100 samples per batch for *Architecture 2*, in which an independent CNN is trained for each channel. The dropout rate was 20%. The number of epochs was adjusted to the rate of convergence for each architecture (300 for *Architecture 1* and 150 for *Architecture 2*).

The average classification accuracy over the 5 folds for each enzymatic class is shown in Table 1 for both fusion schemes, whereas the analytic distribution of samples in each class is shown in the form of confusion matrices in Table 2.

In order to further assess the performance of the deep networks, receiver operating characteristic (ROC) curves and area-under-the-curve (AUC) values were calculated for each class for the selected scheme (based on kNN and *Architecture 2*), as shown in Fig.2). The calculations were performed based on the final decision scores in a one-versus-rest classification scheme. The decision scores for the kNN classifier reflected the ratio of the within-class neighbors over total number of neighbors. The ROC curve represents the true positive rate against the false positive rate and was produced by averaging over the five folds of the cross-validation experiments.

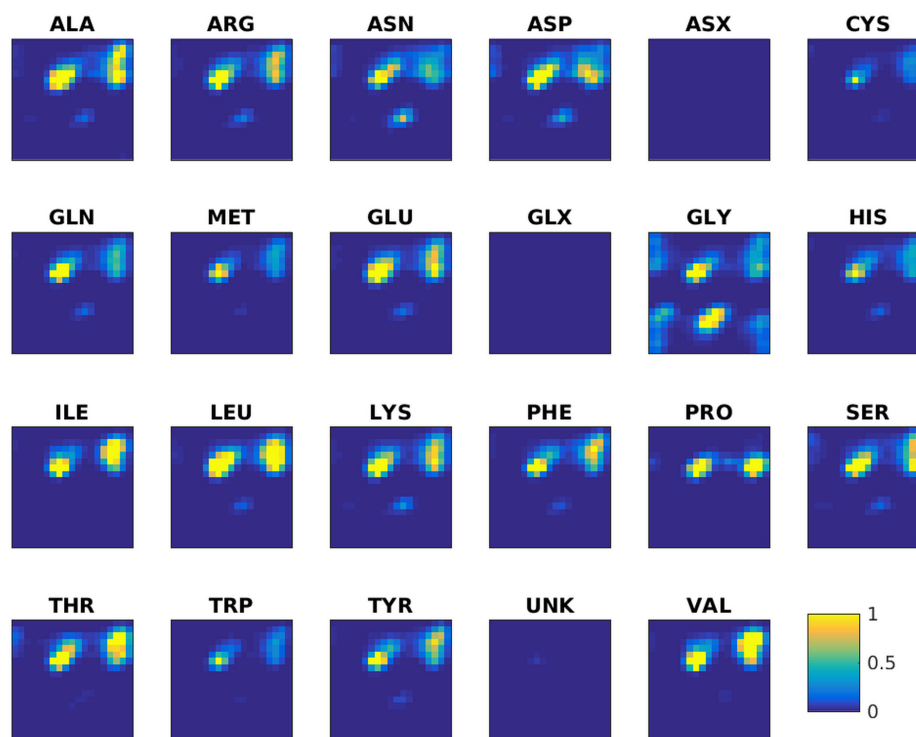


**Figure 2.** ROC curves for each enzymatic class based on kNN and *Architecture 2*

202 **Effect of sequence redundancy and sample size.** Analysis of protein datasets is often performed  
203 after removal of redundancy, such that the remaining entries do not overreach a pre-arranged threshold of  
204 sequence identity. In this particular work the author chose not to employ data filtering strategies, since  
205 the pattern analysis method is based on structure similarity and not sequence similarity. Thus, even if  
206 proteins are present with high sequence identity, the distance metrics during classification do not exploit  
207 it. Based on the (by now) established opinion that structure is far more conserved than sequence in nature  
208 (Illergard2009), the aim was not to jeopardize the dataset by losing reliable structural entries over a  
209 sequence based threshold cutoff. Also, only X-ray crystallography data were used; such data represent  
210 a ‘snapshot’ of a given protein’s 3D structure. In order not to miss the multiple poses that the same  
211 protein may adopt in different crystallography experiments, sequence/threshold metrics were not applied  
212 to remove sequence-redundancy in the presented results.

213 Nevertheless, the performance of the method was also investigated on a non-redundant dataset and the  
214 classification accuracy was compared in respect to the original (redundant) dataset randomly subsampled  
215 to include equal number of proteins. This experiment allows to assess the effect of redundancy under  
216 conditions (number of samples). Since inference in deep networks requires the estimation of a very  
217 large number of parameters, a large amount of training data is required and therefore very strict filtering  
218 strategies could not be applied. A dataset (the *pdbaanr*) pre-compiled by PISCES (Wang and Dunbrack,  
219 2003), was used that includes only non-redundant sequences across all PDB files ( $n = 23242$  proteins, i.e.  
220 half in size of the original dataset). Representative chains are selected based on the highest resolution  
221 structure available and then the best R-values. Non-X-ray structures are considered after X-ray structures.  
222 As a note, the author also explored the Leaf algorithm (Bull et al., 2013) which is especially designed  
223 to maximize the number of retained proteins and has shown improvement over PISCES. However, the  
224 computational cost was too high (possibly due to the large number of samples) and the analysis was not  
225 completed. The classification performance was assessed on *Architecture 2* by using 80% of the samples  
226 for training and 20% of the samples for testing. For the non-redundant dataset the accuracy was 79.3% for  
227 kNN and 75.5% for linear-SVM, whereas for the sub-sampled dataset it was 85.7% for kNN and 83.2%  
228 for linear-SVM. The results show that for the selected classifier (kNN), the accuracy drops 4.4% when the  
229 number of samples are reduced to the half, and it also drops additionally 6.4% if the utilized samples are  
230 non-redundant. Also the decrease in performance is not inconsiderable, the achieved accuracy indicates  
231 that structural similarity is an important criterion for the prediction of enzymatic function.



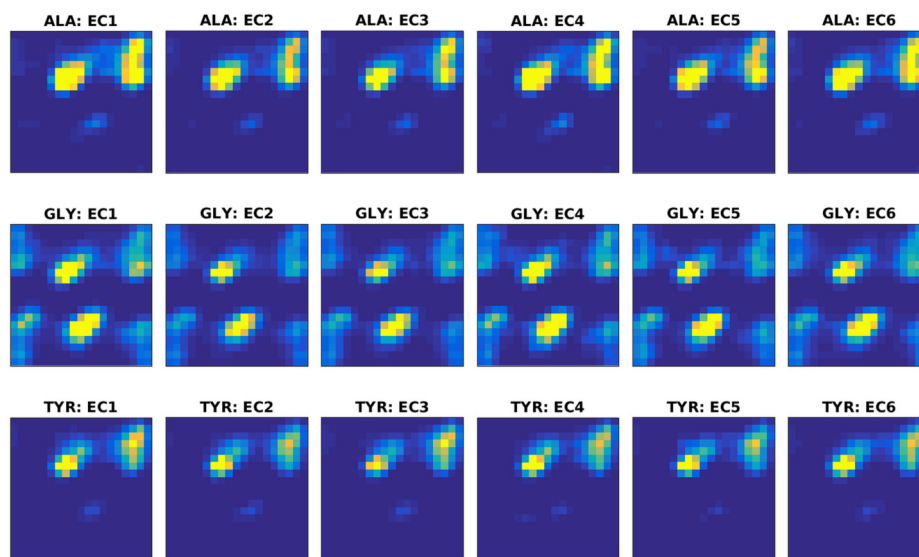


**Figure 3.** Torsion angles density maps (Ramachandran plots) averaged over all samples for each of the 20 standard and 3 non-standard (ASX, GLX, UNK) amino acids. The horizontal and vertical axes at each plot correspond to  $\phi$  and  $\psi$  angles and vary from  $-180^\circ$  (top left) to  $180^\circ$  (right bottom). The color scale (blue to yellow) is in the range  $[0, 1]$ . For an amino acid  $a$ , yellow means that the number of occurrences of the specific value  $(\phi, \psi)$  in all observations of  $a$  (within and across proteins) is at least equal to the number of proteins. On the opposite, blue indicates a small number of occurrences, and is observed for rare amino acids or unfavorable conformations.

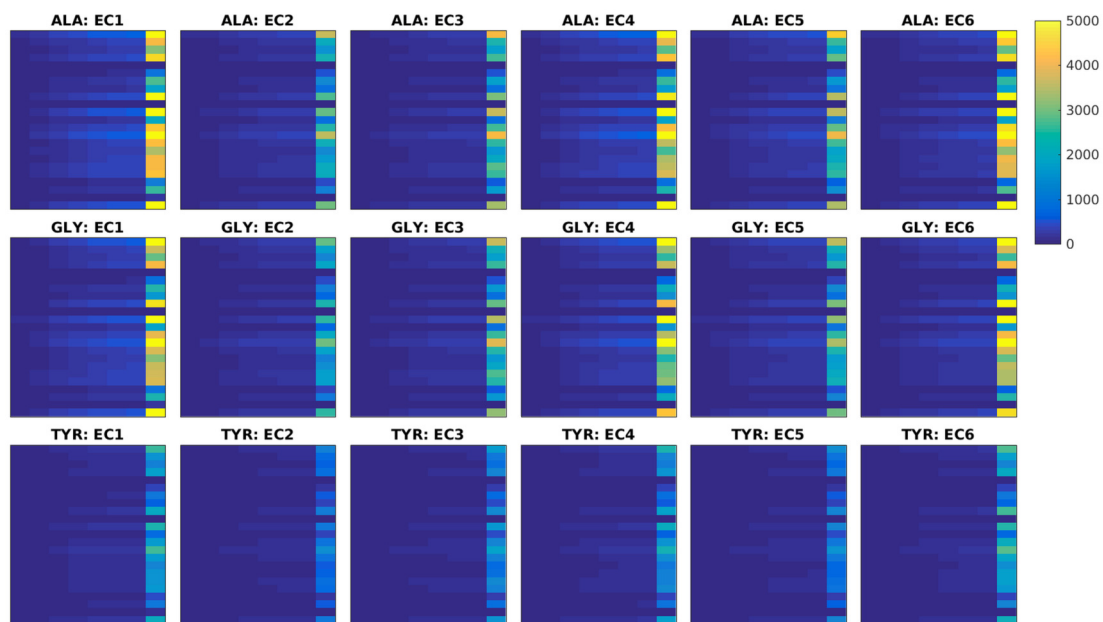
### 232 3.2 Structural representation and complementarity of features

233 Next, some examples of the extracted feature maps are illustrated, in order to provide some insight on the  
 234 representation of protein's 3D structure. The average (over all samples) 2D histogram of torsion angles for  
 235 each amino acid is shown in Fig. 3. The horizontal and vertical axes at each plot represent torsion angles  
 236 (in  $[-180^\circ, 180^\circ]$ ). It can be observed that the non-standard (ASX, GLX, UNK) amino acids are very rare,  
 237 thus their density maps have nearly zero values. The same color scale was used in all plots to make feature  
 238 maps comparable, as "seen" by the deep network. Since the histograms are (on purpose) not normalized  
 239 for each sample, rare amino acids will have few visible features and due to the 'max-pooling operator'  
 240 will not be selected as significant features. The potential of these feature maps to differentiate between  
 241 classes is illustrated in Fig. 4 for three randomly selected amino acids (ALA, GLY, TYR). Overall the  
 242 spatial patterns in each class are distinctive and form a multi-dimensional signature for each sample. As a  
 243 note, before training of the CNN ensemble data standardization is performed by subtracting the mean  
 244 density map. The same map is used to standardize the test sample during assessment.

245 Examples of features maps representing amino acid distances ( $X_D$ ) are illustrated in figures 1 and 5.  
 246 Fig. 1 illustrates an image slice across the 3rd dimension, i.e. one  $[m \times m]$  channel, and as introduced in  
 247 the 2D multichannel CNN, i.e. after mean-centering (over all samples). Fig. 5 illustrates image slices (of  
 248 size  $[m \times h_D]$ ) across the 1st dimension averaged within each class. Fig. 5 has been produced by selecting  
 249 the same amino acids as in Fig. 4 for easiness of comparison of the different feature representations. It  
 250 can be noticed that for all classes most pairwise distances are concentrated in the last bin, corresponding  
 251 to high distances between amino acids. Also, as expected there are differences in quantity of each amino  
 252 acid, e.g. by focusing on the last bin, it can be seen that ALA and GLY have higher values than TYR in  
 253 most classes. Moreover, the feature maps indicate clear differences between samples of different classes.



**Figure 4.** Ramachandran plots averaged across samples within each class. Rows correspond to amino acids and columns to functional classes. Three amino acids (ALA, GLY, TYR) are randomly selected for illustration of class separability. The horizontal and vertical axes at each plot correspond to  $\phi$  and  $\psi$  angles and vary from  $-180^\circ$  (top left) to  $180^\circ$  (right bottom). The color scale (blue to yellow) is in the range  $[0, 1]$  as illustrated in Fig. 3.



**Figure 5.** Histograms of pairwise amino acid distances averaged across samples within each class. The same three amino acids (ALA, GLY, TYR) selected in Fig. 4 are also shown here. The horizontal axis at each plot represents the histogram bins (distance values in the range  $[5, 40]$ ). The vertical axis at each plot corresponds to the 23 amino acids sorted alphabetically from top to bottom (ALA, ARG, ASN, ASP, ASX, CYS, GLN, MET, GLU, GLX, GLY, HIS, ILE, LEU, LYS, PHE, PRO, SER, THR, TRP, TYR, UNK, VAL). Thus each row shows the histogram of distances for a specific pair of the amino acids (the one in the title and the one corresponding to the specific row). The color scale is the same for all plots and shown at the bottom of the figure.

**Table 3.** Cross-validation accuracy (average  $\pm$  standard deviation over 5 folds) for each feature set separately and after fusion of CNN outputs based on *Architecture 2*

Feature sets	linear-SVM	kNN
$X_A$ (angles)	$79.6 \pm 0.5$	$82.4 \pm 0.4$
$X_D$ (distances)	$88.1 \pm 0.4$	$89.8 \pm 0.2$
<b>Ensemble</b>	$88.0 \pm 0.4$	$90.1 \pm 0.2$

254 The discrimination ability and complementary of the extracted features in respect to classification  
255 performance is shown in Table 3. It can be observed that the relative position of amino acids and their  
256 arrangement in space (features  $X_D$ ) predict enzymatic function better than the backbone conformation  
257 (features  $X_A$ ). Also, the fusion of network decisions based on correlation distance outperforms predictions  
258 from either network alone, but the difference is only marginal in respect to the predictions by  $X_D$ . In  
259 all cases the differences in prediction for the performed experiments (during cross validation) was very  
260 small (usually standard deviation  $< 0.5\%$ ), indicating that the method is robust to variations in training  
261 examples.

## 262 4 DISCUSSION

263 A deep CNN ensemble was presented that performs enzymatic function classification through fusion  
264 in feature level and decision level. The method has been applied for the prediction of the primary EC  
265 number and achieved 90.1% accuracy, which is a considerable improvement over the accuracy (73.5%)  
266 achieved in previous work (Amidi et al., 2016) when only structural information was incorporated.

267 Many methods have been proposed in the literature using different features and different classifiers.  
268 Nasibov and Kandemir-Cavas (2009b) obtained 95%-99% accuracy by applying kNN-based classification  
269 on 1200 enzymes based on their amino acid composition. Shen and Chou (2007b) fused results derived  
270 from the functional domain and evolution information and obtained 93.7% average accuracy on 9,832  
271 enzymes. On the same dataset Wang et al. (2011) improved the accuracy (which ranged from 81% to  
272 98% when predicting the first three EC digits) by using sequence encoding and SVM for hierarchy labels.  
273 Kumar and Choudhary (2012b) reported overall accuracy of 87.7% in predicting the main class for 4,731  
274 enzymes using random forests. Volpato et al. (2013) applied neural networks on the full sequence and  
275 achieve 96% correct classification on 6,000 non-redundant proteins. Most of these works have been  
276 applied on a subset of enzymes and have not been tested for large-scale annotation. Also they incorporate  
277 sequence-based features.

278 Assessment of the relationship between function and structure (Todd et al., 2001) revealed 95%  
279 conservation of the fourth EC digit for proteins with up to 30% sequence identity. Similarity, Devos  
280 and Valencia (2000) concluded that enzymatic function is mostly conserved for the first digit of EC  
281 code whereas more detailed functional characteristics are poorly conserved. It is generally believed that  
282 as sequences diverge, 3D protein structure becomes a more reliable predictor than sequence, and that  
283 structure is far more conserved than sequence in nature (Illergård et al., 2009). Thus, the focus of this  
284 study was to explore the predictive ability of 3D structure alone and provide a tool that can generalize in  
285 cases where sequence information is insufficient. Thus the presented results are not directly comparable  
286 to the ones of previous methods which incorporate sequence information. If desired, the current approach  
287 can also be combined with sequence-related features; in such a case it is expected that classification  
288 accuracy would further increase.

289 A possible limitation of the proposed approach is that the extracted features do not capture the  
290 topological properties of the 3D structure. Due to the statistical nature of the implemented descriptors,  
291 calculated by considering the amino acids as elements in Euclidean space, connectivity information is not  
292 strictly retained. The author and colleagues recently started to investigate in parallel the predictive power  
293 of the original 3D structure, represented as a volumetric image, without the extraction of any statistical  
294 features. Since the more detailed representation increased the dimensionality considerably, new ways  
295 are being explored to optimally incorporate the relationship between the structural units (amino-acids) in  
296 order not to impede the learning process.

297 **5 CONCLUSIONS**

298 A method was presented that extracts shape features from the 3D protein geometry that are introduced  
299 into a deep CNN ensemble for enzymatic function prediction. The investigation of protein function  
300 based only on structure reveals relationships hidden at the sequence level and provides the foundation  
301 to build a better understanding of the molecular basis of biological complexity. Overall, the presented  
302 approach can provide quick protein function predictions on extensive datasets opening the path for  
303 relevant applications, such as pharmacological target identification. Future work includes application of  
304 the method for prediction of the hierarchical relation of function subcategories and annotation of enzymes  
305 up to the last digit of the enzyme classification system.

306 **Acknowledgments**

307 The authors want to thank Prof. N. Paragios from the Center for Visual Computing, CentraleSupélec,  
308 Paris, for providing the means to complete this study and Dr. D. Vlachakis from the Multidimensional  
309 Data Analysis and Knowledge Management Laboratory, University of Patras, for useful discussions on  
310 the biological aspects.

311 **REFERENCES**

- 312 (1992). *Enzyme nomenclature 1992. Recommendations of the Nomenclature Committee of the Inter-*  
313 *national Union of Biochemistry and Molecular Biology on the Nomenclature and Classification of*  
314 *Enzymes*. Number Ed. 6. Academic Press.
- 315 Amidi, A., Amidi, S., Vlachakis, D., Paragios, N., and Zacharaki, E. I. (2016). A machine learning  
316 methodology for enzyme functional classification combining structural and protein sequence descriptors.  
317 In *Bioinformatics and Biomedical Engineering*, pages 728–738. Springer.
- 318 Bermejo, G. A., Clore, G. M., and Schwieters, C. D. (2012). Smooth statistical torsion angle potential  
319 derived from a large conformational database via adaptive kernel density estimation improves the  
320 quality of nmr protein structures. *Protein Science*, 21(12):1824–1836.
- 321 Borgwardt, K. M., Ong, C. S., Schönauer, S., Vishwanathan, S., Smola, A. J., and Kriegel, H.-P. (2005).  
322 Protein function prediction via graph kernels. *Bioinformatics*, 21(suppl 1):i47–i56.
- 323 Borro, L. C., Oliveira, S. R., Yamagishi, M. E., Mancini, A. L., Jardine, J. G., Mazoni, I., Santos, E. D.,  
324 Higa, R. H., Kuser, P. R., and Neshich, G. (2006). Predicting enzyme class from protein structure using  
325 bayesian classification. *Genet. Mol. Res*, 5(1):193–202.
- 326 Bull, S. C., Muldoon, M. R., and Doig, A. J. (2013). Maximising the size of non-redundant protein  
327 datasets using graph theory. *PLoS one*, 8(2):e55484.
- 328 Cai, C., Han, L., Ji, Z. L., Chen, X., and Chen, Y. Z. (2003). Svm-prot: web-based support vector machine  
329 software for functional classification of a protein from its primary sequence. *Nucleic acids research*,  
330 31(13):3692–3697.
- 331 Chang, C.-C. and Lin, C.-J. (2011). Libsvm: a library for support vector machines. *ACM Transactions on*  
332 *Intelligent Systems and Technology (TIST)*, 2(3):27.
- 333 Chen, C., Tian, Y.-X., Zou, X.-Y., Cai, P.-X., and Mo, J.-Y. (2006). Using pseudo-amino acid compo-  
334 sition and support vector machine to predict protein structural class. *Journal of Theoretical Biology*,  
335 243(3):444–448.
- 336 Devos, D. and Valencia, A. (2000). Practical limits of function prediction. *Proteins: Structure, Function,*  
337 *and Bioinformatics*, 41(1):98–107.
- 338 Dobson, P. D. and Doig, A. J. (2005). Predicting enzyme class from protein structure without alignments.  
339 *Journal of molecular biology*, 345(1):187–199.
- 340 Godzik, A. (2011). Metagenomics and the protein universe. *Current opinion in structural biology*,  
341 21(3):398–403.
- 342 Han, L., Cai, C., Ji, Z., Cao, Z., Cui, J., and Chen, Y. (2004). Predicting functional family of novel  
343 enzymes irrespective of sequence similarity: a statistical learning approach. *Nucleic acids research*,  
344 32(21):6437–6444.
- 345 Huang, W.-L., Chen, H.-M., Hwang, S.-F., and Ho, S.-Y. (2007). Accurate prediction of enzyme subfamily  
346 class using an adaptive fuzzy k-nearest neighbor method. *Biosystems*, 90(2):405–413.
- 347 Illergård, K., Ardell, D. H., and Elofsson, A. (2009). Structure is three to ten times more conserved  
348 than sequence—a study of structural response in protein cores. *Proteins: Structure, Function, and*  
349 *Bioinformatics*, 77(3):499–508.

- 350 Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional  
351 neural networks. In *Advances in neural information processing systems*, pages 1097–1105.
- 352 Kumar, C. and Choudhary, A. (2012a). A top-down approach to classify enzyme functional classes and  
353 sub-classes using random forest. *EURASIP Journal on Bioinformatics and Systems Biology*, 1:1–14.
- 354 Kumar, C. and Choudhary, A. (2012b). A top-down approach to classify enzyme functional classes  
355 and sub-classes using random forest. *EURASIP Journal on Bioinformatics and Systems Biology*,  
356 2012(1):1–14.
- 357 Lee, B. J., Shin, M. S., Oh, Y. J., Oh, H. S., and Ryu, K. H. (2009). Identification of protein functions  
358 using a machine-learning approach based on sequence-derived properties. *Proteome science*, 7(1):1.
- 359 Li, Y. and Shibuya, T. (2015). Malphite: A convolutional neural network and ensemble learning based  
360 protein secondary structure predictor. In *IEEE Int. Conf. on Bioinformatics and Biomedicine (BIBM)*,  
361 pages 1260–1266.
- 362 Lin, Z., Lanchantin, J., and Qi, Y. (2016). Must-cnn: A multilayer shift-and-stitch deep convolutional  
363 architecture for sequence-based protein structure prediction. In *30th AAAI Conference on Artificial  
364 Intelligence*.
- 365 Lu, L., Qian, Z., Cai, Y.-D., and Li, Y. (2007). Ecs: an automatic enzyme classifier based on functional  
366 domain composition. *Computational biology and chemistry*, 31(3):226–232.
- 367 Nagao, C., Nagano, N., and Mizuguchi, K. (2014). Prediction of detailed enzyme functions and identifica-  
368 tion of specificity determining residues by random forests. *PLoS one*, 9(1):1–12.
- 369 Nasibov, E. and Kandemir-Cavas, C. (2009a). Efficiency analysis of knn and minimum distance-based  
370 classifiers in enzyme family prediction. *Computational biology and chemistry*, 33(6):461–464.
- 371 Nasibov, E. and Kandemir-Cavas, C. (2009b). Efficiency analysis of knn and minimum distance-based  
372 classifiers in enzyme family prediction. *Computational biology and chemistry*, 33(6):461–464.
- 373 Qiu, J.-D., Huang, J.-H., Shi, S.-P., and Liang, R.-P. (2010). Using the concept of chou's pseudo amino  
374 acid composition to predict enzyme family classes: an approach with support vector machine based on  
375 discrete wavelet transform. *Protein and peptide letters*, 17(6):715–722.
- 376 Sharma, M. and Garg, P. (2014). Computational approaches for enzyme functional class prediction: A  
377 review. *Current Proteomics*, 11(1):17–22.
- 378 Shen, H.-B. and Chou, K.-C. (2007a). Ezympred: a top-down approach for predicting enzyme functional  
379 classes and subclasses. *Biochemical and biophysical research communications*, 364(1):53–59.
- 380 Shen, H.-B. and Chou, K.-C. (2007b). Ezympred: a top-down approach for predicting enzyme functional  
381 classes and subclasses. *Biochemical and biophysical research communications*, 364(1):53–59.
- 382 Spencer, M., Eickholt, J., and Cheng, J. (2015). A deep learning network approach to ab initio protein  
383 secondary structure prediction. *IEEE/ACM Trans. on Computational Biology and Bioinformatics  
384 (TCBB)*, 12(1):103–112.
- 385 Todd, A. E., Orengo, C. A., and Thornton, J. M. (2001). Evolution of function in protein superfamilies,  
386 from a structural perspective. *Journal of molecular biology*, 307(4):1113–1143.
- 387 Vedaldi, A. and Lenc, K. (2015). Matconvnet: Convolutional neural networks for matlab. In *Proceedings  
388 of the 23rd ACM international conference on Multimedia*, pages 689–692. ACM.
- 389 Volpato, V., Adelfio, A., and Pollastri, G. (2013). Accurate prediction of protein enzymatic class by n-to-1  
390 neural networks. *BMC bioinformatics*, 14(1):1.
- 391 Wang, G. and Dunbrack, R. L. (2003). Pisces: a protein sequence culling server. *Bioinformatics*,  
392 19(12):1589–1591.
- 393 Wang, Y.-C., Wang, X.-B., Yang, Z.-X., and Deng, N.-Y. (2010). Prediction of enzyme subfamily class  
394 via pseudo amino acid composition by incorporating the conjoint triad feature. *Protein and Peptide  
395 Letters*, 17(11):1441–1449.
- 396 Wang, Y.-C., Wang, Y., Yang, Z.-X., and Deng, N.-Y. (2011). Support vector machine prediction of  
397 enzyme function with conjoint triad feature and hierarchical context. *BMC systems biology*, 5(1):1.
- 398 Yadav, S. K. and Tiwari, A. K. (2015). Classification of enzymes using machine learning based approaches:  
399 a review. *Machine Learning and Applications: An International Journal (MLAIJ)*, 2(3/4).
- 400 Zhou, X.-B., Chen, C., Li, Z.-C., and Zou, X.-Y. (2007). Using chou's amphiphilic pseudo-amino acid  
401 composition and support vector machine for prediction of enzyme subfamily classes. *Journal of  
402 theoretical biology*, 248(3):546–551.