



Contents lists available at ScienceDirect

## Analytical Biochemistry

journal homepage: [www.elsevier.com/locate/yabio](http://www.elsevier.com/locate/yabio)

## Predicting peroxidase subcellular location by hybridizing different descriptors of Chou' pseudo amino acid patterns

Yong-Chun Zuo<sup>a,\*</sup>, Yong Peng<sup>b</sup>, Li Liu<sup>b</sup>, Wei Chen<sup>c</sup>, Lei Yang<sup>d,\*</sup>, Guo-Liang Fan<sup>b,\*</sup><sup>a</sup>The Key Laboratory of Mammalian Reproductive Biology and Biotechnology of the Ministry of Education, Inner Mongolia University, Hohhot 010021, China<sup>b</sup>Laboratory of Theoretical Biophysics, School of Physical Science and Technology, Inner Mongolia University, Hohhot 010021, China<sup>c</sup>Center of Genomics and Computational Biology, College of Sciences, Hebei United University, Tangshan 063000, China<sup>d</sup>College of Bioinformatics Science and Technology, Harbin Medical University, Harbin 150081, China

## ARTICLE INFO

## Article history:

Received 23 February 2014

Received in revised form 22 April 2014

Accepted 25 April 2014

Available online xxxxx

## Keywords:

Peroxidase proteins

Chou' pseudo amino acid patterns

GO-homology annotation

Prediction performance

## ABSTRACT

Peroxidases as universal enzymes are essential for the regulation of reactive oxygen species levels and play major roles in both disease prevention and human pathologies. Automated prediction of functional protein localization is rarely reported and also is important for designing new drugs and drug targets. In this study, we first propose a support vector machine (SVM)-based method to predict peroxidase subcellular localization. Various Chou' pseudo amino acid descriptors and gene ontology (GO)-homology patterns were selected as input features to multiclass SVM. Prediction results showed that the smoothed PSSM encoding pattern performed better than the other approaches. The best overall prediction accuracy was 87.0% in a jackknife test using a PSSM profile of pattern with width = 5. We also demonstrate that the present GO annotation is far from complete or deep enough for annotating proteins with a specific function.

© 2014 Elsevier Inc. All rights reserved.

Peroxidases are ubiquitous enzymes that catalyze a number of oxidative reactions by using various peroxides as electron acceptors [1,2]. These peroxidase proteins are central elements of the antioxidant defense system, which are extremely widespread in almost all microorganisms and higher organisms. They are essential for the regulation of reactive oxygen species levels and for the promotion of various substrates' oxidation [3–5]. There has been increased interest in them over the past few years; for example, the mammalian heme peroxidase enzymes play major roles in both disease prevention and human pathology defense [6,7]. Therefore, knowing the localization of peroxidase proteins will be important for disease prevention and human pathologies.

Proteins in various subcellular locations play distinct roles in biological processes, such as triggering programmed cell death. Protein localization may be used as a starting point for function prediction systems. Knowing a protein's localization is an important step toward understanding its function [8,9]. Experimental and computational methods are two very important methods for annotating protein functional information. During the past 2 decades, a substantial amount of bioinformatics work for predicting

protein subcellular location has been carried out and rapidly developed; significant progress has been achieved with the establishment of various organism-specific benchmark datasets [10–15]. However, to the best of our knowledge, there are few theoretical methods for localization prediction for proteins of specific function.

Therefore, it is becoming crucial to develop a reliable automatic subcellular localizer for identifying the locations of functional proteins. In this study we first attempted to annotate the subcellular localization of a specific oxidoreductase, peroxidase, by using a computational method based on state-of-the-art features. Several different descriptors of the Chou' pseudo amino acid pattern have been discussed for localization prediction [16–21], including amino acid composition (AAC) [22], dipeptide composition (DC) [23,24], split amino acid composition (SAAC) [25], evolutionary information (PSSM) [10,26–28], and gene ontology (GO) of homologous proteins [29–32]. All of the above features were selected as input parameters to establish an automatic subcellular classifier. The best overall prediction accuracy achieved 87.0% in a jackknife test for eight locations by using a PSSM profile with width = 5. The GO-homology annotation with different sequence identities was also discussed; the evaluation results showed the present GO annotation is far from complete or deep enough for accurately annotating the localization of peroxidase proteins.

\* Corresponding authors. Fax: +86 471 5227683.

E-mail addresses: [yczuo@imu.edu.cn](mailto:yczuo@imu.edu.cn) (Y.-C. Zuo), [yanglei\\_hmu@163.com](mailto:yanglei_hmu@163.com) (L. Yang), [eequoliangfan@sina.com](mailto:eequoliangfan@sina.com) (G.-L. Fan).

## 85 Materials and methods

### 86 Benchmark datasets

87 The data of peroxidase proteins used in this research were  
 88 extracted from the PeroxiBase database [33]. PeroxiBase is a  
 89 unique specialized database, which is devoted to established com-  
 90 prehensive peroxidase families and superfamilies from both  
 91 eukaryotes and prokaryotes. More than 10,000 peroxidase-encod-  
 92 ing sequences come from 940 organisms, and each sequence is  
 93 individually annotated in this database. Since the number of mul-  
 94 ti-plex proteins in the existing database is not large enough to con-  
 95 struct a statistically meaningful benchmark dataset for studying a  
 96 case of multiple locations, only the proteins with singleplex  
 97 locations were used in this experiment, and every protein is  
 98 characterized by an expert sequence annotation procedure, with  
 99 manual curation, which is a guarantee of quality necessary for per-  
 100 forming subcellular localization analysis. After the redundant  
 101 sequences were removed using the CD-HIT algorithm [34], 586  
 102 nonredundant peroxidase proteins were obtained. According to  
 103 the annotation information, these defensin sequences can be clas-  
 104 sified into eight subcellular locations: apoplasmic (30), chloroplasmic  
 105 (44), cytosolic (265), mitochondrial (44), peroxisomal (107),  
 106 secreted (23), stromal (37), and thylakoid (37). After measuring  
 107 by the CD-HIT program, most of the protein similarity scores in  
 108 each family were lower than 80%.

### 109 Features and modules

110 Support vector machine (SVM), as a strong machine learning  
 111 technique, is used to evaluate various alternative features of our  
 112 work. SVM is a machine learning algorithm based on statistical  
 113 learning theory, which has been successfully used for classification  
 114 [35]. The basic idea of SVM is to transform the data into a high-  
 115 dimensional feature space and then determine the optimal separ-  
 116 ating hyperplane by using a kernel function. In this work, we used  
 117 the free software LIBSVM to predict peroxidase protein location. A  
 118 radial basis function (RBF) was chosen as the kernel function. For  
 119 multiclassification, SVM uses a one-versus-one strategy and  
 120 constructs  $k \times (k - 1)/2$  classifiers and voting strategy to assign  
 121 the class for an arbitrary protein sequence. Here various features  
 122 of a protein sequence were utilized to perform a comprehensive  
 123 study and achieve maximum accuracy.

### 124 PSSM profile of patterns

125 Evolutionary conservation usually reflects important biological  
 126 function. An amino acid at a conserved site of a protein is preferred  
 127 to locate at a functionally important region [36]. PSI-BLAST is a  
 128 robust measure of residue conservation in a given location. Evolu-  
 129 tionary information on protein sequences like PSSM can be created  
 130 using a PSI-BLAST search. Compared to the compositional informa-  
 131 tion, the PSSM profile provides more important information of  
 132 evolutionary significance about residue conservation at a given  
 133 position in a protein sequence [31,37]. In this study, the PSSM  
 134 was generated using the PSI-BLAST search with a cutoff  $E$  value  
 135 of 0.001 against the Swiss-Prot database.

136 The PSSM provides a matrix of dimension  $L$  rows and 20  
 137 columns for a protein chain with  $L$  amino acid residues, where  
 138 20 columns represent the occurrence/substitution of each type of  
 139 20 amino acids [38]. We summed all of the rows in the PSSM cor-  
 140 responding to the same amino acid in the sequence and then  
 141 divided each element by the length of the sequence. In the predic-  
 142 tion of peroxidase location, we used PSSM profiles with different

similarities to generate 400 dimension ( $20 \times 20$  residue pairs) 143  
 input vectors as parameters. 144

### Composition profile of patterns 145

The aim of calculating the protein composition is to transform 146  
 the variable lengths of the protein sequence to fixed-length vec- 147  
 tors. This is an important and crucial step for protein classification 148  
 using a computational approach because it requires a fixed-length 149  
 pattern. 150

### Amino acid and dipeptide compositions 151

The AAC representation of a given sequence is composed of 20 152  
 different amino acids with a variety of shapes, sizes, and chemical 153  
 properties. A protein can be represented as a 20-dimensional (20D) 154  
 vector according to AAC [22]. DC is the occurrence frequency of 155  
 each of 2 adjacent amino acid residues. It is used to encapsulate 156  
 the global information of each protein sequence, and a protein 157  
 can be represented as a 400D vector by means of DC [39–41]. In 158  
 this study, the AAC and DC of the  $N$ -part split amino acid composi- 159  
 tion were selected as classification vectors. 160

### Split amino acid composition 161

In simple amino acid-, dipeptide-, and pseudo amino acid-based 162  
 compositions, the composition is taken at once for the whole 163  
 sequence, whereas in the split amino acid composition model, 164  
 the protein sequence is divided into different parts and the compo- 165  
 sition of each part is calculated separately [25]. The composition is 166  
 taken independently for the  $N$  parts of the protein sequence [42]. 167  
 Hence, the advantage of SAAC over standard AAC is that it provides 168  
 a greater weight of compositional biasness to proteins that have a 169  
 signal at different sequence regions. In our SAAC model each pro- 170  
 tein is divided into 1 to 10 parts to train the optimal parameter 171  
 combination for the SVM program. 172

### Gene ontology profile of patterns 173

Gene Ontology is one of the databases that describes molecular 174  
 function, and the molecular function of the GO database is corre- 175  
 lated to the subcellular location [43]. Accordingly, protein 176  
 sequences formulated in the GO database space would be clustered 177  
 in a way that better reflects their subcellular locations [29]. How- 178  
 ever, to incorporate more information, instead of using only 0 and 179  
 1 element, as done in Ref. [44], here let us use a different approach 180  
 as described below. 181

182 First, we searched for the homologous proteins of protein **P**  
 183 from the Swiss-Prot database (released on 5 September 2012)  
 184 using the PSI-BLAST method, with the expected value  $E \leq 0.001$   
 185 for the BLAST parameter [31]. Second, we collected those proteins  
 186 that had  $\geq 60\%$  pairwise sequence identity with protein **P** into a  
 187 subset,  $\mathbf{P}^{\text{homo}}$ , called the “homology set” of **P**. All the elements in  
 188  $\mathbf{P}^{\text{homo}}$  could be deemed the “representative proteins” of **P**, sharing  
 189 some similar attributes such as structural conformation and  
 190 biological function. These representative proteins retrieved from  
 191 the Swiss-Prot database must each have their own accession num-  
 192 ber. Third, we searched each of the accession numbers collected in  
 193 the second step against the GO database to find the corresponding  
 194 GO number. Last, we statistically analyzed each coordinate of the  
 195 vector and found that many of the coordinates were equal to 0.  
 196 This denoted that certain GOs did not belong to any protein; these  
 197 GOs were eliminated, and the dimension of the GO feature vector  
 198 was decreased in this manner. 199

Download English Version:

<https://daneshyari.com/en/article/7559187>

Download Persian Version:

<https://daneshyari.com/article/7559187>

[Daneshyari.com](https://daneshyari.com)