

Predict Remote Homology Proteins by Merging Grey Incidence Analysis and Domain Similarity Analysis

Journal:	<i>Bioinformatics</i>
Manuscript ID	Draft
Category:	Original Paper
Date Submitted by the Author:	n/a
Complete List of Authors:	Lin, Weizhong; Jingdezhen Ceramic Institute, School of Information Engineering Xiao, Xuan; Jingdezhen Ceramic Institute Qiu, Wang-Ren; Jingdezhen Ceramic Institute Chou, Kuo-Chen; University of Electronic Science and Technology of China
Keywords:	Protein sequence analysis, Information extraction, Bioinformatics, Remote homology proteins prediction

Predict Remote Homology Proteins by Merging Grey Incidence Analysis and Domain Similarity Analysis

Wei-Zhong Lin¹, Xuan Xiao^{1,2,*}, Wangren Qiu¹, Kuo-Chen Chou^{2,3}

1 Computer Department, Jing-De-Zhen Ceramic Institute, Jing-De-Zhen 333046 China; **2** Gordon Life Science Institute, Boston, MA 02478, USA; **3** Center for Informational Biology, University of Electronic Science and Technology of China, Chengdu, 610054, China

Authors' e-mail address:
Xuan Xiao: jdzxiaoxun@163.com
Wei-Zhong Lin: linweizhongjci@sina.com
Wangren Qiu: qiuone@163.com
Kuo-Chen Chou: kcchou@gordonlifescience.org

*Corresponding author: Xuan Xiao, email: Xuan Xiao: jdzxiaoxun@163.com

Running Title: Remote Homology Protein detection

ABSTRACT

Protein remote homology detection is a challenging problem for drug development. Although there are a couple of methods to deal with this problem, the benchmark datasets based on which the existing models were trained and tested contained many high homologous samples due to the fact that the cutoff threshold was set at 95%. In this study, we reconstructed the benchmark dataset by setting the threshold at 40%, meaning none of the proteins included has more than 40% pairwise sequence identity with any other. Using the new benchmark dataset, we proposed a new method called PHom-GRADSI to detect the remote homologous proteins by integrating various ranking approaches via grey incidence analysis and function domain similarity index. Rigorous cross-validations have indicated that the new predictor is superior to its counterparts in both enhancing successes rates and reducing computational cost. The predictor can be download from <https://github.com/jcilwz/RemoteHomology/tree/master/program>.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1. INTRODUCTION

Detecting remote homology relationship among proteins plays one of the fundamental and central roles in computational proteomics. It is particularly useful for drug development (see, e.g., [1, 2]). With the development of sequencing techniques, the protein sequence data rapidly raise. To find those proteins structure and function is more and more urgent. Although X-ray crystallography is a powerful tool in determining protein 3D structures, it is time-consuming and expensive. Particularly, not all proteins can be successfully crystallized, particularly for membrane proteins. The NMR technique is indeed a very powerful tool in determining the 3D structures for membrane proteins as indicated by a series of recent publications (see, e.g., [3-7]), it is time-consuming and costly. To acquire the structural information in a timely manner, one has to resort to various structural bioinformatics tools based on the sequence similarity principle (see, e.g., [8]). Unfortunately, such principle cannot cover the cases of remote homology proteins. In view of this, considerable efforts [9-14] have been made to detect remote homology proteins.

Although these methods each had their own merits and did play stimulating role in this area, further work is needed. Firstly, the benchmark datasets used in their studies had high similarity. For instance, the benchmark dataset in [9, 12] contains 7329 proteins from 1070 different super families, with pairwise sequence identity cutoff set at 95%. In other words, it would allow those proteins with higher than 80% similarity in the data set. Secondly, the ranking algorithm used in those studies would spend a lot of time to training the learning model. For example, if the training dataset has N proteins, the LambdaMART need to deal with N² proteins pair samples.

The present study was initiated to address the two problems with the aim to develop a more powerful method in this regard.

2. MATERIALS AND METHOD

2.1 Benchmark Dataset

According to Chou's 5-step rules [15] that have been widely and increasingly used by many investigators (see, e.g., [16-32]), the first prerequisite in establishing a new predictor is to construct or select an effective benchmark dataset.

In this study, the benchmark dataset was taken from Liu et al. [12]. It included 7329 proteins from 1070 different super families and 1824 families derived from SCOP database. To reduce the redundancy and homology bias, the program CD-HIT[33] was adopted to cut down those proteins that had ≥40% pairwise sequence identity to any other in the dataset. Furthermore we removed those families that just had one protein sequence. Finally, we obtained 3128 proteins from 540 super-families and 777 families.

2.2 Sample Similarity Analysis

2.2.1 Grey Incidence Analysis of proteins formulated by Grey-PSSM

Given a protein with L amino acid residues, it is usually expressed by

$$\mathbf{P} = R_1 R_2 R_3 \cdots R_i \cdots R_L \quad (1)$$

where R_i ($i = 1, 2, \dots, L$) is the i -th residue in the protein. Since all the existing machine-learning algorithms can only handle vector but not sequence samples [34], one has to convert Eq.1 into a vector model. But a biological sequence expressed as a vector in the discrete framework may lose all the sequence-order or pattern information.

To avoid completely losing this kind of information for proteins, the pseudo amino acid composition (PseAAC)[35, 36] was proposed. Ever since the concept of Chou's PseAAC was proposed, it has been widely used in nearly all the areas of computational proteomics (see, e.g., [37-44] as well as a long list of references cited in [45, 46]). According to the general PseAAC[15], the protein of Eq.1 can be formulated as

$$\mathbf{P} = [\Psi_1 \ \Psi_2 \ \cdots \ \Psi_u \ \cdots \ \Psi_\Omega]^T \quad (2)$$

where \mathbf{T} is the transposing operator, the subscript Ω is an integer, and its value and the components Ψ_u ($u = 1, 2, \dots$) will depend on how to extract the desired features and properties from the protein sequence.

In this study, the model, Grey-PSSM proposed by Lin[47, 48], is adopted. It extracted the sequential evolution information by the Position Specific Scoring Matrix (PSSM). For the concrete procedures, refer to the original papers [47, 48].

After the Grey-PSSM treatment, we have finally got a 60-D PseKNC vector for Eq.2; i.e., its subscript parameter $\Omega = 60$ and each of the 60 components therein has been uniquely defined.

Assume

$$\mathbb{S} = \{P_1, P_2, \dots, P_N\} \quad (3)$$

are the set of protein samples, and P_i ($1 \leq i \leq N$) is the i^{th} protein. According to Equals [6]~[11] in Ref. [49], the distance $\Gamma(P_i, P_j)$ is defined as the grey incidence degree between P_i and P_j . The larger the value of $\Gamma(P_i, P_j)$, the more similar they are.

2.2.2 Domain Similarity Analysis

There are also other models used to formulating proteins except for the PseAAC. Here, we propose a novel method to describe the proteins. For a protein $P_i \in \mathbb{S}$, we describe its functional domains set by the following steps.

Step 1, \mathbb{S}_i^{homo} is the homology set of protein P_i and it is extracted by searching against UniProt release 2018_08 Swiss-Prot FASTA format flatfile by HMMER[50-52]. We just use the top 10 sequences if the search results have

more than 10 sequences. Therefore there are at most 10 proteins in \mathbb{S}_i^{homo} .

Step 2, for a protein in \mathbb{S}_i^{homo} , $h_k^i \in \mathbb{S}_i^{homo}$ ($1 \leq k \leq 10$), we annotate its functional domains by running hmmscan program against Pfam-A database (Pfam release 32.0). The Pfam-A includes 17,929 functional domains and 688 clans. We define the sets \mathbb{F} and \mathbb{C} as following.

$$\begin{aligned}\mathbb{F} &= \{f_1, f_2, \dots, f_{17929}\} \\ \mathbb{C} &= \{c_1, c_2, \dots, c_{688}\}\end{aligned}\quad (3)$$

where f_i ($1 \leq i \leq 17929$) denote the i -th functional domain in \mathbb{F} and c_i ($1 \leq i \leq 688$) the i -th clan in \mathbb{C} . Some functional domains have same clan. For example, the domains of "PF15884" and "PF17050" have the same clan "CL0683". Therefore, the functional domains set of protein h_k^i , the k -th homology protein of protein P_i , is expressed as a set

$$D_k^i = \{f_i | f_i \in \mathbb{F}\} \quad (4)$$

It is mean that all functional domains of h_k^i constitute the set D_k^i .

Step 3, the protein P_i is expressed as a domains set

$$D_i = \bigcup_{k=1}^{10} D_k^i \quad (5)$$

D_i is a set which is unioned together the functional domain set of each homology protein of the protein P_i . We define D_i as the functional domain set formulation of P_i .

As aboved steps, a protein is expressed a set including some functional domains from Pfam-A. For the proteins in same family or clan have similar functional domains, new distance between two proteins, named as Domain Similarity Index (DSI), can be defined based on the functional domains.

The algorithm of distance between P_i, P_j is discribed as follows.

$$1) \quad \text{If } D_i \cap D_j \neq \emptyset, \text{ DSI}(P_i, P_j) = \frac{|D_i \cap D_j|}{|D_i \cup D_j|}$$

$$2) \quad \text{Else If } D_i \cap D_j = \emptyset$$

We define $Clan_i$ to denote the clans set of P_i . $Clan_i$ includes the clans of each element in D_i

2.1) If $Clan_i \cap Clan_j \neq \emptyset$, $\text{DSI}(P_i, P_j) = F$. F is a constant and in this study F is equal to 0.2.

$$2.2) \quad \text{Else, } \text{DSI}(P_i, P_j) = 0$$

Where, $|\bullet|$ means the count of set, \cap is the intersection operator of two sets, and \cup is the union operator of two sets. From above discription, we have $0 \leq \text{DSI}(P_i, P_j) \leq 1$ and the larger the $\text{DSI}(P_i, P_j)$, the more similarity they are.

2.3 Operation Engine or Algorithm

In this study, the Grey Relational Analysis [53, 54] and the Domain Similarity Index was utilized to rank the relationship of proteins. Given a query protein, the system will search it against the benchmark dataset and return the top ranked proteins. The predictor thus formed is called “PHom-GRADSI”. Illustrated in **Figure 1** is a flowchart to show how the proposed predictor is working. In this paper, $w(1)$ and $w(2)$ are equal to 0.5.

3. RESULT AND DISCUSSION

The jackknife test is deemed the least arbitrary and most objective among three cross-validation methods: independent dataset test, K-fold cross-validation test and jackknife test [55]. Because the LambdaMART ranking algorithm used in preview studies [9, 12] consumed more training time and computer memory, as a compromise the 5-fold cross-validation test was adopted there. Now, we employed GRA and DSI to compute the relationship score between the query protein and benchmark dataset proteins, significantly reducing the computing time and memory. Therefore it would be feasible to use the most rigorous jackknife test to examine the prediction quality. The outcome thus obtained are given in **Table 1**, where we can see that PHom-GRADSI achieved the best performance in both the score of ROC1 and the score of ROC50.

In the same time, we used the Jaccard Index (JI) to calculate the similarity of the proteins. The JI is used to measure the similarity of two sets A and B. It is defined as following:

$$JI(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

The proteins were formulated as the functional domain set as Eq.(4). In the method of Jaccard Index, the score of ROC and ROC50 were 0.8196 and 0.8070 respectively. In contrast to this, the score of ROC and ROC50 achieved by DSI were 0.9053 and 0.8454 respectively. It can be concluded that DSI preforms better than Jaccard Index. The main reason is that some different domains have same clan. If there are not same elements in two proteins' functional domain set, the distance of Jaccard Index is zero, but the distance of DSI is greater zero. It is evident that the non zero distance of these proteins is more reasonable because their some functional domain have same clan.

Because not all proteins can be finded their homology protins, these proteins' functional domain set formulation is empty set (see Eq.(4)). For example, there are 23 protiens who can not be formulated as functional domain set in the **Benchmark Dataset**. The distences between these proteins and other proteins are zero according to the definition of JI and DSI . So we can not distinguish the similarity of these proteins who have no homology protiens. This situation is the worst failure DSI. In order to overcome the failure we merged GIA and Jaccard Index and merging GIA and DSI in predicting, respectively. From the

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Tabel 1, it is showed that the values of ROC1 and ROC50 are improved at least 0.04 after merging GIA into Jaccard Index or DSI. We can draw the grey incidence analysis can compensate for the worst failure of DSI.

4. CONCLUSION

Protein remote homology detection is vital for studying protein structures and functions. It is anticipated that the proposed method may become a useful high throughput toll for both basic research and drug design. In this study, a noval method DSI is proposed. It discribes a protein as a functional domain set and measure the distance of two proteins by comparing two proteins’ functional domain set similarity. The work testifies the DSI method is effective. This method formulating proteins and calculating distance between proteins may be used in other fields of predicting protein function or structure. We deliver programs of this noval method in <https://github.com/jcilwz/RemoteHomology/tree/master/program>. Everyone can download programs from this website and the usage is described in ReadMe.txt.

ACKNOWLEDGEMENTS

This work was support by the grants from the National Natural Science Foundation of China (No.61462047, 31560316, 31760315). Natural Science Foundation of Jiangxi Province, China (No. 20171ACB20023), the Department of Education of JiangXi Province (GJJ160866), The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

FIGURE LEGENDS

Figure 1. A flowchart to illustrate how the proposed predictor is working.

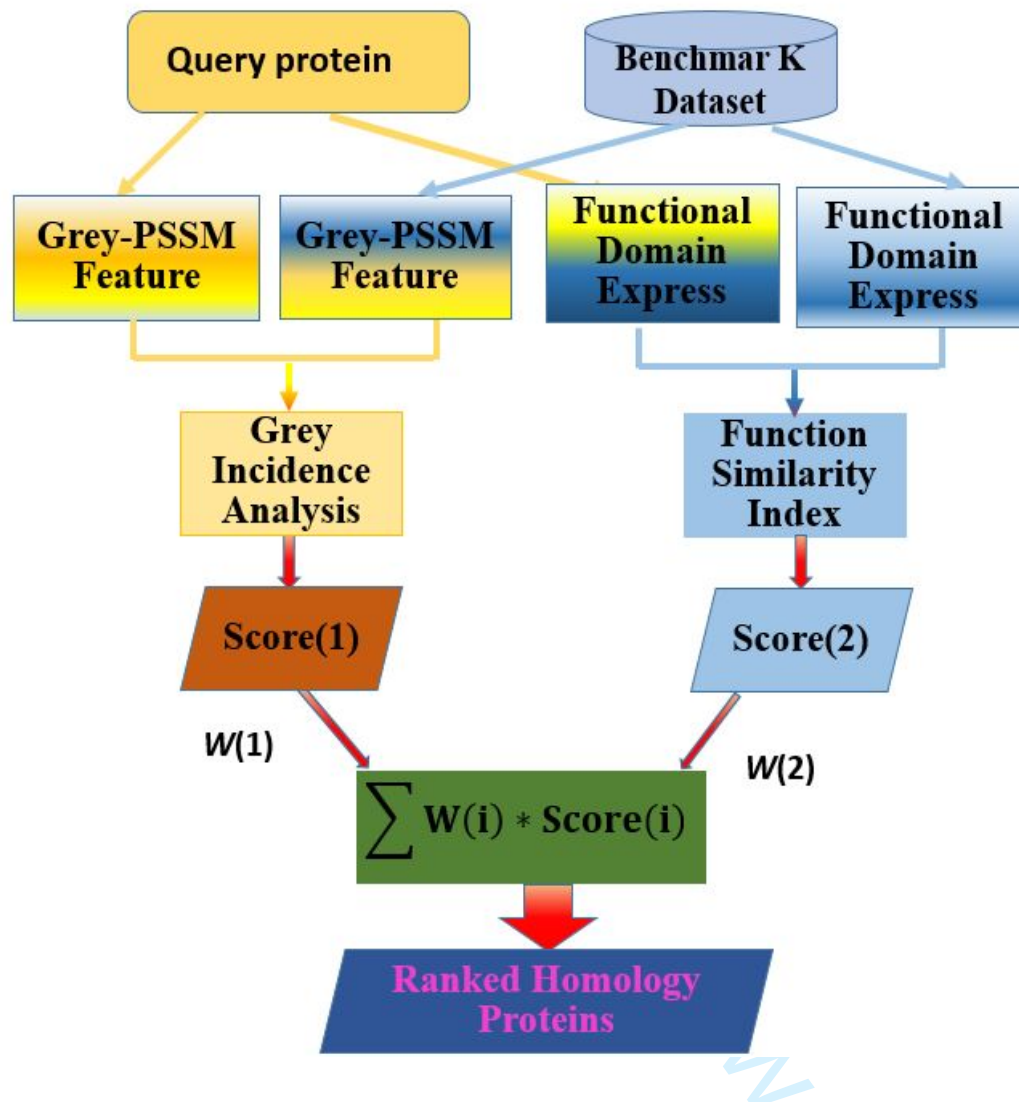


Table 1. A comparison of the jackknife test results for protein remote homology detection on the benchmark dataset

Methods	ROC1	ROC50
PSI-BLAST	0.7113	0.7647
GRA (Grey-PSSM)	0.8937	0.7149
Jaccard Index	0.8196	0.8070
Domain Similarity Index (DSI)	0.9053	0.8454
GRA and Jaccard Index	0.9301	0.8533
PHom-GRADSI	0.9620	0.8861

For Peer Review

REFERENCES

1. Chou, K.-C., K.D. Watenpaugh, and R.L. Heinrikson, *A Model of the Complex between Cyclin-Dependent Kinase 5 and the Activation Domain of Neuronal Cdk5 Activator*. Biochemical and Biophysical Research Communications, 1999. **259**(2): p. 420-428.
2. Zhou, G.P., R.B. Huang, and F.A. Troy, 2nd, *3D structural conformation and functional domains of polysialyltransferase ST8Sia IV required for polysialylation of neural cell adhesion molecules*. Protein Pept Lett, 2015. **22**(2): p. 137-48.
3. Schnell, J.R. and J.J. Chou, *Structure and mechanism of the M2 proton channel of influenza A virus*. Nature, 2008. **451**: p. 591.
4. Berardi, M.J., et al., *Mitochondrial uncoupling protein 2 structure determined by NMR molecular fragment searching*. Nature, 2011. **476**: p. 109.
5. OuYang, B., et al., *Unusual architecture of the p7 channel from hepatitis C virus*. Nature, 2013. **498**: p. 521.
6. Dev, J., et al., *Structural basis for membrane anchoring of HIV-1 envelope spike*. Science, 2016: p. aaf7066.
7. Oxenoid, K., et al., *Architecture of the mitochondrial calcium uniporter*. Nature, 2016. **533**: p. 269.
8. Chou, K.C., *Structural bioinformatics and its impact to biomedical science*. Curr Med Chem, 2004. **11**(16): p. 2105-34.
9. Chen, J., et al., *dRHP-PseRA: detecting remote homology proteins using profile-based pseudo protein sequence and rank aggregation*. Scientific Reports, 2016. **6**: p. 32333.
10. Chen, J., et al., *A comprehensive review and comparison of different computational methods for protein remote homology detection*. Brief Bioinform, 2016.
11. Liu, B., J. Chen, and S. Wang, *Protein Remote Homology Detection by Combining Pseudo Dimer Composition with an Ensemble Learning Method*. Current Proteomics, 2016. **13**(2): p. 86-91.
12. Liu, B., J. Chen, and X. Wang, *Application of learning to rank to protein remote homology detection*. Bioinformatics, 2015. **31**(21): p. 3492-3498.
13. Liu, B., J. Chen, and X. Wang, *Protein remote homology detection by combining Chou's distance-pair pseudo amino acid composition and principal component analysis*. Molecular Genetics and Genomics, 2015. **290**(5): p. 1919-1931.
14. Liu, B., et al., *Combining evolutionary information extracted from frequency profiles with sequence-based kernels for protein remote homology detection*. Bioinformatics, 2014. **30**(4): p. 472-9.
15. Chou, K.C., *Some remarks on protein attribute prediction and pseudo amino acid composition*. J Theor Biol, 2011. **273**(1): p. 236-47.
16. Chen, W., et al., *Using deformation energy to analyze nucleosome positioning in genomes*. Genomics, 2016. **107**(2): p. 69-75.

17. Cheng, X., X. Xiao, and K.-C. Chou, *pLoc-mGneg: Predict subcellular localization of Gram-negative bacterial proteins by deep gene ontology learning via general PseAAC*. Genomics, 2018. **110**(4): p. 231-239.
18. Cheng, X., X. Xiao, and K.C. Chou, *pLoc-mHum: predict subcellular localization of multi-location human proteins via general PseAAC to winnow out the crucial GO information*. Bioinformatics, 2018. **34**(9): p. 1448-1456.
19. Chen, W., et al., *iRNA-AI: identifying the adenosine to inosine editing sites in RNA sequences*. Oncotarget, 2016. **8**(3): p. 4208-4217.
20. Jia, J., et al., *Identification of protein-protein binding sites by incorporating the physicochemical properties and stationary wavelet transforms into pseudo amino acid composition*. Journal of Biomolecular Structure and Dynamics, 2016. **34**(9): p. 1946-1961.
21. Cheng, X., et al., *pLoc-mAnimal: predict subcellular localization of animal proteins with both single and multiple sites*. Bioinformatics, 2017. **33**(22): p. 3524-3531.
22. Jia, J., et al., *pSuc-Lys: Predict lysine succinylation sites in proteins with PseAAC and ensemble random forest approach*. J Theor Biol, 2016. **394**: p. 223-30.
23. Cheng, X., et al., *iATC-mISF: a multi-label classifier for predicting the classes of anatomical therapeutic chemicals*. Bioinformatics, 2017. **33**(3): p. 341-346.
24. Jia, J., et al., *pSumo-CD: predicting sumoylation sites in proteins with covariance discriminant algorithm by incorporating sequence-coupled effects into general PseAAC*. Bioinformatics, 2016. **32**(20): p. 3133-3141.
25. Feng, P., et al., *iRNA-PseColl: Identifying the Occurrence Sites of Different RNA Modifications by Incorporating Collective Effects of Nucleotides into PseKNC*. Molecular Therapy - Nucleic Acids, 2017. **7**: p. 155-163.
26. Liu, B., et al., *iRSpot-EL: identify recombination spots with an ensemble learning approach*. Bioinformatics, 2017. **33**(1): p. 35-41.
27. Liu, B., F. Yang, and K.-C. Chou, *2L-piRNA: A Two-Layer Ensemble Classifier for Identifying Piwi-Interacting RNAs and Their Function*. Molecular Therapy - Nucleic Acids, 2017. **7**: p. 267-277.
28. Cheng, X., X. Xiao, and K.C. Chou, *pLoc-mEuk: Predict subcellular localization of multi-label eukaryotic proteins by extracting the key GO information into general PseAAC*. Genomics, 2018. **110**(1): p. 50-58.
29. Feng, P., et al., *iDNA6mA-PseKNC: Identifying DNA N6-methyladenosine sites by incorporating nucleotide physicochemical properties into PseKNC*. Genomics, 2018.
30. Ehsan, A., et al., *A Novel Modeling in Mathematical Biology for Classification of Signal Peptides*. Scientific Reports, 2018. **8**(1): p. 1039.
31. Liu, B., et al., *iPromoter-2L: a two-layer predictor for identifying promoters and their types by multi-window-based PseKNC*. Bioinformatics, 2018. **34**(1): p. 33-40.
32. Song, J., et al., *PREvall, an integrative approach for inferring catalytic*

- residues using sequence, structural, and network features in a machine-learning framework. *Journal of Theoretical Biology*, 2018. **443**: p. 125-137.
33. Huang, Y., et al., *CD-HIT Suite: a web server for clustering and comparing biological sequences*. *Bioinformatics*, 2010. **26**(5): p. 680-2.
34. Chou, K.C., *Impacts of bioinformatics to medicinal chemistry*. *Med Chem*, 2015. **11**(3): p. 218-34.
35. Chou, K.C., *Prediction of protein cellular attributes using pseudo amino acid composition*. *Proteins: Structure, Function, and Bioinformatics*, 2001(43): p. 246-255.
36. Chou, K.-C., *Using amphiphilic pseudo amino acid composition to predict enzyme subfamily classes*. *Bioinformatics*, 2005. **21**(1): p. 10-19.
37. Dehzangi, A., et al., *Gram-positive and Gram-negative protein subcellular localization by incorporating evolutionary-based descriptors into Chou's general PseAAC*. *Journal of Theoretical Biology*, 2015. **364**: p. 284-294.
38. Behbahani, M., H. Mohabatkar, and M. Nosrati, *Analysis and comparison of lignin peroxidases between fungi and bacteria using three different modes of Chou's general pseudo amino acid composition*. *Journal of Theoretical Biology*, 2016. **411**: p. 1-5.
39. Yu, B., et al., *Accurate prediction of subcellular location of apoptosis proteins combining Chou's PseAAC and PsePSSM based on wavelet denoising*. *Oncotarget*, 2017. **8**(64): p. 107640-107665.
40. Meher, P.K., et al., *Predicting antimicrobial peptides with improved accuracy by incorporating the compositional, physico-chemical and structural features into Chou's general PseAAC*. *Scientific Reports*, 2017. **7**: p. 42362.
41. Rahimi, M., M.R. Bakhtiarizadeh, and A. Mohammadi-Sangcheshmeh, *OOgenesis_Pred: A sequence-based method for predicting oogenesis proteins by six different modes of Chou's pseudo amino acid composition*. *Journal of Theoretical Biology*, 2017. **414**: p. 128-136.
42. Tahir, M., M. Hayat, and M. Kabir, *Sequence based predictor for discrimination of enhancer and their types by applying general form of Chou's trinucleotide composition*. *Computer Methods and Programs in Biomedicine*, 2017. **146**: p. 69-75.
43. Tripathi, P. and P.N. Pandey, *A novel alignment-free method to classify protein folding types by combining spectral graph clustering with Chou's pseudo amino acid composition*. *Journal of Theoretical Biology*, 2017. **424**: p. 49-54.
44. Zhang, S. and X. Duan, *Prediction of protein subcellular localization with oversampling approach and Chou's general PseAAC*. *Journal of Theoretical Biology*, 2018. **437**: p. 239-250.
45. Chou, K.-C., *Pseudo Amino Acid Composition and its Applications in Bioinformatics, Proteomics and System Biology*. *Current Proteomics*, 2009(6): p. 262-274.
46. Chou, K.-C., *An Unprecedented Revolution in Medicinal Chemistry Driven by the Progress of Biological Science*. *Current Topics in Medicinal Chemistry*,

2017. **17**(21): p. 2337-2358.
47. Lin, W.Z., et al., *iLoc-Animal: a multi-label learning classifier for predicting subcellular localization of animal proteins*. Mol Biosyst, 2013. **9**(4): p. 634-644.
48. Lin, W.-Z., et al., *Predicting Secretory Proteins of Malaria Parasite by Incorporating Sequence Evolution Information into Pseudo Amino Acid Composition via Grey System Model*. PLOS ONE, 2012. **7**(11): p. e49040.
49. Lin, W.-Z., X. Xiao, and K.-C. Chou, *GPCR-GIA: a web-server for identifying G-protein coupled receptors and their families with grey incidence analysis*. Protein Engineering, Design and Selection, 2009. **22**(11): p. 699-705.
50. Potter, S.C., et al., *HMMER web server: 2018 update*. Nucleic Acids Res, 2018. **46**(W1): p. W200-W204.
51. Finn, R.D., et al., *HMMER web server: 2015 update*. Nucleic Acids Res, 2015. **43**(W1): p. W30-8.
52. Finn, R.D., J. Clements, and S.R. Eddy, *HMMER web server: interactive sequence similarity searching*. Nucleic Acids Res, 2011. **39**(Web Server issue): p. W29-37.
53. Deng, J.L., *Introduction to Grey System Theory*. The Journal of Grey System, 1989(1): p. 1-24.
54. Liu, S., Z. Fang, and Y. Lin, *A new definition for the degree of grey incidence*. Scientific Inquiry, 2006. **7**(2): p. 111-124.
55. Chou, K.C. and C.T. Zhang, *Prediction of protein structural classes*. Crit Rev Biochem Mol Biol, 1995. **30**(4): p. 275-349.

