

Accepted Manuscript

Prediction of cell-penetrating peptides with feature selection techniques

Hua Tang, Zhen-Dong Su, Huan-Huan Wei, Wei Chen, Hao Lin

PII: S0006-291X(16)30953-6

DOI: [10.1016/j.bbrc.2016.06.035](https://doi.org/10.1016/j.bbrc.2016.06.035)

Reference: YBBRC 35951

To appear in: *Biochemical and Biophysical Research Communications*

Received Date: 30 May 2016

Accepted Date: 8 June 2016

Please cite this article as: H. Tang, Z.-D. Su, H.-H. Wei, W. Chen, H. Lin, Prediction of cell-penetrating peptides with feature selection techniques, *Biochemical and Biophysical Research Communications* (2016), doi: 10.1016/j.bbrc.2016.06.035.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



Prediction of cell-penetrating peptides with feature selection techniques

Hua Tang^{1,*}, Zhen-Dong Su³, Huan-Huan Wei³, Wei Chen^{2,3}, Hao Lin^{3,*}

¹ Department of Pathophysiology, Southwest Medical University, Luzhou 646000, China;

² Department of Physics, School of Sciences, Center for Genomics and Computational Biology, North China University of Science and Technology, Tangshan 063009, China;

³ Key Laboratory for NeuroInformation of Ministry of Education, Center of Bioinformatics and Center for Information in Biomedicine, School of Life Science and Technology, University of Electronic Science and Technology of China, Chengdu 610054, China.

Email:

Hua Tang: tanghua771211@aliyun.com

Hao Lin: hlin@uestc.edu.cn

* Corresponding authors

Abstract: Cell-penetrating peptides are a group of peptides which can transport different types of cargo molecules such as drugs across plasma membrane and have been applied in the treatment of various diseases. Thus, the accurate prediction of cell-penetrating peptides with bioinformatics methods will accelerate the development of drug delivery systems. The study aims to develop a powerful model to accurately identify cell-penetrating peptides. At first, the peptides were translated into a set of vectors with the same dimension by using dipeptide compositions. Secondly, the Analysis of Variance-based technique was used to reduce the dimension of the vector and explore the optimized features. Finally, the support vector machine was utilized to discriminate cell-penetrating peptides from non-cell-penetrating peptides. The five-fold cross-validated results showed that our proposed method could achieve an overall prediction accuracy of 83.6%. The results indicated that our model could provide more precise predictions for new peptides. Based on the proposed model, we constructed a free webserver called C2Pred (<http://lin.uestc.edu.cn/server/C2Pred>).

Keywords: Cell-penetrating peptides; Support vector machine; *g*-gap dipeptide composition; analysis of variance

Introduction

Cell-penetrating peptides (CPPs), a group of short peptides, are able to mediate the intracellular delivery of a series of molecules. They were firstly found in the measurements of the activity of the tat protein from human immunodeficiency virus 1 (HIV-1) in 1988 [1]. Subsequent studies revealed that a positively charged peptide between the 47th and 57th amino acids was in charge of the translocation [2]. Since then, more and more CPPs had been experimentally identified [3,4]. Generally, the CPPs are typically hydrophobic linear arrangements of less than 30 residues and can transport both small and large molecules in vitro and in vivo. Thus, they have great potential in the biological research and medicine development. Recently, active CPPs (ACPPs) have been employed to target cancer cells over-expressing metalloproteinase-2 [5] and treat various inflammatory diseases [6]. Therefore, it is crucial to deeply understand the CPPs function so as to reveal the mechanism involved in membrane translocation. The correct identification of CPPs is the first step for understanding their translocation mechanisms and important for discovering more CPPs, which can be used as drug delivery agent.

Although the biochemical experimental approaches can provide the exact details for CPPs, the wet experimental technique is time-consuming and expensive. With a lot of biological data, it is highly desirable to develop computational methods to identify CPPs. The machine learning-based algorithms allow us to find out CPPs from huge peptide data. In fact, several methods had been proposed to computationally identify CPPs. Dobchev et al. [7] investigated 101 peptides with artificial neural networks (ANN) and principle component analysis (PCA) and obtained the prediction accuracy of 80%-100%. Sanders et al. [8] firstly established a more objective benchmark dataset including 111 experimentally confirmed CPPs and 34 known non-CPPs and then distinguished CPPs from non-CPPs with support vector machine (SVM) and 61 features. The sensitivity (S_n), specificity (S_p) and overall accuracy (Acc) were 75.9%, 23.2% and 75.9%, respectively. Subsequently, Gautam et al. [9] improved the prediction accuracy by considering different features including amino acid and dipeptide compositions as well as physicochemical properties. However, the prediction accuracy for an independent dataset was only 81.31%. Holton et al. [10] developed a web server called CPPpred based on N-to-1 neural network to predict CPPs and obtained the Acc of only 75.86% in 5-fold cross-validation for a non-redundant dataset. Recently, Chen et al. [11] used pseudo amino acid composition (PseAAC) to encode peptides and achieved the ACC of 83.5% in 10-fold cross-validation by using random forest (RF).

The aforementioned methods did yield quite encouraging results. However, the prediction accuracies are still far from satisfactory. Moreover, few online resources of CPPs prediction are available. Therefore, we developed a powerful tool to correctly identify CPPs. In the study, we enhanced the prediction power and quality in identifying CPPs.

According to previous studies [12,13,14], the machine learning method for predicting biological molecules included four steps: benchmark dataset construction, representation of peptide samples, machine learning method selection and webserver construction. In the following sections, we will describe the four steps in detail.

Material and methods

Benchmark datasets

A reliable benchmark dataset is the foundation of an accurate model. Thus, in this study, all the CPPs and non-CPPs were derived from CPPsite2.0 [15]. To avoid any similarity bias which would result in an overestimate of predicted results, we used the CD-HIT program to remove highly similar sequences. Finally, a total of 411 experimentally confirmed CPPs and 411 know non-CPPs were gained. All the data can be freely downloaded from the website (<http://lin.uestc.edu.cn/server/CPPIDen/data>).

Representation of peptide samples

Formulating a given peptide sample \mathbf{P} with a mathematics descriptor is the second step to develop a sequence-based predictor for CPPs prediction. The most straightforward method in the current benchmark dataset is to use its entire amino acid sequence to formulate a peptide \mathbf{P} with L residues as follows

$$\mathbf{P} = R_1 R_2 R_3 R_4 \dots R_L \quad (1)$$

where R_1 , R_2 and R_L respectively denote the 1st, 2nd, and L -th residues of the peptide \mathbf{P} . This sequential model can be statistically analyzed by using BLAST program. Unfortunately, the straightforward and intuitive approach fails when a query peptide sequence has no significant similarity to any known CPP sequence.

To solve this problem, a peptide sequence may be firstly translated into a vector with the same dimension and then machine learning method is used to perform prediction. The peptide samples in the vector format can be more easily handled than those in the sequence format with many existing operation engines. The most simple vector model for a peptide sequence is its amino acid composition (AAC), which has been widely applied in protein classification [16,17,18,19]. However, if AAC was used to represent a peptide sample, all of its sequence order information would be lost. Many studies demonstrated that the residue-order information was very important for peptide structure and function annotation [12,13,14,20]. Thus, in this work, CPPs and non-CPPs were described with dipeptide composition as follows:

$$\mathbf{p} = [f_1, f_2, \dots, f_u, \dots, f_{400}]^T \quad (2)$$

where the f_u is the frequency of the u -th ($u=1, 2, \dots, 400$) dipeptide defined as

$$f_u = \frac{x_u}{\sum_u x_u} \quad (3)$$

where x_u denotes the number of the u -th dipeptide in a peptide.

According to Eqs. (2-3), each sample can be transformed into a 400-dimension vector. However, it is well known that the length of CPPs usually ranges from 12 to 26 residues, suggesting that noise or redundant information is included in this vector model. Generally, garbage information will prevent the proposed model from correctly identifying CPPs. Moreover, the computational time will increase. Thus, it is necessary to pick out the useful features via a feature selection technique. Currently, the technique based on analysis of variance (ANOVA) has been proposed to rank the features and improve the predictive accuracies in protein classification field [13,19,21]. Thus, we also used the feature selection technique to optimize the feature set for improving the predictive performance.

On the basis of the ANOVA theory, the importance of each dipeptide for CPP prediction can be defined as:

$$F(u) = \frac{\sum_{i=1}^2 m_i \left(\frac{\sum_{j=1}^{m_i} f_u(i,j)}{m_i} - \frac{\sum_{i=1}^2 \sum_{j=1}^{m_i} f_u(i,j)}{\sum_{i=1}^2 m_i} \right)^2}{\sum_{i=1}^2 \sum_{j=1}^{m_i} \left(f_u(i,j) - \frac{\sum_{j=1}^{m_i} f_u(i,j)}{m_i} \right)^2} / (m_1 + m_2 - 2) \quad (4)$$

where $f_u(i, j)$ denotes the frequency of the u -th dipeptide of the j -th sample in the i -th group; m_i denotes the number of samples in the i -th group (here $m_1=411$, $m_2=411$). It is obviously that the larger $F(u)$ value means the better discriminative capability of the u -th feature. Thus, we may rank all features according to their F values. Then we investigated the predictive performance of the first feature subset including the feature with the largest F value on CPP prediction. Subsequently, we measured the predictive accuracy of a new feature subset produced by adding a new feature with the second highest F value into the first feature set. This process was repeated from the higher F value to the lower F value until all candidate features were added. The optimal feature subset including u_0 ranked dipeptides which could achieve the highest predictive accuracy was expressed as:

$$\mathbf{p}_{u_0} = [f_1, f_2, \dots, f_{u_0}]^T \quad (5)$$

Based on the feature selection, the high-dimensional data were projected into a low-dimensional space. The final model was built based on the optimal feature subset.

Machine learning method

After the representation of peptide samples, the third step in CPPs prediction is to perform the classification with a machine learning method. With the progress of in mathematical theory, several machine learning methods, such as, fisher discrimination (FD) [22], RF[23], neural network (NN) [7], and k-nearest neighbors (KNN) [24] have been developed and widely applied in bioinformatics. SVM is one of the most powerful and popular methods in protein classification [25,26,27]. Its basic idea is to transform the input vector into a high-dimension Hilbert space and seek a separating hyperplane in this space. Due to its excellent learning ability, especially for small sample size, we also used SVM to perform classification.

In this work, each sample in the benchmark dataset expressed as a vector has a corresponding label $y \in \{+1, -1\}$, where $+1$ and -1 indicate CPPs and non-CPPs, respectively. The SVM projects the input vectors into a high-dimensional feature space for constructing an optimal separating hyperplane with the largest distance between two classes, measured along a line perpendicular to this hyperplane. The decision function of SVM is expressed as:

$$f(\vec{p}) = \text{sgn}(\sum_{i=1}^N y_i \alpha_i \cdot K(\vec{p}_i, \vec{p})) + b \quad (6)$$

where N is the number of samples (here $N=145$); $K(\vec{p}_i, \vec{p}_j)$ is called kernel function, which is an inner product in a high-dimensional feature space. In this work, we used radial basis function (RBF) defined as $K(\vec{p}_i, \vec{p}_j) = \exp(-\gamma \|\vec{p}_i - \vec{p}_j\|^2)$. The coefficients α_i can be solved by the convex Quadratic Programming (QP) problem,

$$\text{Maximize } \{ \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \cdot y_i y_j \cdot K(\vec{p}_i, \vec{p}_j) \} \text{ subject to } 0 < \alpha_i < C \quad (7)$$

where the regularization parameter C can control the tradeoff between margin and misclassification error. These \vec{p}_i are called support vectors only if corresponding $\alpha_i > 0$.

For the convenience of study, a soft package called LibSVM (<https://www.csie.ntu.edu.tw/~cjlin/libsvm/>) was designed to implement SVM. The regularization parameter C and kernel parameter γ were optimized by using a grid search method with cross-validation test. The search spaces for C and γ are $[2^{-15}, 2^5]$ and $[2^{-5}, 2^{-15}]$ with the steps of 2^{-1} and 2, respectively.

Performance evaluation

Three widely used cross-validation methods are widely used in statistical prediction: independent dataset test, sub-sampling (2-, 5- or 10-fold cross-validation) test, and jackknife test [7,8,10,11,12,13,14,18,19,20,21,22,23,25,26,27]. Because jackknife test always yield unique results for a given benchmark dataset, it is often used to evaluate the performance of the proposed methods in practical application [12,13,14,21,22]. However, it is time-consuming to run the test. Thus, in previous studies on CPPs prediction, the n -fold cross-validation was adopted [8,11]. To provide an objective comparison, we also used n -fold cross-validation in this study.

A set of simple methods to measure the prediction quality are introduced as the follows:

$$Sn = \frac{TP}{TP+FN} \quad 0 \leq Sn \leq 1 \quad (8)$$

$$Sp = \frac{TN}{TN+FP} \quad 0 \leq Sp \leq 1 \quad (9)$$

$$Acc = \frac{TP+TN}{TP+TN+FP+FN} \quad 0 \leq Acc \leq 1 \quad (10)$$

where Sn , Sp and Acc are called sensitivity, specificity and overall accuracy, respectively. TP and TN denote the numbers of correctly recognized CPPs and non-CPPs, respectively. FP and FN are the number of the non-CPPs incorrectly predicted as the CPPs and the number of CPPs incorrectly predicted as the non-CPPs, respectively.

Results and Discussion

Feature selection for improving accuracy

Based on the dipeptide composition definition in Eqs. (2-3), each peptide in benchmark dataset may be represented by a 400-dimension vector. To improve the predictive performance, it is necessary to find out the best feature subset which can produce the highest Acc . It is obvious that the optimal feature set can be obtained by investigating the Acc of all combinations of features. However, it is impossible to examine the performance of all feature subsets due to the long computation time. For the amino acid composition including 20 features, there are over 10^6 possible combinations. If the dimension increases to 400, the number of all possible combinations will be greater than 10^{120} , which is beyond the computational capability for most computers.

Thus, in order to reduce computation time, we used the feature selection technique defined in Eq. 4 to optimize features. At first, we used $F(u)$ value to evaluate the importance of each dipeptide for CPP prediction. Secondly, 400 dipeptides were ranked according to their $F(u)$ values. Thirdly, we estimated the Acc of the first feature with the largest $F(u)$ by using SVM. Furthermore, a new feature subset was achieved when the feature with the second highest F value was added. Then the Acc of this feature subset was investigated. We repeated

the process from the large $F(u)$ value to the small $F(u)$ value until the $Accs$ of all candidate features were examined. All examinations were performed by using 5-fold cross-validation to avoid over-fitting.

A large feature set bears more information. However, it will also bring about information redundancy or noise, which will result in the low capability in the generalization of a predictor or reduce the cluster-tolerant capacity so as to lower the cross-validation accuracy. For example, by investigating the accuracy of 400 features for CPPs prediction, we found that 80.7% samples could be correctly predicted in 5-fold cross-validation. However, the low-dimension feature could improve the robustness of a predictor. However, if few features are available, the obtained features are still not the optimal features for prediction because they cannot afford enough information or reflect real characteristics of the CPPs, thus leading to the low predictive accuracy. For instance, 10 dipeptides can only produce the Acc of 75.3% in 5-fold cross-validation.

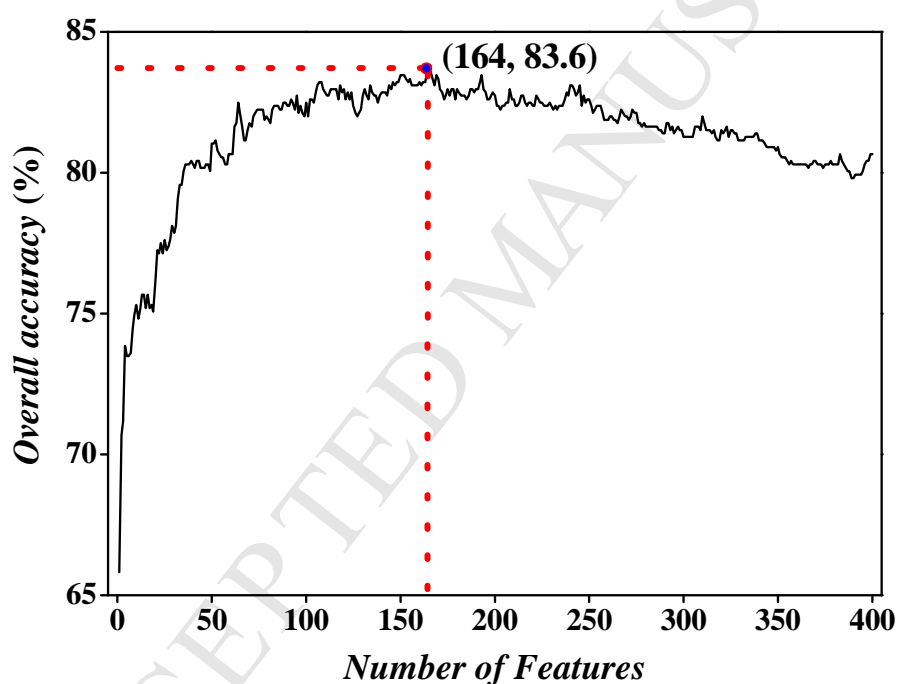


Figure 1. The feature selection results. When the top 164 dipeptides were used to perform prediction, the overall accuracy reached its peak of 83.6%.

To achieve the optimal feature subset and produce the highest Acc , a curve (**Figure 1**) was plotted in a 2D Cartesian coordinate system with the number of features as its abscissa and the Acc as its ordinate. The peak of the curve corresponds to the maximum Acc . As shown in **Fig. 1**, when 164 best dipeptides are used, the maximum Acc of 83.6% is obtained, and the Sn and Sp are 81.5% and 85.6%, respectively.

Comparison with other methods

To demonstrate the advantages of the proposed model, we made a comparison between the proposed method and other published methods. Sanders et al. constructed a

benchmark dataset including 111 experimentally confirmed CPPs and 34 known non-CPPs [8]. We examined our method on this dataset using 10-fold cross-validation and recorded results in **Table 1**. It shows that all index values of published methods are much lower than the corresponding ones achieved by our method. It should be noted that the features used in our model are much less than that in other published methods, suggesting that our model is more robust and ingenious.

Table 1. Comparison with published methods

Methods	<i>Sn</i>	<i>Sp</i>	<i>Acc</i>
Sander et al.'s method [8]	0.759	0.232	0.759
Chen et al.'s method [11]	0.955	0.441	0.835
Our method	0.973	0.765	0.924

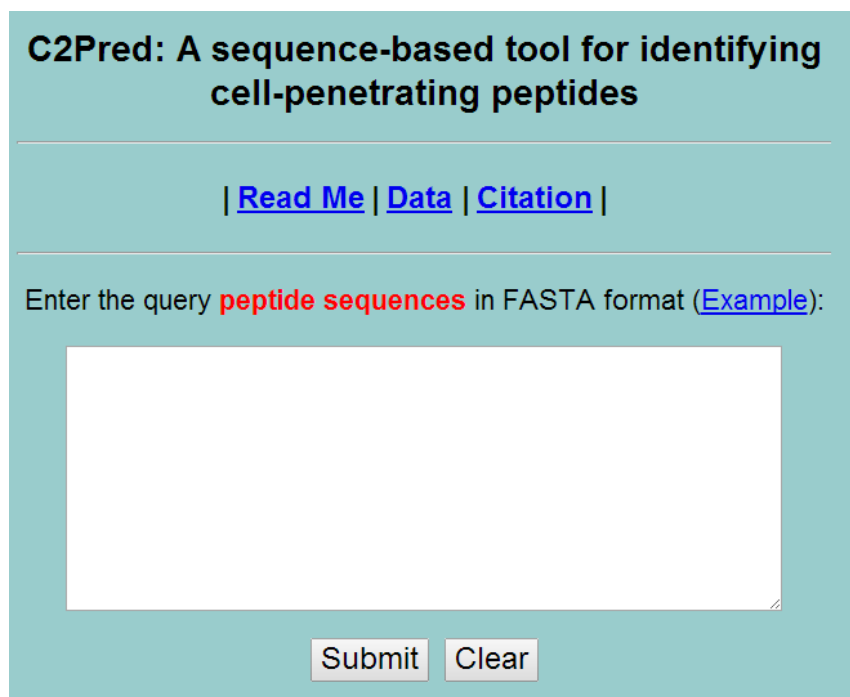
Sander et al. [8] designed a series of experiments to investigate the effect of unbalanced datasets on CPP prediction. Here, we compared the accuracies of our method with Sander's method through two strategies [8]. The first was to yield a set of 111 negative examples randomly repeated produced from the 34 known negatives. Combining with the 111 positive samples, we found that the overall accuracy is 89.92% in 10-fold cross-validation, which is higher than that of Sander's model (88.74%). The second strategy was to yield a set of 34 positive examples randomly selected from the 111 known positive examples. Combining with the 34 known negative examples, our model obtained the overall accuracy of 83.35% in 10-fold cross-validation, which is still higher than that of Sander's model (78.82%).

Gautam et al. [28] constructed a curated database called CPPsite including 843 cell penetrating peptides. To investigate the performance of the proposed method, we randomly selected 5 sets of 34 positive examples from CPPsite. Combining with the 34 known negative examples, we noticed that the averaged overall accuracy is 93.53% in 10-fold cross-validation, suggesting that our method could identify cell-penetrating peptides. Gautam et al. [28] also provided several search pages for users to run a search of a peptide against the CPPsite. However, these servers were based on homologous sequence in the searching dataset. Our model is more flexible because it is just dependent on the residue sequence.

Feature selection by using ANOVA can not only provide deeper insights into the intrinsic properties of peptide sequences, but also economize runtime and computational resource. Furthermore, it is robust to most violations of its assumptions and more intuitive for users to analyze the interaction of the two features. The ANOVA-based method did improve the cross-validated accuracies and robustness of the model. Therefore, our model has the better performance for CPP prediction.

Protein structure and function are correlated with the physiochemical properties of amino acids. Consequently, in the future, we will make our efforts to improve the accuracy by combining the optimal dipeptide composition with physiochemical properties of amino acids.

The positive and negative samples can be found from a number of different experimental techniques (different detection methodologies, and different cell types) [8]. Thus, in the future, we will collect more data from different cell types for the creation of CPPs dataset and the prediction of CPPs.



C2Pred: A sequence-based tool for identifying cell-penetrating peptides

| [Read Me](#) | [Data](#) | [Citation](#) |

Enter the query **peptide sequences** in FASTA format ([Example](#)):

Figure 2. A semi-screenshot to show the top page of the **C2Pred** webserver. Its website address is <http://lin.uestc.edu.cn/server/C2Pred>.

Web-Server Guide

We established a user-friendly web-server called **C2Pred** based on the proposed model to improve the efficiency and avoid a repeated a complicated mathematic process for identifying CPPs.

To facilitate the studies of other scholars, we provided the guideline as follows. One may browse the web server at <http://lin.uestc.edu.cn/server/C2Pred> (top page shown in **Fig.2**). The [Read Me](#) button provides a brief introduction about the predictor and the caveat for use. The [Data](#) button lists a link for downloading the benchmark datasets. The [Citation](#) button gives the relevant paper of **C2Pred**. The [Example](#) button provides example sequences in FASTA format. Users may type or copy/paste the query peptide sequences with FASTA format into the input box at the center of **Fig.2**. After submitting their peptide sequences, results will be shown in a new interface.

We will provide the maintenance of **C2Pred**. The server will be available to the research community for at least five years. With the collection of new CPPs and the development of new algorithms, we will upgrade the server with the higher accuracy and specificity. At this stage, **C2Pred** provides users with a highly practical tool and straightforward web interface for identifying CPPs. We also plan to develop more useful applications on CPP analysis.

Conclusions

In this paper, we developed a novel approach to discriminate CPPs from non-CPPs. In order to improve the prediction capability of the model, we designed a feature selection technique based on the ANOVA. An overall accuracy of 83.6% was achieved in 5-fold cross-validation. By comparing the proposed method with the other existing

methods, we demonstrated that our proposed method was superior to other methods, suggesting that **C2Pred** was a powerful tool for CPPs prediction.

Acknowledgements

This work was supported by the Applied Basic Research Program of Sichuan Province (No. 2015JY0100 and LZ-LY-45), the Scientific Research Foundation of the Education Department of Sichuan Province (11ZB122), the Nature Scientific Foundation of Hebei Province (No. C2013209105), the Fundamental Research Funds for the Central Universities of China (No. ZYGX2015J144) and Program for the Top Young Innovative Talents of Higher Learning Institutions of Hebei Province (No. BJ2014028).

Competing financial interests:

The authors declare that there is no conflict of interests.

References

- [1] Q. Zou, X.B. Li, Y. Jiang, Y.M. Zhao, G.H. Wang, BinMemPredict: a Web Server and Software for Predicting Membrane Protein Types, *Current Proteomics* 10 (2013) 2-9.
- [2] B. Liu, X.L. Wang, Q. Zou, Q.W. Dong, Q.C. Chen, Protein Remote Homology Detection by Combining Chou's Pseudo Amino Acid Composition and Profile-Based Protein Representation, *Molecular Informatics* 32 (2013) 775-782.
- [3] B. Liu, S.Y. Wang, X.L. Wang, DNA binding protein identification by combining pseudo amino acid composition and profile-based protein representation, *Scientific Reports* 5 (2015).
- [4] B. Liu, F.L. Liu, X.L. Wang, J.J. Chen, L.Y. Fang, K.C. Chou, Pse-in-One: a web server for generating various modes of pseudo components of DNA, RNA, and protein sequences, *Nucleic Acids Research* 43 (2015) W65-W71.
- [5] B. Liu, J. Chen, X. Wang, Application of learning to rank to protein remote homology detection, *Bioinformatics* 31 (2015) 3492-3498.
- [6] B. Liu, L. Fang, R. Long, X. Lan, K.C. Chou, iEnhancer-2L: a two-layer predictor for identifying enhancers and their strength by pseudo k-tuple nucleotide composition, *Bioinformatics* (2015).
- [7] B. Liu, L.Y. Fang, F.L. Liu, X.L. Wang, J.J. Chen, K.C. Chou, Identification of Real MicroRNA Precursors with a Pseudo Structure Status Composition Approach, *Plos One* 10 (2015).
- [8] B. Liu, L.Y. Fang, J.J. Chen, F.L. Liu, X.L. Wang, miRNA-dis: microRNA precursor identification based on distance structure status pairs, *Molecular Biosystems* 11 (2015) 1194-1204.
- [9] B. Liu, J.J. Chen, X.L. Wang, Protein remote homology detection by combining Chou's distance-pair pseudo amino acid composition and principal component analysis, *Molecular Genetics And Genomics* 290 (2015) 1919-1931.
- [10] B. Liu, L.Y. Fang, S.Y. Wang, X.L. Wang, H.T. Li, K.C. Chou, Identification of microRNA precursor with the degenerate K-tuple or Kmer strategy, *Journal Of Theoretical Biology* 385 (2015) 153-159.
- [11] B. Liu, F.L. Liu, L.Y. Fang, X.L. Wang, K.C. Chou, repDNA: a Python package to generate various modes of feature vectors for DNA sequences by incorporating user-defined physicochemical properties and sequence-order effects, *Bioinformatics* 31 (2015) 1307-1309.
- [12] H. Lin, W.X. Liu, J. He, X.H. Liu, H. Ding, W. Chen, Predicting cancerlectins by the optimal g-gap dipeptides, *Scientific Reports* 5 (2015).
- [13] H. Lin, E.Z. Deng, H. Ding, W. Chen, K.C. Chou, iPro54-PseKNC: a sequence-based predictor for identifying sigma-54 promoters in prokaryote with pseudo k-tuple nucleotide composition,

- Nucleic Acids Research 42 (2014) 12961-12972.
- [14] P.P. Zhu, W.C. Li, Z.J. Zhong, E.Z. Deng, H. Ding, W. Chen, H. Lin, Predicting the subcellular localization of mycobacterial proteins by incorporating the optimal tripeptides into the general form of pseudo amino acid composition, *Mol Biosyst* 11 (2015) 558-563.
- [15] P. Agrawal, S. Bhalla, S.S. Usmani, S. Singh, K. Chaudhary, G.P. Raghava, A. Gautam, CPPsite 2.0: a repository of experimentally validated cell-penetrating peptides, *Nucleic Acids Res* 44 (2016) D1098-1103.
- [16] T. Parfitt, Georgia: an unlikely stronghold for bacteriophage therapy, *Lancet* 365 (2005) 2166-2167.
- [17] S.H. Guo, E.Z. Deng, L.Q. Xu, H. Ding, H. Lin, W. Chen, K.C. Chou, iNuc-PseKNC: a sequence-based predictor for predicting nucleosome positioning in genomes with pseudo k-tuple nucleotide composition, *Bioinformatics* 30 (2014) 1522-1529.
- [18] H. Lin, H. Ding, F.B. Guo, A.Y. Zhang, J. Huang, Predicting subcellular localization of mycobacterial proteins by using Chou's pseudo amino acid composition, *Protein Pept Lett* 15 (2008) 739-744.
- [19] H. Lin, W. Chen, Prediction of thermophilic proteins using feature selection technique, *J Microbiol Methods* 84 (2011) 67-70.
- [20] E.C. Keen, Phage therapy: concept to cure, *Front Microbiol* 3 (2012) 238.
- [21] A. Pirisi, Phage therapy--advantages over antibiotics?, *Lancet* 356 (2000) 1418.
- [22] K. Kimura, Y. Itoh, Characterization of poly-gamma-glutamate hydrolase encoded by a bacteriophage genome: possible role in phage infection of *Bacillus subtilis* encapsulated with poly-gamma-glutamate, *Appl Environ Microbiol* 69 (2003) 2491-2497.
- [23] L. Rodriguez-Rubio, N. Quiles-Puchalt, B. Martinez, A. Rodriguez, J.R. Penades, P. Garcia, The peptidoglycan hydrolase of *Staphylococcus aureus* bacteriophage 11 plays a structural role in the viral particle, *Appl Environ Microbiol* 79 (2013) 6187-6190.
- [24] C. Verheust, N. Fornelos, J. Mahillon, The *Bacillus thuringiensis* phage GIL01 encodes two enzymes with peptidoglycan hydrolase activity, *FEMS Microbiol Lett* 237 (2004) 289-295.
- [25] W.W. Navarre, H. Ton-That, K.F. Faull, O. Schneewind, Multiple enzymatic activities of the murein hydrolase from staphylococcal phage phi11. Identification of a D-alanyl-glycine endopeptidase activity, *J Biol Chem* 274 (1999) 15847-15856.
- [26] H. Lin, Q.Z. Li, Eukaryotic and prokaryotic promoter prediction using hybrid approach, *Theory Biosci* 130 (2011) 91-100.
- [27] D. Nelson, R. Schuch, P. Chahales, S. Zhu, V.A. Fischetti, PlyC: a multimeric bacteriophage lysin, *Proc Natl Acad Sci U S A* 103 (2006) 10765-10770.
- [28] A. Gautam, H. Singh, A. Tyagi, K. Chaudhary, R. Kumar, P. Kapoor, G.P. Raghava. (2012) CPPsite: a curated database of cell penetrating peptides. Database (Oxford). 2012:bas015.

Table**Table 1. Comparison with published methods**

Methods	<i>Sn</i>	<i>Sp</i>	<i>Acc</i>
Sander et al.'s method [8]	0.759	0.232	0.759
Chen et al.'s method [11]	0.955	0.441	0.835
Our method	0.973	0.765	0.924

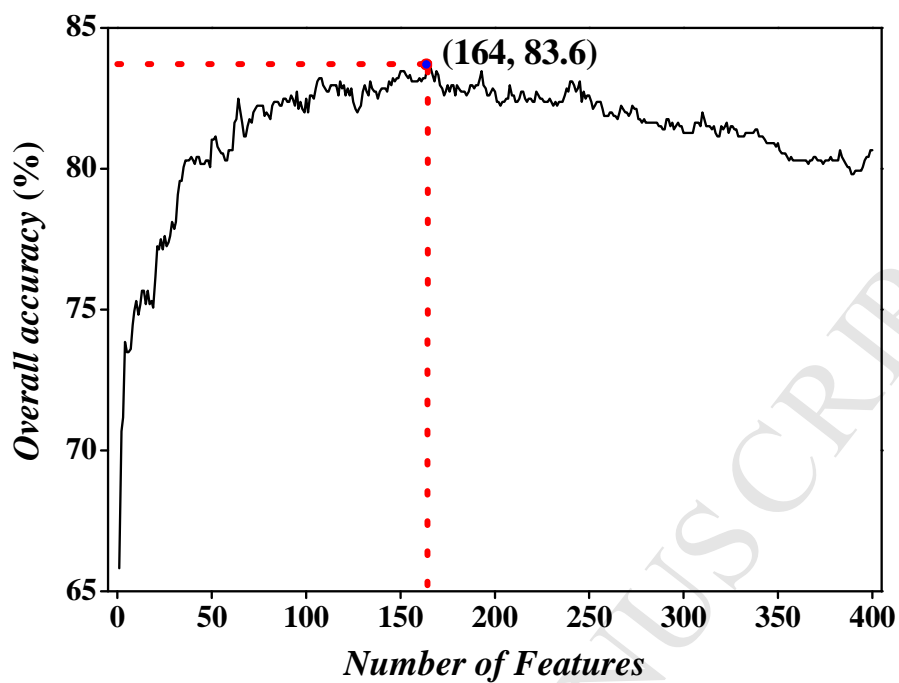
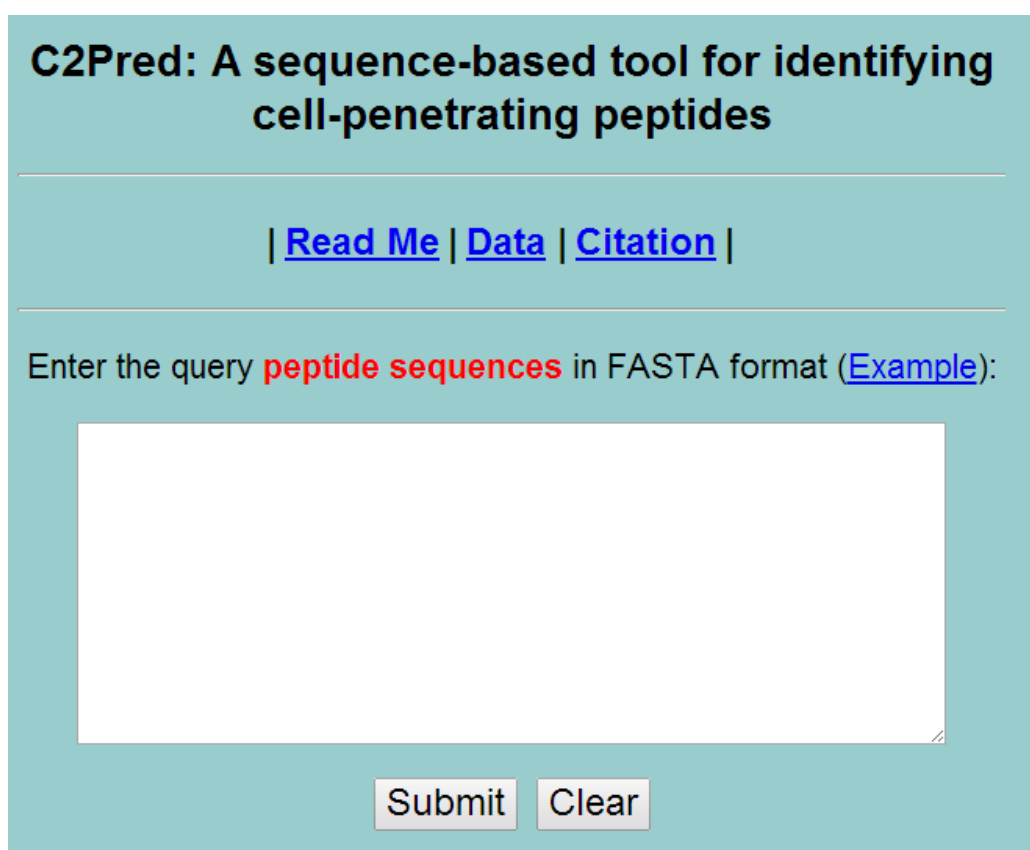


Figure 1. The feature selection results. When the top 164 dipeptides were used to perform prediction, the overall accuracy reached its peak of 83.6%.



C2Pred: A sequence-based tool for identifying cell-penetrating peptides

| [Read Me](#) | [Data](#) | [Citation](#) |

Enter the query **peptide sequences** in FASTA format ([Example](#)):

Figure 2. A semi-screenshot to show the top page of the **C2Pred** webserver. Its website address is <http://lin.uestc.edu.cn/server/C2Pred>.

Highlights

Novel analytical method is developed to predict the cell-penetrating peptides.

A significant feature selection technique is proposed and used to optimize features of proteins

A powerful web server is constructed to identify cell-penetrating peptides.