

# Recent Progresses for Computationally Identifying N<sup>6</sup>-methyladenosine Sites in *Saccharomyces cerevisiae*

Kuo-Chen Chou

Gordon Life Science Institute, Boston, MA 02478, USA.

**How to cite this paper:** Kuo-Chen Chou. (2020) Recent Progresses for Computationally Identifying N<sup>6</sup>-methyladenosine Sites in *Saccharomyces cerevisiae*. *Journal of Applied Mathematics and Computation*, 4(4), 153-173.  
DOI: 10.26855/jamc.2020.12.007

**Received:** October 2, 2020

**Accepted:** October 26, 2020

**Published:** November 10, 2020

\***Corresponding author:** Kuo-Chen Chou, Gordon Life Science Institute, Boston, MA 02478, USA.  
**Email:** kcchou@gordonlifescience.org

## Abstract

N<sup>6</sup>-methyladenosine (m<sup>6</sup>A) plays critical roles in a broad set of biological processes. Knowledge about the precise location of m<sup>6</sup>A site in the transcriptome is vital for deciphering its biological functions. Although experimental techniques have made substantial contributions to identify m<sup>6</sup>A methylations, they are still labor intensive, costly and time consuming. As good complements to experimental methods, in the past few years, a series of computational approaches have been proposed to identify m<sup>6</sup>A sites in *Saccharomyces cerevisiae*. In order to facilitate researchers to select appropriate methods for identifying m<sup>6</sup>A sites, it is necessary to give a comprehensive review and comparison on existing computational methods. In this review, we summarized the current progresses in computational prediction of m<sup>6</sup>A sites and also assessed the performance of computational methods for identifying m<sup>6</sup>A sites on an independent dataset. Finally, challenges and future directions of computationally identifying m<sup>6</sup>A sites were presented as well. Taken together, we anticipate that this review will provide an important guide for future computational analysis of m<sup>6</sup>A and other RNA modifications.

## Keywords

post transcription modification, N<sup>6</sup>-methyladenosine, epitranscriptome, machine learning method, 5-step rules

## 1. Introduction

Among the ~150 kinds of known RNA modifications, the N<sup>6</sup>-methyladenosine (m<sup>6</sup>A) is the most prevalent internal mRNA/lncRNA modifications, which occurs on the sixth nitrogen atom of adenine. As a reversible and dynamic post-transcriptional modification [1], the formation of m<sup>6</sup>A is installed by a multicomponent methyltransferase complex including METTL3, METTL14 and WTAP, while its demethylation is regulated by demethylases FTO and ALKBH5. The biological functions of dynamic m<sup>6</sup>A modification is regulated by m<sup>6</sup>A readers, such as heterogeneous nuclear ribonucleoprotein C (HNRNPC) [2-4], YTH Domain Family proteins 1, 2 and 3 (YTHDF1, YTHDF2, YTHDF3) and YTHDC, etc.

Since discovered in 1970s, m<sup>6</sup>A has been observed in all three kingdoms of life. With the intensive researches on m<sup>6</sup>A methylation in recent years, its functions have been uncovered gradually. It has been found that m<sup>6</sup>A is associated with a broad set of fundamental cellular processes, such as RNA localization and degradation, RNA splicing, circadian rhythm, cell differentiation and reprogramming, immune tolerance and even the occurrence of diseases. However, few of them are currently understood in mechanistic detail. Identifying the precise location of m<sup>6</sup>A site in transcriptomes will be of a great help to investigate its biological mechanisms and functions.

With the development of next-generation sequencing technology, the MeRIP-Seq and m<sup>6</sup>A-seq high-throughput methods have been developed to identify m<sup>6</sup>A sites in *Saccharomyces cerevisiae*, *Homo sapiens*, and *Mus musculus*. However, the resolution of these techniques is low and couldn't identify the exact methylated adenosines. Recently, the miCLIP technique was proposed, which provided the single-nucleotide resolution m<sup>6</sup>A profile of the human transcriptome. Based on these experimental data, several informative databases related with m<sup>6</sup>A modifications have been built. Taken together, these experiments promote the progress of researches on m<sup>6</sup>A modifications. However, experimental techniques are still labor-intensive and expensive for transcriptome-wide detection of m<sup>6</sup>A. Therefore, it is an urgent task to develop effective and low-cost approaches to automatically identify m<sup>6</sup>A sites. As excellent complements to experimental techniques, computational methods are in high demand to accurately detect m<sup>6</sup>A sites.

In 2013, Schwartz *et al.* proposed the first computational model to predict the m<sup>6</sup>A site in the *S. cerevisiae* transcriptome, whose features include relative position in gene, nucleotide composition and predicted secondary structures. Although no public web server or software package was provided for this method, Schwartz *et al.*'s pioneer work provides a new strategy for identifying m<sup>6</sup>A site. Since then, the scientific community witnessed an unprecedented amount of studies considering the application of machine learning method to identify m<sup>6</sup>A sites. For example, a series of machine learning based methods, such as m6Apred, iRNA-Methyl, SRAMP, pRNAm-PC, RAM-ESVM, RAM-NPPS and RNA-MethylPred have been proposed. All these prediction methods were developed in principle by following the guidelines of the Chou's 5-step rule [5] as done in a series of powerful predictors (see, *e.g.*, developed recently for genome or proteome analyses).

Accordingly, these methods also share the advantages: (1) clearer in logic development, (2) more transparent in operation, and (3) more useful in practical application.

To provide the readership with a clear landscape about the recent developments in this important area, in this comprehensive review we are to elaborate their details in observing the Chou's 5-step rule [5]: (1) how to construct or select a valid benchmark dataset to train and test the predictor; (2) how to formulate the biological sequence samples with an effective mathematical expression that can truly reflect their intrinsic correlation with the target to be predicted; (3) how to introduce or develop a powerful algorithm (or engine) to operate the prediction; (4) how to properly perform cross-validation tests to objectively evaluate the anticipated accuracy of the predictor; (5) how to establish a user-friendly web-server for the predictor that is accessible to the public. Moreover, to facilitate users to select appropriate method according to their need, a comparison of existing methods in identifying m<sup>6</sup>A sites is to be performed based on an independent dataset. Finally, the challenges and future perspectives for identifying m<sup>6</sup>A sites are to be discussed.

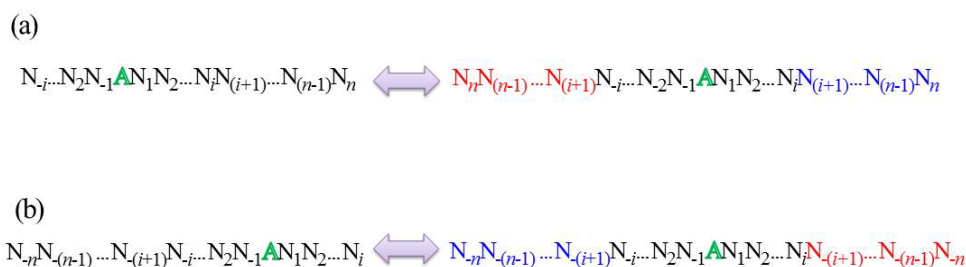
## 2. Benchmark Dataset for Predicting m<sup>6</sup>A Sites

Constructing a valid and reliable benchmark dataset is the critical step to train a computational model with high effectiveness [6, 7]. In literature, the benchmark dataset usually consists of a training dataset and a testing dataset: the former is constructed for the purpose of training a proposed model, while the latter for the purpose of testing it. As pointed out by a comprehensive review [8], however, there is no need to separate a benchmark dataset into a training dataset and a testing dataset for validating a prediction method if it is tested by the jackknife or subsampling (K-fold) cross-validation because the outcome thus obtained is actually from a combination of many different independent dataset tests.

For investigating the m<sup>6</sup>A sites, the benchmark datasets were constructed as follows. The sequence segment surrounding the m<sup>6</sup>A site contains the underlying discrimination information, whose size can be determined with the aid of sliding window scheme. The window length is usually set to  $2n+1$ , whose central element is the experimentally confirmed m<sup>6</sup>A site, with  $n$  flanking nucleotides on both sides of the methylated adenosine. However, there is no uniform standard to set the window size. The determination of  $n$  is always associated with features extraction, prediction method and cross-validation performance.

In 2015, the first publicly available benchmark dataset (called  $\mathcal{S}_1$  here) for the prediction of m<sup>6</sup>A sites was built by Chen *et al.* The positive samples in dataset  $\mathcal{S}_1$  are 832 m<sup>6</sup>A sites with distances to the detected m<sup>6</sup>A-seq peaks less than 10 bp, which were extracted from the 1,307 experimentally confirmed m<sup>6</sup>A sites. The negative samples in dataset  $\mathcal{S}_1$  are the 832 non-m<sup>6</sup>A sites randomly selected from the 33,280 non-methylated adenines. Each sample in the dataset  $\mathcal{S}_1$  is 21-nt long with the m<sup>6</sup>A sites or non-m<sup>6</sup>A site in the center.

In some cases, the m<sup>6</sup>A site locates at the beginning or end of the sequence, which results in that the extracted sequence fragments size is shorter than the given window size. Two strategies are often used to generate fixed window length. The first one is to fill the blank by using the dummy 'X' nucleotide that don't represent any real nucleotide. The second one is to fill the blank using the mirror image method. If the missing nucleotides locate at the beginning (*i.e.* upstream of the m<sup>6</sup>A site), they will be filled by using their mirror images locate at downstream of the m<sup>6</sup>A site, and vice versa (Figure 1). The second approach has been used to construct the benchmark dataset for the prediction of m<sup>6</sup>A sites in *S.cerevisiae*.



**Figure 1. A schematic illustration showing the mirror image for (a) upstream (b) downstream missing nucleotides, respectively. The real RNA segment is colored in blue and its mirror image is colored in red. The methylated A is highlighted in green.**

In 2016, Chen *et al.* built another benchmark dataset (called  $S_2$  here) using the mirror image method, which includes 1,307 m<sup>6</sup>A site containing sequences (positive samples) and the equal number of non-m<sup>6</sup>A containing sequences (negative samples). In dataset  $S_2$ , all the experimentally confirmed m<sup>6</sup>A sites in *S. cerevisiae* were included. The sequences in this dataset are 51-nt long with the sequence similarity less than 85%, with the m<sup>6</sup>A site or non-m<sup>6</sup>A site in the center. Since it has been built, nearly all the computational models for identifying m<sup>6</sup>A site in *S. cerevisiae* were trained and tested on the dataset  $S_2$ .

### 3. Formulation of RNA Samples

The 2<sup>nd</sup> step of the 5-step rules [5] is about the formulation of biological samples. With the explosive growth of biological sequences in the post-genomic era, one of the most important but also most difficult problems in computational biology is how to express a biological sequence with a discrete model or a vector, yet still keep considerable sequence-order information or key pattern characteristic. This is because all the existing machine-learning algorithms (such as “Covariance Discriminant” or “CD” algorithm [9, 10], “Nearest Neighbor” or “NN” algorithm [11, 12], “Support Vector Machine” or “SVM” algorithm [13, 14], and “Random Forest” or “RF” algorithm [15, 16]) can only handle vectors as elaborated in a comprehensive review [17]. However, a vector defined in a discrete model may completely lose all the sequence-pattern information. To avoid completely losing the sequence-pattern information for proteins, the pseudo amino acid composition [18] or PseAAC [19] 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.*, [20-45] as well as a long list of references cited in [46]).

Because it has been widely and increasingly used, four powerful open access soft-wares, called “PseAAC” [47], “PseAAC-Builder” [48], “propy” [49], and “PseAAC-General” [50], were established: the former three are for generating various modes of Chou’s special PseAAC [51]; while the 4th one for those of Chou’s general PseAAC [5], including not only all the special modes of feature vectors for proteins but also the higher level feature vectors such as “Functional Domain” mode (see Eqs.9-10 of [5]), “Gene Ontology” mode (see Eqs.11-12 of [5]), and “Sequential Evolution” or “PSSM” mode (see Eqs.13-14 of [5]). Encouraged by the successes of using PseAAC to deal with protein/peptide sequences, the idea of PseAAC was extended to PseKNC (Pseudo K-tuple Nucleotide Composition) to generate various feature vectors for DNA/RNA sequences [52] that have proved very successful as well [25, 41, 53-60]). Particularly, recently a very powerful web-server called “Pse-in-One” [61] and its updated version “Pse-in-One2.0” [62] have been established that can be used to generate any desired feature vectors for protein/peptide and DNA/RNA sequences according to the users’ need or their own definition.” According to the concept of pseudo components, any RNA sequence sample can be formulated as [53]

$$\mathbf{R} = [\phi_1\phi_2 \cdots \phi_u \cdots \phi_z]^T \tag{1}$$

where

$$\phi_u = \begin{cases} \frac{f_u}{\sum_{i=1}^{4^k} f_{i+w} \sum_{j=1}^{\lambda} \theta_j} & (1 \leq u \leq 4^k) \\ \frac{w\theta_{u-4^k}}{\sum_{i=1}^{4^k} f_{i+w} \sum_{j=1}^{\lambda} \theta_j} & (4^k < u \leq 4^k + \lambda) \end{cases} \tag{2}$$

In Eq.2,  $f_u (u = 1, 2, \dots, 4^k)$  is the normalized occurrence frequency of the  $u$ -th non-overlapping  $k$ -tuple nucleotide in the RNA sequence.  $\lambda$  is the number of the total pseudo components used to reflect the long-range or global sequence effect, and  $w$  is the weight factor.  $\theta_j$  is the  $j$ -th tier correlation factor that reflects the sequence order correlation between all the  $j$ -th most contiguous  $k$ -tuple nucleotide along a  $L$ -nt long RNA sequence as formulated by

$$\theta_j = \frac{1}{L-k-j+1} \sum_{i=1}^{L-k-j+1} C_{i,i+j} \quad (j = 1, 2, \dots, \lambda; \lambda < L - k) \tag{3}$$

where  $C_{i,i+j}$  is the correlation function and is defined by

$$C_{i,i+j} = \frac{1}{\mu} \sum_{g=1}^{\mu} [P_g(R_i R_{i+1} \dots R_{i+k-1}) - P_g(R_{i+j} R_{i+j+1} \dots R_{i+j+k-1})]^2 \tag{4}$$

where  $\mu$  is the number of RNA physicochemical properties considered,  $P_g(R_i R_{i+1} \dots R_{i+k-1})$  is the numerical value of the  $g$ -th ( $g=1, 2, 3, \dots, u$ ) RNA local structural property for the  $k$ -tuple nucleotide  $R_i R_{i+1} \dots R_{i+k-1}$  at position  $i$  and  $P_g(R_{i+j} R_{i+j+1} \dots R_{i+j+k-1})$  the corresponding value for the dinucleotide  $R_i R_{i+j+1} \dots R_{i+j+k-1}$  at position  $i+j$ . The details about PseKNC can be found in our recent review article [53].

### 4. Algorithm or Operation Engine

The 3<sup>rd</sup> step of the 5-step rules [5] is about the operation engine. The two commonly used machine learning algorithms for identifying m6A sites are support vector machine (SVM) and random forest (RF), which were briefly introduced as following.

#### 4.1. Support vector machine (SVM)

SVM is a powerful and popular method for pattern recognition, which has been widely used in the realm of bioinformatics especially very effectively in a series of recent genome analyses (see, e.g., [63-65]). Its basic idea is to transform the input data into a high dimensional feature space and then determine the optimal separating hyperplane. Owing to its effectiveness and speed in training process, the radial basis kernel function (RBF) of SVM was often used to obtain the classification hyperplane. The regularization parameter  $C$  and kernel parameter  $\gamma$  of the SVM operation engine can be optimized in the following ranges  $[2^{-5}, 2^{15}]$  and  $[2^{-15}, 2^{-5}]$  with the steps of 2 and  $2^{-1}$ , respectively. For a brief formulation of SVM and how it works, see the papers [66, 67]. For more details about SVM, see a monograph [68].

#### 4.2. Random forest (RF)

RF is an ensemble of a large number of decision trees. Each tree in the ensemble is trained on a subset of training instances and gives a classification result. The three parameters of RF, namely the number of trees, the number of features randomly selected, and the minimum number of samples required to split an internal node (nsplit) can be determined by using the grid search scheme. The predictive results of RF are based on the ensemble of those decision trees. Since proposed by Breiman in 2001 [69], owing to its advantages in dealing with high-dimensional data, RF has been used in many areas of bioinformatics (see, e.g., [15, 16, 56, 70-83]).

### 5. Performance Evaluation

The 4<sup>th</sup> guideline of the 5-step rules [5] is about how to validate the proposed model. To address this, two issues are needed to be considered. One is what kind of metrics should be used to measure the scores, and the other is what test methods should be adopted to count the scores.

#### 5.1. A set of intuitive metrics

The performance of the computational methods is usually evaluated using the following four metrics [84]: (1) overall accuracy or Acc, (2) Mathew’s correlation coefficient or MCC, (3) sensitivity or Sn, and (4) specificity or Sp. However, the conventional metrics copied from math books are hard to be understood by most experimental scientists due to lacking intuitiveness; especially for the MCC, which is very important to indicate the stability of a predictor. Fortunately, using the symbols introduced by Chou [85] in studying signal peptide cleavage sites, a set of four intuitive metrics were derived [14, 86], as given below

$$\left\{ \begin{array}{l} \text{Sn} = 1 - \frac{N_{+}^{-}}{N_{+}} \qquad 0 \leq \text{Sn} \leq 1 \\ \text{Sp} = 1 - \frac{N_{+}^{-}}{N_{-}} \qquad 0 \leq \text{Sp} \leq 1 \\ \text{Acc} = \Lambda = 1 - \frac{N_{+}^{-} + N_{-}^{+}}{N_{+} + N_{-}} \qquad 0 \leq \text{Acc} \leq 1 \\ \text{MCC} = \frac{1 - \left(\frac{N_{+}^{-}}{N_{+}} + \frac{N_{-}^{+}}{N_{-}}\right)}{\sqrt{\left(1 + \frac{N_{-}^{-} - N_{+}^{-}}{N_{+}}\right)\left(1 + \frac{N_{+}^{-} - N_{-}^{+}}{N_{-}}\right)}} \qquad -1 \leq \text{MCC} \leq 1 \end{array} \right. \tag{5}$$

where  $N^+$  represents the total number of positive samples investigated, while  $N_+^+$  is the number of positive samples incorrectly predicted to be of negative one;  $N^-$  the total number of negative samples investigated, while  $N_+^-$  the number of the negative samples incorrectly predicted to be of positive one. The set of intuitive metrics has been concurred and applauded by a series of recent publications (see, *e.g.*, [14, 16, 57-59, 74, 82, 87-100] [83, 101-114]). It is instructive to point out, however, either the conventional metrics [84] taken from math books or the intuitive metrics of Eq.5 are valid only for single label systems (where each of the constituent samples belong to one, and only one, attribute or class); for the multi-label systems (where a sample may simultaneously belong to several different attributes or classes) whose existence has become more frequent in system biology [6, 7, 29, 115-134], system medicine [135, 136] and biomedicine [78, 137], a completely different set of metrics as defined in [138] is absolutely needed.

## 5.2. Jackknife test

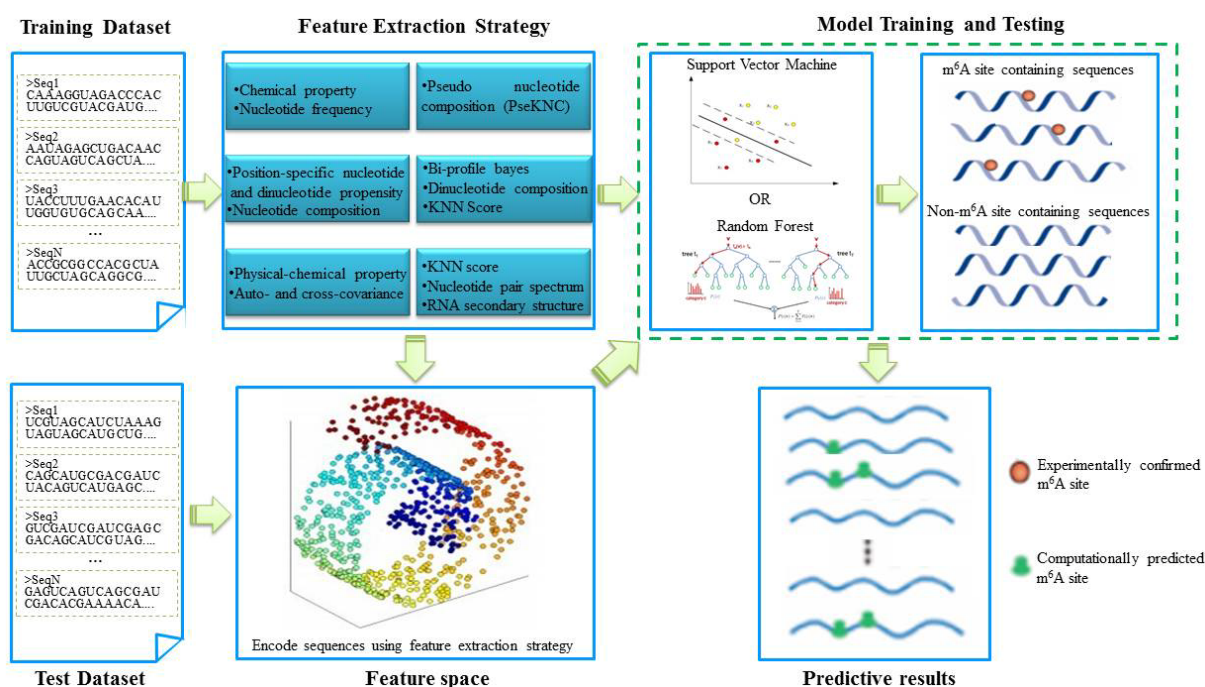
In statistical prediction, the following three cross-validation methods are often used to evaluate the performance of a predictor: independent dataset test, subsampling (or K-fold cross-validation) test, and jackknife test [139]. Among them, however, the jackknife test was deemed the least arbitrary that can always yield a unique result for a given benchmark dataset, as elucidated in [5] and demonstrated by Eqs.28-32 therein. Therefore, the jackknife test has been increasingly recognized and widely adopted by investigators to test the power of various prediction methods (see, *e.g.*, [140-146]). In view of this, the jackknife test was also adopted to evaluate the computational methods in identifying m<sup>6</sup>A sites.

## 6. Web Servers for Detecting m<sup>6</sup>A Sites

The last but not least important step of the Chou's 5-step rules [5] is about the web-server establishment. As pointed out in [147] and demonstrated in a series of recent publications (see, *e.g.*, [55, 57-59, 81, 89, 97, 102, 104, 121-127, 135, 148-160]), user-friendly and publicly accessible web-servers represent the future direction for developing practically more useful prediction methods and computational tools. Actually, many practically useful web-servers have significantly increased the impacts of bioinformatics on medical science [17], driving medicinal chemistry into an unprecedented revolution [46].

## 7. Computational Methods for Detecting m<sup>6</sup>A Sites

Over the past several years, nine different computational methods were proposed to identify m<sup>6</sup>A sites in the *S.cerevisiae* transcriptome. For clarity, their names and web server addresses (if available) are listed in Table 1 according to the chronological order. Show in Figure 2 is the corresponding flowchart.



**Figure 2.** The general framework of computational method for identifying m<sup>6</sup>A sites. The widely used feature extraction strategies and machine learning classifiers were shown in this figure.

**Table 1. List of computational methods for identifying m<sup>6</sup>A sites in *S. cerevisiae***

Methods	Web server address
Schwartz's method (2013)	Not available
m6Apred (2015)	<a href="http://lin-group.cn/server/m6Apred">http://lin-group.cn/server/m6Apred</a>
iRNA-Methyl (2015)	<a href="http://lin-group.cn/server/iRNA-Methyl">http://lin-group.cn/server/iRNA-Methyl</a>
SRAMP (2016)	<a href="http://www.cuilab.cn/sramp/">http://www.cuilab.cn/sramp/</a>
M6A-HPCS (2016)	<a href="http://csbio.njust.edu.cn/bioinf/M6A-HPCS/">http://csbio.njust.edu.cn/bioinf/M6A-HPCS/</a>
pRNAm-PC (2016)	<a href="http://www.jci-bioinfo.cn/pRNAm-PC">http://www.jci-bioinfo.cn/pRNAm-PC</a>
RNA-MethylPred (2016)	Not available
RAM-ESVM (2017)	<a href="http://server.malab.cn/RAM-ESVM/">http://server.malab.cn/RAM-ESVM/</a>
RAM-NPPS (2017)	<a href="http://server.malab.cn/RAM-NPPS/">http://server.malab.cn/RAM-NPPS/</a>
DeepM6APred (2018)	<a href="http://server.malab.cn/DeepM6APred">http://server.malab.cn/DeepM6APred</a>

Using graphic approaches to study biological and medical systems can provide an intuitive vision and useful insights for helping analyze complicated relations therein, as indicated by many previous studies on a series of important biological topics, (see, *e.g.*, [161-174], particularly in enzyme kinetics and protein folding rates [169, 175-177] as well as low-frequency internal motion [178, 179].

Below, we are to make a comparison among the nine different prediction methods via their flowcharts as well.

### 7.1. m<sup>6</sup>A pred

Inspired by Schwartz *et al.*'s pioneer work, a support vector machine (SVM) based method called m6Apred was proposed by Chen *et al.*, which encodes RNA sequences by using both accumulated nucleotide frequency and nucleotide chemical properties (*i.e.*, chemical structure, chemical binding and chemical functionality).

Compared with the classic nucleotide composition, the accumulated nucleotide frequency includes not only the nucleotide frequency information, but also the distribution of each nucleotide in the RNA sequence. The three kinds of nucleotide chemical properties have different impacts on RNA's low-frequency internal motion and its biological function.

Accordingly, each nucleotide in the sequence was represented by a 4-dimensional vector, in which the first element is the accumulated nucleotide frequency and the remaining three elements correspond to the nucleotide chemical properties. For the sequence with a length of  $L$  (where  $L=21$ ), it can be represented by a  $4L$ -dimensional vector and used as the input of SVM. The proposed m6Apred obtained a satisfactory performance for identifying m<sup>6</sup>A site in the *S. cerevisiae* transcriptome based on dataset  $S_1$ .

### 7.2. iRNA-Methyl

It was found that the formation of m<sup>6</sup>A methylation is affected by RNA secondary structure that is highly related with the physicochemical properties of dinucleotide. In view of this, a predictor called "iRNA-Methyl" was proposed by formulating RNA sequences with the pseudo nucleotide composition (PseKNC). By using PseDNC (*i.e.*  $k=2$  in Eq.2), three physicochemical properties, namely enthalpy, entropy, and free energy that can quantify the RNA secondary structures were used to calculate the long-range sequence order effects using the following formula:

$$d_u = \begin{cases} \frac{f_u}{\sum_{i=1}^{16} f_{i+w} \sum_{j=1}^{\lambda} \theta_j} & (1 \leq u \leq 16) \\ \frac{w \theta_{u-16}}{\sum_{i=1}^{16} f_{i+w} \sum_{j=1}^{\lambda} \theta_j} & (16 < u \leq 16 + \lambda) \end{cases} \quad (6)$$

where  $f_u$  ( $u = 1, 2, \dots, 16$ ) is the normalized occurrence frequency of the  $u$ -th non-overlapping dinucleotide in the RNA sequence, and

$$\theta_j = \frac{1}{L-j-1} \sum_{i=1}^{L-j-1} C_{i,i+j} \quad (j = 1, 2, \dots, \lambda; \lambda < L) \quad (7)$$

where  $\theta_j$  is called the  $j$ -tier correlation factor that reflects the sequence order correlation between all the  $j$ -th most contiguous dinucleotide, the coupling factor  $C_{i,i+j}$  is given by

$$C_{i,i+j} = \frac{1}{3} \sum_{g=1}^3 [P_g(D_i) - P_g(D_{i+j})]^2 \quad (8)$$

where  $P_g(D_i)$  ( $g=1, 2, 3$ ) is the normalized value of the above mentioned three RNA physicochemical properties for the  $i$ -th dinucleotide  $D_i$ . By using the 10-fold cross validation test, the optimal values for the parameters  $\lambda$  and  $w$  of PseKNC were obtained (*i.e.*  $\lambda=6$  and  $w=0.9$ ). Accordingly, the samples in dataset  $S_2$  were transferred into a 22-dimensional vector in the iRNA-Methyl method. It was found that iRNA-Methyl obtained an accuracy of 65.59% for identifying m<sup>6</sup>A sites in the rigorous jackknife test. For the convenience of experimental scientists, a web-server for iRNA-Methyl has been established at <http://lin-group.cn/server/iRNA-Methyl>.

### 7.3. pRNAm-PC

In 2016, in order to improve the accuracy of m<sup>6</sup>A site identification, Liu *et al.* proposed the pRNAm-PC method, in which the RNA sequences in dataset  $S_2$  were encoded by using a vector, whose components were derived from a physical-chemical matrix via the auto-covariance and cross-covariance transformations.

Based on the dinucleotide physicochemical properties, a 10×50 dimensional physicochemical property matrix (PC) was defined as following,

$$PC = \begin{bmatrix} P_1(N_1N_2) & P_1(N_2N_3) & \dots & P_1(N_iN_{i+1}) & \dots & P_1(N_{50}N_{51}) \\ P_2(N_1N_2) & P_2(N_2N_3) & \dots & P_2(N_iN_{i+1}) & \dots & P_2(N_{50}N_{51}) \\ \vdots & \vdots & \dots & \vdots & \dots & \vdots \\ P_j(N_1N_2) & P_j(N_2N_3) & \dots & P_j(N_iN_{i+1}) & \dots & P_j(N_{50}N_{51}) \\ \vdots & \vdots & \dots & \vdots & \dots & \vdots \\ P_9(N_1N_2) & P_9(N_2N_3) & \dots & P_9(N_iN_{i+1}) & \dots & P_9(N_{50}N_{51}) \\ P_{10}(N_1N_2) & P_{10}(N_2N_3) & \dots & P_{10}(N_iN_{i+1}) & \dots & P_{10}(N_{50}N_{51}) \end{bmatrix} \quad (9)$$

where  $P_j(N_iN_{i+1})$  is the  $j$ -th ( $j=1, 2, \dots, 10$ ) physicochemical properties value for the dinucleotide  $N_iN_{i+1}$  (namely, AA, AC, AG, AU, CA, ..., or UU) in the RNA sequence. In the pRNAm-PC method, 10 dinucleotide physicochemical properties (*i.e.*  $P_1$ : rise,  $P_2$ : roll,  $P_3$ : shift,  $P_4$ : slide,  $P_5$ : tilt,  $P_6$ : twist,  $P_7$ : enthalpy,  $P_8$ : entropy,  $P_9$ : stack energy, and  $P_{10}$ : free energy) were used.

In order to reflect the correlation of the same and different physicochemical property between two subsequences separated by  $\lambda$  dinucleotides, the auto-covariance (AC) and cross-covariance (CC) method were used to transform the physicochemical property matrix into a length-fixed feature vector and were defined as following.

$$AC(j, \lambda) = \frac{\sum_{i=1}^{50-\lambda} [P_j(N_iN_{i+1}) - \bar{P}_j][P_j(N_{i+\lambda}N_{i+\lambda+1}) - \bar{P}_j]}{50-\lambda} \quad (10)$$

$$CC(j, k, \lambda) = \frac{\sum_{i=1}^{50-\lambda} [P_j(N_iN_{i+1}) - \bar{P}_j][P_k(N_{i+\lambda}N_{i+\lambda+1}) - \bar{P}_k]}{50-\lambda} \quad (j \neq k) \quad (11)$$

By preliminary tests, they found that the best value for  $\lambda$  is 4. Therefore, the RNA sequences in dataset  $S_2$  were encoded by a 400-dimensional vector, of which the 40 elements were deduced from the auto-covariance and the 360 elements from the cross-covariance. Based on this kind of feature, the pRNAm-PC was built and yielded an accuracy of 69.74% for identifying the m<sup>6</sup>A sites in dataset  $S_2$  in the jackknife test, which is ~5% higher than that of iRNA-Methyl. However, the feature dimension of pRNAm-PC was nearly 26 times larger than that of iRNA-Methyl. Moreover, the contributions and biological meanings of the above mentioned 10 physicochemical properties for identifying m<sup>6</sup>A sites were not described at all.

### 7.4. SRAMP

Subsequently, by combining multiple features, Zhou and his colleagues established a random forest based computational predictor, called SRAMP, which is available at <http://www.cuilab.cn/sramp/>. In order to capture more sequence-derived features, the positional nucleotide sequence pattern, K-nearest neighbor information, the position independent nucleotide pair spectrum, and the predicted RNA secondary structure were used to encode the RNA sequences.

For the positional nucleotide sequence pattern, the nucleotide (A, C, G or U) at each position were represented by the binary vector of (1,0,0,0), (0,1,0,0), (0,0,1,0), or (0,0,0,1). For a  $2n+1$  long sequence segment, a  $4 \times (2n+1)$  dimensional vector can be obtained.

In order to measure the extent of how much the flanking window of one query sample resembles those of other m<sup>6</sup>A sites, the K-nearest neighbor information was introduced. Firstly, the flanking window of the query sample was compared with all samples in the training dataset and obtained a pair-wise similarity score,

$$Score = \sum_{i=1}^{2n+1} NUC44(q_i, r_i) \quad (12)$$

where  $q_i$  and  $r_i$  are the nucleotides at the  $i$ -th position of the flanking windows in the query sample and the training samples.  $2n + 1$  is the window size. The NUC44 is a common nucleotide similarity scoring matrix. And then, the fraction of positive samples in the top  $K$  most similar reference samples was taken as the KNN feature. In SRAMP, 30 K values were used (*i.e.*,  $K=50, 100, 150, \dots, 1500$ ).

The sequence context was also reflected by calculating the frequencies of all possible  $d$ -spaced nucleotide pairs, which is defined as

$$\text{Frequency}(np_i) = \frac{C(np_i)}{2n-d-1} \quad (13)$$

where  $C(np_i)$  is the number of  $np_i$  inside a flanking window with a size of  $2n$ ,  $d$  is the space between two nucleotides, and ranged from 0 to 3.

As indicated in their work, the hairpin loop, multiple loop, interior loop, paired and bulged loop from the RNA secondary structure were also used to represent RNA sequences, which were encoded as the binary vectors, namely (1,0,0,0,0), (0,1,0,0,0), (0,0,1,0,0), (0,0,0,1,0), (0,0,0,0,1) and (0,0,0,0,0), respectively. Accordingly, in the structure space, a flanking window with a size of  $2n$  will be converted into a  $2n \times 5$ -dimensional vector.

For each kind of these features, a random forest classifier was built. The final prediction result was the combination of them by using the weighted summing scheme. As indicated by Zhou *et al.*, SRAMP yielded comparable accuracy for identifying m<sup>6</sup>A sites in the *S. cerevisiae* transcriptome to that of iRNA-Methyl. As a big plus, SRAMP can not only identify m<sup>6</sup>A sites in *S. cerevisiae*, but also much more effective for identifying mammalian m<sup>6</sup>A sites.

## 7.5. M6A-HPCS

Later on, with the aim of finding out which physicochemical properties making great contributions for identifying m<sup>6</sup>A sites, Zhang *et al.* proposed a heuristic nucleotide physicochemical property selection algorithm, called M6A-HPCS, to identify m<sup>6</sup>A sites in the *S. cerevisiae* transcriptome. The M6A-HPCS method is based on the iRNA-Methyl and pRNAm-PC methods. However, rather than directly using the physicochemical properties, the relative gain and direct gain methods were used to measure the significance of each of the 23 dinucleotide physicochemical properties for identifying the m<sup>6</sup>A sites. And then, a heuristic algorithm is employed to select the optimal physicochemical properties for the PseKNC (used in iRNA-Methyl) and auto-covariance and cross-covariance (used in pRNAm-PC) encoding schemes, respectively.

For the PseKNC and auto-covariance and cross-covariance encoding scheme, 5 and 13 out of the 23 dinucleotide physicochemical properties were selected out as their optimal candidates to represent the RNA sequences in dataset  $\mathbb{S}_2$ , respectively. In the rigorous jackknife test, the accuracies of 67.33% and 72.38% were obtained for identifying m<sup>6</sup>A sites for both encoding schemes, respectively. Although its predictive accuracy is higher than those of iRNA-Methyl and pRNAm-PC, the shortcomings for M6A-HPCS still exist in the following aspects. First, the biological meanings of using the optimal dinucleotide physicochemical properties are not described at all. Second, although a web-server was developed for M6A-HPCS at <http://csbio.njst.edu.cn/bioinf/M6A-HPCS>, it couldn't be accessed anymore.

## 7.6. RNA-MethylPred

To further improve the accuracy of m<sup>6</sup>A site identification, Jia *et al.* proposed a new computational method called RNA-MethylPred. In this method, three kinds of feature extraction strategies were used to represent the RNA sequences in dataset  $\mathbb{S}_2$ .

Bi-profile bayes vector ( $V$ ) was employed to reflect the posterior probability of positive and negative samples.

$$V = [p_1, p_2, \dots, p_n, p_{n+1}, \dots, p_{2n}] \quad (14)$$

where the first  $n$  components denotes the posterior probability of each nucleotide at the  $i$ -th position in the positive samples, the remaining  $n$  components denotes the posterior probability of each nucleotide at the  $i$ -th position in the negative samples.  $n$  is equal to the length of the RNA sequence (*i.e.*  $n=51$ ).

Two forms of dinucleotide composition were defined to reflect sequence order information.

$$P_{ab} = \frac{N_{ab}}{N_a} \quad (15)$$

$$P'_{ab} = \frac{N_{ab}}{n-1} \quad (16)$$

where  $N_{ab}$  is the number of neighboring dinucleotide ( $a, b$  can be nucleotide A, C, G or U),  $a\bullet$  indicates any adjoining dinucleotides that starting with  $a$ .

K nearest neighbor (KNN) scores were used to measure whether the local sequence similarity. To this end, similarity score  $S(A,B)$  between two sequence fragments A and B was defined as

$$S(A, B) = \sum_{1 \leq i \leq 51} \text{Score}(A[i], B[i]) \quad (17)$$

$A[i]$  indicates the nucleotide at the  $i$ -th position in sequence segment A, and the score of two nucleotides was defined as

$$\text{Score} = \begin{cases} +2 & \text{when two nucleotides matched} \\ -1 & \text{when two nucleotides mismatched} \end{cases} \quad (18)$$

Based on Eqs. (11) and (12), the KNN score was achieved by calculating the percentage of the positive neighbors in its KNNs. In RNA-MethylPred, the 20 considered Ks were 10, 20, ..., 200.

Finally, these features were combined together and used as the input of SVM to perform the prediction. In the jack-knife test, the RNA-MethylPred obtained an accuracy of 76.51% for identifying m<sup>6</sup>A sites in the *S. cerevisiae* transcriptome. Rather than building a web-server, the authors provided a MATLAB package for the RNA-MethylPred method.

### 7.7. RAM-ESVM

As introduced above, various features and predictors have been proposed for identifying m<sup>6</sup>A sites. However, their performances are still not satisfactory. In 2017, Chen *et al.* developed an ensemble classifier called RAM-ESVM, which combines three basic classifiers based on different features including PseKNC, motif features, and optimized K-mer. The first two classifiers (SVM-PseKNC and SVM-motif) were built based on SVM by using PseKNC and motif features as the inputs, respectively. The third one is also a SVM based classifier and its input features are optimized gapped  $k$ -mers, which is achieved by using the GkmSVM software. The three basic classifiers vote for the final result based on the voting score.

$$V_i = \sum_{k=1}^3 f(\text{pre}(C_k), \text{Class}_i) \quad (i = 1, 2; k = 1, 2, 3) \quad (19)$$

where  $V_i$  is the voting score for the RNA sample belonging to the class <sub>$i$</sub>  ( $i=1$ : m<sup>6</sup>A sites;  $i=2$ : non- m<sup>6</sup>A sites), and

$$f(\text{pre}(C_k), \text{Class}_i) = \begin{cases} 1 & \text{if } \text{pre}(C_k) \in \text{Class}_i \\ 0 & \text{otherwise} \end{cases} \quad (20)$$

The final prediction is determined by the argument that maximizes the voting score  $V_i$ ,

$$\text{Sgn}(i) = \arg \max_i \{V_i\} \quad (21)$$

In the jackknife test, the RAM-ESVM produced an accuracy of 78.35% for identifying m<sup>6</sup>A sites in the *S. cerevisiae* transcriptome. The RAM-ESVM can be freely accessed at <http://server.malab.cn/RAM-ESVM/>.

### 7.8. RAM-NPPS

In 2017, another m<sup>6</sup>A site predictor, called RAM-NPPS, was proposed by Xing *et al.*, which is based on multi-interval nucleotide pair position specificity (NPPS).

For a given RNA sequence segment P, it can be represented by  $P = P^+ \cdot P^-$ .  $P^+$  and  $P^-$  can be formulated as

$$P^\xi = p_1^\xi p_2^\xi \dots p_k^\xi \dots p_{51}^\xi \quad (\xi \text{ is } + \text{ or } -) \quad (22)$$

To obtain  $p_k^\xi$ , single nucleotide position matrix  $T_s^\xi$  and dinucleotide position matrix  $T_d^\xi$  are defined as following

$$T_s^\xi = \begin{bmatrix} f_{A,1}^\xi & f_{A,2}^\xi & \dots & f_{A,51}^\xi \\ f_{C,1}^\xi & f_{C,2}^\xi & \dots & f_{C,51}^\xi \\ f_{G,1}^\xi & f_{G,2}^\xi & \dots & f_{G,51}^\xi \\ f_{U,1}^\xi & f_{U,2}^\xi & \dots & f_{U,51}^\xi \end{bmatrix} \quad (23)$$

$$T_d^\xi = \begin{bmatrix} f_{AA,1}^\xi & f_{AA,2}^\xi & \dots & f_{AA,50}^\xi \\ f_{AC,1}^\xi & f_{AC,2}^\xi & \dots & f_{AC,50}^\xi \\ \vdots & \vdots & \dots & \vdots \\ f_{UU,1}^\xi & f_{UU,2}^\xi & \dots & f_{UU,50}^\xi \end{bmatrix} \quad (24)$$

The elements in the two matrices indicate the occurrence probability of the nucleotide in each position and the occurrence probability of the nucleotide pair at position  $i$  and  $i+\eta$ , respectively.

Suppose the dinucleotide between the  $i$ -th nucleotide and  $(i+\eta)$ -th nucleotide is “ab”,  $p_k^\xi$  can be calculated according to the conditional probability,

$$p_k^\xi = \frac{P(a \cap b)}{P(b)} = \frac{f_{ab,i}^\xi}{f_{b,i}^\xi} \quad (25)$$

Accordingly, the RNA sequence can be converted into the feature vector as described in Eq. 18. When the optimal interval value of the two nucleotides is set as  $\eta=5$ , the SVM based computational model RAM-NPPS was built, which is available at <http://server.malab.cn/ram-npps/>. In the jackknife test, RAM-NPPS yielded an accuracy of 79.92% for identifying m<sup>6</sup>A sites in the *S. cerevisiae* transcriptome.

### 7.9. DeepM6APred

More recently, Wei *et al.* proposed a new method called DeepM6APred, which represented the RNA samples by using both the above mentioned NPPS feature and binary string encoding scheme. Different from traditional methods, before directly using these features to make predictions, the deep-belief network was used to automatically learn meaningful feature representations from raw input sequences. Finally, an optimal feature set containing 429 features was obtained, based on which a predictive accuracy of 80.50% was obtained for identifying m<sup>6</sup>A sites in the *S. cerevisiae* transcriptome. To the best of our knowledge, this is the best accuracy for identifying m<sup>6</sup>A sites in the *S. cerevisiae* transcriptome till now. The DeepM6APred can be accessed at <http://server.malab.cn/deepm6apred/>.

## 8. Comparison of Various Prediction Methods

In this section, we performed a comparison on existing methods for identifying m<sup>6</sup>A sites in the *S. cerevisiae* transcriptome. Since SRAMP is a mammalian specific predictor and m6Apred is trained based on dataset  $S_1$ , for a fair comparison, they were not considered here. The predictive accuracies of the other 7 methods for identifying m<sup>6</sup>A sites based on the benchmark dataset  $S_2$  were shown in Figure 3. It was found the performance of DeepM6APred ranks the top.

To further demonstrate the generalization ability of these methods, an independent dataset was built, which includes 239 m<sup>6</sup>A site containing sequences obtained from the RMBase, and the same number of non-m<sup>6</sup>A site containing sequences. All these sequences are 51 nt and independent from the samples in the dataset  $S_2$ , which are available at <https://github.com/chenweimu/m6a>.

It should be point out that the web-server of M6A-HPCS is not accessible anymore as indicated in its homepage, and the pRNAm-PC could not make predictions for these independent sequences. Therefore, the comparisons were performed among the remaining methods (*i.e.* iRNA-Methyl, RAM-ESVM, RNA-MethylPred, RAM-NPPS, and DeepM6APred). Their predictive results for identifying m<sup>6</sup>A sites in the independent dataset were reported in Table 2. It was found that the Sn, Acc and MCC of DeepM6APred are much higher than the other four methods. Although iRNA-Methyl obtained a high Sp, it has lower Sn, Acc and MCC. Thus, we can draw a conclusion that the performance of DeepM6APred is the best, while the performance of iRNA-Methyl is comparable to RNA-MethylPred and RAM-NPPS.

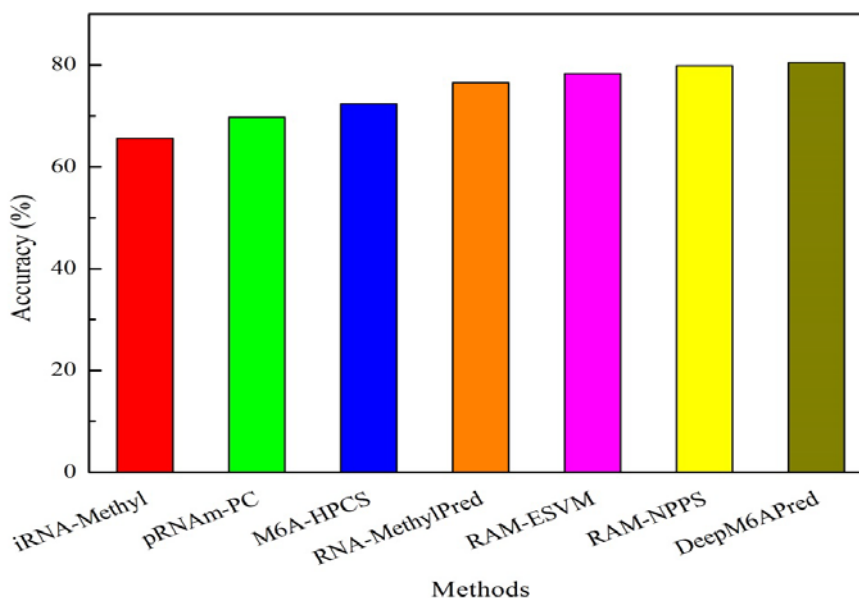
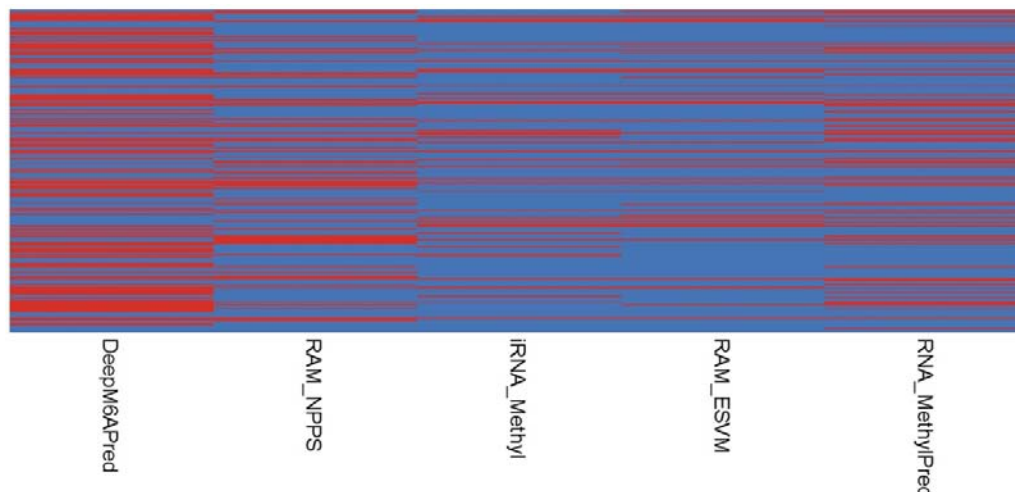


Figure 3. The performance of different methods for identifying m<sup>6</sup>A sites in the benchmark dataset  $S_2$ .

**Table 2. Performance comparisons of different methods for identifying m<sup>6</sup>A sites based on independent dataset. See Eq. 5 for the definition of metrics below.**

Methods	Sn (%)	Sp (%)	Acc (%)	MCC
iRNA-Methyl	19.25	80.75	50.00	0.00
RAM-ESVM	18.83	68.62	43.72	-0.14
RNA-MethylPred	27.20	74.06	50.63	0.01
RAM-NPPS	28.87	71.55	50.21	0.00
DeepM6APred	48.54	70.29	59.41	0.19

**Figure 4. The detail predictive results of DeepM6APred, RAM-NPPS, iRNA-Methyl, RAM-ESVM and RNA-MethylPred based on the independent dataset (Each row is a sample in the independent dataset. The correctly identified m<sup>6</sup>A site containing samples were highlighted in red, and the counterparts were in blue).**

As demonstrated by many previous studies on a series of important biological topics (see, *e.g.*, [161, 162, 166, 168, 170, 173, 180-182]), using image or graphic approaches to study biological systems can provide intuitive insights for helping analyze complicated relations therein, in view of this, the accurately predicted m<sup>6</sup>A sites by the different methods were presented in Figure 4. As we can see from the figure, of the 239 m<sup>6</sup>A site containing sequences, 116 were correctly identified by DeepM6APred, 69 by RAM-NPPS, 65 by RNA-MethylPred, 46 by iRNA-Methyl, and 45 by RAM-ESVM. These results indicate that for users who are interested in identifying m<sup>6</sup>A sites in the *S. cerevisiae* transcriptome, the DeepM6APred predictor should be their first choice, and the other predictors, namely RAM-NPPS, RNA-MethylPred, iRNA-Methyl and RAM-ESVM, may used as complementary tools in this regard.

## 9. Conclusive Remarks and Perspective

In this paper, we comprehensively reviewed the computational methods for identifying m<sup>6</sup>A sites in the *S. cerevisiae* transcriptome and evaluated their performance based on the independent dataset. Although these methods obtained quite good results in identifying m<sup>6</sup>A sites when tested by the benchmark dataset  $S_2$  of Section 2, their exploited or extrapolative effectiveness in practical application [139] was not so ideal as reflected by the fact when tested by the independent dataset.

The poor performance of these methods on the independent dataset is due to the following reason. All these methods were trained based on dataset  $S_2$ , in which both positive and negative samples were obtained by selecting the sequences containing RGAC consensus motif. However, in most cases, the m<sup>6</sup>A site may not locate in the RGAC consensus motif. Thus, the construction of the benchmark dataset in such a way precluded the generalization ability of these methods. In order to improve the performance and generalization ability of the computational methods for identifying m<sup>6</sup>A sites, much more efforts should be made by considering the following aspects.

### 9.1. The performance is dependent on the benchmark dataset

Although several benchmark datasets have been established for training computational models for identifying m<sup>6</sup>A

sites, the challenges still exist in the construction of the benchmark dataset.

Compared with the positive samples, there is no uniform standard to collect negative samples (non-m<sup>6</sup>A samples). The popular strategy of obtaining non-m<sup>6</sup>A samples is to select the adenosines that are not experimentally annotated as being methylated. It indeed raises the possibility that the m<sup>6</sup>A sites are not identified may serve as false negative samples. In addition, in the real case, the number of non-m<sup>6</sup>A sites is significantly higher than that of m<sup>6</sup>A sites. The existing benchmark datasets are all balanced ones that contain roughly equal number of m<sup>6</sup>A samples and the randomly selected non-m<sup>6</sup>A samples. Such randomly sampling of non-m<sup>6</sup>A samples may lead to inadequate learning and the models trained on such a dataset would change when the selected non-m<sup>6</sup>A samples are different. To solve these challenges, more efforts can be made in the following aspects.

First, the one-sided selection (OSS) undersampling and synthetic minority oversampling technique (SMOTE) can be used to balance the non-m<sup>6</sup>A and m<sup>6</sup>A samples and minimize the influences of imbalance issue. The one-sided selection (OSS) undersampling employed the condensed nearest-neighbor to remove redundant negative samples that are far from the boundary of the class and the Tomek links to eliminate borderline samples and samples suffering from class label noise. By doing so, the number of non-m<sup>6</sup>A samples can be decreased. On the other hand, the SMOTE will resample the small class (m<sup>6</sup>A samples) by taking each small class example and introducing synthetic examples along the line segments joining it to the small class nearest neighbors. Accordingly, the positive and negative samples will be balanced.

The second way to deal with such an imbalance problem is to use ensemble techniques, which trains basic classifiers with different sampling data and combines their results to reduce the random sampling bias. The key step of this technique is to select meaningful negative samples to train basic classifiers.

Another strategy is to use cost-sensitive classifiers, such as XGboost (eXtreme Gradient Boosting), which can be trained with all the samples without selecting a subset of negative samples and prevent training model from over-fitting by defining different costs for the misclassified positive and negative samples.

## 9.2. Encode RNA sequences using effective schemes

Feature extraction strategy is another essential step to build computational models for identifying m<sup>6</sup>A sites. The performance of existing models for identifying m<sup>6</sup>A sites depends on how to accurately represent RNA sequences. The encoding schemes are based on the experiences and usually derived from the segments surrounding the m<sup>6</sup>A sites, such as pseudo nucleotide compositions, physicochemical properties, position specific nucleotide/dinucleotide composition, and so on. Although considerable progresses have been achieved, the following aspects should be considered for designing distinguishable feature descriptors in the future work.

Except for Schwartz *et al.*'s and Zhou *et al.*'s works, none of the other existing computational methods represented the RNA samples using RNA secondary structure information. By regulating the interaction of methyltransferase complex with RNA sequences, RNA secondary structure is closely related to the formation of m<sup>6</sup>A. Therefore, it is necessary to integrate this kind of feature when constructing more powerful computational models for identifying m<sup>6</sup>A sites. To this end, the RNAfold tool in ViennaRNA package can be used to predict RNA secondary structure, whose output is dots (indicate unpaired nucleotides) and brackets (indicate paired nucleotides). If encode unpaired nucleotides using 0 and the paired one using 1, a given RNA sample will be transferred into a feature vector with its elements are 0 and 1.

Another shortcoming of existing methods is that existing computational methods directly use the entire features, which may lead to over-fitting problems, reduce the generalization capacity of the model and increase the computational time. In order to alleviate irrelevant features and overcome the above mentioned shortcoming, the feature selection techniques, such as minimal redundancy maximal relevance (mRMR), maximum relevancy maximum distance (MRMD), and analysis of variance (ANOVA), can be used to winnow out the optimal features.

## 9.3. Generalizability of existing computational approaches

Compared with the performance for identifying m<sup>6</sup>A sites in other species, the accuracy for identifying m<sup>6</sup>A sites in the *S. cerevisiae* transcriptome is still far from satisfactory. Therefore, new computational models are still required. Besides support vector machine and random forest, other machine learning methods such as Native Bayes, Logistic Regression, and K-nearest neighbor are all potential candidates to build new computational models for identifying m<sup>6</sup>A sites. With the development of convolutional neural network and deep learning, these advantaged approaches are also suggested to be used in developing computational models. In addition, since most of the existing methods are complementary to each other (Figure 4), it is wise to employ the ensemble classification techniques to develop computational models with high performance.

Besides m<sup>6</sup>A, the pseudouridine, N<sup>1</sup>-Methyladenosine (m<sup>1</sup>A), and 5-methylcytosine (m<sup>5</sup>C) are also frequently observed RNA modifications. However, both the computational models and experimental techniques could not simultaneously identify these different types of RNA modifications. To address such a challenge, more efforts should be made to develop

a platform that can be used to simultaneously detect different types of RNA modifications.

## Acknowledgements

The author thanks Dr Cangzhi Jia for her assistance in running the RNA-MethylPred program.

## Funding

This work was supported by the National Nature Scientific Foundation of China (31771471, 61772119), Natural Science Foundation for Distinguished Young Scholar of Hebei Province (No. C2017209244).

## References

- [1] H. Ding, E. Z. Deng, L. F. Yuan, L. Liu, H. Lin, W. Chen, K. C. Chou. (2014). iCTX-Type: A sequence-based predictor for identifying the types of conotoxins in targeting ion channels, *BioMed Research International (BMRI)*, 2014(2014), 286419.
- [2] G. L. Fan, X. Y. Zhang, Y. L. Liu, Y. Nang, H. Wang. (2015). DSPMP: Discriminating secretory proteins of malaria parasite by hybridizing different descriptors of Chou's pseudo amino acid patterns. *J. Comput. Chem.*, 36(2015), 2317-2327.
- [3] L. Pan, W. Zhao, J. Lai, D. Ding, Q. Zhang, X. Yang, M. Huang, S. Jin, Y. Xu, S. Zeng, J. J. Chou, S. Chen. (2017). Sortase A-Generated Highly Potent Anti-CD20-MMAE Conjugates for Efficient Elimination of B-Lineage Lymphomas, *Small*, 13(2017).
- [4] K. Oxenoid, Y. S. Dong, C. Cao, T. Cui, Y. Sancak, A. L. Markhard, Z. Grabarek, L. Kong, Z. Liu, B. Ouyang, Y. Cong, V. K. Mootha, J. J. Chou. (2016). Architecture of the Mitochondrial Calcium Uniporter, *Nature*, 533(2016), 269-273.
- [5] K. C. Chou. (2011). Some remarks on protein attribute prediction and pseudo amino acid composition (50th Anniversary Year Review, 5-steps rule). *J. Theor. Biol.*, 273(2011), 236-247.
- [6] K. C. Chou, H. B. Shen. (2008). Cell-PLoc: A package of Web servers for predicting subcellular localization of proteins in various organisms. *Nature Protocols*, 3(2008), 153-162.
- [7] K. C. Chou, H. B. Shen. (2010). Cell-PLoc 2.0: An improved package of web-servers for predicting subcellular localization of proteins in various organisms. *Natural Science*, 2(2010), 1090-1103.
- [8] K. C. Chou, H. B. Shen. (2007). Recent progresses in protein subcellular location prediction. *Anal. Biochem.*, 370(2007), 1-16.
- [9] K. C. Chou, D. W. Elrod. (2002). Bioinformatical analysis of G-protein-coupled receptors. *Journal of Proteome Research*, 1(2002), 429-433.
- [10] W. Chen, H. Lin, P. M. Feng, C. Ding, Y. C. Zuo, K. C. Chou. (2012). iNuc-PhysChem: A Sequence-Based Predictor for Identifying Nucleosomes via Physicochemical Properties. *PLoS ONE*, 7(2012), e47843.
- [11] Y. D. Cai, K. C. Chou. (2004). Predicting subcellular localization of proteins in a hybridization space. *Bioinformatics*, 20(2004), 1151-1156.
- [12] K. C. Chou, Y. D. Cai. (2006). Prediction of protease types in a hybridization space. *Biochem Biophys Res Comm (BBRC)*, 339(2006), 1015-1020.
- [13] P. M. Feng, W. Chen, H. Lin, K. C. Chou. (2013). iHSP-PseRAAAC: Identifying the heat shock protein families using pseudo reduced amino acid alphabet composition. *Anal. Biochem.*, 442(2013), 118-125.
- [14] W. Chen, P. M. Feng, H. Lin, K. C. Chou. (2013). iRSpot-PseDNC: identify recombination spots with pseudo dinucleotide composition. *Nucleic Acids Research*, 41(2013), e68.
- [15] W. Z. Lin, J. A. Fang, X. Xiao, K. C. Chou. (2011). iDNA-Prot: Identification of DNA Binding Proteins Using Random Forest with Grey Model. *PLoS ONE*, 6(2011), e24756.
- [16] J. Jia, Z. Liu, X. Xiao, B. Liu, K. C. Chou. (2016). pSuc-Lys: Predict lysine succinylation sites in proteins with PseAAC and ensemble random forest approach. *J. Theor. Biol.*, 394(2016), 223-230.
- [17] K. C. Chou. (2015). Impacts of bioinformatics to medicinal chemistry. *Medicinal Chemistry*, 11(2015), 218-234.

- [18] K. C. Chou. (2001). Prediction of protein cellular attributes using pseudo amino acid composition, *PROTEINS: Structure, Function, and Genetics* (Erratum: *ibid.*, 2001, Vol. 44, 60), 43(2001), 246-255.
- [19] K. C. Chou. (2005). Using amphiphilic pseudo amino acid composition to predict enzyme subfamily classes. *Bioinformatics*, 21(2005), 10-19.
- [20] K. C. Chou, Y. D. Cai. (2003). Prediction and classification of protein subcellular location: sequence-order effect and pseudo amino acid composition, *Journal of Cellular Biochemistry* (Addendum, *ibid.* 2004, 91, 1085), 90(2003) 1250-1260.
- [21] M. Arif, M. Hayat, Z. Jan. (2018). iMem-2LSAAC: A two-level model for discrimination of membrane proteins and their types by extending the notion of SAAC into Chou's pseudo amino acid composition. *J. Theor. Biol.*, 442(2018), 11-21.
- [22] J. Mei, J. Zhao. (2018). Prediction of HIV-1 and HIV-2 proteins by using Chou's pseudo amino acid compositions and different classifiers. *Sci Rep*, 8(2018), 2359.
- [23] J. Mei, J. Zhao. (2018). Analysis and prediction of presynaptic and postsynaptic neurotoxins by Chou's general pseudo amino acid composition and motif features. *J. Theor. Biol.*, 427(2018), 147-153.
- [24] M. S. Krishnan. (2018). Using Chou's general PseAAC to analyze the evolutionary relationship of receptor associated proteins (RAP) with various folding patterns of protein domains. *J. Theor. Biol.*, 445(2018), 62-74.
- [25] L. Zhang, L. Kong. (2018). iRSpot-ADPM: Identify recombination spots by incorporating the associated dinucleotide product model into Chou's pseudo components. *J. Theor. Biol.*, 441(2018), 1-8.
- [26] S. Zhang, X. Duan. (2018). Prediction of protein subcellular localization with oversampling approach and Chou's general PseAAC. *J. Theor. Biol.*, 437(2018), 239-250.
- [27] A. H. Butt, N. Rasool, Y. D. Khan. (2018). Predicting membrane proteins and their types by extracting various sequence features into Chou's general PseAAC. *Molecular biology reports*, 18(2018), 39-58.
- [28] E. Contreras-Torres. (2018). Predicting structural classes of proteins by incorporating their global and local physicochemical and conformational properties into general Chou's PseAAC. *J. Theor. Biol.*, 454(2018), 139-145.
- [29] F. Javed, M. Hayat. (2018). Predicting subcellular localizations of multi-label proteins by incorporating the sequence features into Chou's PseAAC. *Genomics*, 17(2018), 793-821.
- [30] Z. Ju, S. Y. Wang. (2018). Prediction of citrullination sites by incorporating k-spaced amino acid pairs into Chou's general pseudo amino acid composition. *Gene*, 664(2018), 78-83.
- [31] Y. Liang, S. Zhang. (2018). Identify Gram-negative bacterial secreted protein types by incorporating different modes of PSSM into Chou's general PseAAC via Kullback-Leibler divergence. *J. Theor. Biol.*, 454(2018), 22-29.
- [32] J. Mei, Y. Fu, J. Zhao. (2018). Analysis and prediction of ion channel inhibitors by using feature selection and Chou's general pseudo amino acid composition. *J. Theor. Biol.*, 456(2018), 41-48.
- [33] M. Mousavizadegan, H. Mohabatkar. (2018). Computational prediction of antifungal peptides via Chou's PseAAC and SVM. *Journal of bioinformatics and computational biology*, (2018), 1850016.
- [34] W. Qiu, S. Li, X. Cui, Z. Yu, M. Wang, J. Du, Y. Peng, B. Yu. (2018). Predicting protein submitochondrial locations by incorporating the pseudo-position specific scoring matrix into the general Chou's pseudo-amino acid composition. *J. Theor. Biol.*, 450(2018), 86-103.
- [35] S. M. Rahman, S. Shatabda, S. Saha, M. Kaykobad, M. Sohel Rahman. (2018). DPP-PseAAC: A DNA-binding Protein Prediction model using Chou's general PseAAC. *J. Theor. Biol.*, 452(2018), 22-34.
- [36] E. S. Sankari, D. D. Manimegalai. (2018). Predicting membrane protein types by incorporating a novel feature set into Chou's general PseAAC. *J. Theor. Biol.*, 455(2018), 319-328.
- [37] A. Srivastava, R. Kumar, M. Kumar. (2018). BlaPred: predicting and classifying beta-lactamase using a 3-tier prediction system via Chou's general PseAAC. *J. Theor. Biol.*, 457(2018), 29-36.
- [38] S. Zhang, Y. Liang. (2018). Predicting apoptosis protein subcellular localization by integrating auto-cross correlation and PSSM into Chou's PseAAC. *J. Theor. Biol.*, 457(2018), 163-169.
- [39] W. Zhao, L. Wang, T. X. Zhang, Z. N. Zhao, P. F. Du. (2018). A brief review on software tools in generating Chou's pseudo-factor representations for all types of biological sequences. *Protein Pept Lett*, 25(2018), 822-829.

- [40] S. Akbar, M. Hayat. (2018). iMethyl-STTNC: Identification of N(6)-methyladenosine sites by extending the Idea of SAAC into Chou's PseAAC to formulate RNA sequences. *J. Theor. Biol.*, 455(2018), 205-211.
- [41] M. A. Al Maruf, S. Shatabda. (2018). iRSpot-SF: Prediction of recombination hotspots by incorporating sequence based features into Chou's Pseudo components. *Genomics*, 18(2018), 63-82.
- [42] Y. Pan, S. Wang, Q. Zhang, Q. Lu, D. Su, Y. Zuo, L. Yang. (2019). Analysis and prediction of animal toxins by various Chou's pseudo components and reduced amino acid compositions. *J. Theor. Biol.*, 462(2019), 221-229.
- [43] M. Tahir, M. Hayat, S. A. Khan. (2019). iNuc-ext-PseTNC: an efficient ensemble model for identification of nucleosome positioning by extending the concept of Chou's PseAAC to pseudo-tri-nucleotide composition, *Molecular genetics and genomics: MGG*, 294(2019), 199-210.
- [44] M. Tahir, H. Tayara, K. T. Chong. (2019). iRNA-PseKNC(2methyl): Identify RNA 2'-O-methylation sites by convolution neural network and Chou's pseudo components. *J. Theor. Biol.*, 465(2019), 1-6.
- [45] B. Tian, X. Wu, C. Chen, W. Qiu, Q. Ma, B. Yu. (2019). Predicting protein-protein interactions by fusing various Chou's pseudo components and using wavelet denoising approach. *J. Theor. Biol.*, 462(2019), 329-346.
- [46] K. C. Chou. (2017). An unprecedented revolution in medicinal chemistry driven by the progress of biological science. *Current Topics in Medicinal Chemistry*, 17(2017), 2337-2358.
- [47] H. B. Shen, K. C. Chou. (2008). PseAAC: a flexible web-server for generating various kinds of protein pseudo amino acid composition. *Anal. Biochem.*, 373(2008), 386-388.
- [48] P. Du, X. Wang, C. Xu, Y. Gao. (2012). PseAAC-Builder: A cross-platform stand-alone program for generating various special Chou's pseudo amino acid compositions, *Anal. Biochem.*, 425(2012), 117-119.
- [49] D. S. Cao, Q. S. Xu, Y. Z. Liang. (2013). Propy: a tool to generate various modes of Chou's PseAAC. *Bioinformatics*, 29(2013), 960-962.
- [50] P. Du, S. Gu, Y. Jiao. (2014). PseAAC-General: Fast building various modes of general form of Chou's pseudo amino acid composition for large-scale protein datasets. *International Journal of Molecular Sciences*, 15(2014), 3495-3506.
- [51] K. C. Chou. (2009). Pseudo amino acid composition and its applications in bioinformatics, proteomics and system biology. *Current Proteomics*, 6(2009), 262-274.
- [52] W. Chen, T. Y. Lei, D. C. Jin, H. Lin, K. C. Chou. (2014). PseKNC: a flexible web-server for generating pseudo K-tuple nucleotide composition, *Anal. Biochem.*, 456(2014), 53-60.
- [53] W. Chen, H. Lin, K. C. Chou. (2015). Pseudo nucleotide composition or PseKNC: an effective formulation for analyzing genomic sequences, *Mol BioSyst*, 11(2015), 2620-2634.
- [54] W. Chen, H. Tang, J. Ye, H. Lin, K. C. Chou. (2016). iRNA-PseU: Identifying RNA pseudouridine sites *Molecular Therapy - Nucleic Acids*, 5(2016), e332.
- [55] B. Liu, L. Fang, R. Long, X. Lan, K. C. Chou. (2016). iEnhancer-2L: a two-layer predictor for identifying enhancers and their strength by pseudo k-tuple nucleotide composition. *Bioinformatics*, 32(2016), 362-369.
- [56] B. Liu, R. Long, K. C. Chou. (2016). iDHS-EL: Identifying DNase I hypersensitive sites by fusing three different modes of pseudo nucleotide composition into an ensemble learning framework. *Bioinformatics*, 32(2016), 2411-2418.
- [57] P. Feng, H. Ding, H. Yang, W. Chen, H. Lin, K. C. Chou. (2017). iRNA-PseColl: Identifying the occurrence sites of different RNA modifications by incorporating collective effects of nucleotides into PseKNC. *Molecular Therapy - Nucleic Acids*, 7(2017), 155-163.
- [58] B. Liu, S. Wang, R. Long, K. C. Chou. (2017). iRSpot-EL: identify recombination spots with an ensemble learning approach. *Bioinformatics*, 33(2017), 35-41.
- [59] B. Liu, F. Yang, K. C. Chou. (2017). 2L-piRNA: A two-layer ensemble classifier for identifying piwi-interacting RNAs and their function. *Molecular Therapy - Nucleic Acids*, 7(2017), 267-277.
- [60] M. F. Sabooh, N. Iqbal, M. Khan, M. Khan, H. F. Maqbool. (2018). Identifying 5-methylcytosine sites in RNA sequence using composite encoding feature into Chou's PseKNC. *J. Theor. Biol.*, 452(2018), 1-9.
- [61] B. Liu, F. Liu, X. Wang, J. Chen, L. Fang, K. C. Chou. (2015). Pse-in-One: a web server for generating various modes of pseudo components of DNA, RNA, and protein sequences. *Nucleic Acids Res.*, 43(2015), W65-W71.

- [62] B. Liu, H. Wu, K. C. Chou. (2017). Pse-in-One 2.0: An improved package of web servers for generating various modes of pseudo components of DNA, RNA, and protein sequences. *Natural Science*, 9(2017), 67-91.
- [63] Z. D. Su, Y. Huang, Z. Y. Zhang, Y. W. Zhao, D. Wang, W. Chen, K. C. Chou, H. Lin. (2018). iLoc-lncRNA: predict the subcellular location of lncRNAs by incorporating octamer composition into general PseKNC. *Bioinformatics*, 34(2018), 4196-4204.
- [64] W. Chen, P. Feng, H. Yang, H. Ding, H. Lin, K.C. Chou. (2018). iRNA-3typeA: identifying 3-types of modification at RNA's adenosine sites. *Molecular Therapy: Nucleic Acid*, 11(2018), 468-474.
- [65] H. Yang, W. R. Qiu, G. Liu, F. B. Guo, W. Chen, K.C. Chou, H. Lin. (2018). iRSpot-Pse6NC: Identifying recombination spots in *Saccharomyces cerevisiae* by incorporating hexamer composition into general PseKNC. *International Journal of Biological Sciences*, 14(2018), 883-891.
- [66] K. C. Chou, Y. D. Cai. (2002). Using functional domain composition and support vector machines for prediction of protein subcellular location. *J. Biol. Chem.*, 277(2002), 45765-45769.
- [67] Y. D. Cai, G. P. Zhou, K. C. Chou. (2003). Support vector machines for predicting membrane protein types by using functional domain composition. *Biophys. J.*, 84(2003), 3257-3263.
- [68] N. Cristianini, J. Shawe-Taylor. (2000). *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*, Chapter 3, Cambridge University Press.
- [69] L. Breiman, *Random Forests*. (2001). *Machine learning*, 45(2001), 5-32.
- [70] K. K. Kandaswamy, K. C. Chou, T. Martinetz, S. Moller, P. N. Suganthan, S. Sridharan, G. Pugalenth. (2011). AFP-Pred: A random forest approach for predicting antifreeze proteins from sequence-derived properties. *J. Theor. Biol.*, 270(2011), 56-62.
- [71] G. Pugalenth, K. K. Kandaswamy, K. C. Chou, S. Vivekanandan, P. Kolatkar. (2012). RSARF: Prediction of Residue Solvent Accessibility from Protein Sequence Using Random Forest Method. *Protein & Peptide Letters*, 19(2012), 50-56.
- [72] Y. Xu, J. Ding, L. Y. Wu, K. C. Chou. (2013). iSNO-PseAAC: Predict cysteine S-nitrosylation sites in proteins by incorporating position specific amino acid propensity into pseudo amino acid composition PLoS ONE, 8(2013), e55844.
- [73] J. Jia, Z. Liu, X. Xiao, K. C. Chou. (2015). iPPI-Esml: an ensemble classifier for identifying the interactions of proteins by incorporating their physicochemical properties and wavelet transforms into PseAAC. *J. Theor. Biol.*, 377(2015), 47-56.
- [74] J. Jia, Z. Liu, X. Xiao, B. Liu, K. C. Chou. (2016). Identification of protein-protein binding sites by incorporating the physicochemical properties and stationary wavelet transforms into pseudo amino acid composition (iPPBS-PseAAC). *J Biomol Struct Dyn (JBSD)*, 34(2016), 1946-1961.
- [75] J. Jia, Z. Liu, X. Xiao, B. Liu, K. C. Chou. (2016). iSuc-PseOpt: Identifying lysine succinylation sites in proteins by incorporating sequence-coupling effects into pseudo components and optimizing imbalanced training dataset. *Anal. Biochem.*, 497(2016), 48-56.
- [76] J. Jia, Z. Liu, X. Xiao, B. Liu, K. C. Chou. (2016). iCar-PseCp: identify carbonylation sites in proteins by Monto Carlo sampling and incorporating sequence coupled effects into general PseAAC. *Oncotarget*, 7(2016), 34558-34570.
- [77] W. R. Qiu, B. Q. Sun, X. Xiao, Z. C. Xu, K. C. Chou. (2016). iHyd-PseCp: Identify hydroxyproline and hydroxylysine in proteins by incorporating sequence-coupled effects into general PseAAC. *Oncotarget*, 7(2016), 44310-44321.
- [78] W. R. Qiu, B. Q. Sun, X. Xiao, Z. C. Xu, K. C. Chou. (2016). iPTM-mLys: identifying multiple lysine PTM sites and their different types. *Bioinformatics*, 32(2016), 3116-3123.
- [79] W. R. Qiu, X. Xiao, Z. C. Xu, K. C. Chou. (2016). iPhos-PseEn: identifying phosphorylation sites in proteins by fusing different pseudo components into an ensemble classifier. *Oncotarget*, 7(2016), 51270-51283.
- [80] X. Xiao, H. X. Ye, Z. Liu, J. H. Jia, K. C. Chou. (2016). iROS-gPseKNC: predicting replication origin sites in DNA by incorporating dinucleotide position-specific propensity into general pseudo nucleotide composition. *Oncotarget*, 7(2016), 34180-34189.

- [81] W. R. Qiu, B. Q. Sun, X. Xiao, D. Xu, K. C. Chou. (2017). iPhos-PseEvo: Identifying human phosphorylated proteins by incorporating evolutionary information into general PseAAC via grey system theory. *Molecular Informatics*, 36(2017), UNSP 1600010.
- [82] B. Liu, F. Yang, D. S. Huang, K. C. Chou. (2018). iPromoter-2L: a two-layer predictor for identifying promoters and their types by multi-window-based PseKNC. *Bioinformatics*, 34(2018), 33-40.
- [83] J. Jia, X. Li, W. Qiu, X. Xiao, K. C. Chou. (2019). iPPI-PseAAC(CGR): Identify protein-protein interactions by incorporating chaos game representation into PseAAC. *Journal of Theoretical Biology*, 460(2019), 195-203.
- [84] J. Chen, H. Liu, J. Yang, K. C. Chou. (2007). Prediction of linear B-cell epitopes using amino acid pair antigenicity scale. *Amino Acids*, 33(2007), 423-428.
- [85] K. C. Chou. (2001). Prediction of signal peptides using scaled window, *Peptides*, 22(2001), 1973-1979.
- [86] Y. Xu, X. J. Shao, L. Y. Wu, N. Y. Deng, K. C. Chou. (2013). iSNO-AAPair: incorporating amino acid pairwise coupling into PseAAC for predicting cysteine S-nitrosylation sites in proteins, *PeerJ*, 1(2013), e171.
- [87] Y. Xu, X. Wen, X. J. Shao, N. Y. Deng, K. C. Chou. (2014). iHyd-PseAAC: Predicting hydroxyproline and hydroxylysine in proteins by incorporating dipeptide position-specific propensity into pseudo amino acid composition, *Int. J. Mol. Sci.*, 15(2014), 7594-7610.
- [88] W. R. Qiu, X. Xiao, W. Z. Lin, K. C. Chou. (2014). iMethyl-PseAAC: Identification of Protein Methylation Sites via a Pseudo Amino Acid Composition Approach, *Biomed Res Int (BMRI)*, 2014(2014), 947416.
- [89] H. Lin, E. Z. Deng, H. Ding, W. Chen, K. C. Chou. (2014). iPro54-PseKNC: a sequence-based predictor for identifying sigma-54 promoters in prokaryote with pseudo k-tuple nucleotide composition, *Nucleic Acids Res.*, 42 (2014) 12961-12972.
- [90] Y. Xu, X. Wen, L. S. Wen, L. Y. Wu, N. Y. Deng, K. C. Chou. (2014). iNitro-Tyr: Prediction of nitrotyrosine sites in proteins with general pseudo amino acid composition, *PLoS ONE*, 9(2014), e105018.
- [91] X. Xiao, J. L. Min, W. Z. Lin, Z. Liu, X. Cheng, K. C. Chou. (2015). iDrug-Target: predicting the interactions between drug compounds and target proteins in cellular networking via the benchmark dataset optimization approach, *J Biomol Struct Dyn (JBSD)*, 33(2015), 2221-2233.
- [92] C. J. Zhang, H. Tang, W. C. Li, H. Lin, W. Chen, K. C. Chou. (2016). iOri-Human: identify human origin of replication by incorporating dinucleotide physicochemical properties into pseudo nucleotide composition, *Oncotarget*, 7(2016), 69783-69793.
- [93] Z. Liu, X. Xiao, D.J. Yu, J. Jia, W.R. Qiu, K.C. Chou, pRNAm-PC: Predicting N-methyladenosine sites in RNA sequences via physical-chemical properties, *Anal. Biochem.*, 497 (2016), 60-67.
- [94] W. Chen, H. Ding, P. Feng, H. Lin, K.C. Chou, iACP: a sequence-based tool for identifying anticancer peptides, *Oncotarget*, 7 (2016), 16895-16909.
- [95] J. Jia, Z. Liu, X. Xiao, B. Liu, K. C. Chou. (2016). iPPBS-Opt: A Sequence-Based Ensemble Classifier for Identifying Protein-Protein Binding Sites by Optimizing Imbalanced Training Datasets, *Molecules*, 21(2016), E95.
- [96] W. R. Qiu, S. Y. Jiang, B. Q. Sun, X. Xiao, X. Cheng, K. C. Chou. (2017). iRNA-2methyl: identify RNA 2'-O-methylation sites by incorporating sequence-coupled effects into general PseKNC and ensemble classifier, *Medicinal Chemistry*, 13(2017), 734-743.
- [97] W. Chen, P. Feng, H. Yang, H. Ding, H. Lin, K. C. Chou. (2017). iRNA-AI: identifying the adenosine to inosine editing sites in RNA sequences, *Oncotarget*, 8(2017), 4208-4217.
- [98] P. K. Meher, T. K. Sahu, V. Saini, A. R. Rao. (2017). Predicting antimicrobial peptides with improved accuracy by incorporating the compositional, physico-chemical and structural features into Chou's general PseAAC, *Sci Rep*, 7(2017), 42362.
- [99] J. Jia, L. Zhang, Z. Liu, X. Xiao, K. C. Chou. (2016). pSumo-CD: Predicting sumoylation sites in proteins with covariance discriminant algorithm by incorporating sequence-coupled effects into general PseAAC, *Bioinformatics*, 32(2016), 3133-3141.
- [100] A. Ehsan, K. Mahmood, Y. D. Khan, S. A. Khan, K. C. Chou. (2018). A Novel Modeling in Mathematical Biology for Classification of Signal Peptides, *Scientific Reports*, 8(2018), 1039.
- [101] J. Wang, J. Li, B. Yang, R. Xie, T. T. Marquez-Lago, A. Leier, M. Hayashida, T. Akutsu, Y. Zhang, K. C. Chou, J. Selkrig, T. Zhou, J. Song, T. Lithgow. (2018). Bastion3: a two-layer approach for identifying type III secreted ef-

- factors using ensemble learning, *Bioinformatics*, 35 (2018), 2017-2028.
- [102] F. Li, C. Li, T. T. Marquez-Lago, A. Leier, T. Akutsu, A. W. Purcell, A. I. Smith, T. Lightow, R. J. Daly, J. Song, K. C. Chou. (2018). Quokka: a comprehensive tool for rapid and accurate prediction of kinase family-specific phosphorylation sites in the human proteome, *Bioinformatics*, 34(2018), 4223-4231.
- [103] F. Li, Y. Wang, C. Li, T. T. Marquez-Lago, A. Leier, N. D. Rawlings, G. Haffari, J. Revote, T. Akutsu, K. C. Chou, A. W. Purcell, R. N. Pike, G. I. Webb, A. Ian Smith, T. Lithgow, R. J. Daly, J. C. Whisstock, J. Song. (2018). Twenty years of bioinformatics research for protease-specific substrate and cleavage site prediction: a comprehensive revisit and benchmarking of existing methods, *Brief in Bioinform*, doi:10.1093/bib/bby077 (2018).
- [104] J. Song, Y. Wang, F. Li, T. Akutsu, N. D. Rawlings, G. I. Webb, K. C. Chou. (2018). iProt-Sub: a comprehensive package for accurately mapping and predicting protease-specific substrates and cleavage sites, *Brief in Bioinform*, 20(2018), 638-658.
- [105] W. Chen, H. Ding, X. Zhou, H. Lin, K. C. Chou. (2018). iRNA(m6A)-PseDNC: Identifying N6-methyladenosine sites using pseudo dinucleotide composition, *Anal. Biochem.*, 561-562(2018), 59-65.
- [106] Y. D. Khan, M. Jamil, W. Hussain, N. Rasool, S. A. Khan, K. C. Chou. (2019). pSSbond-PseAAC: Prediction of disulfide bonding sites by integration of PseAAC and statistical moments, *J. Theor. Biol.*, 463(2019), 47-55.
- [107] X. Cheng, W. Z. Lin, X. Xiao, K. C. Chou. (2019). pLoc\_bal-mAnimal: predict subcellular localization of animal proteins by balancing training dataset and PseAAC, *Bioinformatics*, 35(2019), 398-406.
- [108] K. C. Chou. (2019). Advance in predicting subcellular localization of multi-label proteins and its implication for developing multi-target drugs, *Current Medicinal Chemistry*, 26(2019), 4918-4943.
- [109] K. C. Chou, X. Cheng, X. Xiao. (2019). pLoc\_bal-mEuk: predict subcellular localization of eukaryotic proteins by general PseAAC and quasi-balancing training dataset, *Med Chem*, 15(2019), 472-485.
- [110] A. Ehsan, M.K. Mahmood, Y. D. Khan, O. M. Barukab, S. A. Khan, K. C. Chou. (2019). iHyd-PseAAC (EPSV): Identify hydroxylation sites in proteins by extracting enhanced position and sequence variant feature via Chou's 5-step rule and general pseudo amino acid composition, *Current Genomics*, 20(2019), 124-133.
- [111] P. Feng, H. Yang, H. Ding, H. Lin, W. Chen, K. C. Chou. (2019). iDNA6mA-PseKNC: Identifying DNA N(6)-methyladenosine sites by incorporating nucleotide physicochemical properties into PseKNC, *Genomics*, 111(2019), 96-102.
- [112] W. Hussain, S. D. Khan, N. Rasool, S. A. Khan, K. C. Chou. (2019). SPalmitoylC-PseAAC: A sequence-based model developed via Chou's 5-steps rule and general PseAAC for identifying S-palmitoylation sites in proteins, *Anal. Biochem.*, 568(2019), 14-23.
- [113] W. Hussain, Y. D. Khan, N. Rasool, S. A. Khan, K. C. Chou. (2019). SPrenylC-PseAAC: A sequence-based model developed via Chou's 5-steps rule and general PseAAC for identifying S-prenylation sites in proteins, *J. Theor. Biol.*, 468(2019), 1-11.
- [114] Y. Lu, S. Wang, J. Wang, G. Zhou, Q. Zhang, X. Zhou, B. Niu, Q. Chen, K. C. Chou. (2019). An Epidemic Avian Influenza Prediction Model Based on Google Trends, *Letters in Organic Chemistry*, 16(2019), 303-310.
- [115] H. B. Shen, K. C. Chou. (2010). Gneg-mPLoc: A top-down strategy to enhance the quality of predicting subcellular localization of Gram-negative bacterial proteins, *Journal of Theoretical Biology*, 264(2010), 326-333.
- [116] K. C. Chou, H. B. Shen. (2010). Plant-mPLoc: A Top-Down Strategy to Augment the Power for Predicting Plant Protein Subcellular Localization, *PLoS ONE*, 5(2010), e11335.
- [117] H. B. Shen, K. C. Chou. (2009). A top-down approach to enhance the power of predicting human protein subcellular localization: Hum-mPLoc 2.0, *Anal. Biochem.*, 394(2009), 269-274.
- [118] Z. C. Wu, X. Xiao, K. C. Chou. (2012). iLoc-Gpos: A Multi-Layer Classifier for Predicting the Subcellular Localization of Singleplex and Multiplex Gram-Positive Bacterial Proteins, *Protein & Peptide Letters*, 19(2012), 4-14.
- [119] W. Z. Lin, J. A. Fang, X. Xiao, K. C. Chou. (2013). iLoc-Animal: A multi-label learning classifier for predicting subcellular localization of animal proteins *Molecular BioSystems*, 9(2013), 634-644.
- [120] Z. C. Wu, X. Xiao, K. C. Chou. (2011). iLoc-Plant: a multi-label classifier for predicting the subcellular localization of plant proteins with both single and multiple sites, *Molecular BioSystems*, 7(2011), 3287-3297.
- [121] X. Cheng, X. Xiao, K. C. Chou. (2017). pLoc-mPlant: predict subcellular localization of multi-location plant proteins via incorporating the optimal GO information into general PseAAC, *Molecular BioSystems*, 13(2017),

1722-1727.

- [122] X. Cheng, X. Xiao, K. C. Chou. (2017). pLoc-mVirus: predict subcellular localization of multi-location virus proteins via incorporating the optimal GO information into general PseAAC, *Gene (Erratum: ibid., 2018, Vol.644, 156-156)*, 628(2017), 315-321.
- [123] X. Cheng, S. G. Zhao, W. Z. Lin, X. Xiao, K. C. Chou. (2017). pLoc-mAnimal: predict subcellular localization of animal proteins with both single and multiple sites, *Bioinformatics*, 33(2017), 3524-3531.
- [124] X. Xiao, X. Cheng, S. Su, Q. Nao, K. C. Chou. (2017). pLoc-mGpos: Incorporate key gene ontology information into general PseAAC for predicting subcellular localization of Gram-positive bacterial proteins, *Natural Science*, 9(2017), 330-349.
- [125] X. Cheng, X. Xiao, K. C. Chou. (2018). pLoc-mEuk: Predict subcellular localization of multi-label eukaryotic proteins by extracting the key GO information into general PseAAC, *Genomics*, 110(2018), 50-58.
- [126] X. Cheng, X. Xiao, K. C. Chou. (2018). pLoc-mGneg: Predict subcellular localization of Gram-negative bacterial proteins by deep gene ontology learning via general PseAAC, *Genomics*, 110(2018), 231-239.
- [127] X. Cheng, X. Xiao, K. C. Chou. (2018). pLoc-mHum: predict subcellular localization of multi-location human proteins via general PseAAC to winnow out the crucial GO information, *Bioinformatics*, 34(2018), 1448-1456.
- [128] E. Pacharawongsakda, T. Theeramunkong. (2013). Predict Subcellular Locations of Singleplex and Multiplex Proteins by Semi-Supervised Learning and Dimension-Reducing General Mode of Chou's PseAAC, *IEEE Transactions on Nanobioscience*, 12(2013), 311-320.
- [129] J. Z. Cao, W. Q. Liu, H. Gu. (2012). Predicting Viral Protein Subcellular Localization with Chou's Pseudo Amino Acid Composition and Imbalance-Weighted Multi-Label K-Nearest Neighbor Algorithm, *Protein and Peptide Letters*, 19(2012), 1163-1169.
- [130] L. Q. Li, Y. Zhang, L. Y. Zou, Y. Zhou, X. Q. Zheng. (2012). Prediction of Protein Subcellular Multi-Localization Based on the General form of Chou's Pseudo Amino Acid Composition, *Protein & Peptide Letters*, 19(2012), 375-387.
- [131] S. Mei. (2012). Predicting plant protein subcellular multi-localization by Chou's PseAAC formulation based multi-label homolog knowledge transfer learning, *J. Theor. Biol.*, 310(2012), 80-87.
- [132] C. Huang, J. Yuan. (2013). Using radial basis function on the general form of Chou's pseudo amino acid composition and PSSM to predict subcellular locations of proteins with both single and multiple sites, *Biosystems*, 113(2013), 50-57.
- [133] X. Wang, G. Z. Li, W. C. Lu. (2013). Virus-ECC-mPLoc: a multi-label predictor for predicting the subcellular localization of virus proteins with both single and multiple sites based on a general form of Chou's pseudo amino acid composition, *Protein & Peptide Letters*, 20(2013), 309-317.
- [134] M. Mandal, A. Mukhopadhyay, U. Maulik. (2015). Prediction of protein subcellular localization by incorporating multiobjective PSO-based feature subset selection into the general form of Chou's PseAAC, *Medical & biological engineering & computing*, 53(2015), 331-344.
- [135] X. Cheng, S. G. Zhao, X. Xiao, K. C. Chou. (2017). iATC-mISF: a multi-label classifier for predicting the classes of anatomical therapeutic chemicals, *Bioinformatics (Corrigendum, ibid., 2017, Vol. 33, 2610)*, 33(2017), 341-346.
- [136] X. Cheng, S. G. Zhao, X. Xiao, K. C. Chou. (2017). iATC-mHyb: a hybrid multi-label classifier for predicting the classification of anatomical therapeutic chemicals, *Oncotarget*, 8(2017), 58494-58503.
- [137] X. Xiao, P. Wang, W. Z. Lin, J. H. Jia, K. C. Chou. (2013). iAMP-2L: A two-level multi-label classifier for identifying antimicrobial peptides and their functional types, *Anal. Biochem.*, 436(2013), 168-177.
- [138] K. C. Chou. (2013). Some remarks on predicting multi-label attributes in molecular biosystems, *Molecular Biosystems*, 9(2013), 1092-1100.
- [139] K. C. Chou, C. T. Zhang. (1995). Review: Prediction of protein structural classes, *Crit. Rev. Biochem. Mol. Biol.*, 30(1995), 275-349.
- [140] H. Mohabatkar. (2010). Prediction of cyclin proteins using Chou's pseudo amino acid composition, *Protein & Peptide Letters*, 17(2010), 1207-1214.
- [141] G. P. Zhou, K. Doctor. (2003). Subcellular location prediction of apoptosis proteins, *Proteins: Struct., Funct., Ge-*

- net., 50(2003), 44-48.
- [142] S. S. Sahu, G. Panda. (2010). A novel feature representation method based on Chou's pseudo amino acid composition for protein structural class prediction, *Computational Biology and Chemistry*, 34(2010), 320-327.
- [143] Zia-ur-Rehman, A. Khan. (2012). Identifying GPCRs and their Types with Chou's Pseudo Amino Acid Composition: An Approach from Multi-scale Energy Representation and Position Specific Scoring Matrix, *Protein & Peptide Letters*, 19(2012), 890-903.
- [144] G. L. Fan, Q. Z. Li. (2013). Discriminating bioluminescent proteins by incorporating average chemical shift and evolutionary information into the general form of Chou's pseudo amino acid composition, *J. Theor. Biol.*, 334(2013), 45-51.
- [145] C. Huang, J. Q. Yuan. (2013). Predicting protein subchloroplast locations with both single and multiple sites via three different modes of Chou's pseudo amino acid compositions, *J. Theor. Biol.*, 335(2013), 205-212.
- [146] Z. Hajisharifi, M. Piryaiee, M. Mohammad Beigi, M. Behbahani, H. Mohabatkar. (2014). Predicting anticancer peptides with Chou's pseudo amino acid composition and investigating their mutagenicity via Ames test, *J. Theor. Biol.*, 341(2014), 34-40.
- [147] K. C. Chou, H. B. Shen. (2009). Recent advances in developing web-servers for predicting protein attributes *Natural Science*, 1(2009), 63-92.
- [148] H. B. Shen, K. C. Chou. (2008). HIVcleave: a web-server for predicting HIV protease cleavage sites in proteins, *Anal. Biochem.*, 375(2008), 388-390.
- [149] W. R. Qiu, X. Xiao, K. C. Chou. (2014). iRSpot-TNCPseAAC: Identify recombination spots with trinucleotide composition and pseudo amino acid components, *Int J Mol Sci (IJMS)*, 15(2014), 1746-1766.
- [150] B. Liu, L. Fang, F. Liu, X. Wang, J. Chen, K. C. Chou. (2015). Identification of real microRNA precursors with a pseudo structure status composition approach, *PLoS ONE*, 10(2015), e0121501.
- [151] J. Wang, B. Yang, J. Revote, A. Leier, T. T. Marquez-Lago, G. Webb, J. Song, K. C. Chou, T. Lithgow. (2017). POSSUM: a bioinformatics toolkit for generating numerical sequence feature descriptors based on PSSM profiles, *Bioinformatics*, 33(2017), 2756-2758.
- [152] Z. Chen, P. Y. Zhao, F. Li, Leier A, T. T. Marquez-Lago, Y. Wang, G. I. Webb, A. I. Smith, R. J. Daly, K. C. Chou, J. Song. (2018). iFeature: a python package and web server for features extraction and selection from protein and peptide sequences, *Bioinformatics*, 34(2018), 2499-2502.
- [153] J. Song, F. Li, A. Leier, T. T. Marquez-Lago, T. Akutsu, G. Haffari, K. C. Chou, G. I. Webb, R. N. Pike. (2018). PROSPERous: high-throughput prediction of substrate cleavage sites for 90 proteases with improved accuracy, *Bioinformatics*, 34(2018), 684-687.
- [154] J. Song, F. Li, K. Takemoto, G. Haffari, T. Akutsu, K. C. Chou, G. I. Webb. (2018). PREvalL, an integrative approach for inferring catalytic residues using sequence, structural and network features in a machine learning framework, *Journal of Theoretical Biology*, 443(2018), 125-137.
- [155] W. R. Qiu, B. Q. Sun, X. Xiao, Z. C. Xu, J. H. Jia, K. C. Chou. (2018). iKcr-PseEns: Identify lysine crotonylation sites in histone proteins with pseudo components and ensemble classifier, *Genomics*, 110(2018), 239-246.
- [156] L. M. Liu, Y. Xu, K. C. Chou, iPGK-PseAAC: identify lysine phosphoglycerylation sites in proteins by incorporating four different tiers of amino acid pairwise coupling information into the general PseAAC, *Med Chem*, 13(2017), 552-559.
- [157] W. R. Qiu, S. Y. Jiang, Z. C. Xu, X. Xiao, K. C. Chou. (2017). iRNAm5C-PseDNC: identifying RNA 5-methylcytosine sites by incorporating physical-chemical properties into pseudo dinucleotide composition, *Oncotarget*, 8(2017), 41178-41188.
- [158] J. Wang, B. Yang, A. Leier, T. T. Marquez-Lago, M. Hayashida, A. Rocker, Z. Yanju, T. Akutsu, K. C. Chou, R. A. Strugnell, J. Song, T. Lithgow. (2018). Bastion6: a bioinformatics approach for accurate prediction of type VI secreted effectors, *Bioinformatics*, 34(2018), 2546-2555.
- [159] Y. Xu, C. Li, K. C. Chou. (2017). iPreny-PseAAC: identify C-terminal cysteine prenylation sites in proteins by incorporating two tiers of sequence couplings into PseAAC, *Med Chem*, 13(2017), 544-551.
- [160] B. Liu, H. Wu, D. Zhang, X. Wang, K. C. Chou. (2017). Pse-Analysis: a python package for DNA/RNA and protein/peptide sequence analysis based on pseudo components and kernel methods, *Oncotarget*, 8(2017),

13338-13343.

- [161] K. C. Chou, S. P. Jiang, W. M. Liu, C. H. Fee. (1979). Graph theory of enzyme kinetics: 1. Steady-state reaction system, *Scientia Sinica*, 22(1979), 341-358.
- [162] K. C. Chou, S. Forsen. (1980). Graphical rules for enzyme-catalyzed rate laws, *Biochem. J.*, 187(1980), 829-835.
- [163] K. C. Chou, S. Forsen, G. Q. Zhou. (1980). Three schematic rules for deriving apparent rate constants, *Chemica Scripta*, 16(1980), 109-113.
- [164] K. C. Chou, R. E. Carter, S. Forsen. (1981). A new graphical method for deriving rate equations for complicated mechanisms, *Chemica Scripta*, 18(1981), 82-86.
- [165] K. C. Chou, S. Forsen. (1981). Graphical rules of steady-state reaction systems, *Can. J. Chem.*, 59(1981), 737-755.
- [166] G. P. Zhou, M. H. Deng. (1984). An extension of Chou's graphic rules for deriving enzyme kinetic equations to systems involving parallel reaction pathways, *Biochem. J.*, 222(1984), 169-176.
- [167] K. C. Chou. (1989). Graphic rules in steady and non-steady enzyme kinetics, *J. Biol. Chem.*, 264(1989), 12074-12079.
- [168] I. W. Althaus, J. J. Chou, A. J. Gonzales, M. R. Diebel, K. C. Chou, F. J. Kezdy, D. L. Romero, P. A. Aristoff, W. G. Tarpley, F. Reusser. (1993). Steady-state kinetic studies with the non-nucleoside HIV-1 reverse transcriptase inhibitor U-87201E, *J. Biol. Chem.*, 268(1993), 6119-6124.
- [169] K. C. Chou. (1990). Review: Applications of graph theory to enzyme kinetics and protein folding kinetics. Steady and non-steady state systems, *Biophysical Chemistry*, 35(1990), 1-24.
- [170] I. W. Althaus, A. J. Gonzales, J. J. Chou, M. R. Diebel, K. C. Chou, F. J. Kezdy, D. L. Romero, P. A. Aristoff, W. G. Tarpley, F. Reusser. (1993). The quinoline U-78036 is a potent inhibitor of HIV-1 reverse transcriptase, *J. Biol. Chem.*, 268(1993), 14875-14880.
- [171] K. C. Chou. (2010). Graphic rule for drug metabolism systems, *Current Drug Metabolism*, 11(2010), 369-378.
- [172] G. P. Zhou. (2011). The disposition of the LZCC protein residues in wenxiang diagram provides new insights into the protein-protein interaction mechanism, *J. Theor. Biol.*, 284(2011), 142-148.
- [173] I. W. Althaus, J. J. Chou, A. J. Gonzales, M. R. Diebel, K. C. Chou, F. J. Kezdy, D. L. Romero, P. A. Aristoff, W. G. Tarpley, F. Reusser. (1993). Kinetic studies with the nonnucleoside HIV-1 reverse transcriptase inhibitor U-88204E, *Biochemistry*, 32(1993), 6548-6554.
- [174] K. C. Chou, W. Z. Lin, X. Xiao. (2011). Wenxiang: a web-server for drawing wenxiang diagrams *Natural Science*, 3(2011), 862-865
- [175] K. C. Chou, S. Forsen. (1980). Diffusion-controlled effects in reversible enzymatic fast reaction system: Critical spherical shell and proximity rate constants, *Biophysical Chemistry*, 12(1980), 255-263.
- [176] K. C. Chou, T. T. Li, S. Forsen. (1980). The critical spherical shell in enzymatic fast reaction systems, *Biophysical Chemistry*, 12(1980), 265-269.
- [177] H. B. Shen, J. N. Song, K. C. Chou. (2009). Prediction of protein folding rates from primary sequence by fusing multiple sequential features *Journal of Biomedical Science and Engineering (JBISE)*, 2(2009), 136-143.
- [178] K. C. Chou, N. Y. Chen, S. Forsen. (1981). The biological functions of low-frequency phonons: 2. Cooperative effects, *Chemica Scripta*, 18(1981), 126-132.
- [179] K. C. Chou. (1988). Review: Low-frequency collective motion in biomacromolecules and its biological functions, *Biophysical Chemistry*, 30(1988), 3-48.
- [180] K. C. Chou, F. J. Kezdy, F. Reusser. (1994). Review: Kinetics of processive nucleic acid polymerases and nucleases, *Anal. Biochem.*, 221(1994), 217-230.
- [181] X. Xiao, S. Shao, Y. Ding, Z. Huang, X. Chen, K. C. Chou. (2005). Using cellular automata to generate Image representation for biological sequences, *Amino Acids*, 28(2005), 29-35.
- [182] Z. C. Wu, X. Xiao, K. C. Chou. (2010). 2D-MH: A web-server for generating graphic representation of protein sequences based on the physicochemical properties of their constituent amino acids, *J. Theor. Biol.*, 267(2010), 29-34.