



DeepKhib: A Deep-Learning Framework for Lysine 2-Hydroxyisobutyrylation Sites Prediction

Luna Zhang^{1†}, Yang Zou^{2†}, Ningning He², Yu Chen¹, Zhen Chen^{3,4*} and Lei Li^{1,2*}

¹ School of Data Science and Software Engineering, Qingdao University, Qingdao, China, ² School of Basic Medicine, Qingdao University, Qingdao, China, ³ Collaborative Innovation Center of Henan Grain Crops, Henan Agricultural University, Zhengzhou, China, ⁴ Key Laboratory of Rice Biology in Henan Province, Henan Agricultural University, Zhengzhou, China

OPEN ACCESS

Edited by:

Jian Ren,
Sun Yat-sen University, China

Reviewed by:

Santosh Panjikar,
Australian Synchrotron, Australia
Annamaria Tonazzi,
National Research Council (CNR), Italy

*Correspondence:

Zhen Chen
chenzhen-win2009@163.com

Lei Li
leili@qdu.edu.cn;
lileime@hotmail.com

[†] These authors have contributed
equally to this work

Specialty section:

This article was submitted to
Cellular Biochemistry,
a section of the journal
Frontiers in Cell and Developmental
Biology

Received: 05 July 2020

Accepted: 17 August 2020

Published: 09 September 2020

Citation:

Zhang L, Zou Y, He N, Chen Y,
Chen Z and Li L (2020) DeepKhib:
A Deep-Learning Framework
for Lysine 2-Hydroxyisobutyrylation
Sites Prediction.
Front. Cell Dev. Biol. 8:580217.
doi: 10.3389/fcell.2020.580217

As a novel type of post-translational modification, lysine 2-Hydroxyisobutyrylation (K_{hib}) plays an important role in gene transcription and signal transduction. In order to understand its regulatory mechanism, the essential step is the recognition of K_{hib} sites. Thousands of K_{hib} sites have been experimentally verified across five different species. However, there are only a couple traditional machine-learning algorithms developed to predict K_{hib} sites for limited species, lacking a general prediction algorithm. We constructed a deep-learning algorithm based on convolutional neural network with the one-hot encoding approach, dubbed CNN_{OH} . It performs favorably to the traditional machine-learning models and other deep-learning models across different species, in terms of cross-validation and independent test. The area under the ROC curve (AUC) values for CNN_{OH} ranged from 0.82 to 0.87 for different organisms, which is superior to the currently available K_{hib} predictors. Moreover, we developed the general model based on the integrated data from multiple species and it showed great universality and effectiveness with the AUC values in the range of 0.79–0.87. Accordingly, we constructed the on-line prediction tool dubbed DeepKhib for easily identifying K_{hib} sites, which includes both species-specific and general models. DeepKhib is available at <http://www.bioinfo.org/DeepKhib>.

Keywords: post-translational modification, lysine 2-hydroxyisobutyrylation, deep learning, modification site prediction, machine learning

INTRODUCTION

Protein post-translational modification (PTM) is a key mechanism to regulate cellular functions through covalent modification and enzyme modification, which dynamically regulates a variety of biological events (Beltrao et al., 2013; Skelly et al., 2016). Recently, an evolutionarily conserved short-chain lysine acylation modification dubbed lysine 2-hydroxyisobutyrylation (K_{hib}) has been reported, which introduces a steric bulk with a mass shift of +86.03 Da (**Supplementary Figure 1A**) and neutralize the positive charge of lysine (Dai et al., 2014; Xiao et al., 2015). It involves various biological functions including biosynthesis of amino acids, starch biosynthesis, carbon metabolism, glycolysis / gluconeogenesis and transcription (Dai et al., 2014; Huang et al., 2017, 2018a; Li et al., 2017; Meng et al., 2017; Yu et al., 2017; Wu et al., 2018; Yin et al., 2019). For instance, the decrease

of this modification on K281 of glycolytic enzyme ENO1 reduces its catalytic activity (Huang et al., 2018b). The three-dimension structure of the peptide containing K281 in the center was shown as **Supplementary Figure 1B**.

Thousands of K_{hib} sites have been identified in different species including humans, plants and prokaryotes through large-scale experimental approaches (Dai et al., 2014; Huang et al., 2018a), which is summarized in **Supplementary Table 1**. The experimental methods, however, are time-consuming and expensive and thus the development of prediction algorithms *in silico* is necessary for the high-throughput recognition of K_{hib} sites. Two classifiers (i.e., iLys-Khib and Khibpred) have been reported for predicting the K_{hib} sites in a few species (Ju and Wang, 2019; Wang et al., 2020). As many different organisms have been investigated and the number of K_{hib} sites has increased, it is indispensable to compare the characteristics of this modification in different species and investigate whether it is suitable to develop a general model with high confidence. Additionally, the reported models were based on traditional machine-learning (ML) algorithms (e.g., Random Forest (RF)). Recently, the deep learning (DL) algorithms, as the modern ML architecture, have demonstrated superior prediction performance in the field of bioinformatics, such as the prediction of modification sites on DNA, RNA and proteins (Wang et al., 2017; Huang et al., 2018c; Long et al., 2018; Tahir et al., 2019; Tian et al., 2019). We have developed a few DL approaches for the prediction of PTM sites and they all demonstrate their superiority over conventional ML algorithms (Chen et al., 2018a, 2019; Zhao et al., 2020). Therefore, we attempted to compare the DL models with the traditional ML models for the prediction of K_{hib} sites.

In this study, we constructed a convolutional neural network (CNN)-based architecture with one-hot encoding approach, named as CNN_{OH} . This model performed favorably to the traditional ML models and other DL models across different species, in terms of cross-validation and independent test. It is also superior to the documented K_{hib} predictors. Furthermore, we constructed a general model based on the integrated data from multiple species and it demonstrated great generality and effectiveness. Finally, we shared both species-specific models and the general model as the on-line prediction tool DeepKhib for easily identifying K_{hib} sites.

MATERIALS AND METHODS

Dataset Collection

The experimentally identified K_{hib} sites from five different organisms including *Homo sapiens* (human), *Oryza sativa* (rice), *Physcomitrella patens* (moss) and two one-celled eukaryotes *Toxoplasma gondii* and *Saccharomyces cerevisiae*. The data of the species were pre-processed and the related procedure was exemplified using the human data, as listed below (**Supplementary Figure 2**).

We collected 12,166 K_{hib} sites from 3,055 human proteins (Wu et al., 2018). These proteins were classified into 2,466 clusters using CD-HIT with the threshold of 40% according to the previous studies (Li and Godzik, 2006; Huang et al., 2010).

In each cluster, the protein with the most K_{hib} sites was selected as the representative of the cluster. On the 2,466 representatives, 9,473 K_{hib} sites were considered positives whereas the remaining K sites were taken as negatives. We further estimated the potential redundancy of the positive sites by extracting the peptide segment of seven residues with the K_{hib} site in the center and count the number of unique segments (Chen et al., 2018a; Xie et al., 2018). The number (9,444) of the unique segments is 99.7% of the total segments, suggesting considerable diversity of the positive segments. The number of the negative sites (103,987) is 11 times larger than that of the positive sites. To avoid the potential impact of biased data on model construction, we referred to previous studies and balanced positives and negatives by randomly selecting the same number of negative sites (Huang et al., 2018c; Tahir et al., 2019). These positives and negatives composed the whole human dataset.

To determine the optimal sequence window for model construction, we tested different sequence window sizes ranging from 21 to 41, referring to the previous PTM studies where the optimal window sizes are between 31 and 39 (Wang et al., 2017; Chen et al., 2018a; Huang et al., 2018b). The window size of 37 corresponded to the largest area under the ROC curve (AUC) through 10-fold cross-validation (**Supplementary Figure 3**) and was therefore selected in this study. It should be noted that if the central lysine residue is located near the N-terminus or C-terminus of the protein sequence, the symbol "X" is added at the related terminus to ensure the same window size of the sequences.

Figure 1 showed the flowcharts for all the species. The dataset of each species was randomly separated into five groups of which four were used for 10-fold cross-validation and the rest for independent test. Each group contained the same number of positives and negatives. Specifically, the cross-validation datasets included 15,156/15,464/10,204/12,354 samples for *H. sapiens/T. gondii/O. sativa/P. patens*, respectively. Accordingly, the independent test sets comprised 3,790/3,866/2,552/3,090 samples for these organisms, separately. These datasets are available at <http://www.bioinfo.org/DeepKhib>.

Feature Encodings

The ZSCALE Encoding

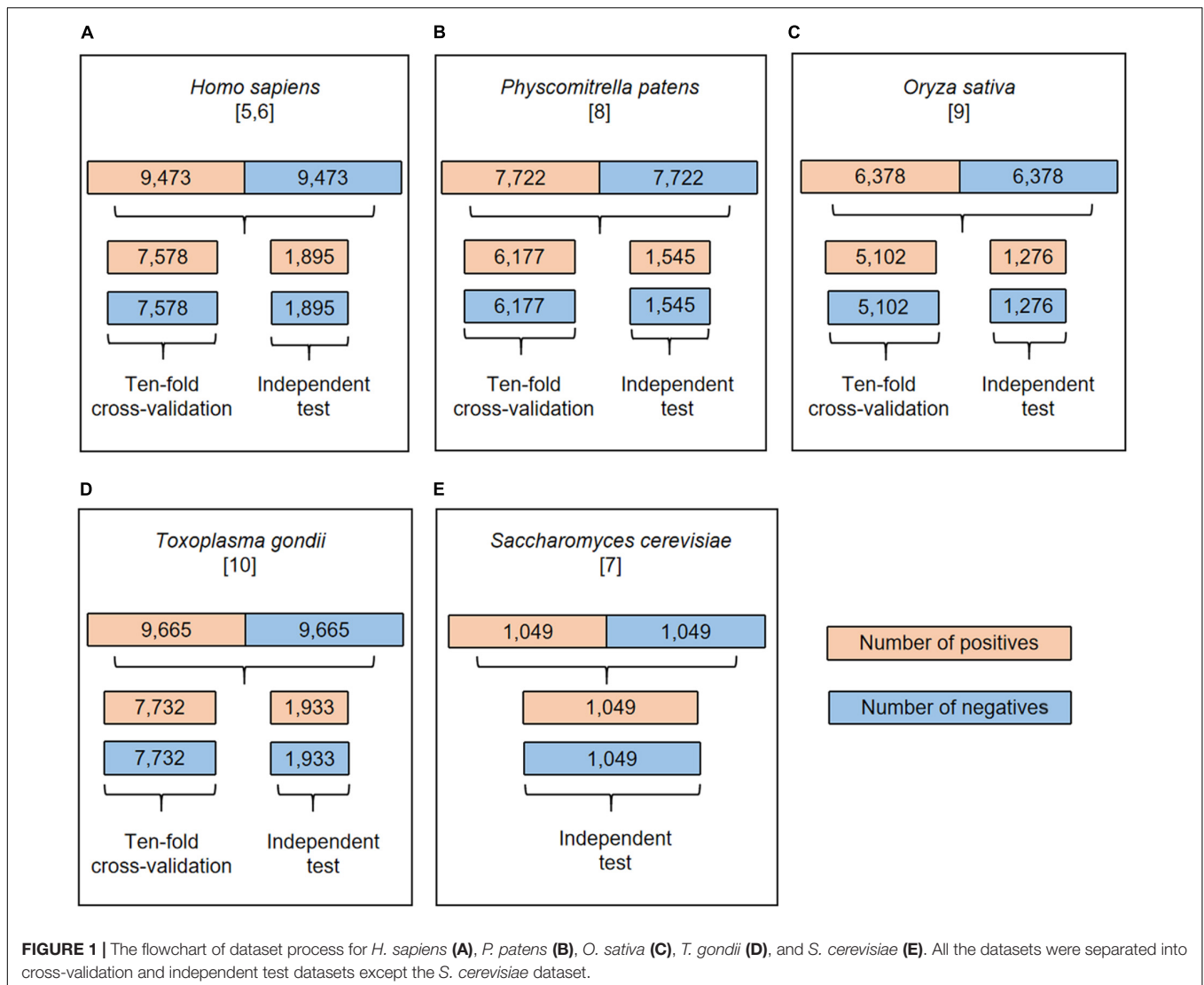
Each amino acid is characterized by five physiochemical descriptor variables (Sandberg et al., 1998; Chen et al., 2012).

The Encoding of Extended Amino Acid Composition (EAAC) Encoding

The EAAC encoding is based on the calculation of the amino acid composition (AAC) that indicates the amino acid frequencies for every position in the sequence window. EAAC is calculated by continuously sliding using a fixed-length sequence window (the default is 5) from the N-terminus to the C-terminus of each peptide (Chen et al., 2018b). The related formula is listed below:

$$f(t, win) = \frac{N(t, win)}{N(win)}, t \in \{A, C, D, \dots, Y\},$$

$$win \in \{window1, window2, \dots, window37\} \quad (1)$$



where $N(t, win)$ is the number of amino acid t in the sliding window win , and $N(win)$ is the size of the sliding window win .

The Enhanced Grouped Amino Acids Content (EGAAC) Encoding

The EGAAC feature (Zhao et al., 2020) is developed based on the grouped amino acids content (GAAC) feature (Chen et al., 2018b, 2020). In the GAAC feature, the 20 amino acid types are categorized into five groups (g1: GAVLMI, g2: FYW, g3: KRH, g4: DE and g5: STCPNQ) according to their physicochemical properties and the frequencies of the groups are calculated for every position in the sequence window. For the EGAAC feature, the GAAC values are calculated in the window of fixed length (the default as 5) continuously sliding from the N- to C-terminal of each peptide sequence.

The One-Hot Encoding

The one-hot encoding is represented by the conversion of the 20 types of amino acids to 20 binary bits. By considering the

complemented symbol “X,” a vector of size $(20+1)$ bits is used to represent a single position in the peptide sequence. For example, the amino acid “A” is represented by “10000000000000000000,” “Y” is represented by “000000000000000000010,” and the symbol “X” is represented by “000000000000000000001.”

Architecture of the Machine-Learning Models

The CNN Model With One-Hot Encoding

The CNN algorithm (Fukushima, 1980) decomposes an overall pattern into many sub-patterns (features) through a neurocognitive machine, and then enters the hierarchically connected feature plane for processing. The architecture of the CNN model with one-hot encoding (called as CNN_{OH}) contained four layers as follows (Figure 2A).

- (i) The first layer was the input layer where peptide sequences were represented using the one-hot encoding approach.

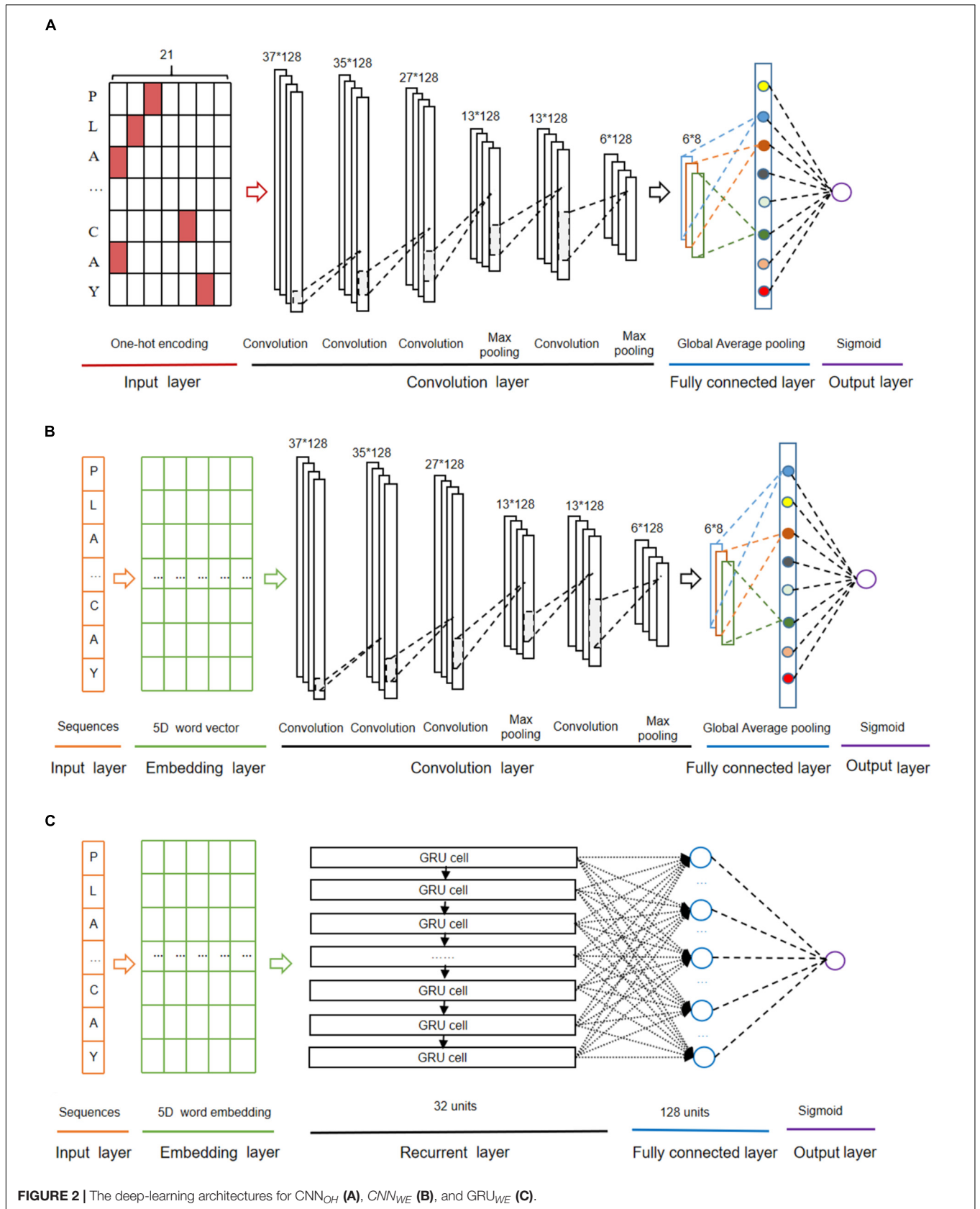
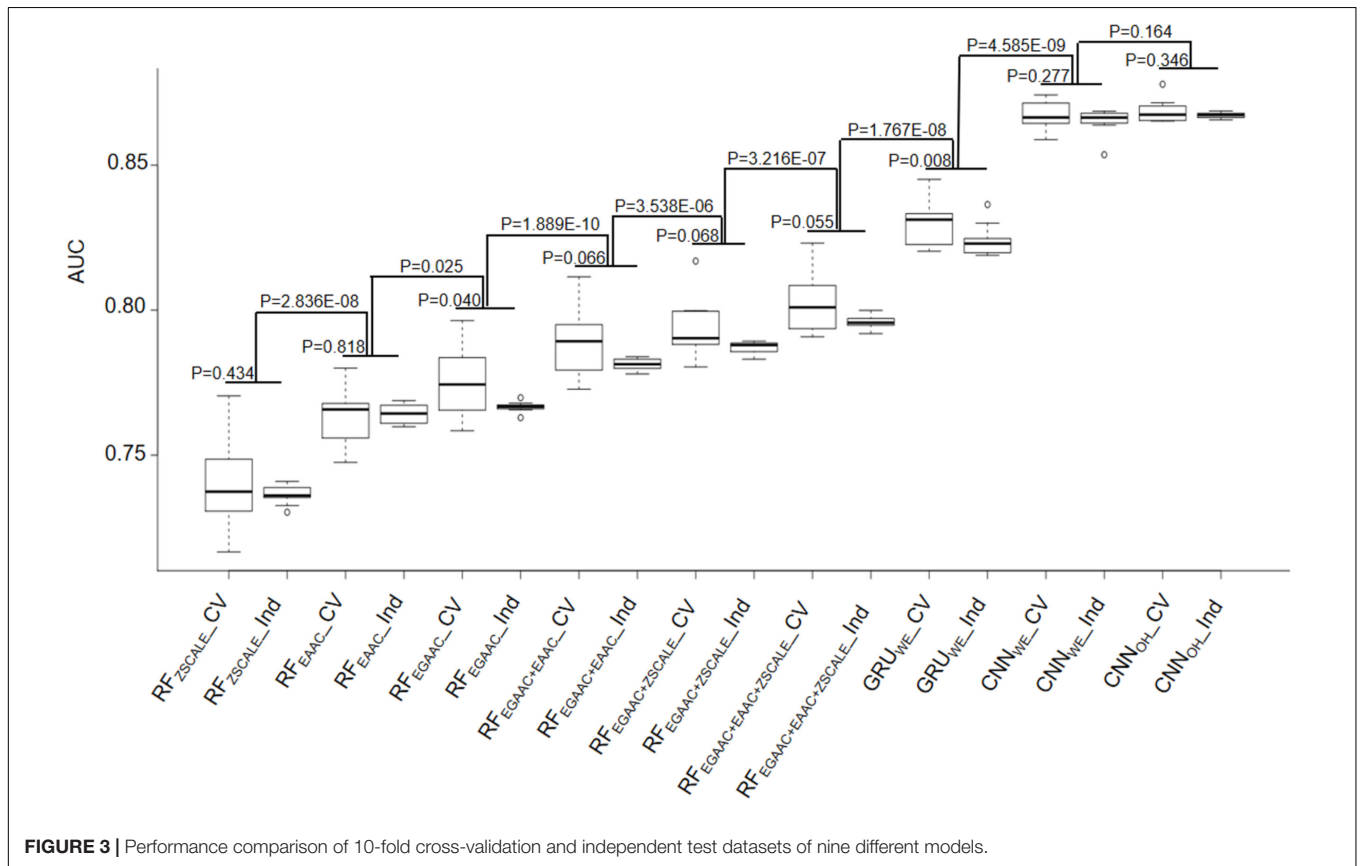


TABLE 1 | Performances comparison of the different classifiers for human K_{hib} prediction.

	Classifier	Sn	Sp	Acc	MCC	AUC
10-fold cross-validation	RF _{EGAAC}	0.727 ± 0.015	0.682 ± 0.017	0.704 ± 0.011	0.409 ± 0.022	0.775 ± 0.011
	RF _{EAAc}	0.744 ± 0.025	0.645 ± 0.023	0.695 ± 0.010	0.391 ± 0.020	0.763 ± 0.008
	RF _{ZSCALE}	0.681 ± 0.016	0.662 ± 0.018	0.672 ± 0.011	0.344 ± 0.023	0.740 ± 0.014
	RF _{EGAAC+EAAC}	0.748 ± 0.019	0.691 ± 0.023	0.719 ± 0.012	0.439 ± 0.025	0.789 ± 0.011
	RF _{EGAAC+ZSCALE}	0.726 ± 0.019	0.707 ± 0.015	0.716 ± 0.012	0.433 ± 0.025	0.794 ± 0.010
	RF _{EGAAC+EAAC+ZSCALE}	0.751 ± 0.016	0.702 ± 0.022	0.727 ± 0.013	0.454 ± 0.026	0.802 ± 0.010
	GRU _{WE}	0.821 ± 0.024	0.683 ± 0.033	0.752 ± 0.009	0.509 ± 0.018	0.830 ± 0.007
	CNN _{WE}	0.849 ± 0.035	0.722 ± 0.042	0.786 ± 0.007	0.578 ± 0.012	0.867 ± 0.005
	CNN _{OH}	0.876 ± 0.025	0.700 ± 0.026	0.788 ± 0.007	0.586 ± 0.014	0.868 ± 0.004
	Independent test	RF _{EGAAC}	0.719 ± 0.006	0.676 ± 0.007	0.698 ± 0.002	0.395 ± 0.004
RF _{EAAc}		0.755 ± 0.003	0.638 ± 0.007	0.697 ± 0.003	0.396 ± 0.006	0.764 ± 0.003
RF _{ZSCALE}		0.680 ± 0.008	0.658 ± 0.009	0.669 ± 0.005	0.337 ± 0.011	0.736 ± 0.003
RF _{EGAAC+EAAC}		0.740 ± 0.006	0.678 ± 0.005	0.709 ± 0.002	0.419 ± 0.005	0.781 ± 0.002
RF _{EGAAC+ZSCALE}		0.728 ± 0.006	0.692 ± 0.006	0.710 ± 0.002	0.420 ± 0.005	0.787 ± 0.002
RF _{EGAAC+EAAC+ZSCALE}		0.752 ± 0.005	0.693 ± 0.004	0.723 ± 0.002	0.446 ± 0.005	0.796 ± 0.002
GRU _{WE}		0.806 ± 0.015	0.692 ± 0.029	0.749 ± 0.004	0.501 ± 0.007	0.824 ± 0.005
CNN _{WE}		0.846 ± 0.035	0.719 ± 0.042	0.783 ± 0.006	0.572 ± 0.009	0.865 ± 0.004
CNN _{OH}		0.874 ± 0.026	0.690 ± 0.035	0.782 ± 0.005	0.575 ± 0.005	0.871 ± 0.001

The data sets for 10-fold cross-validation and an independent test were described in the section "Materials and Methods." The RF classifier with the different encoding approach was named as RF_{EGAAC}, RF_{EAAc}, RF_{ZSCALE}, RF_{EGAAC+EAAC}, RF_{EGAAC+ZSCALE}, and RF_{EGAAC+EAAC+ZSCALE}. The RNN/CNN classifier with the word embedding encoding approach was named as GRU_{WE} /CNN_{WE}, respectively. The CNN classifier with one-hot encoding was named as CNN_{OH}. Ten models were constructed in the 10-fold cross validation and evaluated using the ten different validation datasets and the same independent dataset. Accordingly, the value Sn, Sp, Acc, MCC, and AUC were represented by average ± standard deviation.



- (ii) The second layer was the convolution layer that consisted of four convolution sublayers and two max pooling sublayers. The convolution sublayers, each sublayer uses 128 convolution filters, the length of which are 1, 3, 9, and 10, respectively. The two max pooling sublayers followed the third and fourth convolution sublayers, individually.
- (iii) The third layer contained the fully connected sublayer, which contained a fully connected sublayer with eight neuron units without flattening, and a global average pooling sublayer, which was adopted to correlate the feature mapping with category output in order to reduce training parameters and avoid over-fitting.
- (iv) The last layer was the output layer that included a single unit outputting the probability score of the modification, calculated using the "Sigmoid" function. If the probability score is greater than a specified threshold (e.g., 0.5), the peptide is predicted to be positive.

The "ReLU" function (Hahnloser et al., 2000) was used as the activation function of the convolution sublayers and fully connected sublayers of the above layers to avoid gradient dispersion in the training process. The Adam optimizer (Kingma and Jimmy, 2014) was used to optimize the hyper-parameters of this model, which include batch size, maximum epoch, learning rate and dropout rate. The maximum training period was set as 1000 epochs to ensure the convergence of the loss function values. In each epoch, the training data set was separated and iterated in a batch size of 1024. To avoid over-fitting, the dropout of neurons units in each convolution sublayer of the second layer was set 70% and that in the full connection sublayer of the third layer was set 30% (Nitish et al., 2014), the early stop strategy was adopted and the best model was saved.

The CNN Algorithm With Word Embedding

The CNN algorithm with word embedding (CNN_{WE}) contained five layers (Figure 2B). The input layer receives the sequence of window size 37 and each residue is transformed into a five-dimensional word vector in the embedding layer. The rest layers are the same as the corresponding layers in CNN_{OH} .

The GRU Algorithm With Word Embedding

The GRU algorithm (Cho et al., 2014) includes an update gate and a reset gate. The former is used to control the extent to which the state information at the previous moment is brought into the current state, whereas the latter is used to control the extent to which the state information at the previous moment is ignored. The GRU algorithm with word embedding (GRU_{WE}) contained

five layers (Figure 2C). The first, the second and the last layers are the same as the corresponding layers in CNN_{WE} . The third layer is the recurrent layer where each word vector from the previous layer was sequentially inputted into the related GRU unit that contains 32 hidden neuron units. The fourth layer was the fully connected layer that contains 128 neuron units with "ReLU" as the activation function.

The RF Algorithms With Different Features

The Random Forest algorithm (Breiman, 2001) contains multiple decision trees, which remain unchanged under the scaling of feature values and various other transformations, and the output category is determined by the mode of the category output by the individual tree. Each tree depends on the values of a random vector sampled independently with the same distribution for all trees in the forest. The number of decision trees was set 140. This classifier was developed based on the Python module "sklearn."

Cross-Validation and Performance Evaluation

To evaluate the performance of K_{hib} sites prediction, we adopted four statistical measurement methods. They included sensitivity (Sn), specificity (Sp), accuracy (ACC), and Matthew's correlation coefficient (MCC), listed as follows:

$$Sn = \frac{TP}{TP + FN} \quad (2)$$

$$Sp = \frac{TN}{TN + FP} \quad (3)$$

$$Acc = \frac{TP + TN}{TP + FP + TN + FN} \quad (4)$$

$$MCC = \frac{TP \times TN - FN \times FP}{\sqrt{(TP + FN) \times (TN + FP) \times (TP + FP) \times (TN + FN)}} \quad (5)$$

In the above equations, TP is true positives, FP is false positives, TN is true negatives, FN is false negatives. In addition, the area under the receiver operating characteristic (ROC) curve (AUC) values was calculated to evaluate the performance of the prediction model.

TABLE 2 | The AUC values of the CNN_{OH} model constructed for *O. sativa*, *P. patens*, *T. gondii*, and *H. sapiens*, respectively.

Species	10-fold cross-validation	Independent test
<i>O. sativa</i>	0.823	0.818
<i>P. patens</i>	0.830	0.831
<i>T. gondii</i>	0.862	0.865
<i>H. sapiens</i>	0.868	0.871

TABLE 3 | The AUC values of different CNN_{OH} models in terms of independent test for five distinct organisms.

Prediction models	Independent data sets				
	<i>O. sativa</i>	<i>P. patens</i>	<i>T. gondii</i>	<i>H. sapiens</i>	<i>S. cerevisiae</i>
<i>O. sativa</i>	0.818	0.788	0.782	0.803	0.721
<i>P. patens</i>	0.761	0.831	0.812	0.837	0.806
<i>T. gondii</i>	0.781	0.813	0.865	0.827	0.776
<i>H. sapiens</i>	0.778	0.818	0.832	0.871	0.785
General	0.802	0.840	0.860	0.868	0.789

The top two models with best performance are bold.

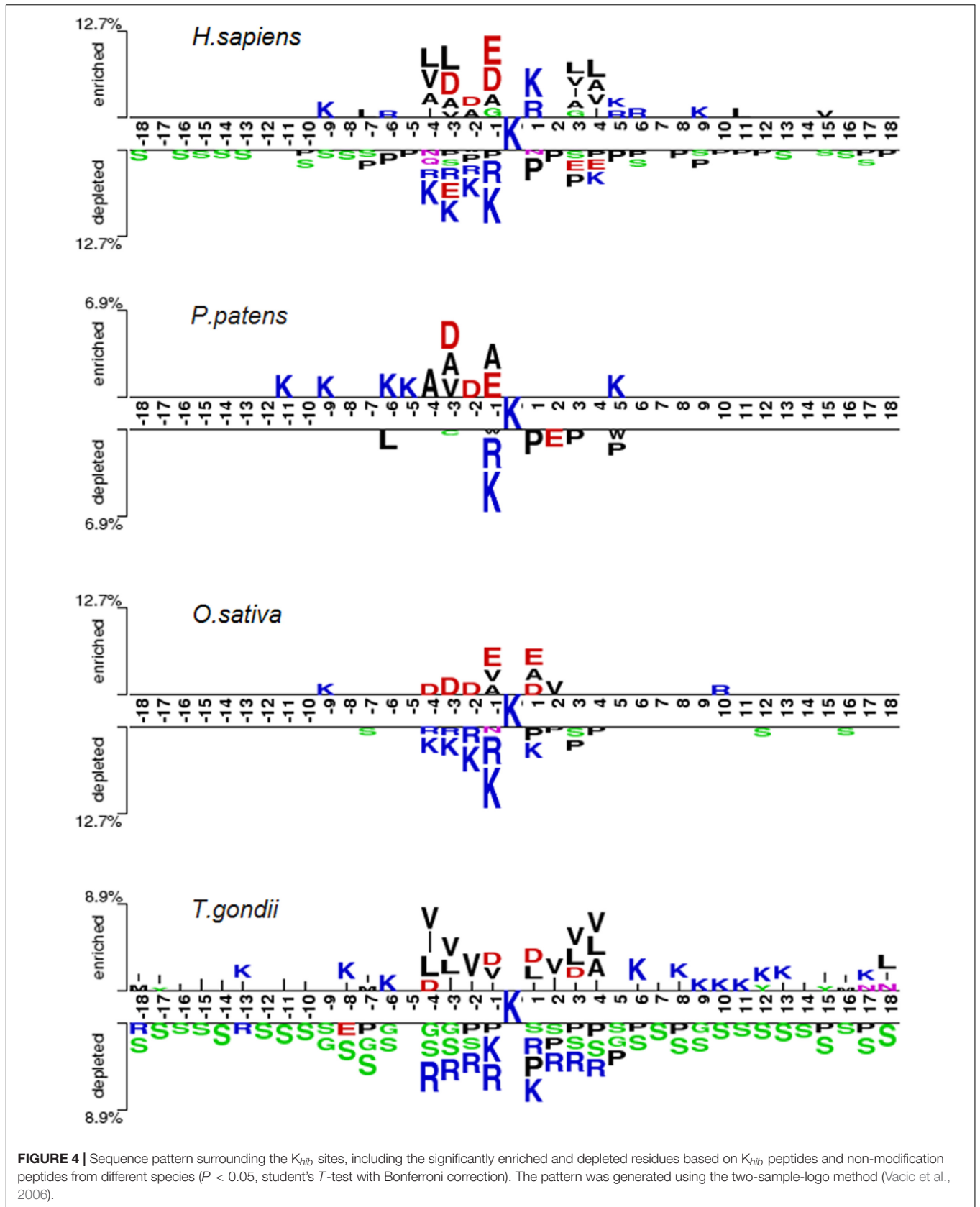


FIGURE 4 | Sequence pattern surrounding the K_{hib} sites, including the significantly enriched and depleted residues based on K_{hib} peptides and non-modification peptides from different species ($P < 0.05$, student's *T*-test with Bonferroni correction). The pattern was generated using the two-sample-logo method (Vacic et al., 2006).

Statistical Methods

The paired student's *t*-test was used to test the significant difference between the mean values of the two paired populations. As for multiple comparisons, the adjusted *P* value with the Benjamini-Hochberg (BH) method was adopted.

RESULTS AND DISCUSSION

A couple of computational approaches has been developed for the prediction of K_{hib} sites (Ju and Wang, 2019; Wang et al., 2020). Recently, this modification has been investigated across five different species, ranging from single-celled organisms to multiple-celled organisms and from plants to mammals. Additionally, the number of reported sites has been significantly increased. These raised our interest to develop novel prediction algorithms and explore the characteristics of this modification. We pre-processed the data from different species and separated them into the cross-validation dataset and the independent test set (see section "Materials and Methods" for detail; **Figure 1**). We first took the human data as the representative to compare different models and then applied the model with the best performance to other species. The human cross-validation dataset contained 15,156 samples and the independent test set covered 3,790 samples, in each of which half were positives and half were negatives.

CNN_{OH} Showed Superior Performance

We constructed nine models, divided into two categories: six traditional ML models and three DL models (see section "Materials and Methods" for details). The traditional ML models were based on the RF algorithm combined with different encoding schemes. The DL models included a Gated Recurrent Unit (GRU) model with the word-embedding encoding approach dubbed GRU_{WE} and two CNN models with the one-hot and word-embedding encoding approaches named CNN_{OH} and CNN_{WE}, respectively. Both encoding methods are common in the DL algorithms (Chen et al., 2018a; Xie et al., 2018).

The RF-based models were developed with different common encoding schemes, including EAAC, EGAAC and ZSCALE. Among these encoding schemes, EGAAC had the best performance followed by EAAC whereas ZSCALE was the worst in terms of AUC and ACC for both 10-fold cross-validation and the independent test (**Table 1** and **Figure 3**). For instance, EGAAC corresponded to the average AUC value as 0.775, EAAC had the value as 0.763 and ZSCALE had the value as 0.740 for cross validation. Because different encodings represent distinct characteristics of K_{hib} -containing peptides, we evaluated the combinations of the encoding schemes. The combinations showed better performances than individual scheme and the combination of all the three was the best for both cross-validation and the independent test, in terms of AUC, MCC, and ACC (**Table 1** and **Figure 3**). Therefore, the K_{hib} prediction accuracy could be improved by the integration of different encoding schemes.

As the DL algorithms showed superior to the traditional ML algorithms for a few PTM predictions in our previous studies (Chen et al., 2019; Zhao et al., 2020), we examined the DL

algorithms for the K_{hib} prediction. Traditionally, CNN is popular for image prediction with spatial invariant features while RNN is ideal for text prediction with sequence features. However, many cases demonstrate that CNN also has good performance when applied to sequence data (Sainath et al., 2013; Tahir et al., 2019). Accordingly, we developed both RNN and CNN models for the K_{hib} prediction with two common encoding approaches: one-hot and word-embedding. Expectedly, all three DL models were significantly better than the traditional ML models constructed above in the cross-validation and independent test (**Table 1** and **Figure 3**). For instance, the average AUC values of the DL models were above 0.824 whereas those of the ML models were less than 0.802.

In these DL models, two CNN models CNN_{OH} and CNN_{WE} had similar performances and both compared favorably to GRU_{WE} (**Table 1** and **Figure 3**). CNN_{OH} had the AUC value as 0.868 for the cross-validation and its values of SN, SP, ACC and MCC were 0.876, 0.700, 0.788, and 0.586, respectively. Here, we chose CNN_{OH} as the 2-Hydroxyisobutyrylation predictor. We evaluated the robustness of our models by comparing their performances between the cross-validation and independent tests. As their performances between these two tests had no statistically different ($P > 0.01$), we concluded that our constructed models were robust and neither over-fitting nor under-fitting.

Construction and Comparison of Predictors for Other Species

We constructed nine models for the human organism and chose CNN_{OH} as the final prediction model. We applied the CNN_{OH} architecture to the other three organisms (i.e., *T. gondii*, *O. sativa*, and *P. patens*). For each organism, we separated the dataset as the cross-validation set and the independent set. Similar to the human species, the CNN_{OH} models for these species had similar performances between cross-validation and independent test and their AUC values were larger than 0.818 (**Table 2**). It indicates that these constructed models are effective and robust.

As lysine 2-Hydroxyisobutyrylation is conserved across different types of species, we hypothesized that the model built for one species may be used to predict K_{hib} sites for other species. To test this hypothesis, we compared the performances of the CNN_{OH} models in terms of the independent data sets of individual species. Additionally, we built a general CNN_{OH} model based on the training datasets integrated from all the four species. **Table 3** shows that the AUC values of these predictions were larger than 0.761, suggesting that the cross-species prediction had reliable performances. Specifically, given a species, the best prediction performances were derived from

TABLE 4 | The prediction performance of CNN_{OH} compared to iLys-Khib in terms of the same cross-validation and independent test datasets.

Dataset	Model	Sn	Sp	Acc	MCC	AUC
10-fold cross-validation	iLys-Khib	0.745	0.658	0.701	0.404	0.770
	CNN _{OH}	0.830	0.713	0.772	0.547	0.847
Independent test	iLys-Khib	0.725	0.643	0.648	0.186	0.756
	CNN _{OH}	0.861	0.685	0.696	0.281	0.860

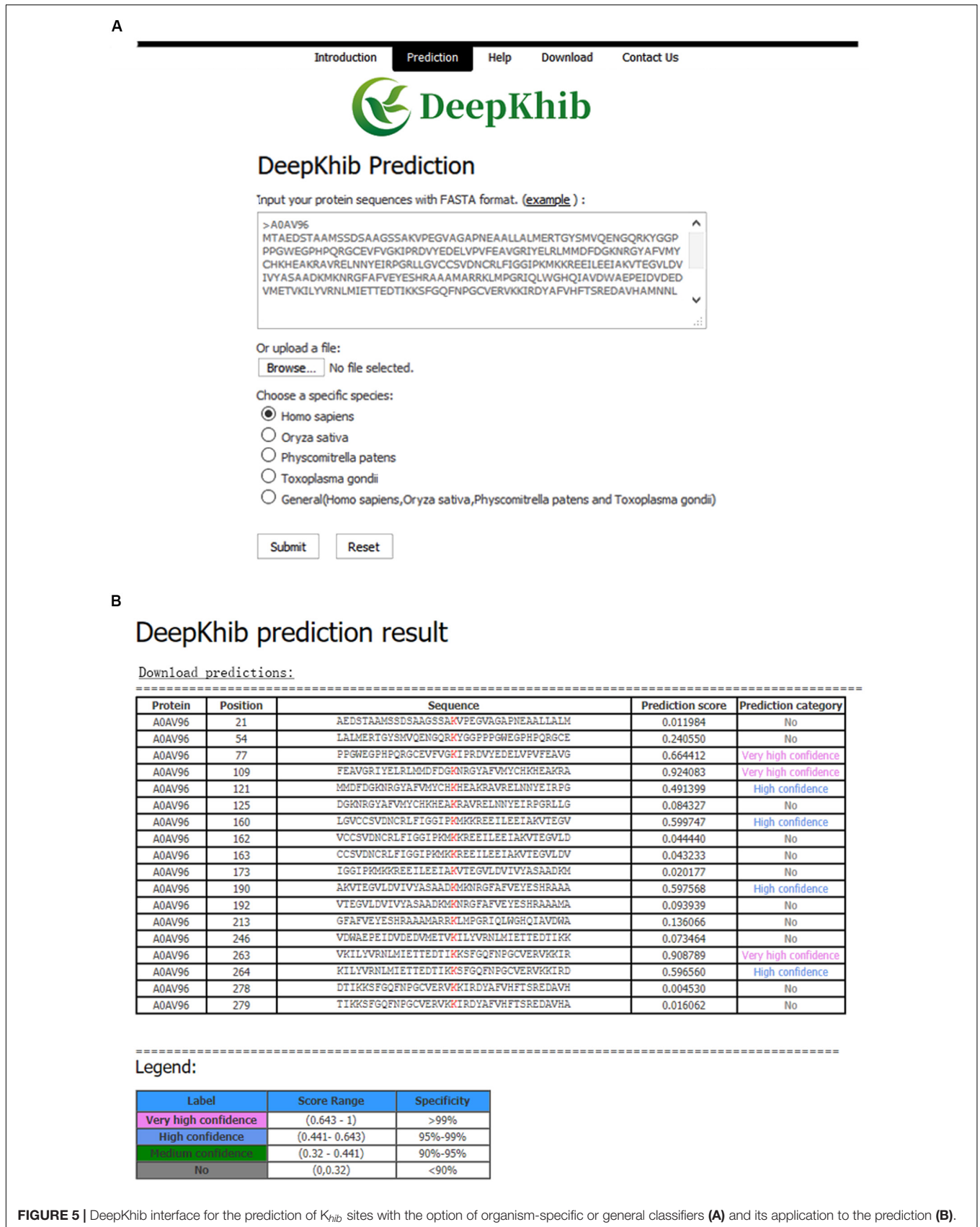


FIGURE 5 | DeepKhib interface for the prediction of K_{hib} sites with the option of organism-specific or general classifiers (A) and its application to the prediction (B).

the general model and the model developed specifically for this species. For instance, the human CNN_{OH} model had the best performance followed by the general model in terms of the human independent test whereas the general model had the best accuracy followed by the moss-specific model for the moss independent test. These suggest that on one hand, lysine 2-Hydroxyisobutyrylation of each species has its own characteristics; on the other hand, these modifications across different species share strong commonalities. Therefore, the general model may be effectively applied to any species. Furthermore, we evaluated the generality of the general CNN_{OH} model using the dataset of *S. cerevisiae* that contained 1,049 positive and 1,049 negative samples, which may not be enough for build an effective DL predictor (Chen et al., 2018a). The general model got the AUC value as 0.789, indicating the generality of this model. In other words, the general model is effective to predict K_{hib} sites for any organism.

We identified and compared the significant patterns and conserved motifs between K_{hib} and non- K_{hib} sequences across the different organisms using the two-sample-logo program with *t*-test ($P < 0.05$) with Bonferroni correction (Vacic et al., 2006). **Figure 4** shows the similarities and differences between the species. For instance, the residues R and K at the -1 position (i.e., R&K@P-1) and P at $+1$ position (i.e., P@P+1) are significantly depleted across the species. On the contrary, K&R@P+1 tend to be enriched for *H. sapiens* but depleted for *T. gondii* whereas both species have the depleted residue Serine across the positions ranging from P-18 to P+18. These similarities between the organisms may result in the generality and effectiveness of the general CNN_{OH} model.

Comparison of CNN_{OH} With the Reported Predictors

We assessed the performance of CNN_{OH} by comparing it with the existing K_{hib} predictors KhibPred (Wang et al., 2020) and iLys-Khib (Ju and Wang, 2019). First, we compared CNN_{OH} with KhibPred for individual species in terms of 10-fold cross-validation (Wang et al., 2020). The average AUC values of CNN_{OH} were 0.868/0.830/0.823 for *H. sapiens*/*P. patens*/*O. sativa*, respectively (**Table 2**). On the contrary, the corresponding values of KhibPred were 0.831/0.781/0.825 (Wang et al., 2020). Thus, CNN_{OH} compares favorably to KhibPred. Second, the model iLys-Khib was constructed and tested using 9,318 human samples as the 10-fold cross-validation data set and 4,219 human samples as the independent test set. We used the same datasets to construct CNN_{OH} and compared it with iLys-Khib. CNN_{OH} outperformed iLys-Khib in terms of all the measurements of performance (e.g., Sn, Sp, Acc, MCC, and AUC) for both 10-fold cross-validation and independent test (**Table 4**). For instance, the AUC value of CNN_{OH} was 0.860 for the independent test whereas that of iLys-Khib was 0.756. In summary, CNN_{OH} is a competitive predictor.

Construction of the On-Line K_{hib} Predictor

We developed an easy-to-use Web tool for the prediction of K_{hib} sites, dubbed as DeepKhib. It contains five CNN_{OH} models,

including one general model and four models specific to the species (i.e., *H. sapiens*, *O. sativa*, *P. patens*, and *T. gondii*). Given a species of interest, users could select the suitable model (e.g., the general model or the model specific to an organism) for prediction (**Figure 5A**). After the protein sequences as the fasta file format are uploaded, the prediction results will be shown with five columns: Protein, Position, Sequence, Prediction score, and Prediction category (**Figure 5B**). The prediction category covered four types according to the prediction scores: no (0–0.320), medium confidence (0.320–0.441), high confidence (0.441–0.643), and very high confidence (0.643–1).

CONCLUSION

The common PTM classifiers are mainly based on the traditional ML algorithms that require the pre-defined informative features. Here, we applied the advanced DL algorithm CNN_{OH} for predicting K_{hib} sites. CNN_{OH} shows its superior performance, because of the capability of the multi-layer CNN algorithm to extract complex features and learn sparse representation in a self-taught manner. Moreover, the general CNN_{OH} model demonstrates great generality and effectiveness, due to the conservation of K_{hib} modification from single-cell to multiple-cell organisms. The outstanding performance of DL in the prediction of K_{hib} sites suggests that DL may be applied broadly to predicting other types of modification sites.

DATA AVAILABILITY STATEMENT

All datasets presented in this study are accessible through <http://www.bioinfogo.org/DeepKhib/download.php>.

AUTHOR CONTRIBUTIONS

LL conceived this project. LZ and YZ constructed the algorithms under the supervision of LL and ZC. LZ and NH analyzed the data. LL, YZ, YC, and LZ wrote the manuscript. All authors read and approved the final manuscript.

FUNDING

This work was supported in part by funds from the Young Scientists Fund of the National Natural Science Foundation of China (Grant No. 31701142 to ZC) and the National Natural Science Foundation of China (Grant No. 31770821 to LL). LL was supported by the “Distinguished Expert of Overseas Tai Shan Scholar” program. YZ was supported by the Qingdao Applied Research Project.

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fcell.2020.580217/full#supplementary-material>

REFERENCES

- Beltrao, P., Bork, P., Krogan, N. J., and Van Noort, V. (2013). Evolution and functional cross-talk of protein post-translational modifications. *Mol. Syst. Biol.* 9:714. doi: 10.1002/msb.201304521
- Breiman, L. (2001). Random Forests. *Machine Learn.* 45, 5–32. doi: 10.1023/A:1010933404324
- Chen, Y. Z., Chen, Z., Gong, Y. A., Ying, G., and Parkinson, J. (2012). Sumohydro: a novel method for the prediction of sumoylation sites based on hydrophobic properties. *PLoS One* 7:e39195. doi: 10.1371/journal.pone.0039195
- Chen, Z., He, N., Huang, Y., Qin, W. T., Liu, X., and Li, L. (2018a). Integration of A Deep Learning Classifier with A Random Forest Approach for Predicting Malonylation Sites. *Genom. Proteom. Bioinform.* 16, 451–459.
- Chen, Z., Zhao, P., Li, F., Leier, A., Marquez-Lago, T. T., and Wang, Y. (2018b). iFeature: a python package and web server for features extraction and selection from protein and peptide sequences. *Bioinformatics* 34, 2499–2502.
- Chen, Z., Liu, X., Li, F., Li, C., Marquez-Lago, T., Leier, A., et al. (2019). Large-scale comparative assessment of computational predictors for lysine post-translational modification sites. *Brief. Bioinform.* 20, 2267–2290.
- Chen, Z., Zhao, P., Li, F., Marquez-Lago, T. T., Leier, A., et al. (2020). iLearn: an integrated platform and meta-learner for feature engineering, machine-learning analysis and modeling of DNA, RNA and protein sequence data. *Brief. Bioinform.* 21, 1047–1057.
- Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., et al. (2014). Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. *Comput. Ence.* 2014, 1724–1734. doi: 10.3115/v1/D14-1179
- Dai, L., Peng, C., Montellier, E., Lu, Z., Chen, Y., Ishii, H., et al. (2014). Lysine 2-hydroxyisobutyrylation is a widely distributed active histone mark. *Nat. Chem. Biol.* 10, 365–370.
- Fukushima, K. (1980). Neocognitron: a self organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biol. Cybern.* 36, 193–202. doi: 10.1007/BF00344251
- Hahnloser, R. H., Sarpeshkar, R., Mahowald, M. A., Douglas, R. J., and Seung, H. S. (2000). Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit. *Nature* 405, 947–951. doi: 10.1038/35016072
- Huang, H., Luo, Z., Qi, S., Huang, J., Xu, P., Wang, X., et al. (2018a). Landscape of the regulatory elements for lysine 2-hydroxyisobutyrylation pathway. *Cell Res.* 28, 111–125.
- Huang, H., Tang, S., Ji, M., Tang, Z., Shimada, M., Liu, X., et al. (2018b). p300-Mediated Lysine 2-Hydroxyisobutyrylation Regulates Glycolysis. *Mol. Cell* 70, 663–678.e.
- Huang, Y., He, N., Chen, Y., Chen, Z., and Li, L. (2018c). BERMP: a cross-species classifier for predicting mA sites by integrating a deep learning algorithm and a random forest approach. *Int. J. Biol. Sci.* 14, 1669–1677.
- Huang, J., Luo, Z., Ying, W., Cao, Q., and Dai, J. (2017). 2-hydroxyisobutyrylation on histone h4k8 is regulated by glucose homeostasis in *saccharomyces cerevisiae*. *Proc. Natl. Acad. Sci. U.S.A.* 114, 8782–8787. doi: 10.1073/pnas.1700796114
- Huang, Y., Niu, B., Gao, Y., Fu, L., and Li, W. (2010). CD-HIT Suite: a web server for clustering and comparing biological sequences. *Bioinformatics* 26, 680–682. doi: 10.1093/bioinformatics/btq003
- Ju, Z., and Wang, S.-Y. (2019). iLys-Khib: Identify lysine 2-Hydroxyisobutyrylation sites using mRMR feature selection and fuzzy SVM algorithm. *Chemometr. Intell. Laborat. Syst.* 191, 96–102. doi: 10.1016/j.chemolab.2019.06.009
- Kingma, D. P., and Jimmy, B. (2014). *Adam: A Method for Stochastic Optimization*. New York, NY: Cornell University.
- Li, Q. Q., Hao, J. J., Zhang, Z., Krane, L. S., Hammerich, K. H., Sanford, T., et al. (2017). Proteomic analysis of proteome and histone post-translational modifications in heat shock protein 90 inhibition-mediated bladder cancer therapeutics. *Sci. Rep.* 7:201.
- Li, W., and Godzik, A. (2006). Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* 22, 1658–1659. doi: 10.1093/bioinformatics/btl158
- Long, H., Liao, B., Xu, X., and Yang, J. (2018). A Hybrid Deep Learning Model for Predicting Protein Hydroxylation Sites. *Int. J. Mol. Sci.* 19:2817.
- Meng, X., Xing, S., Perez, L. M., Peng, X., Zhao, Q., Redona, E. D., et al. (2017). Proteome-wide Analysis of Lysine 2-hydroxyisobutyrylation in Developing Rice (*Oryza sativa*) Seeds. *Sci. Rep.* 7:17486.
- Nitish, S., Geoffrey, H., Alex, K., Ilya, S., and Ruslan, S. (2014). Dropout: a simple way to prevent neural networks from overfitting. *J. Machine Learn. Res.* 15, 1929–1958.
- Sainath, T. N., Mohamed, A., Kingsbury, B., and Ramabhadran, B. (2013). *Deep convolutional neural networks for LVCSR*. New Jersey: IEEE. doi: 10.1109/ICASSP.2013.6639347
- Sandberg, M., Eriksson, L., and Jonsson, J. (1998). New chemical descriptors relevant for the design of biologically active peptides. *a multivariate characterization of 87 amino acids*. *J. Med. Chem.* 41, 2481–2491. doi: 10.1021/jm9700575
- Skelly, M. J., Frungillo, L., and Spoel, S. H. (2016). Transcriptional regulation by complex interplay between post-translational modifications. *Curr. Opin. Plant Biol.* 33, 126–132. doi: 10.1016/j.pbi.2016.07.004
- Tahir, M., Tayara, H., and Chong, K. T. (2019). iPseU-CNN: Identifying RNA Pseudouridine Sites Using Convolutional Neural Networks. *Mol. Ther. Nucl. Acids* 16, 463–470. doi: 10.1016/j.omtn.2019.03.010
- Tian, Q., Zou, J., Tang, J., Fang, Y., Yu, Z., and Fan, S. (2019). MRCNN: a deep learning model for regression of genome-wide DNA methylation. *BMC Genomics* 20:192. doi: 10.1186/s12864-019-5488-5
- Vacic, V., Iakoucheva, L. M., and Radivojac, P. (2006). Two Sample Logo: a graphical representation of the differences between two sets of sequence alignments. *Bioinformatics* 22, 1536–1537. doi: 10.1093/bioinformatics/btl151
- Wang, D., Zeng, S., Xu, C., Qiu, W., Liang, Y., Joshi, T., et al. (2017). Musitedeep: a deep-learning framework for general and kinase-specific phosphorylation site prediction. *Bioinformatics* 33, 3909–3916.
- Wang, Y. G., Huang, S. Y., Wang, L. N., Zhou, Z. Y., and Qiu, J. D. (2020). Accurate prediction of species-specific 2-hydroxyisobutyrylation sites based on machine learning frameworks. *Anal. Biochem.* 602:113793. doi: 10.1016/j.ab.2020.113793
- Wu, Q., Ke, L., Wang, C., Fan, P., Wu, Z., and Xu, X. (2018). Global Analysis of Lysine 2-Hydroxyisobutyrylation upon SAHA Treatment and Its Relationship with Acetylation and Crotonylation. *J. Proteome Res.* 17, 3176–3183.
- Xiao, H., Xuan, W., Shao, S., Liu, T., and Schultz, P. G. (2015). Genetic Incorporation of epsilon-N-2-Hydroxyisobutyryl-lysine into Recombinant Histones. *ACS Chem. Biol.* 10, 1599–1603. doi: 10.1021/cb501055h
- Xie, Y., Luo, X., Li, Y., Chen, L., Ma, W., Huang, J., et al. (2018). DeepNitro: Prediction of Protein Nitration and Nitrosylation Sites by Deep Learning. *Genom. Proteom. Bioinform.* 16, 294–306.
- Yin, D., Zhang, Y., Wang, D., Sang, X., Feng, Y., Chen, R., et al. (2019). Global Lysine Crotonylation and 2-Hydroxyisobutyrylation in Phenotypically Different *Toxoplasma gondii* Parasites. *Mole. Cell. Proteom.* 18, 2207–2224.
- Yu, Z., Ni, J., Sheng, W., Wang, Z., and Wu, Y. (2017). Proteome-wide identification of lysine 2-hydroxyisobutyrylation reveals conserved and novel histone modifications in *Physcomitrella patens*. *Sci. Rep.* 7:15553. doi: 10.1038/s41598-017-15854-z
- Zhao, Y., He, N., Chen, Z., and Li, L. (2020). Identification of Protein Lysine Crotonylation Sites by a Deep Learning Framework With Convolutional Neural Networks. *IEEE Access* 8, 14244–14252. doi: 10.1109/ACCESS.2020.2966592

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2020 Zhang, Zou, He, Chen, Chen and Li. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.