

1 **To improve the predictions of binding residues with DNA, RNA,**  
2 **carbohydrate, and peptide via multiple-task deep neural networks**

3 **Zhe Sun<sup>1</sup>, Shuangjia Zheng<sup>1</sup>, Huiying Zhao<sup>2</sup>, Zhangming Niu<sup>3</sup>, Yutong Lu<sup>1</sup>, Yi Pan<sup>5</sup> and**  
4 **Yuedong Yang<sup>1,4\*</sup>**

5 <sup>1</sup>School of Data and Computer Science, Sun Yat-sen University, Guangzhou 510000, China

6 <sup>2</sup>Sun Yat-sen Memorial Hospital, Sun Yat-sen University, Guangzhou 510000, China

7 <sup>3</sup>Aladdin Healthcare Technologies Ltd. 24-26 Baltic Street West, London EC1Y 0UR, UK

8 <sup>4</sup>Key Laboratory of Machine Intelligence and Advanced Computing (Sun Yat-sen University), Ministry of  
9 Education, China

10 <sup>5</sup>Department of Computer Science, Georgia State University, USA

11

12 **Abstract**

13 **Motivation:** The interactions of proteins with DNA, RNA, peptide, and carbohydrate play key roles in  
14 various biological processes. The studies of uncharacterized protein–molecules interactions could be  
15 aided by accurate predictions of residues that bind with partner molecules. However, the existing  
16 methods for predicting binding residues on proteins remain of relatively low accuracies due to the  
17 limited number of complex structures in databases. As different types of molecules partially share  
18 chemical mechanisms, the predictions for each molecular type should benefit from the binding  
19 information with other molecules types.

20 **Results:** In this study, we employed a multiple task deep learning strategy to develop a new  
21 sequence-based method for simultaneously predicting binding residues/sites with multiple important  
22 molecule types named MTDsite. By combining four training sets for DNA, RNA, peptide, and  
23 carbohydrate-binding proteins, our method yielded accurate and robust predictions with AUC values of  
24 0.852, 0.836, 0.758, and 0.776 on their respective independent test sets, which are 0.52 to 6.6% better  
25 than other state-of-the-art methods. More importantly, this study provides a new strategy to improve  
26 predictions by combining multiple similar tasks.

27 **Availability:** <http://biomed.nscg-gz.cn/server/MTDsite/>

28 **Contact:** [yangyd25@mail.sysu.edu.cn](mailto:yangyd25@mail.sysu.edu.cn)

29

## 1 **1. Introduction**

2 Predicting proteins interactions with other molecules is critical for understanding biological processes  
3 and discovering drugs. The most important molecules controlling biological activities include DNA,  
4 RNA , peptides and carbohydrates (CBH) (Hanson, et al., 2019). For example, interactions between  
5 protein and nucleic acids are central to many of the vital processes in molecular biology (Jia, et al.,  
6 2019) such as transcription, translation, post-transcriptional modification and regulation. Additionally,  
7 RNA-binding proteins can modulate or stabilize RNA structures to make RNA catalytically active (Pan  
8 and Shen, 2017). The carbohydrate-binding proteins act as important biomarkers for cell  
9 communication, cell adhesion, fertilization, development and differentiation (Lu and Pieters, 2019).  
10 Peptide interactions with protein domains occur in many cell processes particular signaling pathways  
11 (Petsalaki and Russell, 2008). In order to better understand the molecular mechanisms, we need to  
12 know the binding residues of these proteins interacting with their respective binding partner.  
13 Traditional experimental techniques such as X-ray and NMR experiments are robust but are expensive  
14 and time-consuming (London, et al., 2010) (Su, et al., 2018). With the exponentially increasing protein  
15 sequences, it is demanding to make predictions of binding residues from sequences.

16  
17 Many bioinformatics methods combining physicochemical and evolutionary features (Wang L, 2010)  
18 have been developed in the past decades. For DNA-binding residue predictions, representative methods  
19 include TargetDNA (Jun Hu, 2017) and HMMBinder (Rianon Zaman, 2017) based on SVM, CNNsites  
20 (Wang, 2016) based on CNN network, DRNAPred (Jing Yan, 2017) for accurately predicting and  
21 discriminating between DNA- and RNA-binding residues, and SPOT-DNA-Seq (Zhao, et al., 2014)  
22 based on alignments with known DNA-binding proteins. For RNA-binding residue predictions, there  
23 are RNAProSite (Meijian Sun, 2016) based on the random forest classifier and PredRBR (Yongjun  
24 Tang, 2017) based on the gradient tree boosting. For predictions of carbohydrate binding residues, the  
25 common idea is to find residues frequently observed on the sugar interface(Ghazaleh Taherzadeh,  
26 2016) (Sujatha M S, 2004). For predictions of peptide-binding residues, several methods have been  
27 developed based on machine learning techniques, e.g. SPRINT(Zhou, 2016) and SPOT-Peptide(Zhou,  
28 2019).

29  
30 Although many successful methods have been proposed, most of them suffer from low accuracies  
31 due to small training sample sizes. The underlying reason is that the complex structures of proteins are  
32 difficult to obtain by experiments. To be worse, the predictions of binding residues with different types  
33 of molecules were traditionally treated as different problems (Miao and Westhof, 2015), and the  
34 prediction tasks were usually performed separately. Thus, the small sizes of individual binding data  
35 sets prevent the applications of deep learning techniques. In fact, most biological molecules are organic  
36 molecules, and the similarities in physic-chemical properties enable the sharing of interaction patterns.  
37 For example, the combined inputs of DNA- and RNA-binding residues into the same learning system  
38 have been proven to improve the predictions through the support vector machine techniques (Zhang X,  
39 2016) (Su, et al., 2019) and the artificial neural network (Zhang, et al., 2012). However, these studies  
40 have ignored the differences between DNA and RNA molecules, and no study has yet been performed  
41 to combine binding information with other molecular types. For this purpose, the multi-task learning  
42 provides a promising framework that learns shared information through common networks while

1 retaining the task-related output layers (Rich Caruana, 1997). The shared networks between multiple  
2 similar tasks enable a bigger training set that could maintain a larger network.

3 In this study, we designed a multi-task network architecture (namely MTDsite) to simultaneously  
4 predict respective binding residues with DNA, RNA, carbohydrate, and peptide molecules. The shared  
5 networks among all tasks can help learn common representations and thereby obtain relatively strong  
6 abstracting capabilities, and we used the LSTM as our Shared network to collect the information of  
7 long-range residues in the protein chain. At the same time, four small specific sub-networks were  
8 respectively trained for four individual types to extract individual properties. The benchmark tests  
9 indicated the employing of multi-task learning leads to averagely 3.6% improvements over  
10 state-of-the-art methods when measured by the area under the receiver operating characteristic curve  
11 (AUC).

## 12 2. Methods

### 13 2.1 Benchmark Datasets

14 We evaluated our method by using the previously curated training and testing datasets. The datasets  
15 include the protein binding with DNA, RNA, peptide, and carbohydrate, where a residue was defined  
16 as a binding residue if it contains at least one atom within 3.5Å from its binding partner. The sequence  
17 identities between the proteins in the training and test sets are less than 30% according to  
18 BLASTCLUST(Johnson, 2008). Table 1 lists the datasets, with details as:

19

20 **The DNA and RNA datasets:** The datasets were collected from a recent study (Yan and Kurgan,  
21 2017), where the training and testing datasets are 309 and 47 chains for DNA, and 157 and 17 chains  
22 for RNA, respectively.

23

24 **The peptide dataset:** The dataset was downloaded from a recent study (Tahezadeh G, et al., 2016),  
25 where protein-peptide complex structures were extracted from the BioLip protein-ligands database  
26 (Yang, et al., 2013) with peptides as the ligands derived from the Protein Data Bank (PDB). The  
27 dataset includes 1115 proteins as the training set, and 125 proteins as the independent test set that were  
28 randomly split by the previous study.

29

30 **The carbohydrate dataset (CBH):** The dataset was downloaded from a recent study (Zhao, 2014)  
31 that includes 157 chains for training and 17 chains for test. The dataset was originally derived from  
32 the PROCARB (Malik, et al., 2010) by keeping only high resolution (<3Å) complex structures.

33

### 34 2.2 Input Features

35 Our input includes totally 54 features that are composed of G-PSSM (20 features), G-HHM (20), and  
36 G\_SPD3 (14), as detailed below:

37

38 **G-PSSM:** Evolutionarily conserved residues are considered to be the same or similar residues  
39 maintained between species by natural selection. They have important functional roles like acting as  
40 binding sites. In this study, we employed the position specific scoring matrix (PSSM) which is a 20\*L  
41 dimensional matrix (where L is protein length) generated from PSI-BLAST with E-value threshold of  
42 0.001 in three iterations.

1

2 **G-HHM:** The hidden Markov models (HMMs) have been successfully used in protein structure  
3 prediction (Remmert, 2012) that assumed a Markov process with unobserved states, and the profile  
4 HMMs accomplish the protein structure prediction task well based on the HMM-HMM alignment. The  
5 sequence alignment generated by HHblits has been found with higher accuracy than by PSI-BLAST  
6 (Stephen F. Altschul, 1997). Here, HMM profile was generated by HHblits that compared to the query  
7 sequence with the proteins in the uniprot20\_2015\_06 database

8

9 **G-SPD3:** Predicted structural properties by SPIDER3

10 The structural properties were predicted by SPIDER 3.0 (Rhys Heffernan, 2018) , including:

11

12 **ASA (2):** The accessible surface area (ASA) means the surface area of a biomolecule accessible to a  
13 solvent, which reflects the functional importance of residues. Here, ASA of each residue was obtained  
14 by SPIDER 3.0, and the ASA was predicted as relative ASA named as rASA. We also computed the  
15 average rASA of the residue and its four adjacent residues.

16

17 **Torsional angles (8):** The backbone torsional angles are composed of  $\Phi$ ,  $\psi$ , and  $\omega$  that are used to  
18 describe local backbone structure of a protein. The  $\omega$  was not used here because it is usually at  $180^\circ$   
19 due to the planarity of the pep-tide bond. The angles between neighboring residues include  $\theta$  and  $\tau$ .  
20 Here, we let  $\theta$  be the angle between  $C_{ai-1}-C_{ai}-C_{ai+1}$ , and  $\tau$  be the dihedral angle rotated about the  
21  $C_{ai}-C_{ai+1}$  vector to compute the feature. We extracted the features by using cosine and sine values of  
22 the four angles, totally 8 features.

23

24 **CN (1):** CN is the number of residues within a distance cutoff to a given residue in  
25 three-dimensional space.

26

27 **HSE (3):** Half-Sphere Exposure (HSE) is an extension of CN, which considers directions of two  
28 residues in a top and bottom half of the sphere. Two methods were used to define the plane separating  
29 the upper and lower hemisphere, which included  $HSE\alpha$  based on the neighboring  $C\alpha-C\alpha$ directional  
30 vector and  $HSE\beta$  based on the  $C\alpha-C\beta$  directional vector. In this study, we used the  $HSE\alpha$  to describe  
31 the feature. For CN and HSE, residue distance is defined as the distance between  $C\alpha$  atoms with a  $13 \text{ \AA}$   
32 cutoff.

33

### 34 **2.3 Performance evaluation**

35 The performance of methods was measured by the number of correct classified and the number of  
36 misclassified instances using the terms below:

37 **TP:** number of actual binding residues predicted as binding sites

38 **TN:** number of actual non-binding residues predicted as non-binding sites

39 **FP:** number of actual non-binding residues incorrectly predicted as binding sites

40 **FN:** number of actual binding residues incorrectly predicted as non-binding sites

41 We evaluated the performance of our proposed prediction method in terms of the Matthews'  
42 Correlation Coefficient (MCC), Accuracy, Receiver Operating Characteristic (ROC) curve as:

43

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

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$Sensitivity = \frac{TP}{POS}$$

$$Specificity = \frac{TN}{NEG}$$

$$Precision = \frac{TP}{TP + FP}$$

1 , where POS is the number of known binding residues, and NEG is the number of known non-binding  
2 residues. MCC varies between 0 and 1, with 1 representing correct prediction for all residues, and 0  
3 by random. Additionally, the Area Under the Curve (AUC) was adopted as our primary evaluation  
4 index because of our unbalanced datasets, i.e., the much less numbers of positive than negative  
5 samples.

6

#### 7 **2.4 MTDSite architecture**

8 As shown in Fig 1, the deep learning network in this method consists of two parts. The first part is a  
9 shared Bi-directional Long Short-term Memory (BiLSTM) network called shared network (referred to  
10 shared BiLSTM), which was used to extract common information from different binding molecules.  
11 The second part is four individual networks composed of full connection layer. The predictive fraction  
12 of protein-molecules binding residues can be obtained from this part.

13

#### 14 **Bi-LSTM network:**

15 Long short-term memory (LSTM) is a recurrent neural network (RNN) architecture (Hochreiter and  
16 Schmidhuber, 1997). The LSTM network was shown to have better performances than traditional RNN  
17 in processing, classifying, and predicting time series when there are indefinitely long separations  
18 between important events (Zhiheng Huang, 2015). This is the main reason why LSTM outperforms  
19 alternative Hidden Markov Models and other sequence learning methods in numerous applications  
20 (Yequan Wang, 2016).

21

22 In the prediction of binding residues, extracting information between long range residue pairs is  
23 important for constructing an accurate model. Traditional machine learning methods, such as Xgboost  
24 and SVM, extract feature information of adjacent residues by creating slide-window. In comparison,  
25 LSTM collects information on adjacent residues by establishing various 'gates', which saves time on  
26 training and adjusting parameters.

27

28 BiLSTM is a combination network of BRNN and LSTM by stacking multiple BiLSTM-RNN hidden  
29 layers together to build a deep bidirectional LSTM-RNN. The output of one layer is used as the input

30 of the next layer. The hidden state sequence,  $h^n$ , consist of forward and backward sequence  $\vec{h}^n$  and

31  $\overleftarrow{h}^n$ , iteratively computed from  $n = 1$  to N and  $t = 1$  to T as follows:

32

$$\vec{h}^n = \mathfrak{N} \left( W_{h^{n-1}h^n} \vec{h}_t^{n-1} + W_{h^n h^n} \vec{h}_{t-1}^n + b_h^n \right)$$

$$\overleftarrow{h}^n = \mathfrak{N} \left( W_{h^{n-1}h^n} \overleftarrow{h}_t^{n-1} + W_{h^n h^n} \overleftarrow{h}_{t-1}^n + b_h^n \right)$$

$$y_t = \mathfrak{N} \left( W_{hNy} \vec{h}_t^N + W_{hNy} \overleftarrow{h}_t^N + b_y \right)$$

1 where  $y = (y_1, y_2, y_3, \dots, y_t, \dots, y_T)$  is the output,  $\mathfrak{N}$  is the activation function

2

### 3 **Multi-task learning:**

4 A neural network is a non-linear classifier that performs repeated linear and non-linear transformations  
5 on the input. Let  $x_i$  represent the input of  $i$ -th layer of the network (where  $x_0$  is just the feature vector).

6 The transformation performed as:

$$X_{i+1} = \sigma(W_i X_i + b_i)$$

7 where  $W_i$  and  $b_i$  are the weight matrix and bias value for the  $i$ -th layer respectively, and  $\sigma$  is a  
8 nonlinear function. After  $N$  conversions, the final layer of the network  $X^N$  is then fed to a simple  
9 linear classifier (the softmax in our method) to predict the probability of input  $x_0$  in label  $j$ :

$$P(y = j | x_0) = \frac{e^{(w^j)^T \times N}}{\sum_{m=1}^M e^{(w^m)^T \times N}}$$

10 where  $M$  is the number of possible labels (here  $M = 4$ ) and  $w^1, w^2, \dots, w^M$  are weight vectors. The  
11 above parameters  $W_i, b_i$  and  $w^m$  were parameters learned during training by the back-propagation  
12 algorithm. A multitask network attaches  $N$  SoftMax classifiers, one for each task, to the final layer  $x_L$ .  
13 Here, we defined a “task” as classifier associated with a particular task, and one particular binding data  
14 set in our collection.

15

### 16 **MTDsite networks:**

17 The networks consist of two parts, the shared network and four specific networks for individual tasks.  
18 The input of shared network is a  $54 \times L$  ( $L$  is the length of a protein and 54 is the number of the  
19 features) feature matrix. Through two shared LSTM hidden layers, a  $2 \times L$  scoring matrix can be  
20 obtained, which is used as input for the next part.

21

22 As different tasks have specific properties, the second part is four independent fully-connected  
23 networks specifically trained for four different tasks. The four specific networks have identical  
24 structures (layer and neuron sizes). For each task, only the corresponding specific network will be  
25 updated, with the left three specific networks unchanged.

26

### 27 **2.5 Cross-Validation and Independent Test**

28 The cross validations and independent tests were employed to evaluate the robustness and performance  
29 of the method. The training data set was evenly divided into ten pieces (folds) at random. In each  
30 round, nine folds were employed for training and the remaining fold was used for test. This process  
31 was repeated for 10 times so that each fold has been tested once, and all outputs were collected to  
32 compute the performances. Based on the optimal hyper-parameters, a model trained by using the whole  
33 training set was then tested on the independent test set.

### 34 **2.6 Model selection and Parameters optimization**

1 During the optimization of the MTDsite models and hyperparameters, we only randomly selected 1/10  
2 of the training samples, and selected the optimal parameters with the highest AUC value. We didn't use  
3 the 10-fold cross validation to select as it will take 10 times of training costs. Finally, we used the ELU  
4 and Cross-Entropy as the network activation and loss functions in order to improve convergence speed  
5 and accuracy. The final optimal parameters of EPOCH and LR were 21 and 0.001, respectively, the  
6 number of hidden layers is 2, and the hidden nodes for shared network were both set as 128. The  
7 specific networks were all set as simple, single-layer, fully connected networks with 64 hidden nodes.  
8 These hyper-parameters (LR, the number of hidden layers and the number of hidden nodes) were  
9 determined by the GridSearchCV () function in python sklearn library. Due to the different lengths of  
10 all protein sequences, batch-size was set to 1 to avoid the decrease in accuracy caused by padding,  
11 though it reduces the training speed of the network. After the optimal parameters were decided, the  
12 models were evaluated by the cross-validations on the training sets, and independent tests on the test  
13 sets.

## 14 **2.7 MTDsite-single Models**

15 For a direct comparison, we trained MTDsite-single models like traditional ways by independently  
16 inputting single binding type of training data, and thus four models were obtained for four binding  
17 types, respectively. If not specifically mentioned, the model trained on one binding type will be  
18 employed to test the corresponding type. The architecture of MTDsite-single is a two-layer BiLSTM  
19 network, which is the same as the shared network of MTDsite. Similarly, the hyperparameters were  
20 also optimized on the training sets. The final EPOCH parameters are 17, 19, 26, and 28 for DNA,  
21 RNA, peptide and carbohydrate-binding models, and the learning rate is 0.001.

## 22 **3. Results**

### 23 **3.1 Performances of MTDsite on the 10-fold cross validations and independent tests**

24 Table 2 shows the performances of MTDsite measured by AUC, MCC, sensitivity, and specificity in  
25 prediction of DNA-, RNA-, peptide-, and carbohydrate-binding residues. By the 10-fold cross  
26 validation, MTDsite achieved AUC values of 0.866, 0.857, 0.760, and 0.779 for the four binding types,  
27 respectively. The predictions of DNA- and RNA-binding residues have relatively greater AUC values,  
28 likely because DNA and RNA are negatively charged and thus the binding residues are relatively easier  
29 to predict. The independent tests have obtained essentially the same AUC values with differences of  
30 only 0.003~0.024 for the four tests, indicating the robustness of our models. Fig 2 shows the ROC  
31 curves by the cross validations and independent tests on four binding types, respectively. The  
32 performances were also confirmed by the consistent maximum MCC values between the 10-fold cross  
33 validations and independent tests, with differences of 0.003~0.021. At the thresholds with the  
34 maximum MCC values, our models have high specificities while relatively low sensitivities due to the  
35 much greater numbers of non-binding residues (negatives) than the binding residues (positives).

36  
37 We evaluated the contributions of individual feature group by using only single feature group or  
38 excluding one feature group from all features. As shown in Table 3, when individual feature group was  
39 used in the prediction, G-PSSM, the evolution features produced by PSIBLAST, yielded the greatest  
40 values in regard with the average values of both AUC and MCC. G-HHM, another feature group  
41 produced by HHblits, yielded slightly lower AUC and MCC values. G-SPD3, the structural feature  
42 group produced by SPIDER3 package, yielded significantly lower AUC values in average. These

1 results suggest the importance of evolution information for protein binding, consistent with previous  
2 findings (Hong Su, et al, 2018). When excluding individual feature group, the removal of G-PSSM  
3 caused the largest decreases in the average values of both AUC and MCC, again indicating its most  
4 important role. Though the removal of the G-SPD3 feature group caused the smallest decrease, the  
5 difference is significant ( $P=0.001$ ) according to the paired t-test. The decreases are small likely because  
6 the G\_SPD3 features were derived from the PSSM and HHM profiles, and our neural networks could  
7 partly catch the structural information from the two profiles.

### 8 9 **3.2 Contributions by the shared networks**

10 To evaluate the contributions of the networks shared by different binding types, we trained four  
11 MTDsite-single models for comparison, each with single type of binding data. As shown in Fig 2, the  
12 ROC curves of MTDsite-single are consistently below those of MTDsite by independent tests on four  
13 binding types. As a result, the AUC values by MTDsite-single are 3.6%, 3.2%, 3.3%, and 4.0% lower  
14 than those by MTDsite for DNA, RNA, peptide, and carbohydrate, respectively (Table 4). The  
15 difference of ROC curves between MTDsite and MTDsite-single reflect the improvement affected by  
16 the data set fusion and multi-task learning. When measured by MCC, the MTDsite-single are 4.7%,  
17 36.7%, 37.2%, and 15.4% worse than the MTDsite.

18  
19 We further performed cross tests by using the classifier trained from one binding type to test other  
20 types. As shown in Table 4, on prediction of DNA and RNA binding residues, similar performances  
21 have been achieved by the MTD-single (DNA) and MTD-single (RNA), indicating a similarity  
22 between DNA- and RNA-binding sites. This also explains why predictions could be improved by a  
23 simple combination of DNA and RNA binding residues in the previous studies. By comparison, there  
24 are significant decreases on other cross tests, among which the MTD-single (PEP) made almost random  
25 predictions on other three binding types although the peptide-binding dataset used for training has the  
26 biggest sample size. Interestingly, the MTD-single (DNA) and MTD-single (RNA) models produced  
27 good predictions on the carbohydrate-binding dataset: only 3.9% and 5.2% lower AUC than  
28 MTD-single (CBH). This is likely because carbohydrate has common properties with the deoxyribose  
29 and ribose respectively contained in DNA and RNA molecules. By comparison in the MTDsite,  
30 although the specific networks have always produced the best predictions for their respective binding  
31 types, other specific networks in the MTDsite could also achieve reasonable predictions. For example,  
32 on the prediction of peptide-binding residues, MTDsite (CBH) achieved essentially the same AUC  
33 value as the MTDsite (PEP), indicating its potential generality to other binding types.

34  
35 Fig 3 shows a direct comparison of MTDsite and MTDsite-single. Among the 237 protein chains  
36 from four independent datasets, MTDsite significantly outperforms MTDsite-single with P-value of  
37 0.003 according to the paired t-test, where MTDsite have greater AUC value for 162 (68%) chains.

### 38 39 **3.3 Comparisons with other methods**

40 MTDsite was compared with other methods on the independent test sets respectively for four binding  
41 types. As shown in Table 5, MTDsite achieved the highest AUC values through all four types of  
42 independent datasets, which are 2.2%, 4.4%, 6.6%, and 0.52% higher than other best methods for  
43 DNA, RNA, peptide, and carbohydrate datasets, respectively. The improvements are mostly  
44 contributed by the multi-task learning as the MTDsite-single models without using multi-task learning



1 yielded close and mostly lower performances than other state-of-the-art methods. Fig 2 also indicated  
2 that all other methods are mostly below the ROC curves by MTDsite in all four tests. The only  
3 exception is the SPRINT-CBH method, previously developed by our group for the  
4 carbohydrate-binding prediction, which is marginally above the ROC by MTDsite, though MTDsite  
5 achieved a higher AUC value. Here, the lower MCC by MTDsite is likely because MTDsite has been  
6 optimized for the greatest AUC values while SPRINT-CBH was optimized for MCC values.  
7 Additionally, the carbohydrate dataset shows a much larger ratio between the number of negatives and  
8 the positives in the test set (26.0 on the carbohydrate compared to 8.5, 6.5, and 16.8 on the DNA, RNA,  
9 peptide datasets, respectively), and thus the shared network wasn't well optimized for the carbohydrate  
10 dataset.

11 As a result, the MCC values of MTDsite are 11.8%, 1.4%, and 7.2% percent higher than other  
12 methods for DNA, RNA, and peptide, respectively. At the same time, MTDsite achieved the highest  
13 accuracy, sensitivity, and precision in most cases. It should also be noticed that PepSite and Peptimap  
14 are structure-based methods. Though MTDsite used sequence-based information only, our method still  
15 achieved better performances.

16

### 17 **3.4 Case study**

18 As an example, we demonstrated the prediction on a restriction-modification controller DNA-binding  
19 protein (PDBid: 3s8qB). As shown in Fig 4, MTDsite predicted 15 binding residues that include 14  
20 residues are truly binding. In comparison, MTDsite-single predicted 14 binding residues including 13  
21 truly binding residues. Thus, total accuracies of two methods are 97% and 94%, respectively. The  
22 difference of MTD-single in prediction performance on different data sets showed that our method  
23 could distinguish the binding sites of different molecules on proteins, and also proved that  
24 cross-prediction had less impact on MTDsites.

## 25 **4. Conclusion and Discussion**

26 We have developed a new architecture to use multiple-task learning for predicting binding residues. To  
27 our best knowledge, this is the first attempt of the strategy for prediction of binding residues. The  
28 sharing between DNA, RNA, peptide, and carbohydrate-binding information was indicated to improve  
29 the predictions of binding residues for all four types. Our method was indicated robust by the consistent  
30 performances between the cross validations and independent tests. The method was proven to  
31 outperform state-of-the-state methods. Such strategy was expected to extend to other tasks like  
32 predictions of ligand- and lipid-binding residues. The framework is also promising for other similar  
33 tasks, such as protein function prediction, modification sites of protein or DNA/RNA, and predicting  
34 properties for chemical compounds.

35

36 As this study focused on predicting binding residues from sequences, the prediction from protein  
37 structures (from experiments or structural modeling) was expected to further advance the predictions.  
38 This has been proved by our previous studies (Taherzadeh, et al., 2017).

39

40 The current sharing of information was simply based on a port of shared network. The  
41 performances were hurt by the differences between the binding types, for example, different ratios of  
42 positives and negatives, the preference of positive-charge residues in the DNA/RNA binding. Recently,

1 there is a significant advance in the transferred learning techniques (José Juan Almagro Armenteros,  
2 2019). We will employ these recent techniques in future studies.

### 3 **5. Acknowledgement**

4 This study was supported in part by the National Key R&D Program of China (2018YFC0910500),  
5 National Natural Science Foundation of China (U1611261, 61772566, and 81801132), Guangdong  
6 Frontier & Key Tech Innovation Program (2018B010109006, 2019B020228001) and Introducing  
7 Innovative and Entrepreneurial Teams (2016ZT06D211).

### 8 **6. References**

- 9 Taherzadeh G, Zhou Y, Liew A W C, et al. Sequence-based prediction of protein-carbohydrate binding sites using support  
10 vector machines[J]. *Journal of chemical information and modeling*, 2016, 56(10): 2115-2122.
- 11 Hanson J, Litfin T, Paliwal K, et al. Identifying molecular recognition features in intrinsically disordered regions of proteins by  
12 transfer learning[J]. *Bioinformatics*, 2019.
- 13 Jia J, Li X, Qiu W, et al. iPPI-PseAAC (CGR): Identify protein-protein interactions by incorporating chaos game representation  
14 into PseAAC[J]. *Journal of theoretical biology*, 2019, 460: 195-203.
- 15 Yan J, Kurgan L. DRNApred, fast sequence-based method that accurately predicts and discriminates DNA-and RNA-binding  
16 residues[J]. *Nucleic acids research*, 2017, 45(10): e84-e84.
- 17 Johnson M, Zaretskaya I, Raytselis Y, et al. NCBI BLAST: a better web interface[J]. *Nucleic acids research*, 2008, 36(suppl\_2):  
18 W5-W9.
- 19 Armenteros J J A, Tsirigos K D, Sønderby C K, et al. SignalP 5.0 improves signal peptide predictions using deep neural  
20 networks[J]. *Nature biotechnology*, 2019, 37(4): 420.
- 21 Hu J, Li Y, Zhang M, et al. Predicting protein-DNA binding residues by weightedly combining sequence-based features and  
22 boosting multiple SVMs[J]. *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, 2017, 14(6):  
23 1389-1398.
- 24 London N, Movshovitz-Attias D, Schueler-Furman O. The structural basis of peptide-protein binding strategies[J]. *Structure*,  
25 2010, 18(2): 188-199.
- 26 Lu W, Pieters R J. Carbohydrate-protein interactions and multivalency: Implications for the inhibition of influenza A virus  
27 infections[J]. *Expert opinion on drug discovery*, 2019, 14(4): 387-395.
- 28 Malik A, Firoz A, Jha V, et al. PROCARB: a database of known and modelled carbohydrate-binding protein structures with  
29 sequence-based prediction tools[J]. *Advances in bioinformatics*, 2010, 2010.
- 30 Sun M, Wang X, Zou C, et al. Accurate prediction of RNA-binding protein residues with two discriminative structural  
31 descriptors[J]. *BMC bioinformatics*, 2016, 17(1): 231.
- 32 Miao Z, Westhof E. Prediction of nucleic acid binding probability in proteins: a neighboring residue network based score[J].  
33 *Nucleic acids research*, 2015, 43(11): 5340-5351.
- 34 Pan X, Shen H B. RNA-protein binding motifs mining with a new hybrid deep learning based cross-domain knowledge  
35 integration approach[J]. *BMC bioinformatics*, 2017, 18(1): 136.
- 36 Petsalaki E, Russell R B. Peptide-mediated interactions in biological systems: new discoveries and applications[J]. *Current*  
37 *opinion in biotechnology*, 2008, 19(4): 344-350.
- 38 Remmert M, Biegert A, Hauser A, et al. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM  
39 alignment[J]. *Nature methods*, 2012, 9(2): 173.
- 40 Heffernan R, Paliwal K, Lyons J, et al. Single-sequence-based prediction of protein secondary structures and solvent  
41 accessibility by deep whole-sequence learning[J]. *Journal of computational chemistry*, 2018, 39(26): 2210-2216.

- 1 Zaman R, Chowdhury S Y, Rashid M A, et al. Hmmbinder: Dna-binding protein prediction using hmm profile based features[J].
- 2 BioMed research international, 2017, 2017.
- 3 Altschul S F, Madden T L, Schäffer A A, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search
- 4 programs[J]. Nucleic acids research, 1997, 25(17): 3389-3402.
- 5 Su, H., et al. Improving the prediction of protein–nucleic acids binding residues via multiple sequence profiles and the consensus
- 6 of complementary methods. Bioinformatics 2018;35(6):930-936.
- 7 Sujatha M S, Balaji P V. Identification of common structural features of binding sites in galactose-specific proteins[J]. Proteins:
- 8 Structure, Function, and Bioinformatics, 2004, 55(1): 44-65.
- 9 Sun T, Shao Y, Li X, et al. Learning Sparse Sharing Architectures for Multiple Tasks[J]. arXiv preprint arXiv:1911.05034, 2019.
- 10 Taherzadeh G, Zhou Y, Liew A W C, et al. Structure-based prediction of protein–peptide binding regions using Random
- 11 Forest[J]. Bioinformatics, 2017, 34(3): 477-484.
- 12 Zhou J, Lu Q, Xu R, et al. Cnnsite: Prediction of dna-binding residues in proteins using convolutional neural network with
- 13 sequence features[C]//2016 IEEE International Conference on Bioinformatics and Biomedicine (BIBM). IEEE, 2016: 78-85.
- 14 Wang L, Huang C, Yang M Q, et al. BindN+ for accurate prediction of DNA and RNA-binding residues from protein sequence
- 15 features[J]. BMC Systems Biology, 2010, 4(1): S3.
- 16 Yan J, Kurgan L. DRNAPred, fast sequence-based method that accurately predicts and discriminates DNA-and RNA-binding
- 17 residues[J]. Nucleic acids research, 2017, 45(10): e84-e84.
- 18 Yang J, Roy A, Zhang Y. BioLiP: a semi-manually curated database for biologically relevant ligand–protein interactions[J].
- 19 Nucleic acids research, 2012, 41(D1): D1096-D1103.
- 20 Wang Y, Huang M, Zhao L. Attention-based LSTM for aspect-level sentiment classification[C]//Proceedings of the 2016
- 21 conference on empirical methods in natural language processing. 2016: 606-615.
- 22 Tang Y, Liu D, Wang Z, et al. A boosting approach for prediction of protein-RNA binding residues[J]. BMC bioinformatics,
- 23 2017, 18(13): 465.
- 24 Zhang T, Faraggi E, Xue B, et al. SPINE-D: accurate prediction of short and long disordered regions by a single neural-network
- 25 based method[J]. Journal of Biomolecular Structure and Dynamics, 2012, 29(4): 799-813.
- 26 Zhang X, Liu S. RBPPred: predicting RNA-binding proteins from sequence using SVM[J]. Bioinformatics, 2016, 33(6):
- 27 854-862.
- 28 Zhao H, Wang J, Zhou Y, et al. Predicting DNA-binding proteins and binding residues by complex structure prediction and
- 29 application to human proteome[J]. PloS one, 2014, 9(5): e96694.
- 30 Taherzadeh G, Yang Y, Zhang T, et al. Sequence-based prediction of protein–peptide binding sites using support vector
- 31 machine[J]. Journal of computational chemistry, 2016, 37(13): 1223-1229.
- 32 Zhao H, Yang Y, Von Itzstein M, et al. Carbohydrate-binding protein identification by coupling structural similarity searching
- 33 with binding affinity prediction[J]. Journal of computational chemistry, 2014, 35(30): 2177-2183.
- 34 Huang Z, Xu W, Yu K. Bidirectional LSTM-CRF models for sequence tagging[J]. arXiv preprint arXiv:1508.01991, 2015.
- 35 Taherzadeh G, Yang Y, Zhang T, et al. Sequence-based prediction of protein–peptide binding sites using support vector
- 36 machine[J]. Journal of computational chemistry, 2016, 37(13): 1223-1229.
- 37 Zhou J, Lu Q, Gui L, et al. MTTFsite: Cross-cell-type TF Binding Site Prediction by using Multi-task Learning[J].
- 38 Bioinformatics, 2019.
- 39 Caruana R. Multitask learning[J]. Machine learning, 1997, 28(1): 41-75.
- 40 Litfin T, Yang Y, Zhou Y. SPOT-peptide: Template-based Prediction of Peptide-binding Proteins and Peptide-binding Sites[J].
- 41 Journal of chemical information and modeling, 2019, 59(2): 924-930.
- 42

1 **Table 1.** Summary of the benchmark datasets

<i>Types</i>	<i>Traning</i>		<i>Test</i>		<i>References</i>
	Chains	BD/ nonBD	Chains	BD/ nonBD	
<i>DNA</i>	309	6832/58270	47	875/8231	(Su, et al., 2019)
<i>RNA</i>	157	4627/30052	17	409/5448	(Su, et al., 2019)
<i>Peptide</i>	1115	14953/251708	125	1716/29154	(Zhou, 2016)
<i>CBH</i>	100	1028/25958	49	508/13230	(Ghazaleh Taherzadeh, 2016)

2

3 **Table 2.** Method performance on 10-fold cross validation and independent test

<b>Dataset</b>		<b>AUC</b>	<b>MCC</b>	<b>Sen</b>	<b>Spe</b>
<b>DNA</b>	CV	0.866	0.401	0.691	0.877
	Test	0.852	0.397	0.684	0.873
<b>RNA</b>	CV	0.857	0.369	0.635	0.936
	Test	0.836	0.361	0.612	0.932
<b>Peptide</b>	CV	0.760	0.304	0.346	0.956
	Test	0.758	0.299	0.344	0.953
<b>CBH</b>	CV	0.779	0.276	0.428	0.957
	Test	0.776	0.255	0.414	0.954

4

5 **Table 3.** The average AUC and MCC values on four independent datasets by employing or excluding  
6 individual feature group from the MTDsite model.

<b>Feature Groups<sup>a</sup></b>	<b>AUC</b>	<b>MCC</b>	<b>Feature Groups<sup>b</sup></b>	<b>AUC</b>	<b>MCC</b>
-	-	-	MTDsite	0.806	0.328
G-PSSM	0.742	0.247	-G-PSSM	0.757	0.281
G-HHM	0.739	0.244	-G-HHM	0.767	0.302
G-SPD3	0.61	0.201	-G-SPD3	0.794	0.320

7

<sup>a</sup> Performances based on individual feature groups;

8

<sup>b</sup> Performances by excluding individual feature groups.

9

10 **Table 4.** Comparison of AUC values of different MTDsite models for different binding molecules on  
11 Independent test

<b>Models</b>	<b>DNA</b>	<b>RNA</b>	<b>Peptide</b>	<b>CBH</b>
<b>MTD-DNA</b>	<b>0.852</b>	0.833	0.747	0.756
<b>MTD-RNA</b>	0.829	<b>0.836</b>	0.738	0.742
<b>MTD-PEP</b>	0.732	0.748	<b>0.758</b>	0.731
<b>MTD-CBH</b>	0.768	0.761	0.757	<b>0.776</b>
<b>MTD-Single</b> ( DNA )	0.822	0.801	0.601	0.718
<b>MTD-Single</b> ( RNA )	0.801	0.810	0.600	0.709
<b>MTD-Single</b> ( PEP )	0.531	0.520	0.734	0.51

MTD-Single      0.632      0.614      0.552      0.746  
 ( CBH )

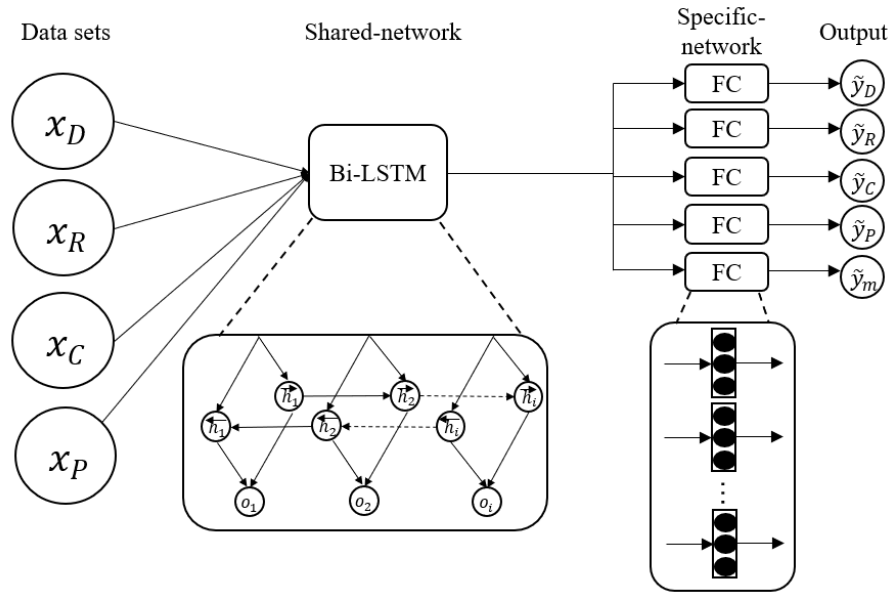
1

2 **Table 5.** Comparison with other methods on the independent test sets based on AUC, ACC, MCC.

3 Sensitivity and Precision

	Method	AUC	MCC	ACC	Sen	Pre
<b>DNA</b>	BindN+	0.806	0.256	NA	NA	NA
	SVMnuc	0.833	0.336	NA	0.26	0.51
	NucBind	0.834	0.355	NA	0.29	0.50
	DRNApred	0.757	0.232	NA	0.32	0.25
	COACH-D	0.810	0.266	NA	NA	0.70
	MTDsite-single	0.822	0.379	0.901	0.562	0.628
	MTDsite	<b>0.852</b>	<b>0.397</b>	<b>0.916</b>	<b>0.684</b>	<b>0.734</b>
<b>RNA</b>	BindN+	0.738	0.219	NA	NA	NA
	SVMnuc	0.784	0.205	NA	0.19	0.26
	NucBind	0.801	0.210	NA	0.18	0.29
	DRNApred	0.687	0.056	NA	0.09	0.08
	COACH-D	0.711	0.356	NA	0.30	<b>0.45</b>
	MTDsite-single	0.810	0.264	0.945	0.455	0.345
	MTDsite	<b>0.836</b>	<b>0.361</b>	<b>0.931</b>	<b>0.612</b>	0.351
<b>Peptide</b>	VisGrid	NA	0.224	0.609	0.38	NA
	PinUp	NA	0.0071	0.547	0.207	NA
	PepSite	NA	0.141	0.598	0.221	NA
	Peptimap	NA	0.262	0.629	0.436	NA
	SPRINT	0.711	0.279	0.662	<b>0.642</b>	NA
	MTDsite-single	0.734	0.218	0.947	0.32	0.29
	MTDsite	<b>0.758</b>	<b>0.299</b>	<b>0.956</b>	0.344	<b>0.366</b>
<b>CBH</b>	SPRINT-CBH	0.772	<b>0.285</b>	0.961	<b>0.223</b>	NA
	MTDsite-single	0.746	0.221	0.958	0.374	0.177
	MTDsite	<b>0.776</b>	0.255	<b>0.963</b>	0.209	0.171

1

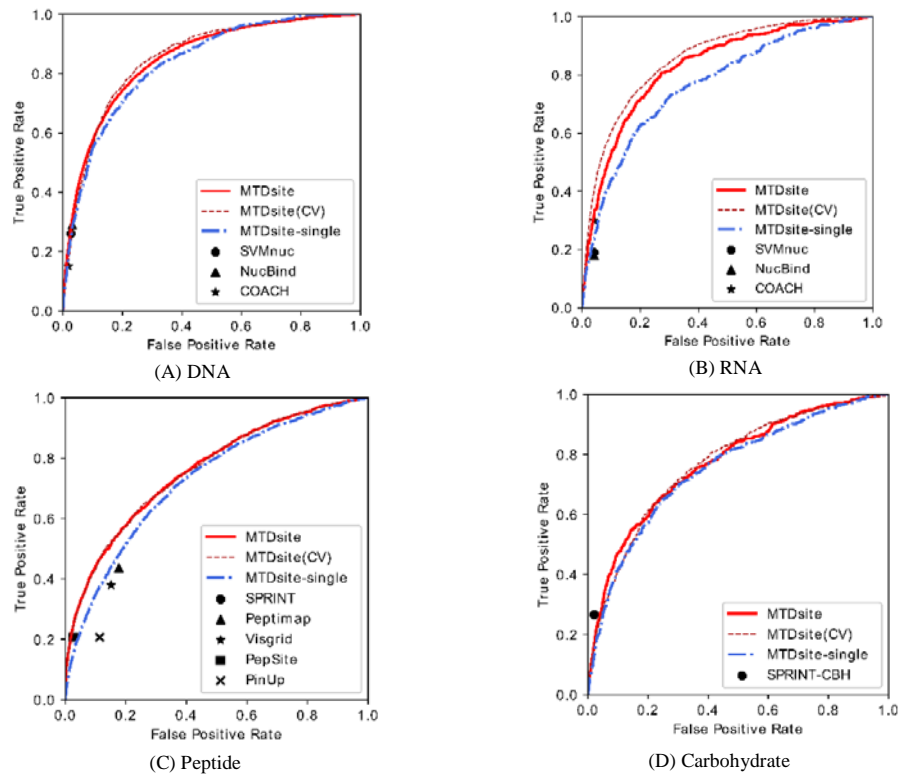


2

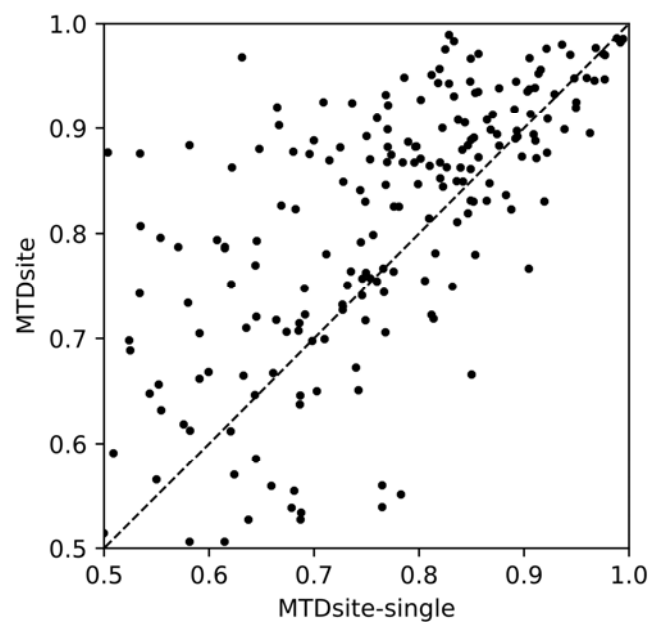
**Fig 1** The general architecture of the MTDsite

3

4



**Figure 2.** The receiver of characteristic curves for MTDsite (10 fold cross-validations and independent tests) and MTDsite-single on the independent tests for (A) DNA, (B) RNA, (C) peptide, and (D) carbohydrate datasets. The reported results by other methods were labelled on respective plots.

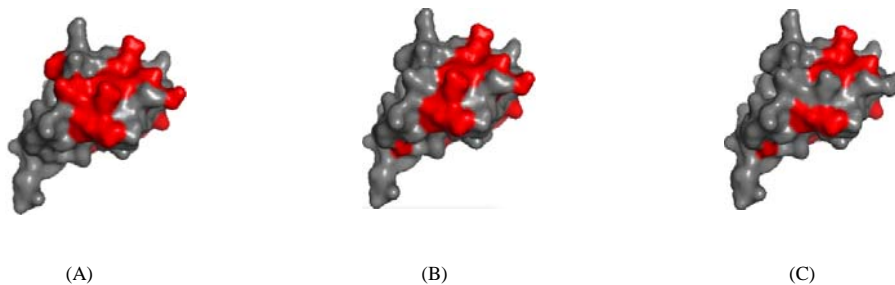


1

**Figure 3.** The comparison of AUC values by MTDsite and MTDsite-single on independent dataset, each point representing one protein chain.

2

3



4

5

**Figure 4.** (A) The actual binding residues, and the predicted binding residues by (B) MTDsite and (C) MTDsite-single for a Restriction-Modification controller DNA-binding protein (PDBid: 3s8qB).