

# RPI-Pred: Predicting ncRNA-protein interaction using sequence and structural information

V Suresh<sup>1\*</sup>, Liang Liu<sup>1\*</sup>, Donald Adjeroh<sup>2</sup>, and Xiaobo Zhou<sup>1\*\*</sup>

<sup>1</sup>Department of Radiology, Wake Forest University Health Science, Medical Center Boulevard, Winston-Salem, NC 27157, USA

<sup>2</sup>Lane Department of Computer Science and Electrical Engineering, West Virginia University, Morgantown, WV 26505

\* Authors contributed equally to this work

\*\* Corresponding author (Xiaobo Zhou): Department of Radiology, Wake Forest School of Medicine, Medical Center Boulevard, Winston-Salem, NC 27157, USA. Tel: +1 336-713-1879; Fax: +1 336-713-5891; Email: xizhou@wakehealth.edu

## ABSTRACT

RNA-protein complexes are essential in mediating important fundamental cellular processes such as, transport and localization. In particular, ncRNA-protein interactions play an important role in post-transcriptional gene regulation like mRNA localization, mRNA stabilization, poly-adenylation, splicing and translation. The experimental methods to solve RNA-protein interaction prediction problem remain expensive and time-consuming. Here, we present the RPI-Pred (RNA-protein interaction predictor), a new support-vector machine (SVM) based method, to predict protein-RNA interaction pairs, based on both the sequences and structures. The results show that RPI-Pred can correctly predict RNA-protein interaction pairs with ~94% prediction accuracy when using sequence and experimentally determined protein and RNA structures, and with ~83% when using sequences and predicted protein and RNA structures. Further, our proposed method RPI-Pred was superior to other existing ones by predicting more experimentally validated ncRNA-protein interaction pairs from different organisms. Motivated by the improved performance of RPI-Pred, we further applied our method for reliable construction of ncRNA-protein interaction networks. The RPI-Pred is publically available at: <http://ctsb.is.wfubmc.edu/projects/rpi-pred>.

## INTRODUCTION

RNA-protein interactions (RPI) play a crucial role in fundamental cellular processes such as human diseases (1), viral replication and transcription (2,3) and pathogen resistance in plants (4-6). Recent high throughput techniques produce remarkable evidences to prove that protein can interact with RNA to mediate different kinds of cellular functions. During the post-transcriptional regulation process, RPI complex interacts with targeted mRNAs and/or non-coding RNAs (ncRNAs) to regulate cellular functions such as RNA splicing, RNA transport, RNA stability and RNA translation (7-9). Experimental studies on RPI reveal that many functional ncRNAs play pivotal roles in gene expression and regulation (10-16). Although a few individual ncRNAs have been well studied, e.g., HOTAIR (17), MALAT-1 (18), and Xist

1  
2  
3 (19), the majority are still not well understood. Over 30,000 ncRNAs have been identified and this number  
4 is expected to increase every year (14,15,20). Currently, NPInter (21) is the only database, which  
5 provides the functional information for all the experimentally validated ncRNA-protein interactions. The  
6 experimental techniques are generally time consuming and expensive. Our understanding of function of  
7 individual ncRNAs is far outpaced by the sheer volume and diversity of the available data. Furthermore,  
8 our understanding of ncRNA-protein interactions in gene regulatory networks is very limited, especially  
9 when compared to the regulatory roles of protein-protein and DNA-protein complexes. This is because  
10 the advances in genomics and proteomics techniques have resulted in tremendous amounts of data on  
11 protein-protein and protein-DNA interactions (22-24); however, much less information is available on  
12 ncRNA-protein interactions.  
13  
14  
15  
16  
17  
18

19 In despite of the increasing amount (~400) of successfully identified RNA binding proteins (RBP) in the  
20 human genome (25,26), we still lack a complete understanding of RPI complexes and their roles in post-  
21 transcriptional regulatory networks (7,27). Although the sequence-homology based approaches, such as  
22 BLAST (28-30) and PFAM (31-33), helped in detecting the functional regions (binding domains) of  
23 proteins and therefore the possible functions, these approaches lack the ability to identify the interacting  
24 partners (RNAs) for a given protein, or determine whether a given pair of protein and RNA can form  
25 interaction or not. To our knowledge, currently very few computational approaches are available to predict  
26 RNA-protein interactions. One of the first computational methods for predicting ncRNA-protein interaction  
27 was reported in 2011 by Pancaldi and Bähler (34). They trained random forest (RF) and support vector  
28 machine (SVM) classifiers using more than 100 features extracted from protein secondary structure and  
29 localization, protein and gene physical properties, and untranslated regions (UTRs). Thereafter,  
30 catRAPID (35) was developed by exploiting the physicochemical properties including secondary  
31 structure, hydrogen bonding and van der Waals propensities. Next, Muppirala et al. (36) introduced a  
32 method called RPISeq, which was constructed by using the features derived from protein and RNA  
33 sequences. They also trained RF and SVM classifiers using 3-mer and 4-mer conjoint triad features for  
34 amino acid and nucleotide sequences, respectively (37). Wang et al. (38) proposed an approach based  
35 on Naïve Bayes (NB) and Extended Naïve Bayes (ENB) classifiers using the same datasets and similar  
36 triad features reported in Muppirala et al's work. More recently, Lu et al. (39) proposed a method called  
37 'IncPro' for predicting ncRNA-protein associations, using Fisher linear discriminant approach. His training  
38 features were three types of classical protein secondary structures, hydrogen-bond and Van der Waals  
39 propensities, as well as six types of RNA secondary structures.  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50

51 Muppirala et al. (36) and Wang et. al. (38) proposed their methods based on sequence features to predict  
52 RPI interactions. Other methods (34,35) have also been proposed by combining sequence and structural  
53 features. IncPro (39) method used protein and RNA secondary structures, hydrogen-bond and Van der  
54 Waals propensities. However, none of above methods used the high-order three dimensional protein and  
55 RNA structure features, which are known to be the key of their possible functions (40).  
56  
57  
58  
59  
60

1  
2  
3 In the present work, we presented a computational approach to predicting protein-RNA interaction pairs,  
4 and/or identify the binding partners of a given protein or RNA from candidates. In addition to sequence  
5 features, we combine the high-order structures of both proteins and RNAs, for a comprehensive  
6 understanding of RPI interactions. We consider the protein structures in terms of 16 structural fragments  
7 called protein blocks (41). The protein blocks (PB) provide a more accurate representation of known  
8 protein structures than classical three state protein secondary structures ( $\alpha$ -helix,  $\beta$ -sheet, and coil), and  
9 have been applied in many protein structure-based analysis (42,43). For the RNA high-order structure,  
10 we considered five classes of RNA secondary structures (RSS), namely stem, hairpin, loop, bulges, and  
11 internal loop. These PB and RSS were combined with their corresponding amino acid and nucleotide  
12 sequences. Using these features, we developed a support vector machine (SVM) based machine  
13 learning approach, RPI-Pred, to predict protein-RNA interactions. Our training database was constructed  
14 using sequence and experimentally validated structures of proteins and RNAs from the Protein Data Bank  
15 (PDB) (44). We also used sequence and predicted structures to test our RPI-Pred on different datasets,  
16 such as RPI369 and RPI2241, and ncRNA-protein interaction datasets such as RPI367, RPI13243 and  
17 NPInter10412 (21). We extended our analysis to construct an *in silico* network to study potential  
18 interactions between proteins and ncRNAs, which can help us in further understanding of ncRNA's  
19 functions. Finally, a web server for this proposed method was also developed and freely accessed at  
20 <http://ctsb.is.wfubmc.edu/projects/rpi-pred>.  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30

## 31 MATERIALS AND METHODS

### 32 Work flow

33  
34  
35 Figure 1 shows the work flow for the development of RPI-Pred method. The proposed method includes  
36 three steps; (1) extraction of sequence and structure features for protein and RNA to develop the RPI-  
37 Pred prediction method, (2) prediction of ncRNA-protein interactions, and (3) construction of the *in silico*-  
38 based biological network on the predicted results in step 2. Step 1 includes various processes, such as;  
39 construction of the training dataset, removal of redundant RNA-protein pairs, feature extraction from  
40 sequences and structures in the training dataset, and development of the 'RPI-Pred' model. Step 2  
41 includes the feature extraction in terms of primary sequence and predicted structure for given protein  
42 and/or ncRNA and ncRNA-protein interaction prediction using RPI-Pred. Step 3 consists of construction  
43 of interaction networks based on RPI-Pred interaction predictions. More detailed descriptions for each  
44 step are given below.  
45  
46  
47  
48  
49  
50

### 51 Training datasets

52  
53 To develop the RPI-Pred method, first, we build a non-redundant training dataset of RNA-protein  
54 interaction complexes by parsing the Nucleic Acid Database (NDB) (39) and the protein-RNA interface  
55 database (PRIDB) (45). The former provides data for RNA-protein complexes, whereas the latter provides  
56  
57  
58  
59  
60

1  
2  
3 atomic-interfaces for RNA-protein interacting pairs. A total of 1,560 RPI complexes available in NDB (as  
4 of Feb. 1, 2014) were used in this study. We extracted the atomic and chain interfaces for 1,336  
5 complexes from PRIDB, resulting in 13,163 protein and 2,715 RNA chains. These 1,336 complexes were  
6 further used to construct our training dataset, which consist of both possible positive and negative protein-  
7 RNA pairs.  
8  
9

10  
11 The procedure for constructing the training dataset included removing redundant protein/RNA pairs  
12 through sequence similarity criteria, as follows. For instance, the RNA-protein interaction complex with the  
13 PDB id '1a9n' has four protein chains (A, B, C, D) and two RNA chains (Q, R), respectively. We obtained  
14 the possible interaction pairs from PRIDB as A-Q, B-Q, C-R and D-R. Then the homologous RNA-protein  
15 pairs (i.e. similar protein chains interacting with similar RNA chains) were removed by searching the  
16 sequence similarity between protein (RNA) sequences. In this study, we used EMBOSS needle program  
17 (46) with the standard sequence identity cut-off  $\geq 30\%$  to remove proteins (and RNA chains) with a high  
18 sequence similarity. In the current example with '1a9n', the protein chains A&C, B&D, and RNA chains  
19 Q&R are 100% sequence-similar, therefore we removed the redundant protein pairs. Finally, A-Q and B-  
20 Q were identified as non-redundant RNA-protein pairs.  
21  
22

23  
24 The above selected non-redundant pairs were further tested for atomic interactions with a distance  
25 threshold (3.40 Å). This distance threshold helped to strengthen positive pairs in the training dataset by  
26 including only strongly interacting RPI pairs. Different thresholds have been used to distinguish the  
27 binding RNA-protein pairs from non-binding ones (35,36,38,39,47). We used the threshold 3.40 Å (47).  
28 The threshold 3.4 Å is reasonable and sufficient to cover 'strong' and 'moderate' hydrogen bonds and  
29 energy-rich van der Waals contacts (48). Therefore we set the threshold (3.4 Å) to distinguish the strongly  
30 interacting protein-RNA pairs (positive pairs) and weakly interacting protein-RNA pairs (negative pairs). In  
31 the above given example, the pair B-Q, which had at least two atoms, one from protein and another from  
32 RNA, with distance  $\leq 3.40$  Å, was considered as a positive pair. The pair A-Q, which had no atom-atom  
33 distance within the threshold, was considered as a negative pair. This procedure was applied to all 1,336  
34 RNA-protein complexes to identify the positive and negative pairs. Further, the peptides (protein with  
35 sequence length  $< 25$  amino acids) and small RNA (with sequence lengths  $< 15$  nucleotides) were  
36 excluded from these positive and negative datasets.  
37  
38

39  
40 As a result, we obtained a training dataset, namely RPI1807, with 1,807 positive pairs (consisting of 1,807  
41 protein and 1,078 RNA chains) and 1,436 negative pairs (with 1,436 protein and 493 RNA chains). The  
42 positive and negative pairs of RPI1807 are shown in Supplementary Table S1.  
43  
44

### 45 46 47 48 49 50 51 52 53 **Test datasets**

54  
55 The RPI-Pred was tested with different datasets, including four datasets from previous studies and the  
56 new dataset constructed in this work. The first three datasets were obtained from (36) and denoted as  
57 RPI369, RPI2241 and RPI13243 based on number of protein-RNA pairs (369, 2,241 and 13,243),  
58  
59

1  
2  
3 respectively. In (36), the first two datasets were used as training datasets to develop the classifier and the  
4 third was used to evaluate the classifier. RPI13243 consists of 13,243 RNA-protein interactions, which  
5 includes all 5,166 protein-mRNA interactions published by Hogan et. al., (49). The fourth dataset  
6 (denoted as RPI367) consists of 367 protein-ncRNA interactions, constructed by Wang et. al., (38) from  
7 the NPInter database (21).  
8  
9

10  
11 The pairs in training dataset RPI1807 were also used to construct the fifth test dataset. In this case, the  
12 RPI-Pred was applied to predict RNA-protein interactions by using sequence and predicted structures for  
13 both protein and RNA, instead of the experimentally determined structures obtained from PDB (44). The  
14 sixth test dataset was extracted from the NPInter database (21), namely RPI10412, including 10,412  
15 ncRNA-protein interaction pairs from six different model organisms. These ncRNA and protein pairs had  
16 been experimentally determined to have physical associations and listed in the 'ncRNA binding protein'  
17 category.  
18  
19

### 20 21 22 **Protein blocks and RNA secondary structure**

23  
24 In addition to primary sequences, we used structures, obtained from experimental determinations  
25 (available in PDB) or theoretical predictions, of both proteins and RNAs in our RPI-Pred. A protein three-  
26 dimensional (3D) structure can be represented by 16-letter one-dimensional structural fragments, called  
27 protein blocks (PB) (41). The PDB-2-PB database (50) provides the PB information based on the  
28 experimentally solved protein structures available in PDB (44). We used the PDB-2-PB to retrieve the 16-  
29 letter PB structure features for each protein in our training dataset.  
30  
31

32  
33 We used the 3DNA suite (51) to extract the RNA secondary structures (RSS) from the corresponding 3D  
34 structures (44). We used five category of RSS, namely Stem (S), Hairpin (H), Loop (L), Bulges (B) and  
35 Internal loop (I) to construct our RPI-Pred method. In this study, the pseudoknot RNA structures were not  
36 considered, because they were less numbers in our training dataset. These PB and RSS were further  
37 combined with corresponding protein and RNA sequences to represent proteins and RNAs, respectively.  
38 These combined sequence-structure features were used to develop our RPI-Pred method for predicting  
39 potential ncRNA-protein interaction pairs.  
40  
41  
42  
43  
44

### 45 46 **Representation of sequence and structural features**

47  
48 The sequence and structural features of protein and RNA used in this work were represented as follows.  
49 The protein sequence of 20 amino acids were classified into 7 groups (7-letter reduced sequence  
50 alphabets) according to their dipole moments and side chain volume: {A,G,V}, {I,L,F,P}, {Y,M,T,S},  
51 {H,N,Q,W}, {R,K}, {D,E} and {C}. Then, we combined these 7-letter sequence features with the 16-letter  
52 PB structure representations, resulting in 112 (7x16) possible combinations. The normalized frequencies  
53 of 112 combinations formed the 112 protein vectors. Similarly, RNA sequence and RSS representations  
54 resulted in 4x5 possible combinations (4 for the nucleotide types; A, U, C and G and 5 for the RSS), and  
55  
56  
57  
58  
59  
60

1  
2  
3 normalized frequencies of these 20 combinations resulted in 20 RNA vectors. The labels of sequence and  
4 structural features obtained from the proteins and RNA are given in Supplementary Table S2. In  
5 summary, to construct the RPI-Pred method, we used the sequence and structure features of 132  
6 vectors, in which the first 112 vectors represented proteins and the remaining 20 vectors represented  
7 RNAs.  
8  
9

### 10 11 **Support Vector Machine (SVM) classifier**

12  
13  
14 The support vector machine (SVM) approach is a popular supervised machine learning technique used  
15 for many classification and regression problems (52). Here, we applied a well-known SVM classifier,  
16 LIBSVM-3.17 package (53) implemented as a standalone in-house program, to perform RPI prediction.  
17 We constructed our RPI-Pred method using 132-feature vectors that represent protein and RNA  
18 sequences and structures. RPI-Pred was optimized using different kernel functions with their  
19 corresponding parameters. We selected the 'polynomial' kernel function, which gave better prediction  
20 accuracy than others. RPI-Pred was trained to efficiently predict protein and RNA interaction pairs with  
21 the following optimized parameters:  $C=1000$ ,  $\gamma=1$ ,  $\text{cofe0}=1$  and  $\text{degree}=4$ .  
22  
23  
24  
25

### 26 **Predicting protein blocks and RNA secondary structure**

27  
28  
29 Since many proteins and RNAs have not been experimentally solved, we must use theoretical  
30 approaches to predict their structures. Many research groups have proposed PB prediction methods  
31 (42,54). In this work, we used the PB-kPRED method (55) to predict the PB structures for proteins  
32 included in all the test datasets.  
33  
34

35  
36 Likewise, RSS can also be predicted with available RNA structure prediction methods (56-64). Here we  
37 selected RNAfold from the Vienna package - an in-house standalone program (64) to predict the RSSs  
38 for RNAs in our test datasets. The predicted PB and RSS were combined with the corresponding amino  
39 acid and nucleotide sequences, respectively, and used in our RPI-Pred method.  
40  
41

### 42 **Performance evaluation**

43  
44  
45 The performance of RPI-Pred was evaluated using 10-fold cross validation (10-fold CV) approach. To  
46 perform this test, the training dataset was divided into 10 subsets of equal size. Each subset was used for  
47 testing, while the remaining nine subsets were used for training. This process was repeated 10 times to  
48 cover all possibilities. Finally, we recorded the average performance over all ten testing subsets. We  
49 evaluated the prediction performance by using Precision (PRE), Recall (REC), F-Score (FSC) and  
50 Accuracy (ACC), defined as follows:  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

$$Precision = \frac{tp}{tp + fp}$$

$$Recall = \frac{tp}{tp + fn}$$

$$F - Measure = \frac{2 * Precision * Recall}{Precision + Recall}$$

$$Accuracy = \frac{tp + tn}{tp + fp + fn + tn} * 100$$

Where,  $tp$  and  $tn$  denote the number of correctly predicted positive and negative pairs, respectively, and  $fp$  and  $fn$  denote the number of wrongly predicted positive and negative pairs, respectively. The area under curve (AUC) of the receiver operation characteristic (ROC) curve was calculated using a 10-fold cross validation (CV). The AUC ranges from 0 to 1, with 1 indicating the best prediction.

## RESULTS AND DISCUSSION

We built the RPI-Pred method for identifying binding partners of proteins or RNAs. In this section, we tested the performance of RPI-Pred on different test datasets, including RPI1807, RPI2241 and RPI369, and compared with previous methods. We also applied the RPI-Pred method to a large ncRNA-protein interaction dataset and the predicted results were compared with other existing approaches.

### Performance of RPI-Pred with experimentally determined structures

The performance of our RPI-Pred method was evaluated using the 10-fold CV on RPI1807, RPI2241 and RPI369 datasets. The experimentally validated structures were extracted from PDB database (44). The performance of RPI-Pred was evaluated by calculating ACC, AUC, PRE, REC and FSC for each dataset. Our RPI-Pred successfully predicted the RNA-protein pairs on RPI1807 dataset with prediction accuracy (ACC) of 93%. The other measurements (AUC, PRE, REC and FSC) were observed as 0.97, 0.94, 0.95 and 0.95, respectively. The high prediction accuracy indicated that our method based on sequence and structure was reliably predicted RNA-protein interaction.

Similarly, the performance of RPI-Pred was evaluated using the positive and negative pairs of RPI2241 and RPI369 datasets. The positive pairs were directly adopted from RPI2241 and RPI369 datasets and their corresponding negative pairs were generated by following the steps reported in (38). Then, the RPI-Pred was applied on these datasets, to correctly predict all positive and negative pairs. The RPI-pred reached the prediction accuracy (ACC) of ~84% for the RPI2241 dataset. The AUC, PRE, REC and FSC were also observed as 0.89, 0.88, 0.78 and 0.83, respectively. Applying RPI-pred on the RPI369 dataset resulted in a prediction accuracy of ~92%, and AUC, PRE, REC and FSC of 0.95, 0.89, 0.89 and 0.89, respectively. The results in our newly constructed dataset RPI1807 showed ~10% and 2% increase in accuracy over RPI2241 and RPI369 results, respectively.

### Comparison of RPI-Pred with existing methods

We compared the performance of the RPI-Pred with Muppirala's method (36) on the RPI2241 and RPI369 datasets, respectively, using 10-fold CV. This comparison shows the prediction performance of the RPI-Pred method and the importance of structures in the prediction of RNA-protein interaction.

We compared our RPI-Pred results obtained from RPI2241 and RPI369 datasets, (RPI2241-RPI-Pred and RPI369-RPI-Pred, respectively) using AUC, PRE, REC, FSC and ACC measurements. The comparison results are shown in Table 1. We denoted Muppirala et al's results as RPI2241-SVM, RPI369-SVM, RPI2241-RF and RPI369-RF, based on the two classifiers, SVM and RF, and the used training databases.

As shown in the Table 1, our RPI2241-RPI-Pred result showed a prediction accuracy of 84%. This is ~3% and ~5% less than the results from the RPI2241-SVM (~87%) and RPI2241-RF (~89%) results, respectively. On the other hand, the RPI369-RPI-Pred result showed a prediction accuracy of 92%, implying an increase of ~24%, and ~18% over RPI369-SVM (~72%) and RPI369-RF (~76%) results, respectively.

These results illustrate that our RPI369-RPI-Pred could outperform RPI369-SVM and RPI369-RF classifiers in predicting the pairs of non-ribosomal RNA interacting with protein. Also inclusion of structure features can improve the RNA-protein interaction prediction (36). Slightly lower prediction accuracy was observed for the RPI2241 dataset, which contained more ribosomal RNAs paired with proteins. Ribosomal RNA structures are more likely to contain pseudoknot structures (40,65-67). However, the RNAfold which was used in this work does not have the ability to predict such structures, and thus the proposed RPI-Pred cannot consider pseudoknot structures. This may affect the correct prediction of ribosomal RNAs interacting with proteins. Further, Muppirala et. al., (36) used 3-mer sequence features while our RPI-Pred uses 1-mer features of sequence and structure to perform RNA-protein interaction prediction. This leads to an increased dimensionality in feature space, which could lead to an improved prediction. However, this also results in a more complex model, and a significantly longer processing time.

Recently, Wang et. al., (38) developed RNA-protein interaction prediction method using Naive Bayes (NB) and Extended Naive Bayes (ENB) classifiers on RPI2241 and RPI369 datasets. We also compared the performance of our prediction of RPI-Pred method on these two datasets with Wang et. al., (38) reported results. For this comparison, we grouped the results in (38) into four categories: RPI2241-NB, RPI369-NB, RPI2241-ENB and RPI369-ENB, based on the dataset and classifiers used. The RPI-Pred method had an increased prediction accuracy of ~9% and ~10% over RPI2241-NB (75.7%) and RPI2241-ENB (74.0%), respectively. Our RPI-Pred method on RPI369 dataset also showed an increased prediction accuracy of ~15% and ~17% over RPI369-NB (77.7) and RPI369-ENB (75.0%), respectively.



### RPI-Pred performance with sequences and predicted structures

We further tested the RPI-Pred method by using sequences and predicted structures, instead of experimentally determined structures. This experiment was necessary due to the lack of experimentally validated structures for many RNAs, especially ncRNAs, and proteins. The objective was to understand to what extent RPI-Pred performance might be affected by using predicted (rather than known) structures. To perform this analysis, we used the RPI-model constructed based on the RPI1807 dataset and tested within the same dataset. We observed a prediction accuracy of ~83%. For the remaining measurements AUC, PRE, REC and FSC, the performance was 0.89, 0.79, 0.94 and 0.86, respectively. Compare with the results obtained using known structures as reported earlier (0.97, 0.94, 0.95 and 0.95 for AUC, PRE, REC and FSC, respectively). In particular, the prediction accuracy (ACC) decreased by nearly 10% compared with the performance of RPI-Pred on RPI1807 dataset with known structures. We can observe similar decreases in the other performance measures. Expectedly, precision was significantly reduced when using predicted structures, while there was little or no impact on precision. Since, there is no experimental structural features were available for the rest of our test datasets (i.e., RPI367, RPI13243 and NPInter10412) we used the predicted protein blocks and RNA secondary structures in order to perform the RPI prediction.

### Performance of RPI-Pred on predicting ncRNA-protein interaction pairs

Although most of the DNA transcripts are non-coding RNAs (ncRNAs), very few have known functions. The ncRNA function can be predicted by identifying the different interacting partners such as DNA, RNA and protein. It is currently believed that ncRNAs interact with proteins and then perform their regulatory functions such as chromatin remodeling, to enhance or suppress gene expression (17-20). Therefore, studying ncRNA-protein interactions can reveal the importance of ncRNA in the post-transcriptional regulatory process. Very few computational studies (34-39) have been developed to predict the binding partner either for protein or ncRNA using both sequence and structural information. Here we investigated the performance of our RPI-Pred in terms of predicting the binding partner for a given protein or ncRNA using both sequence and high-order structural information. We tested RPI-Pred method on RPI367, RPI13243 and NPInter10412 datasets, which contain ncRNA-protein interaction pairs. The results obtained from our RPI-Pred method for these three datasets were further compared with those results obtained by other exiting approaches.

Our RPI-Pred method was first tested using small RPI367 dataset (38), consisting of 367 ncRNA-protein interaction pairs across six different model organisms; *Caenorhabditis elegans*, *Drosophila melanogaster*, *Escherichia coli*, *Homo sapiens*, *Mus musculus* and *Saccharomyces cerevisiae*. The RPI-Pred performances for the above six model organisms were given in Table 2. We compared the RPI-Pred prediction results with Wang et al's four classifiers (RPI369-NB (62%), RPI369-NB (77%), RPI2241-ENB (66%) and RPI2241-ENB (79%) results on the RPI367 dataset. Our RPI-Pred method outperformed with

1  
2  
3 a prediction accuracy of 89% (328 out of 367 pairs were correctly predicted), and none of Wang et al's  
4 classifiers performed at more than 80% accuracy (38). When compared prediction results from Wang et  
5 al's classifiers, RPI369-NB (62%), RPI369-NB (77%), RPI2241-ENB (66%) and RPI2241-ENB (79%), our  
6 RPI-Pred result showed greater improvement in prediction accuracy by 27%, 22%, 23% and 10%,  
7 respectively. Especially, our RPI-Pred predicted more ncRNA-protein interaction pairs (with the prediction  
8 accuracies of 100%, 92%, 96% and 91% for *C. elegans*, *D. melanogaster*, *E. coli* and *H. sapiens*) than  
9 each of Wang et al's classifiers.  
10  
11

12  
13  
14 Our next ncRNA-protein interaction prediction analysis was performed on a larger dataset used by  
15 Muppirla (36) , which contains 13,243 ncRNA-protein interaction pairs. Our RPI-Pred method correctly  
16 predicted 12,240 out of 13,243 interaction pairs of this dataset with the prediction accuracy of ~92%. Our  
17 method showed ~27% and ~14% increases in accuracy, compared with Muppirla reported accuracies  
18 (65% and 78% with the SVM and RF classifiers, respectively).  
19  
20  
21

22  
23 Finally, we tested the ability of our RPI-Pred method to predict ncRNA-protein interaction pairs in the  
24 currently available NPInter database (version 2.0). NPInter database (21) is the only the resource that  
25 provides the experimentally verified ncRNA-protein interaction pairs for different model organisms. Our  
26 new NPInter10412 dataset consists of 10,412 ncRNA-protein interaction pairs of six model organisms  
27 from NPInter database. Since there are no experimentally validated negative pairs available in NPInter  
28 database, we randomly shuffled (ie. by keeping the RNA fixed and reordered the proteins) all the positive  
29 pairs in NPInter10412 dataset to make the negative dataset. The performance of our RPI-Pred on  
30 NPInter10412 dataset was evaluated by predicting correct positive and negative ncRNA-protein  
31 interaction pairs. The RPI-Pred had a prediction accuracy of ~87% on NPInter10412 dataset. The  
32 remaining measurements (PRE, REC and FSC) were observed as 0.85, 0.90 and 0.87, respectively.  
33  
34  
35  
36  
37

38 Among the tested 10,412 positive ncRNA-protein interaction pairs, our RPI-Pred method correctly  
39 predicted 9,335 interaction pairs with the accuracy of ~90%. The RPI-Pred method predicted fewer  
40 ncRNA-protein pairs for *C. elegans*, *D. melanogaster* and *E. Coli* with the prediction accuracies of ~78%,  
41 ~77% and ~76%, respectively. The RPI-Pred prediction accuracies for *H. sapiens*, *M. musculus* and  
42 *S.cerevisiae* were ~89%, ~97% and ~82%, respectively. Table 3 shows the total number of positive  
43 ncRNA-protein pairs tested for each organism and the total number of pairs correctly predicted with our  
44 RPI-Pred method. The RPI-PRed prediction results in the NPInter10412 dataset for each organism are  
45 shown in Supplemental Table S3.  
46  
47  
48  
49  
50

51 We further analyzed incorrect the prediction results of some specific complexes in each organism. We  
52 found a few cases in four organisms (*D. melanogaster*, *E. coli*, *H. sapiens*, and *M. musculus*. Most of the  
53 cases, these incorrect predictions were observed for same protein that interacts with different ncRNAs.  
54 Among 21 false-negative ncRNA-protein pairs in *D. melanogaster*, 16 involved three proteins (Uniprot  
55 ID's: P17133, P26017 and Q9V3W7). Similarly, among 48 incorrect predictions in *E. coli*, 18 involved  
56  
57  
58  
59  
60

1  
2  
3 three proteins (P0AFZ3, P0C077 and P10121). The RPI-Pred method also failed to predict the pairs  
4 involving in two proteins (P19338 and P62312) in *H. Sapiens*. In the above few mentioned protein-RNA  
5 interactions, one protein interacts with multiple ncRNAs. This is due to wrongly predicted protein  
6 structures. In this proposed approach, our RPI-Pred method uses the features extracted from predicted  
7 structures of proteins and ncRNAs. Hence, our RPI-Pred prediction performance will be strongly  
8 influenced by the protein or RNA structure prediction approaches.  
9  
10  
11

### 12 **Comparison of RPI-Pred with RPISeq for ncRNA-protein interaction prediction**

13  
14  
15 Performance of our RPI-Pred method on NPInter10412 dataset was further compared with the results  
16 obtained from existing approaches. We tested the NPInter10412 dataset with RPISeq (36) and then the  
17 predicted results were compared with our RPI-Pred results. The stand-alone and locally implemented  
18 RPISeq program was obtained from the developers. Here, the RPISeq models, developed based on RF  
19 and SVM classifiers with RPI2241 and RPI369 datasets, were used to measure the RPI prediction  
20 performance on the NPInter10412 dataset. To perform this test we used both the NPInter10412 positive  
21 pairs and the corresponding shuffled negative pairs. As previously reported (36), interactions with  
22 probability score  $\geq 0.5$  from RPISeq were considered as correct prediction. We further compared the  
23 predicted results of RPI2241-SVM, RPI369-SVM, RPI2241-RF and RPI369-RF models and the  
24 comparisons are given in Table 4.  
25  
26  
27  
28  
29  
30

31 The prediction performance of the RPI2241-RF model on NPInter10412 dataset was 0.50, 0.97 0.66 and  
32 0.50 for PRE, REC, FSC and ACC, respectively. The accuracy of RPISeq on NPInter10412 dataset was  
33 just ~50%. This performance is very low, when compared to our RPI-Pred accuracy ~87%. However,  
34 RPI2241-RF correctly predicted more positive pairs as true positives (10157 out of 10412) and very few  
35 interactions were predicted as true negatives (288 out of 10412). The other measurements PRE, REC  
36 and FSC were observed as 0.50, 0.98 and 0.66, respectively. We also observed the similar trend in the  
37 analysis of RPI2241-SVM results. The performance of RPI2241-SVM was ~49% with more true positives  
38 (9682 out of 10412 positive pairs were predicted) and fewer true negatives (730 out of 10412 negative  
39 pairs were predicted). The other measurements PRE, REC and FSC were observed as 0.50, 0.93 and  
40 0.65, respectively.  
41  
42  
43  
44  
45

46 Similarly, we analyzed the results obtained by RPI369-RF and RPI369-SVM models. Table 4 shows that  
47 the RPI369-RF model predicted very few interactions as true positives (3972 out of 10412), with nearly  
48 half of the interactions were predicted as true negatives (5125 out of 10412). Therefore, the overall  
49 prediction accuracy was just ~44%. The PRE, REC and FSC scores were observed at 0.43, 0.38 and  
50 0.40, respectively. Similarly the RPI369-SVM model had an overall prediction accuracy of only ~39%. This  
51 model predicted more than half of the positive interactions as true positives (6271 out of 10412) and  
52 many fewer negative interactions as true negatives (1851 out of 10412). The remaining measurements  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 PRE, REC and FSC were observed at 0.42, 0.60 and 0.50, respectively. The RPISeq prediction scores  
4 for each ncRNA-protein pair is given in Supplementary Table S3.  
5  
6

### 7 **Application of RPI-Pred for ncRNA-protein network construction**

8  
9 We further extended RPI-Pred method for *in silico* construction of ncRNA-Protein interaction networks.  
10 Nacher and Araki (68) were among the first to study computational construction of ncRNA-protein  
11 interaction networks. They built interaction networks based on ncRNA-protein interactions of various  
12 model organisms, available in NPInter databases. Following the approach, Muppurala et. al., (36) also  
13 used results from their proposed RPISeq method for the construction of ncRNA-protein interactions  
14 networks. Here, we extended our RPI-Pred approach to construct the the ncRNA-protein interaction  
15 networks to further study the important functions of ncRNAs. We evaluated our performance in predicting  
16 ncRNA-protein interactions in the NPInter database.  
17  
18  
19  
20  
21

22 In Figure 2, we show the interaction networks for 91 ncRNA-protein interactions of *D. melanogaster* that  
23 were obtained from NPInter database. Among the 91 positive interactions, the RPI-Pred method  
24 successfully predicted 70 interactions. The ncRNA-protein interactions of *D. melanogaster* contain both  
25 protein hubs (one protein interacting with multiple RNAs), and RNA hubs (one RNA interacting with  
26 multiple proteins). The *in silico*-based network construction helps to understand how many interactions  
27 were correctly predicted by our RPI-Pred method in the same protein or RNA hubs, and the reliability of  
28 our model in derive new ncRNA-protein interactions and construct new biologic networks.  
29  
30  
31  
32

### 33 **CONCLUSION**

34  
35 Lots of the important fundamental cellular processes are mediated by protein and RNA interactions  
36 (RPIs); therefore the study of RPI is valuable for the understanding of their functions. In the recent years,  
37 the high-throughput sequencing methods have led to the discovery of enormous amount of non-coding  
38 RNAs, which also interact with protein and regulate gene expression. Hence it is very important to  
39 understand their function by studying the correct interaction partners. However, the experimental methods  
40 to determine correct interacting partner(s) for ncRNA are expensive and labor-intensive. In this case,  
41 computational approaches were highly relied to predict the interacting partner for ncRNA molecules. To  
42 our knowledge, very few studies have been reported for RPI prediction, and, none of the methods was  
43 considered the high-order protein and RNA structures, which are known to be vital to their functions.  
44  
45  
46  
47  
48  
49

50 In this work, we have developed a computational method, RPI-Pred, to address the prediction of RNA-  
51 protein interaction, and identification interacting partners of any given proteins or RNAs, using both  
52 sequences and structures of proteins and RNAs. Our proposed approach considered high-order structural  
53 features namely; protein blocks and RNA secondary structures, combined with their corresponding  
54 primary sequences for the investigation of RNA-protein interactions. Both experimental and predicted  
55 structures were used for RPI-Pred training and testing purposes. We tested the RPI-Pred method with a  
56  
57  
58  
59  
60

1  
2  
3 set of (nc)RNA-protein interaction datasets, and the results indicated that the proposed RPI-Pred was  
4 able to identify (nc)RNA-protein interactions with higher accuracy, when compared with other existing  
5 approaches. Therefore, our method is reliable to be applied to identify the binding partner(s) either for a  
6 protein or RNA. We further applied the method to *in silico* construction of ncRNA-protein networks. In  
7 addition, the proposed RPI-Pred method can also be extended to determine the binding partners (RNAs)  
8 for other types of proteins, such as transcription factors, which are able to interact with both DNA and  
9 RNA (69). A web server for the RPI-Pred can be freely accessed at [http://ctsb.is.wfubmc.edu/projects/rpi-](http://ctsb.is.wfubmc.edu/projects/rpi-pred)  
10 [pred](http://ctsb.is.wfubmc.edu/projects/rpi-pred).  
11  
12  
13  
14  
15  
16  
17

## 18 ACKNOWLEDGEMENT

19  
20 The authors would like to thank the members in the Bioinformatics group in our lab, especially, Dr. Hua  
21 Tan for valuable discussions. We also thank Dr. Muppirla for providing stand-alone version of RPISeq  
22 methods for the results comparison.  
23  
24  
25

## 26 FUNDING

27  
28 This work was partially supported by funds from the National Institutes of Health [1R01LM010185,  
29 1U01CA166886, and 1U01HL111560 to X.Z.].  
30  
31  
32

## 33 REFERENCES

- 34  
35  
36 1. Khalil, A.M. and Rinn, J.L. (2011) RNA-protein interactions in human health and disease. *Seminars in cell & developmental biology*, **22**, 359-365.  
37  
38 2. Li, Z. and Nagy, P.D. (2011) Diverse roles of host RNA binding proteins in RNA virus replication. *RNA biology*, **8**, 305-315.  
39  
40 3. Sola, I., Mateos-Gomez, P.A., Almazan, F., Zuniga, S. and Enjuanes, L. (2011) RNA-RNA and RNA-  
41 protein interactions in coronavirus replication and transcription. *RNA biology*, **8**, 237-248.  
42  
43 4. Barkan, A. (2009) Genome-wide analysis of RNA-protein interactions in plants. *Methods in*  
44 *molecular biology*, **553**, 13-37.  
45  
46 5. Kim, M.Y., Hur, J. and Jeong, S. (2009) Emerging roles of RNA and RNA-binding protein network  
47 in cancer cells. *BMB reports*, **42**, 125-130.  
48  
49 6. Zvereva, A.S. and Pooggin, M.M. (2012) Silencing and innate immunity in plant defense against  
50 viral and non-viral pathogens. *Viruses*, **4**, 2578-2597.  
51  
52 7. Kishore, S., Lubner, S. and Zavolan, M. (2010) Deciphering the role of RNA-binding proteins in the  
53 post-transcriptional control of gene expression. *Briefings in functional genomics*, **9**, 391-404.  
54  
55 8. Licatalosi, D.D. and Darnell, R.B. (2010) RNA processing and its regulation: global insights into  
56 biological networks. *Nature reviews. Genetics*, **11**, 75-87.  
57  
58 9. Singh, R. (2002) RNA-protein interactions that regulate pre-mRNA splicing. *Gene Expression*, **10**,  
59 79-92.  
60  
61 10. Hafner, M., Landthaler, M., Burger, L., Khorshid, M., Hausser, J., Berninger, P., Rothballer, A.,  
Ascano, M., Jungkamp, A.C., Munschauer, M. *et al.* (2010) PAR-CLIP--a method to identify

- transcriptome-wide the binding sites of RNA binding proteins. *Journal of visualized experiments : JoVE*, **41**.
11. Keene, J.D., Komisarow, J.M. and Friedersdorf, M.B. (2006) RIP-Chip: the isolation and identification of mRNAs, microRNAs and protein components of ribonucleoprotein complexes from cell extracts. *Nature protocols*, **1**, 302-307.
  12. Ray, D., Kazan, H., Chan, E.T., Pena Castillo, L., Chaudhry, S., Talukder, S., Blencowe, B.J., Morris, Q. and Hughes, T.R. (2009) Rapid and systematic analysis of the RNA recognition specificities of RNA-binding proteins. *Nature biotechnology*, **27**, 667-670.
  13. Djebali, S., Davis, C.A., Merkel, A., Dobin, A., Lassmann, T., Mortazavi, A., Tanzer, A., Lagarde, J., Lin, W., Schlesinger, F. *et al.* (2012) Landscape of transcription in human cells. *Nature*, **489**, 101-108.
  14. Derrien, T., Johnson, R., Bussotti, G., Tanzer, A., Djebali, S., Tilgner, H., Guernec, G., Martin, D., Merkel, A., Knowles, D.G. *et al.* (2012) The GENCODE v7 catalog of human long noncoding RNAs: analysis of their gene structure, evolution, and expression. *Genome research*, **22**, 1775-1789.
  15. Hattori, M. (2005) [Finishing the euchromatic sequence of the human genome]. *Tanpakushitsu kakusan koso. Protein, nucleic acid, enzyme*, **50**, 162-168.
  16. International Human Genome Sequencing, C. (2004) Finishing the euchromatic sequence of the human genome. *Nature*, **431**, 931-945.
  17. Rinn, J.L., Kertesz, M., Wang, J.K., Squazzo, S.L., Xu, X., Brugmann, S.A., Goodnough, L.H., Helms, J.A., Farnham, P.J., Segal, E. *et al.* (2007) Functional demarcation of active and silent chromatin domains in human HOX loci by noncoding RNAs. *Cell*, **129**, 1311-1323.
  18. Tripathi, V., Ellis, J.D., Shen, Z., Song, D.Y., Pan, Q., Watt, A.T., Freier, S.M., Bennett, C.F., Sharma, A., Bubulya, P.A. *et al.* (2010) The nuclear-retained noncoding RNA MALAT1 regulates alternative splicing by modulating SR splicing factor phosphorylation. *Molecular cell*, **39**, 925-938.
  19. Kohlmaier, A., Savarese, F., Lachner, M., Martens, J., Jenuwein, T. and Wutz, A. (2004) A chromosomal memory triggered by Xist regulates histone methylation in X inactivation. *PLoS biology*, **2**, E171.
  20. Mercer, T.R. and Mattick, J.S. (2013) Structure and function of long noncoding RNAs in epigenetic regulation. *Nature structural & molecular biology*, **20**, 300-307.
  21. Yuan, J., Wu, W., Xie, C., Zhao, G., Zhao, Y. and Chen, R. (2014) NPInter v2.0: an updated database of ncRNA interactions. *Nucleic acids research*, **42**, D104-108.
  22. Luscombe, N.M., Laskowski, R.A. and Thornton, J.M. (2001) Amino acid-base interactions: a three-dimensional analysis of protein-DNA interactions at an atomic level. *Nucleic acids research*, **29**, 2860-2874.
  23. Jones, S., van Heyningen, P., Berman, H.M. and Thornton, J.M. (1999) Protein-DNA interactions: A structural analysis. *J Mol Biol*, **287**, 877-896.
  24. Nadassy, K., Wodak, S.J. and Janin, J. (1999) Structural features of protein-nucleic acid recognition sites. *Biochemistry-U.S.*, **38**, 1999-2017.
  25. Cook, K.B., Kazan, H., Zuberi, K., Morris, Q. and Hughes, T.R. (2011) RBPDB: a database of RNA-binding specificities. *Nucleic acids research*, **39**, D301-D308.
  26. Ray, D., Kazan, H., Cook, K.B., Weirauch, M.T., Najafabadi, H.S., Li, X., Gueroussov, S., Albu, M., Zheng, H., Yang, A. *et al.* (2013) A compendium of RNA-binding motifs for decoding gene regulation. *Nature*, **499**, 172-177.
  27. Mittal, N., Roy, N., Babu, M.M. and Janga, S.C. (2009) Dissecting the expression dynamics of RNA-binding proteins in posttranscriptional regulatory networks. *P Natl Acad Sci USA*, **106**, 20300-20305.

- 1
  - 2
  - 3
  - 4
  - 5
  - 6
  - 7
  - 8
  - 9
  - 10
  - 11
  - 12
  - 13
  - 14
  - 15
  - 16
  - 17
  - 18
  - 19
  - 20
  - 21
  - 22
  - 23
  - 24
  - 25
  - 26
  - 27
  - 28
  - 29
  - 30
  - 31
  - 32
  - 33
  - 34
  - 35
  - 36
  - 37
  - 38
  - 39
  - 40
  - 41
  - 42
  - 43
  - 44
  - 45
  - 46
  - 47
  - 48
  - 49
  - 50
  - 51
  - 52
  - 53
  - 54
  - 55
  - 56
  - 57
  - 58
  - 59
  - 60
28. Zhang, Z., Schwartz, S., Wagner, L. and Miller, W. (2000) A greedy algorithm for aligning DNA sequences. *J Comput Biol*, **7**, 203-214.
29. Altschul, S.F., Gish, W., Miller, W., Myers, E.W. and Lipman, D.J. (1990) Basic local alignment search tool. *J Mol Biol*, **215**, 403-410.
30. Gish, W. and States, D.J. (1993) Identification of protein coding regions by database similarity search. *Nat Genet*, **3**, 266-272.
31. Bateman, A., Birney, E., Durbin, R., Eddy, S.R., Howe, K.L. and Sonnhammer, E.L. (2000) The Pfam protein families database. *Nucleic acids research*, **28**, 263-266.
32. Bateman, A., Birney, E., Cerruti, L., Durbin, R., Etwiller, L., Eddy, S.R., Griffiths-Jones, S., Howe, K.L., Marshall, M. and Sonnhammer, E.L.L. (2002) The Pfam protein families database. *Nucleic acids research*, **30**, 276-280.
33. Finn, R.D., Bateman, A., Clements, J., Coggill, P., Eberhardt, R.Y., Eddy, S.R., Heger, A., Hetherington, K., Holm, L., Mistry, J. *et al.* (2014) Pfam: the protein families database. *Nucleic acids research*, **42**, D222-230.
34. Pancaldi, V. and Bahler, J. (2011) In silico characterization and prediction of global protein-mRNA interactions in yeast. *Nucleic acids research*, **39**, 5826-5836.
35. Bellucci, M., Agostini, F., Masin, M. and Tartaglia, G.G. (2011) Predicting protein associations with long noncoding RNAs. *Nature methods*, **8**, 444-445.
36. Muppurala, U.K., Honavar, V.G. and Dobbs, D. (2011) Predicting RNA-Protein Interactions Using Only Sequence Information. *Bmc Bioinformatics*, **12**, 489.
37. Shen, J., Zhang, J., Luo, X., Zhu, W., Yu, K., Chen, K., Li, Y. and Jiang, H. (2007) Predicting protein-protein interactions based only on sequences information. *Proc Natl Acad Sci U S A*, **104**, 4337-4341.
38. Wang, Y., Chen, X., Liu, Z.P., Huang, Q., Wang, Y., Xu, D., Zhang, X.S., Chen, R. and Chen, L. (2013) De novo prediction of RNA-protein interactions from sequence information. *Molecular bioSystems*, **9**, 133-142.
39. Lu, Q.S., Ren, S.J., Lu, M., Zhang, Y., Zhu, D.H., Zhang, X.G. and Li, T.T. (2013) Computational prediction of associations between long non-coding RNAs and proteins. *Bmc Genomics*, **14**, 651.
40. Chen, S.J. (2008) RNA folding: conformational statistics, folding kinetics, and ion electrostatics. *Annual review of biophysics*, **37**, 197-214.
41. de Brevern, A.G., Etchebest, C. and Hazout, S. (2000) Bayesian probabilistic approach for predicting backbone structures in terms of protein blocks. *Proteins*, **41**, 271-287.
42. Joseph, A.P., Agarwal, G., Mahajan, S., Gelly, J.C., Swapna, L.S., Offmann, B., Cadet, F., Bornot, A., Tyagi, M., Valadie, H. *et al.* (2010) A short survey on protein blocks. *Biophysical reviews*, **2**, 137-147.
43. Suresh, V., Ganesan, K. and Parthasarathy, S. (2013) A protein block based fold recognition method for the annotation of twilight zone sequences. *Protein and peptide letters*, **20**, 249-254.
44. Rose, P.W., Beran, B., Bi, C., Bluhm, W.F., Dimitropoulos, D., Goodsell, D.S., Prlic, A., Quesada, M., Quinn, G.B., Westbrook, J.D. *et al.* (2011) The RCSB Protein Data Bank: redesigned web site and web services. *Nucleic acids research*, **39**, D392-401.
45. Lewis, B.A., Walia, R.R., Terribilini, M., Ferguson, J., Zheng, C., Honavar, V. and Dobbs, D. (2011) PRIDB: a Protein-RNA interface database. *Nucleic acids research*, **39**, D277-282.
46. Rice, P., Longden, I. and Bleasby, A. (2000) EMBOSS: the European Molecular Biology Open Software Suite. *Trends in genetics : TIG*, **16**, 276-277.
47. Schneider, B., Cerny, J., Svozil, D., Cech, P., Gelly, J.C. and de Brevern, A.G. (2014) Bioinformatic analysis of the protein/DNA interface. *Nucleic acids research*, **42**, 3381-3394.
48. Rajagopal, S. and Vishveshwara, S. (2005) Short hydrogen bonds in proteins. *The FEBS journal*, **272**, 1819-1832.

- 1  
2  
3 49. Hogan, D.J., Riordan, D.P., Gerber, A.P., Herschlag, D. and Brown, P.O. (2008) Diverse RNA-binding proteins interact with functionally related sets of RNAs, suggesting an extensive regulatory system. *PLoS biology*, **6**, e255.
- 4  
5  
6  
7 50. Suresh, V., Ganesan, K. and Parthasarathy, S. (2012) PDB-2-PB: a curated online protein block sequence database. *J Appl Crystallogr*, **45**, 127-129.
- 8  
9  
10 51. Zheng, G.H., Lu, X.J. and Olson, W.K. (2009) Web 3DNA-a web server for the analysis, reconstruction, and visualization of three-dimensional nucleic-acid structures. *Nucleic acids research*, **37**, W240-W246.
- 11  
12  
13 52. Vapnik, V.N. (1999) An overview of statistical learning theory. *Ieee T Neural Networ*, **10**, 988-999.
- 14  
15  
16 53. Chang, C.C. and Lin, C.J. (2011) LIBSVM: A Library for Support Vector Machines. *Acm T Intel Syst Tec*, **2**, 27.
- 17  
18  
19 54. Suresh, V. and Parthasarathy, S. (2014) SVM-PB-Pred: SVM based protein block prediction method using sequence profiles and secondary structures. *Protein and peptide letters*, **21**, 736-742.
- 20  
21  
22 55. Offmann, B., Tyagi, M. and de Brevern, A.G. (2007) Local protein structures. *Curr Bioinform*, **2**, 165-202.
- 23  
24  
25 56. Ding, F., Sharma, S., Chalasani, P., Demidov, V.V., Broude, N.E. and Dokholyan, N.V. (2008) Ab initio RNA folding by discrete molecular dynamics: from structure prediction to folding mechanisms. *Rna*, **14**, 1164-1173.
- 26  
27  
28 57. Zhang, J., Dundas, J., Lin, M., Chen, R., Wang, W. and Liang, J. (2009) Prediction of geometrically feasible three-dimensional structures of pseudoknotted RNA through free energy estimation. *Rna*, **15**, 2248-2263.
- 29  
30  
31 58. Hajdin, C.E., Ding, F., Dokholyan, N.V. and Weeks, K.M. (2010) On the significance of an RNA tertiary structure prediction. *Rna-a Publication of the Rna Society*, **16**, 1340-1349.
- 32  
33  
34 59. Jonikas, M.A., Radmer, R.J., Laederach, A., Das, R., Pearlman, S., Herschlag, D. and Altman, R.B. (2009) Coarse-grained modeling of large RNA molecules with knowledge-based potentials and structural filters. *Rna*, **15**, 189-199.
- 35  
36  
37 60. Cao, S. and Chen, S.J. (2011) Physics-based de novo prediction of RNA 3D structures. *The journal of physical chemistry. B*, **115**, 4216-4226.
- 38  
39  
40 61. Liu, L. and Chen, S.J. (2010) Computing the conformational entropy for RNA folds. *The Journal of chemical physics*, **132**, 235104.
- 41  
42  
43 62. Liu, L. and Chen, S.J. (2012) Coarse-grained prediction of RNA loop structures. *PloS one*, **7**, e48460.
- 44  
45  
46 63. Zhao, Y., Huang, Y., Gong, Z., Wang, Y., Man, J. and Xiao, Y. (2012) Automated and fast building of three-dimensional RNA structures. *Scientific reports*, **2**, 734.
- 47  
48  
49 64. Gruber, A.R., Lorenz, R., Bernhart, S.H., Neubock, R. and Hofacker, I.L. (2008) The Vienna RNA websuite. *Nucleic acids research*, **36**, W70-74.
- 50  
51  
52 65. Chen, G., Chang, K.Y., Chou, M.Y., Bustamante, C. and Tinoco, I., Jr. (2009) Triplex structures in an RNA pseudoknot enhance mechanical stability and increase efficiency of -1 ribosomal frameshifting. *Proc Natl Acad Sci U S A*, **106**, 12706-12711.
- 53  
54  
55 66. Giedroc, D.P. and Cornish, P.V. (2009) Frameshifting RNA pseudoknots: Structure and mechanism. *Virus Res*, **139**, 193-208.
- 56  
57  
58 67. Yu, C.H., Teulade-Fichou, M.P. and Olsthoorn, R.C. (2014) Stimulation of ribosomal frameshifting by RNA G-quadruplex structures. *Nucleic acids research*, **42**, 1887-1892.
- 59  
60 68. Nacher, J.C. and Araki, N. (2010) Structural characterization and modeling of ncRNA-protein interactions. *Bio Systems*, **101**, 10-19.



- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60
69. Cassidy, L.A. and Maher, L.J., 3rd. (2002) Having it both ways: transcription factors that bind DNA and RNA. *Nucleic acids research*, **30**, 4118-4126.

**TABLE LEGENDS**

Table 1. Performance of RPI-Pred using 10-fold cross validation on RPI1807, RPI2241 and RPI369 datasets.

Table 2. Comparison of RPI-Pred and Wang et al.'s classifiers (38) on the RPI367 dataset.

Table 3. Performance of RPI-Pred on the NPIner10412 dataset, for different organisms.

Table 4. Comparison of RPI-Pred and RPISeq models (36) on the NPInter10412 dataset.

**FIGURES LEGENDS**

Figure 1. Step-wise work-flow for the proposed RPI-Pred method.

Figure 2. The ncRNA-protein interaction networks constructed based on interaction pairs predicted by RPI-Pred, for *D. melanogaster*. The ncRNA and proteins are shown in green (square) and yellow (oval/circular) nodes, respectively, while the correctly- and wrongly-predicted ncRNA-protein interactions are shown as blue and red edges, respectively.

Table 1.

Measurements	RPI2241			RPI369		
	RPI-Pred	RPISeq-SVM	RPISeq- RF	RPI-Pred	RPISeq-SVM	RPISeq- RF
AUC	0.89	0.97	0.92	0.95	0.81	0.81
PRE	0.88	0.87	0.89	0.89	0.73	0.75
REC	0.78	0.88	0.90	0.89	0.73	0.78
FSC	0.83	0.87	0.90	0.89	0.73	0.77
ACC	84.0	87.1	89.6	92.0	72.8	76.2

Table 2.

Organism	Total RNA-Protein pairs	RPI-Pred method (%)	RPI369-NB classifier (%)	RPI369-ENB classifier (%)	RP2241-NB classifier (%)	RP2241-ENB classifier (%)
<i>C. elegans</i>	3	3 (100%)	3 (100%)	1 (33%)	1 (33%)	1 (33%)
<i>D.melanogaster</i>	26	24 (92%)	13 (50%)	19 (74%)	23 (88%)	25 (96%)
<i>E. coli</i>	25	24 (96%)	13 (52%)	17 (68%)	12 (48%)	15 (60%)
<i>H. sapiens</i>	148	135 (91%)	93 (63%)	77 (52%)	84 (57%)	91 (62%)
<i>M. musculus</i>	46	37 (80%)	30 (65%)	40 (87%)	28 (61%)	37 (80%)
<i>S. cerevisiae</i>	119	105 (88%)	76 (64%)	89 (75%)	94 (79%)	118 (99%)
Total	367	328 (89%)	228 (62%)	243 (67%)	242 (66%)	287 (79%)

Table 3.

Organism	Total ncRNA-protein pairs in NPInter10412	RPI-PRed performance (%)
<i>Caenorhabditis elegans</i>	36	28 (78%)
<i>Drosophila melanogaster</i>	91	70 (77%)
<i>Escherichia coli</i>	202	154 (76%)
<i>Homo sapiens</i>	6,975	6,193 (89%)
<i>Mus musculus</i>	2,198	2,147 (98%)
<i>Saccharomyces cerevisiae</i>	910	743 (81%)
Total	10,412	9,335 (90%)

Table 4.

	RPI-Pred	RPI2241-RF	RPI2241-SVM	RPI369-RF	RPI369-SVM
PRE	0.85	0.50	0.50	0.43	0.42
REC	0.90	0.98	0.93	0.38	0.60
FSC	0.87	0.66	0.65	0.40	0.50
ACC	86.9	50.2	49.2	43.8	39.0

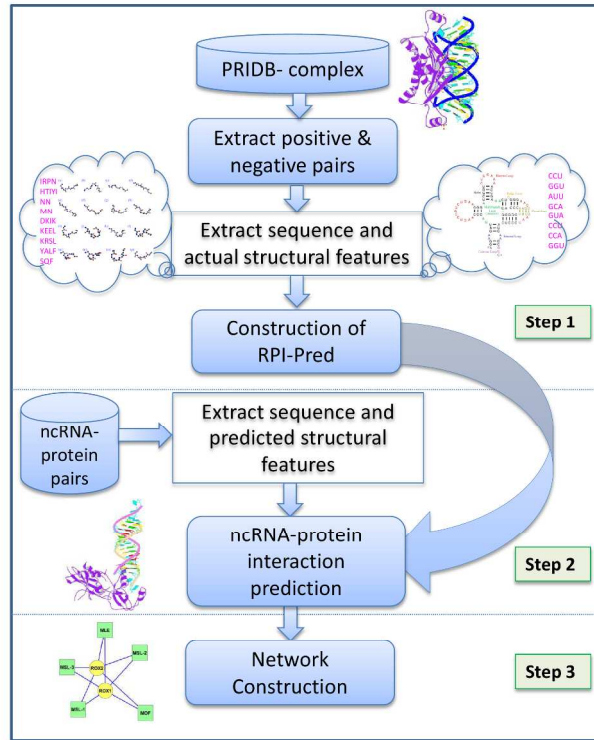


Figure 1

Figure 1. Step-wise work-flow for the proposed RPI-Pred method.  
254x190mm (300 x 300 DPI)

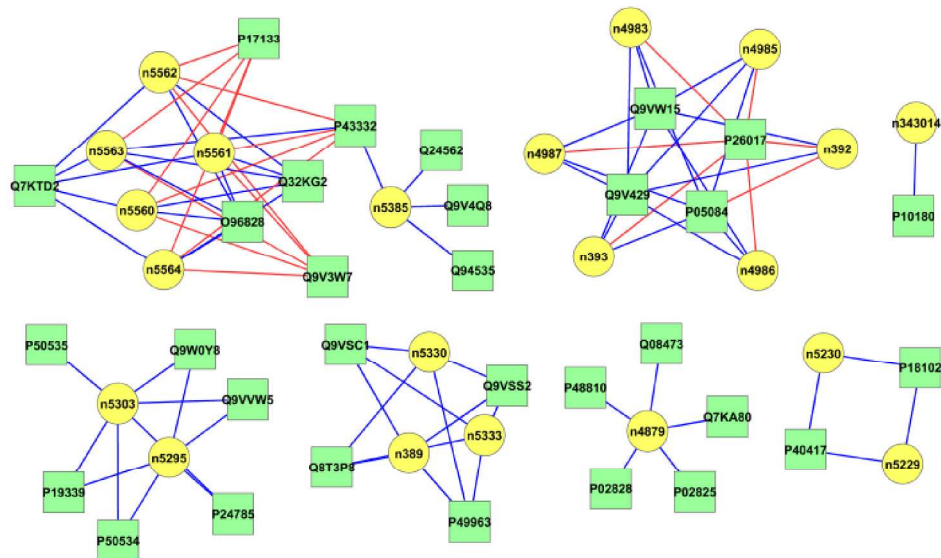


Figure 2

Figure 2. The ncRNA-protein interaction networks constructed based on interaction pairs predicted by RPI-Pred, for *D. melanogaster*. The ncRNA and proteins are shown in green (square) and yellow (oval/circular) nodes, respectively, while the correctly- and wrongly-predicted ncRNA-protein interactions are shown as blue and red edges, respectively.  
254x190mm (300 x 300 DPI)

**Supplemental data:**

Table S1: Training dataset (RNA-protein positive and negative pairs)

Table S2: The sequence and structural feature labels for protein and RNA

Table S3: The RPI-Pred prediction results compared with RPISeq – all four models (30) on NPInter10412 dataset