

Deep learning based DNA:RNA triplex forming potential prediction

Yu ZHANG

Nanyang Technological University <https://orcid.org/0000-0003-4952-8095>

Yahui Long

Hunan University

Chee Keong Kwoh (✉ asckwoh@ntu.edu.sg)

Nanyang Technological University <https://orcid.org/0000-0002-8547-6387>

Software

Keywords: long noncoding RNAs, DNA:RNA triplex, deep learning

Posted Date: November 5th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-41662/v3>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published on November 12th, 2020. See the published version at <https://doi.org/10.1186/s12859-020-03864-0>.

1 Title page

2 **Deep learning based DNA:RNA triplex forming potential**

3 **prediction**

4

5 Yu Zhang¹, Yahui Long² and Chee Keong Kwoh^{1,*}

6 ¹School of Computer Science and Engineering, Nanyang Technological University, 639798,

7 Singapore, ²College of Computer Science and Electronic Engineering, Hunan University,

8 Changsha 410000, China.

9

10 * Correspondence: asckkwoh@ntu.edu.sg

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

1 **Abstract**

2 **Background:** Long non-coding RNAs (lncRNAs) can exert functions via forming triplex
3 with DNA. The current methods in predicting the triplex formation mainly rely on
4 mathematic statistic according to the base paring rules. However, these methods have two
5 main limitations: i) they identify a large number of triplex-forming lncRNAs, but the limited
6 number of experimentally verified triplex-forming lncRNA indicates that maybe not all of
7 them can form triplex in practice, and ii) their predictions only consider the theoretical
8 relationship while lacking the features from the experimentally verified data.

9 **Results:** In this work, we develop an integrated program named TriplexFPP (Triplex
10 Forming Potential Prediction), which is the first machine learning model in DNA:RNA
11 triplex prediction. TriplexFPP predicts the most likely triplex-forming lncRNAs and DNA
12 sites based on the experimentally verified data, where the high-level features are learned by
13 the convolutional neural networks. In the 5-fold cross validation, the average values of Area
14 Under the ROC curves and PRC curves for removed redundancy triplex-forming lncRNA
15 dataset with threshold 0.8 are 0.9649 and 0.9996, and these two values for triplex DNA sites
16 prediction are 0.8705 and 0.9671, respectively. Besides, we also briefly summarize the *cis*
17 and *trans* targeting of triplexes lncRNAs.

18 **Conclusions:** The TriplexFPP is able to predict the most likely triplex-forming lncRNAs
19 from all the lncRNAs with computationally defined triplex forming capacities and the
20 potential of a DNA site to become a triplex. It may provide insights to the exploration of
21 lncRNA functions.

22 **Keywords:** long noncoding RNAs, DNA:RNA triplex, deep learning

23

24 **Background**

1 The advances in sequencing technologies enable the discovery of the vast amount of Long
2 non-coding RNAs (lncRNAs). lncRNAs can serve as signals, decoys, guides, and scaffolds to
3 carry out functions like chromatin states modulation and gene expression regulation. They act
4 via the interactions with DNA, protein, and other RNA, in the way of coordinating regulatory
5 proteins, localizing to target loci, shaping three-dimensional (3D) nuclear organization [1-3] ,
6 *etc.*

7 One way for lncRNA to interact with DNA is to form triplex structures [4]. Triplex is a
8 kind of direct RNA-DNA interaction mechanism, it is formed through the binding of RNA
9 sites and purine rich strand of duplex DNA under the forward or reverse Hoogsteen base-
10 pairing rule. Some lncRNAs are proved to execute functions via forming DNA:RNA
11 triplexes, for example, promoter associated lncRNA interacts with TTF-I to repress
12 transcription of rRNA [5], FENDRR increases PRC2 occupancy at the triplex formation sites
13 [6], MEG3 forms DNA-lncRNA triplex with TGF-β gene to modulate the gene activity [7],
14 PARTICLE binds to MAT2A promoter CpG island as triplex to contribute to gene-silencing
15 machineries [8], KHPS1 interacts with SPHK1 to anchor the lncRNA and associated effector
16 proteins to the gene promoter [9], HOTAIR forms triplex with PCDH7 and HOXB2 to
17 regulate adipogenic differentiation [10], MIR100HG acts via triplex formation to regulate
18 p27 [11], and promoter and pre-rRNA antisense guides associated CHD4/ NuRD to the
19 rDNA promoter [12].

20 Although the recent development of high throughput techniques, such as Chromatin
21 Isolation by RNA Purification (ChIRP-seq) [13], capture hybridization analysis of RNA
22 targets (CHART-seq) [14], RNA Antisense Purification (RAP-seq) [15], and chromatin oligo
23 affinity precipitation (ChOP-seq) [7], has helped to generate the genome-wide map of
24 lncRNA chromatin interactions for specific lncRNAs via deciphering their binding sites, most
25 of them are implemented in crosslinked chromatin which contain RNA associated to DNA

1 binding proteins. Therefore, they cannot provide reliable references for the studying of
2 DNA:RNA triplex formation. To reveal the existence of the DNA:RNA triplex interactions in
3 vivo, Cetin *et al.* developed a method to map the genome-wide DNA:RNA triplexes by
4 excluding the chromatin crosslinking [16]. This method proved the physiological relevance of
5 DNA:RNA triplex structures.

6 Currently, the prediction of DNA:RNA triplex mainly relies on the base paring rules-
7 related mathematic statistics. Triplexator is proposed to systematically identify the potential
8 triplex forming sites of RNA and the targeting sites on DNA by taking the Hoogsteen and
9 reverse-Hoogsteen base-pairing into account [17], Triplex-Inspector is designed to select
10 sequence-specific ligands and targets by considering the gene location and genomic
11 architecture [18], LongTarget is presented to detect motifs and binding sites in forming
12 triplex by considering non-canonical rules [19], and Triplex Domain Finder (TDF) is
13 developed to predict triplexes and characterize lncRNA and the corresponding DNA targets
14 [20].

15 Although the above methods can identify potential triplexes according to the canonical
16 rules, they predict a large population of lncRNAs with triplex forming potential. However,
17 the limited number of experimentally identified triplex-forming lncRNAs indicates that
18 maybe not all of them can form triplex in practice. Besides, the computational methods only
19 calculate the theoretical triplex potential, while do not consider any *in vivo* or *in vitro* assays
20 verified data.

21 In this work, we have the following two aims: i) predicting the most likely triplex-forming
22 lncRNAs in practice from the lncRNAs owning triplex forming capabilities calculated by the
23 computational methods, and ii) predicting the potential of DNA sites in forming triplex based
24 on the experimentally verified data. For these purposes, we develop TriplexFPP (Triplex
25 Forming Potential Prediction). It is the first machine learning program in DNA:RNA triplex

1 forming potential prediction according to our knowledge. In triplex lncRNA prediction, the
2 average values of Area Under the ROC (AUROC) and Area Under the PRC (AUPRC) for 5-
3 fold cross validation on removed redundancy dataset with threshold 0.8 are 0.9649 and
4 0.9996 separately. Besides, the average cross-validation AUROC and AUPRC values for the
5 triplex DNA sites potential prediction are 0.8705 and 0.9671 separately. The general good
6 performances of TriplexFPP illustrate its effectiveness in triplex forming potential prediction,
7 and could provide references to the lncRNA function exploration. Furthermore, we also
8 briefly summarize the *cis* and *trans* targeting of triplexes lncRNAs, which may provide some
9 insights to the exploration of lncRNA binding mechanisms.

10

11 **Implementation**

12 **Dataset**

13 ***Triplex lncRNA prediction dataset***

14 The positive data for triplex lncRNA prediction is collected in 2 ways. On one hand, we
15 extracted the lncRNAs according to the TriplexRNA regions (DNA:RNA triplex forming
16 peaks in RNA) reported in the work of Sentürk *et al.* [16] by considering both Solid Phase
17 Reversible Immobilization-based paramagnetic bead size selection and immunopurification
18 with anti-DNA antibody RNA separation in Hela S3 cell. We used GENCODE release 33
19 lncRNA annotation [21] to extract the lncRNAs that cover the TriplexRNA regions and
20 obtained 476 unique samples in this way. We named these lncRNAs as triplexlncRNA. On
21 the other hand, we also collected lncRNAs that are verified by either *in vivo* or *in vitro* assays
22 to from triplexes with DNA from the peer-reviewed publications. These lncRNAs are named as
23 reported triplex lncRNA, including MEG3 [7], PARTICL [8], MIR100HG [11], FENDRR
24 [22], and HOTAIR [10]. All the variants of the reported triplex lncRNA were taken into
25 consideration. The total number of reported triplex lncRNA is 159.

1 Since our goal is to predict the most likely triplex-forming lncRNAs in practice from the
2 lncRNAs owning triplex forming capacities predicted by computational methods, we used
3 TDF [20] to further filter the data. When evaluating the triplex forming potential of the above
4 635 lncRNAs with the whole gene promoters (except for chromosome Y and M) with TDF
5 (default parameters), 104 of them do not contain DNA Binding Domains (DBDs) with
6 powerful Triplex Forming Oligonucleotide (TFO) support. There are two possible
7 explanations for the phenomenon of low triplex forming capacity obtained from TDF in
8 collected positive data: i) for reported triplex lncRNA, each lncRNA gene may have multiple
9 transcripts with splice variants, but maybe not all of their variants have the abilities to form
10 the triplex with DNA, and ii) for triplexlncRNA, the overlap regions between R-Loop and
11 TriplexRNA regions cannot be confirmed as triplex forming or not [16]. To assure the data
12 reliability, we deleted the 104 lncRNAs with low triplex forming capacity and finally got 531
13 samples in the positive dataset.

14 The negative samples are collected following the same filter rules as that of positive. After
15 removing the lncRNAs in our original positive dataset from GENCODE annotation, we
16 evaluated the triplex forming potential of all remaining lncRNAs with the whole gene
17 promoters (except for chrY and chrM) by TDF. We only kept the lncRNAs with at least one
18 powerful TFO supported DBD and at least 123 DNA Binding Sites (DBSs) (the smallest
19 number of DBSs is 123 in positive data). From the qualified lncRNAs, we further removed
20 one lncRNA with letter ‘N’ in its sequence and the variants of lncRNA MALAT1 which is
21 reported to form RNA-RNA triplex [23]. Finally, the negative dataset contains 36021
22 lncRNAs which have comparable triplex forming abilities as the positive data.

23 We also prepared two more datasets with removed redundancies. We used the CD-HIT
24 [24] to remove the redundancy in each class of the original dataset with threshold of 0.9 and

1 0.8 separately. The positive and negative data amounts in two removed redundancy datasets
2 are 384 and 28012, and 286 and 22681, respectively.

3 ***Triplex DNA sites potential prediction dataset***

4 To predict triplex-forming sites in DNA on the basis of experimental data, we adopted the
5 TriplexDNA regions (DNA:RNA triplex forming peaks in DNA) obtained from [16] as
6 positive data. These RNA-associated DNAs are enriched by an unbiased approach. After
7 removing 5 samples containing the letter 'N' in their sequences, the final positive dataset size
8 is 2542. The negative data is selected as the random regions in promoters. We downloaded all
9 ensembl annotated promoters [25], from which we generated 12735 regions (5 times amount
10 of originally positive data). These regions were obtained by randomly selecting chromosomes
11 (except chrY and chrM) and DNA regions with the same lengths as TriplexDNA. The
12 sequence data for both positive and negative is extracted from the DNA minus strand.

13

14 **Feature extraction**

15 Two types of sequence-related feature extraction strategies are considered here: k-mer and
16 kmerscore. Both two types of strategies have been successfully applied to classification
17 problems in RNA [26, 27]. K-mer is a popular method to transform a sequence into a vector,
18 it counts the frequencies of single or multiple nucleotide compositions in a sequence and
19 represents the sequence into a 4^k dimensional vector. K-mer features can be calculated as

20
$$kmer(i) = \frac{\text{Total number of } k \text{ neighboring nucleic acids}(i)}{n-k+1} \quad (1)$$

21 Where $kmer(i)$ is the frequency of the i th nucleotide composition in all 4^k possibilities, and
22 the denominator $n - k + 1$ represents the total number of all possible k neighboring nucleic
23 acids in a sequence with length n . For example, in the 3-mer circumstance, $k = 3$,
24 considering sequence $S = AAAAC$, whose $n = 5$, its frequencies under the nucleotide

1 composition of *AAA* and *AAC* are 2/3 and 1/3 separately, while the frequencies for other 3-
2 mer compositions like *AAG*, *et al.* are 0.
3 kmerscore is an overall measure of the k-mer nucleotide composition bias in a sequence, it
4 is obtained from k-mer features. To calculate it, firstly, the k-mer features for all sequences
5 need to be calculated, then the mean k-mer vectors from the positive and negative dataset are
6 obtained from the corresponding k-mer features separately, which are represented as
7 $M_{pos}(h_i)$ and $M_{neg}(h_i)$, where $i = 1, 2, \dots, 4^k$. Finally, for a nucleotide sequence $S =$
8 $s_1 s_2 \dots s_n$ with k-mer sequence $S = h_1 h_2 \dots h_{n-k+1}$, the kmerscore can be represented as

$$kmerscore = \frac{1}{n-k+1} \sum_{i=1}^{n-k+1} \log \frac{M_{neg}(h_i)}{M_{pos}(h_i)} \quad (2)$$

10

11 Model construction

12 We developed an integrated machine learning program called TriplexFPP (Triplex Forming
13 Potential Prediction) in triplex forming potential prediction. It consists of two individual
14 models, including triplex lncRNA prediction model and triplex DNA sites potential
15 prediction model.

16 We adopted the 2-layer Convolutional Neural Network to construct the models, which can
17 effectively learn the high-level features. The detailed description of the structure and
18 parameters of the model can be found at section of Initial training of TriplexFPP below. As
19 the positive dataset sizes and negative dataset sizes in triplex lncRNA prediction and triplex
20 DNA sites prediction are imbalanced with ratios around 1:68 and 1:5 separately, to avoid the
21 model bias, we applied the random down-sampling technique during the training process. The
22 negative training data were randomly selected as the same amount as positive training data.
23 Due to the extremely small positive dataset size in triplex lncRNA prediction, *i.e.* 531, we
24 also applied oversampling on this dataset. Because the positive data in triplex lncRNA
25 prediction are collected from two sources whose amounts are in a ratio around 2.5 to 1

1 (triplexlncRNA to reported triplex lncRNA), to force the model to learn the global features
2 for all positive data rather than the common features for the data from the majority type, *i.e.*
3 triplexlncRNA, we assigned more weights to the weak type data, *i.e.* reported triplex
4 lncRNA, during augmenting the positive data in the training process, and we named this
5 practice as the weighted bagging strategy.

6

7 Model evaluation

8 To demonstrate the model performances, the evaluation criteria of accuracy (Acc), sensitivity
9 (Sn), specificity (Sp), and AUROC are used. Besides, criteria in evaluating imbalance data
10 are also adopted, including AUPRC, F1-score, and harmonic mean (Hm). The equations for
11 calculating the above criteria are listed below.

$$12 \quad Accuracy (Acc) = \frac{TP+TN}{TP+TN+FP+FN} \quad (3)$$

$$13 \quad Sensitivity (Sn) = \frac{TP}{TP+FN} \quad (4)$$

$$14 \quad Specificity (Sp) = \frac{TN}{TN+FP} \quad (5)$$

$$15 \quad Harmonic mean of Sn and Sp (Hm) = \frac{2 \times Sn \times Sp}{Sn + Sp} \quad (6)$$

$$16 \quad F1 - score = \frac{2 \times PRE \times Sn}{PRE + Sn}, \text{ where } PRE = \frac{TP}{TP + FP} \quad (7)$$

17 Where FN, FP, TN, and TP denote the number of false negative, false positive, true
18 negative, and true positive, respectively.

19

20 Results

21 Initial training of TriplexFPP

22 Each input sequence is represented into a fix-length 90-dim vector by considering both k-mer
23 ($k = 1, 2, 3, \text{number of feature} = 84$) and kmerscore ($k = 1 - 6, \text{number of feature} =$

1 6) features. The influences of the k-mer and kmerscore features to TriplexFPP can be found
2 at Additional file 1: Figure S1. We use the mean k-mer feature values from training data to
3 calculate the kmerscore features of both training and test data in the corresponding split to
4 exclude any information from test data.

5 The parameters in TriplexFPP, such as the number of convolutional layer, the kernel size,
6 the activation function, *etc.* are determined according to the corresponding random split
7 datasets in two individual models. Each time we change the value of one parameter while
8 keeping other parameters fixed, and then select the one that achieves the highest value of S_n
9 as the final choice for that parameter. The setting of the parameters and the corresponding
10 performances for triplex lncRNA prediction model and triplex DNA sites prediction model
11 are shown in Additional file 1: Figure S2 and S3 separately. The detailed architectures and
12 parameters for TriplexFPP can be found in Fig. 1.

13

14 **Evaluation of triplex lncRNA prediction**

15 In this section, we evaluate the triplex lncRNA prediction model in TriplexFPP regard to its
16 ability in predicting the most likely triplex-forming lncRNA from the lncRNAs owning
17 triplex forming capability predicted by the computation methods.

18 We first visualize the nucleotide compositions of lncRNA sequences in our positive
19 dataset (triplexlncRNA and reported triplex lncRNA) and negative dataset (Additional file 1:
20 Figure S4). The nucleotide compositions for the reported triplex lncRNAs and the negative
21 lncRNAs are more consistent, from which nucleotide A and T taking up heavier percentages;
22 whereas the triplexlncRNA follows a different pattern, whose sequences are mainly CG rich.
23 Because the amount of triplexlncRNA is larger than that of reported triplex lncRNA in our
24 positive dataset, to ensure the model to learn the high-level features for all positive data rather
25 than the sequence composition features of triplexlncRNA, we assign more weights to the

1 reported triplex lncRNA when augmenting the positive data in the training process. We triple
2 the positive training data with the weighted bagging strategy, where two-thirds of them is
3 bagged from the original positive training data directly and one-third of them is extra bagged
4 from the reported triplex lncRNA. We compare our method with baseline models like Deep
5 Neural Network (NN), Support Vector Machine (SVM), Random Forest (RF), and Gradient
6 Boosting. The parameters for the NN model are determined as the ones with the best Sn value
7 among several choices on random split training and test datasets, and the parameters for other
8 baseline models are determined as the optimal ones by cross-validated grid-search over a
9 parameter grid based on the criteria of Recall. The candidate parameters and the final
10 determined parameters for each baseline model are recorded in Additional file 1: Table S1.
11 The box and whisker plots for 5-fold cross validation are demonstrated for all models in Fig.
12 2a, where the negative training data are randomly selected as the same amount of augmented
13 positive training data from the 4 training folds in each validation.

14 With the mission of finding the most likely triplex-forming lncRNAs in practice, although
15 Gradient Boosting method realizes the best Acc and F1-score values, its Sn values mainly
16 concentrate between 87.74% and 92.45%. Conversely, the Sn values for CNN model range
17 from 93.40% to 97.17%, concentrating at a high-value region, which indicates its superiority
18 to baseline models on the model performance. Besides, the overall high values of other
19 evaluation metrics in CNN model, e.g. 98.35% of average Acc, 0.9926 of average AUROC,
20 0.9999 of average AUPRC, 0.992 of average F1-score, and 0.969 of average Hm, further
21 illustrate the effectiveness of our CNN model in TriplexFPP in the lncRNA triplex forming
22 potential prediction.

23 However, one issue for the above 5-fold cross validation is that the high performance may
24 be caused by the high data similarity between training and test. In our positive dataset, the
25 531 samples relate to 135 genes (Additional file 1: Figure S5), where 57 genes owning two or

1 more variants. The different variants from the same lncRNA gene can share high sequence
2 composition similarities, thereby leading to a good prediction result. To evaluate whether the
3 model is powerful enough in predicting lncRNA triplex forming potential, we further execute
4 5-fold cross validation on the datasets with removed redundancy and carry out the leave-out
5 validation.

6 We compare the performances of CNN model on datasets without removing redundancy,
7 removing redundancy with threshold of 0.9, and removing redundancy with threshold of 0.8.
8 The results are demonstrated in Fig. 2b. Although the average performances on the data with
9 removed redundancy are slightly lower than that of without removing redundancy, the values
10 in evaluation matrices on the removed redundancy datasets remain at high levels. For
11 example, in removed redundancy datasets with threshold of 0.9, the AUROC and AUPRC
12 values range from 0.9637 to 0.9880 and 0.9993 to 0.9998 separately; in removed redundancy
13 datasets with threshold of 0.8, the AUROC and AUPRC values range from 0.9497 to 0.9809
14 and 0. 9994 to 0.9998 separately.

15 For the leave-out validation, we select four lncRNAs with the most amounts of variants as
16 the test data, including MIR100HG, PVT1, LINC00963, and MEG3. Their variants amounts
17 are 87, 73, 52, and 46, respectively. The four lncRNAs follow different data sources, PVT1
18 and LINC00963 belong to triplexlncRNA, and MIR100HG and MEG3 belong to reported
19 triplex lncRNA. In each leave-out validation, we select one lncRNA and use all of its variants
20 as positive test data, while using all the remaining positive lncRNAs as the training data. The
21 training and test process are repeated 5 times, each time the negative training data are
22 randomly selected as the same amount with positive training data from one of the 5-fold cross
23 validation split above, and the negative test data are the above 5-fold cross validation test
24 data.

1 When leaving lncRNA PVT1 and LINC00963 out, our model predicts all their variants
2 correctly as positive. The average AUROC values for PVT1 and LINC00963 are 0.9996 and
3 0.9968 separately. However, when leaving lncRNA MIR100HG and MEG3 out, their
4 average AUROC values are 0.6594 and 0.3220 separately, which are a bit low. One possible
5 reason for the different performances between triplexlncRNA and reported triplex lncRNAs
6 is that, in triplexlncRNA, we adopt the variants which are overlapped with experimentally
7 verified triplex forming regions; whereas in the reported triplex lncRNA, we adopt all the
8 variants of that gene, however, maybe not all of these variants could form triplex in practice.
9 Interestingly, when we only use six kmerscore features to train leave-out model for MEG3,
10 its average AUROC value could increase to 0.7610, but this phenomenon is not found when
11 leaving MIR100HG out.

12

13 **Evaluation of triplex DNA sites potential prediction**

14 We use the 5-fold cross validation to evaluate the performance of triplex DNA sites potential
15 prediction model in TriplexFPP, and compare it with other baseline models (Fig. 3a). The
16 determination of the parameters for baseline models in predicting triplex DNA sites follows
17 the same procedure as that of triplex lncRNA prediction. The candidate parameters and the
18 final determined parameters for each baseline model are recorded in Additional file 1: Table
19 S2. The overall performance of CNN is better than that of baseline models, whose average
20 AUROC and AUPRC values are 0.8705 and 0.9671 separately; whereas for baselines
21 methods, the average values of AUROC located in the scope of 0.8635 to 0.8667, and the
22 average values of AUPRC are from 0.9642 to 0.9660, respectively. We then visualize the
23 predicted probability scores of CNN model in each fold (Fig. 3b and Additional file 1: Figure
24 S6). Although some samples are wrongly predicted, the predicted probability scores for most

1 positive data concentrate at around value of 1. This phenomenon indicates that our model can
2 predict most data correctly with high confidence.

3 Besides, the region of 649-708 in HOTAIR sequence is verified to form DNA:RNA
4 triplex [10], our model correctly predicts this site as the triplex forming type. Overall, with
5 the limited data, our results demonstrate that TriplexFPP can effectively distinguish the *in*
6 *vivo* assay defined triplex-forming DNA sites from those background sites with only
7 nucleotide sequence features as the input.

8

9 **TriplexFPP model interpretation**

10 Take the first fold validation in two models in TriplexFPP as examples, we plot the average
11 feature values for each class in the format of heatmap before training (original features),
12 trained after one CNN layer, and trained after two CNN layers in Fig. 4 and Additional file 1:
13 Figure S7. The kmerscore features (the first six features in original features) show obvious
14 differences in two classes, which indicate that the nucleotide compositions have different
15 preferences with regard to the positions in positive and negative data. The kmerscore features
16 also lead to different convolution values after trained with only one CNN layer. However, for
17 the remaining features, the differences of their convolutional values between two classes do
18 not show noticeable differences until trained after 2 CNN layers.

19

20 ***in cis / in trans* triplex-forming lncRNA: an exploration**

21 The lncRNA can form triplex structures with DNA both *in cis* and *in trans*, but if differences
22 exist between *cis* and *trans* targeting of triplexes lncRNAs remains unknown [28]. In this
23 work, we explore the *cis* and *trans* targeting of triplexes lncRNAs following two data
24 sources.

1 From the TDF results between each triplexlncRNA and TriplexDNA, 238 lncRNAs in
2 triplexlncRNAs show both *in cis* and *in trans* interactions with DNA, whose *in cis* binding
3 numbers range from 1 to 2450; whereas the other 141 lncRNAs only show *in trans*
4 interactions with DNA (Additional file 1: Figure S8). Moreover, for a certain lncRNA gene,
5 its variants may show different binding patterns. Among all 130 genes related to
6 triplexlncRNAs, 26 genes contain variants either belong to the type of *in cis* and *in trans*
7 interactions or only *in trans* interactions (Additional file 1: Figure S9 and S10).

8 Besides, for those reported triplex lncRNA, we collect their binding information from the
9 published work in Additional file 1: Table S3 [7, 8, 10, 11, 23, 29]. According to the
10 corresponding experiments, the lncRNA HOTAIR, MEG3, and MIR100HG show *in trans*
11 binding, and PARTICLE and FENDRR show *in cis* binding.

12

13 Discussion

14 LncRNA can exert functions via interacting with DNA. Among all kinds of interactions,
15 DNA:RNA triplex formation is still less understood to us due to the limited number of
16 validation assays. Although varieties of canonical rule-based computational methods have
17 been developed to predict the triplex forming potential for lncRNA and DNA sites, they
18 identify a large number of lncRNAs which can form triplex. However, the limited number of
19 experimentally verified data indicates that maybe not all of them can form triplex in practice.
20 Besides, those computational methods only theoretically calculate the triplex potential, while
21 do not consider any *in vivo* and *in vitro* verified data.

22 Trained with the data obtained from *in vitro* and *in vivo* assays, our newly developed
23 program, namely TriplexFPP, exhibits good prediction performances. Its triplex lncRNA
24 prediction model works effectively by achieving high average scores of evaluation matrices
25 in the 5-fold cross validation. For example, in the removed redundancy datasets with

1 threshold 0.8, the average cross fold validation value of Acc, AUROC, AUPRC, f1-score,
2 and Hm are 95.28%, 0.9649, 0.9996, 0.976, and 0.904, respectively. Besides, the triplex
3 DNA sites potential prediction model in TriplexFPP also works effectively. In the 5-fold
4 cross validation, its average AUROC and AUPRC values are 0.8705 and 0.9671 separately.
5 And most data are predicted correctly with high confidences.

6 We also summarized the *cis* and *trans* targeting of triplexes lncRNAs following the
7 different data sources collected in this work, which may provide some insights to the
8 exploration of lncRNA *cis* and *trans* binding mechanisms.

9 However, one limitation for this work is that the positive data amount is small. And also,
10 some data in our negative class may belong to the positive but are not yet verified. Therefore,
11 we expect more data to be explored to help implementing this tool. Besides, a small fraction
12 lncRNA in our collected data may belong to the R-Loop forming type, which may influence
13 the results somehow.

14

15 Conclusion

16 We proposed a deep learning based program in DNA:RNA triplex formation prediction,
17 namely TriplexFPP. TriplexFPP predicts the most likely triplex-forming lncRNAs from all
18 the lncRNAs with computationally defined triplex forming capacities, and it also predicts the
19 potential of a DNA site to become a triplex. TriplexFPP narrows the scope of possible
20 lncRNAs in forming triplex compared to those mathematic statistic methods. We expect the
21 TriplexFPP can provide insights and references to help to decipher the codes of the lncRNA
22 functions.

23

24 Availability and requirements

- 1 **Project name:** TriplexFPP
- 2 **Project home page:** <https://github.com/yuuuuzhang/TriplexFPP>
- 3 **Operating system(s):** Platform independent
- 4 **Programming language:** Python
- 5 **Other requirements:** python3, tensorflow 1 (≥ 1.12)
- 6 **License:** GNU General Public License v3.0
- 7 **Any restrictions to use by non-academics:** Not applicable
- 8

9 **List of abbreviations**

- 10 TriplexFPP: Triplex Forming Potential Prediction; AUROC: Area Under the ROC; AUPRC: Area Under the PRC; Acc: accuracy; Sn: sensitivity; Sp: specificity; Hm: harmonic mean.
- 11 DBD: DNA Binding Domain; TFO: Triplex Forming Oligonucleotide; DBS: DNA Binding Site; NN: Deep Neural Network; SVM: Support Vector Machine; RF: Random Forest.

14

15 **Declarations**

16 **Ethics approval and consent to participate**

17 Not applicable.

18 **Consent for publication**

19 Not applicable.

20 **Availability of data and materials**

21 The datasets used in this work can be found at

22 https://github.com/yuuuuzhang/TriplexFPP_data, and the code can be found at
23 <https://github.com/yuuuuzhang/TriplexFPP>.

24 **Competing interests**

25 The authors declare that they have no competing interests.

1 **Funding**

2 Publication costs are founded by A*STAR-NTU-SUTD AI Partnership [RGANS1905],
3 Singapore Ministry of Education Academic Research Fund Tier 1 [2020-T1-001-
4 130(RG15/20)], and Singapore Ministry of Education Academic Research Fund Tier 2
5 [MOE2019-T2-2-175]. The funding bodies played no role in the design of the study and
6 collection, analysis, and interpretation of data and in writing the manuscript.

7 **Authors' contributions**

8 ZY, LY, and CKK designed the experiments and wrote the manuscript. ZY implemented the
9 algorithm and do the analyses. All authors have read and approved the final manuscript.

10 **Acknowledgements**

11 We would like to thank Prof. Ivan G. Costa and Chao-Chung Kuo for clarifying our questions
12 regarding to their work.

13

14 **References**

- 15 1. Mercer TR, Mattick JS. Structure and function of long noncoding RNAs in
16 epigenetic regulation. *Nature structural & molecular biology*. 2013 Mar;20(3):300.
- 17 2. Guttman M, Rinn JL. Modular regulatory principles of large non-coding RNAs.
18 *Nature*. 2012 Feb;482(7385):339-46.
- 19 3. Engreitz JM, Ollikainen N, Guttman M. Long non-coding RNAs: spatial amplifiers
20 that control nuclear structure and gene expression. *Nature reviews Molecular cell
21 biology*. 2016 Dec;17(12):756.
- 22 4. Antonov I, Medvedeva YA. Purine-rich low complexity regions are potential RNA
23 binding hubs in the human genome. *F1000Research*. 2018;7.
- 24 5. Schmitz KM, Mayer C, Postepska A, Grummt I. Interaction of noncoding RNA with
25 the rDNA promoter mediates recruitment of DNMT3b and silencing of rRNA genes.
26 *Genes & development*. 2010 Oct 15;24(20):2264-9.
- 27 6. Grote P, Herrmann BG. The long non-coding RNA Fendrr links epigenetic control
28 mechanisms to gene regulatory networks in mammalian embryogenesis. *RNA
29 biology*. 2013 Oct 1;10(10):1579-85.
- 30 7. Mondal T, Subhash S, Vaid R, Enroth S, Uday S, Reinius B, Mitra S, Mohammed
31 A, James AR, Hoberg E, Moustakas A. MEG3 long noncoding RNA regulates the
32 TGF-β pathway genes through formation of RNA–DNA triplex structures. *Nature
33 communications*. 2015 Jul 24;6:7743.

- 1 8. O'Leary VB, Ovsepian SV, Carrascosa LG, Buske FA, Radulovic V, Niyazi M,
2 Moertl S, Trau M, Atkinson MJ, Anastasov N. PARTICLE, a triplex-forming long
3 ncRNA, regulates locus-specific methylation in response to low-dose irradiation.
4 Cell reports. 2015 Apr 21;11(3):474-85.

5 9. Postepska-Igielska A, Giwojna A, Gasri-Plotnitsky L, Schmitt N, Dold A, Ginsberg
6 D, Grummt I. LncRNA Khps1 regulates expression of the proto-oncogene SPHK1
7 via triplex-mediated changes in chromatin structure. Molecular cell. 2015 Nov
8 19;60(4):626-36.

9 10. Kalwa M, Hänelmann S, Otto S, Kuo CC, Franzen J, Joussen S, Fernandez-Rebolledo
10 E, Rath B, Koch C, Hofmann A, Lee SH. The lncRNA HOTAIR impacts on
11 mesenchymal stem cells via triple helix formation. Nucleic acids research. 2016 Dec
12 15;44(22):10631-43.

13 11. Wang S, Ke H, Zhang H, Ma Y, Ao L, Zou L, Yang Q, Zhu H, Nie J, Wu C, Jiao B.
14 LncRNA MIR100HG promotes cell proliferation in triple-negative breast cancer
15 through triplex formation with p27 loci. Cell death & disease. 2018 Jul 24;9(8):1-1.

16 12. Zhao Z, Sentürk N, Song C, Grummt I. lncRNA PAPAS tethered to the rDNA
17 enhancer recruits hypophosphorylated CHD4/NuRD to repress rRNA synthesis at
18 elevated temperatures. Genes & development. 2018 Jun 1;32(11-12):836-48.

19 13. Chu C, Qu K, Zhong FL, Artandi SE, Chang HY. Genomic maps of long noncoding
20 RNA occupancy reveal principles of RNA-chromatin interactions. Molecular cell.
21 2011 Nov 18;44(4):667-78.

22 14. Simon MD, Wang CI, Kharchenko PV, West JA, Chapman BA, Alekseyenko AA,
23 Borowsky ML, Kuroda MI, Kingston RE. The genomic binding sites of a noncoding
24 RNA. Proceedings of the National Academy of Sciences. 2011 Dec
25 20;108(51):20497-502.

26 15. Engreitz JM, Pandya-Jones A, McDonel P, Shishkin A, Sirokman K, Surka C, Kadri
27 S, Xing J, Goren A, Lander ES, Plath K. The Xist lncRNA exploits three-
28 dimensional genome architecture to spread across the X chromosome. Science. 2013
29 Aug 16;341(6147):1237973.

30 16. Sentürk Cetin N, Kuo CC, Ribarska T, Li R, Costa IG, Grummt I. Isolation and
31 genome-wide characterization of cellular DNA: RNA triplex structures. Nucleic
32 acids research. 2019 Mar 18;47(5):2306-21.

33 17. Buske FA, Bauer DC, Mattick JS, Bailey TL. Triplexator: detecting nucleic acid
34 triple helices in genomic and transcriptomic data. Genome research. 2012 Jul
35 1;22(7):1372-81.

36 18. Buske FA, Bauer DC, Mattick JS, Bailey TL. Triplex-Inspector: an analysis tool for
37 triplex-mediated targeting of genomic loci. Bioinformatics. 2013 Aug
38 1;29(15):1895-7.

39 19. He S, Zhang H, Liu H, Zhu H. LongTarget: a tool to predict lncRNA DNA-binding
40 motifs and binding sites via Hoogsteen base-pairing analysis. Bioinformatics. 2015
41 Jan 15;31(2):178-86.

- 1 20. Kuo CC, Hänelmann S, Sentürk Cetin N, Frank S, Zajzon B, Derks JP, Akhade
2 VS, Ahuja G, Kanduri C, Grummt I, Kurian L. Detection of RNA–DNA binding
3 sites in long noncoding RNAs. *Nucleic acids research*. 2019 Apr 8;47(6):e32-.
- 4 21. Frankish A, Diekhans M, Ferreira AM, Johnson R, Jungreis I, Loveland J, Mudge
5 JM, Sisu C, Wright J, Armstrong J, Barnes I. GENCODE reference annotation for
6 the human and mouse genomes. *Nucleic acids research*. 2019 Jan 8;47(D1):D766-
7 73.
- 8 22. Navarro C, Cano C, Cuadros M, Herrera-Merchan A, Molina M, Blanco A. A
9 mechanistic study of lncRNA Fendrr regulation of FoxF1 lung cancer tumor
10 suppressor. InInternational Conference on Bioinformatics and Biomedical
11 Engineering 2016 Apr 20 (pp. 781-789). Springer, Cham.
- 12 23. Ageeli AA, McGovern-Gooch KR, Kaminska MM, Baird NJ. Finely tuned
13 conformational dynamics regulate the protective function of the lncRNA MALAT1
14 triple helix. *Nucleic acids research*. 2019 Feb 20;47(3):1468-81.
- 15 24. Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next-
16 generation sequencing data. *Bioinformatics*. 2012 Dec 1;28(23):3150-2.
- 17 25. Hunt SE, McLaren W, Gil L, Thormann A, Schuilenburg H, Sheppard D, Parton A,
18 Armean IM, Trevanion SJ, Flück P, Cunningham F. Ensembl variation resources.
19 Database. 2018 Jan 1;2018.
- 20 26. Zhang Y, Jia C, Fullwood MJ, Kwoh CK. DeepCPP: a deep neural network based
21 on nucleotide bias information and minimum distribution similarity feature selection
22 for RNA coding potential prediction. *Briefings in Bioinformatics*. 2020 Mar. doi:
23 <https://doi.org/10.1093/bib/bbaa039>.
- 24 27. Zhang Y, Jia C, Kwoh CK. Predicting the interaction biomolecule types for
25 lncRNA: an ensemble deep learning approach. *Briefings in Bioinformatics*. 2020
26 Oct. doi: <https://doi.org/10.1093/bib/bbaa228>.
- 27 28. Mishra K, Kanduri C. Understanding Long Noncoding RNA and Chromatin
28 Interactions: What We Know So Far. *Non-Coding RNA*. 2019 Dec;5(4):54.
- 29 29. Li Y, Syed J, Sugiyama H. RNA-DNA triplex formation by long noncoding RNAs.
30 *Cell chemical biology*. 2016 Nov 17;23(11):1325-33.
- 31
- 32
- 33
- 34
- 35

36 **Figure Legends**

- 37 **Fig. 1.** The architecture of TriplexFPP. TriplexFPP is composed of two models, the
38 corresponding model architecture and parameters are shown.
- 39 **Fig. 2.** Evaluation of lncRNA triplex prediction model in TriplexFPP. **a** The box and whisker
40 plot of the 5-fold cross validation for CNN and 4 baseline modes (SVM, RF, Gradient

1 Boosting, and NN). **b** The comparisons of the 5-fold cross validation performances among
2 data without removing redundancy, removing redundancy with threshold 0.9, and removing
3 redundancy with threshold 0.8.

4 **Fig. 3.** Evaluation of triplex DNA sites potential prediction model in TriplexFPP. **a** The box
5 and whisker plot of 5-fold cross validation for CNN and baseline models. **b** The visualization
6 of the distribution for predicted probability score of the first fold validation data.

7 **Fig. 4.** The average feature values in each class of triplex lncRNA prediction model. Top:
8 original features (the 90-dim features are reshaped to 9*10), middle: features after trained
9 with one CNN layer (x-axis: filter, y-axis: convolution values the 1st to the 15th), and
10 bottom: features after trained with two CNN layers (x-axis: filter, y-axis: convolution values
11 the 1st to the 15th); left: positive data, right: negative data.

12

13 **Supplementary information**

14 **Supplementary information** accompanies this paper at

15 **Additional file 1:** Supplementary materials.

16

Figures

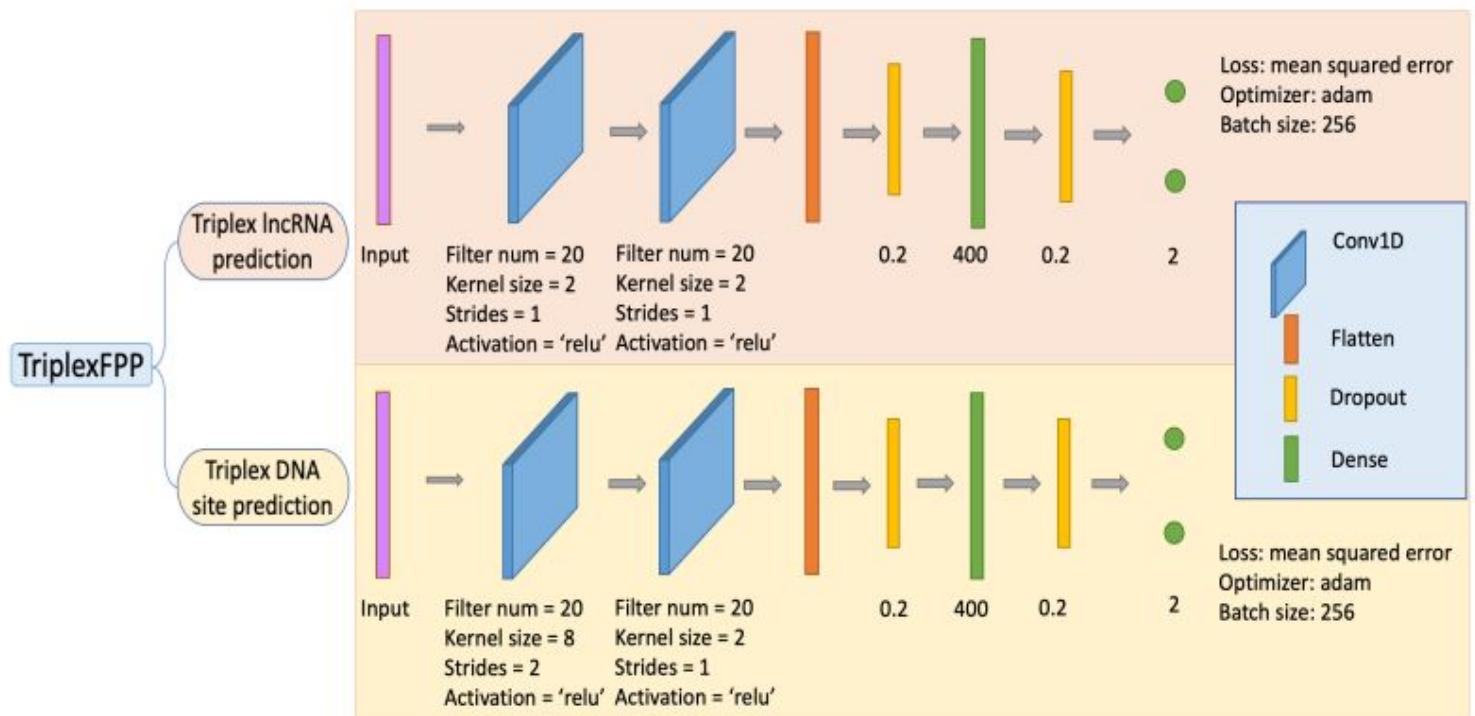


Figure 1

The architecture of TriplexFPP. TriplexFPP is composed of two models, the corresponding model architecture and parameters are shown.

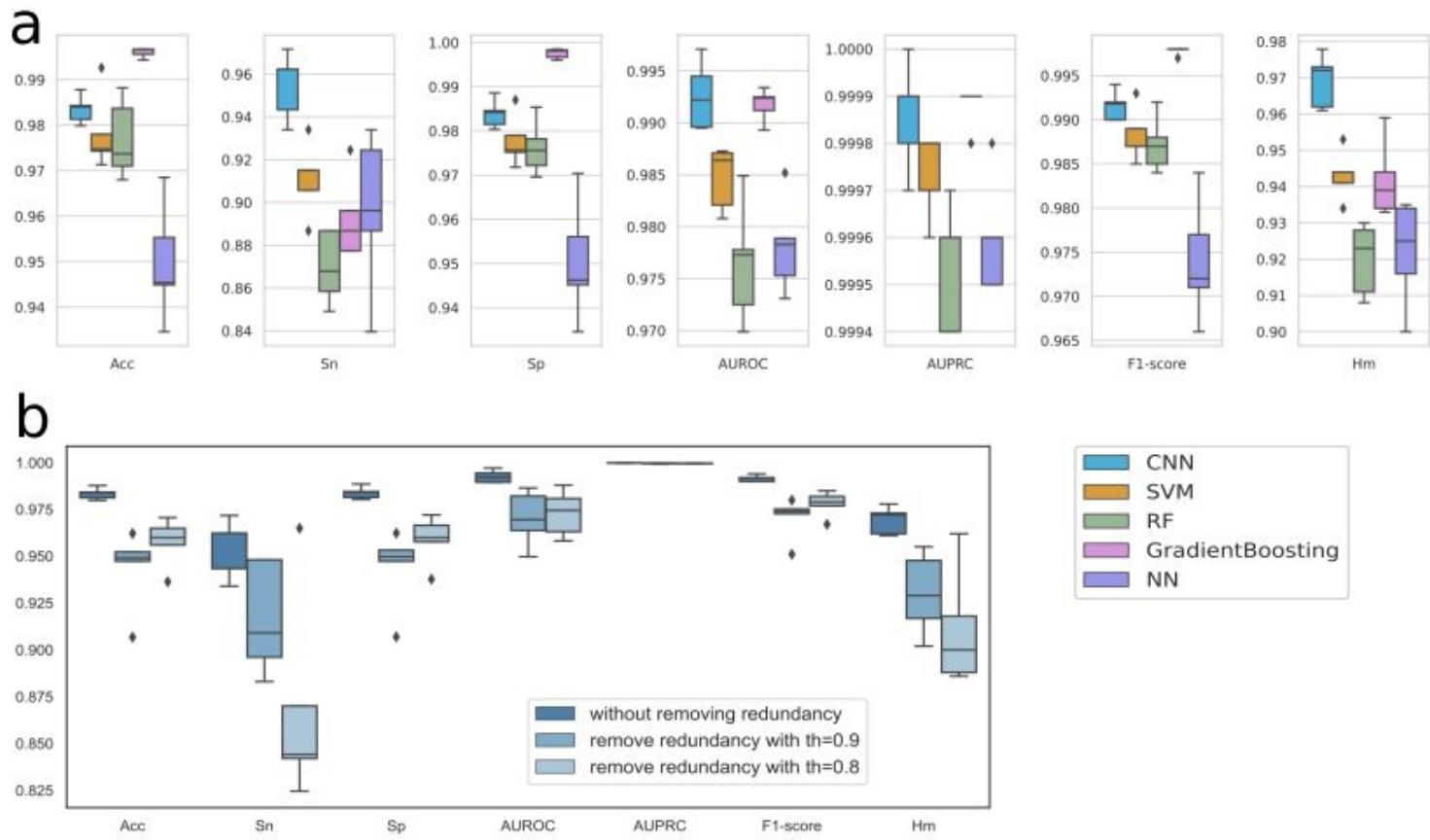


Figure 2

Evaluation of lncRNA triplex prediction model in TriplexFPP. a The box and whisker plot of the 5-fold cross validation for CNN and 4 baseline modes (SVM, RF, Gradient Boosting, and NN). b The comparisons of the 5-fold cross validation performances among data without removing redundancy, removing redundancy with threshold 0.9, and removing redundancy with threshold 0.8.

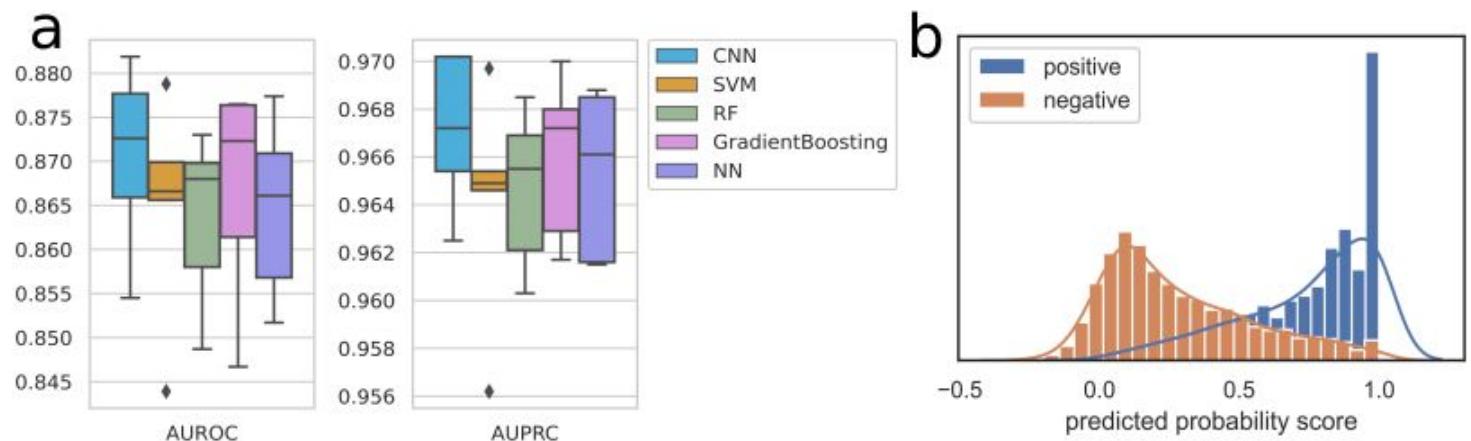


Figure 3

Evaluation of triplex DNA sites potential prediction model in TriplexFPP. a The box and whisker plot of 5-fold cross validation for CNN and baseline models. b The visualization of the distribution for predicted

probability score of the first fold validation data.

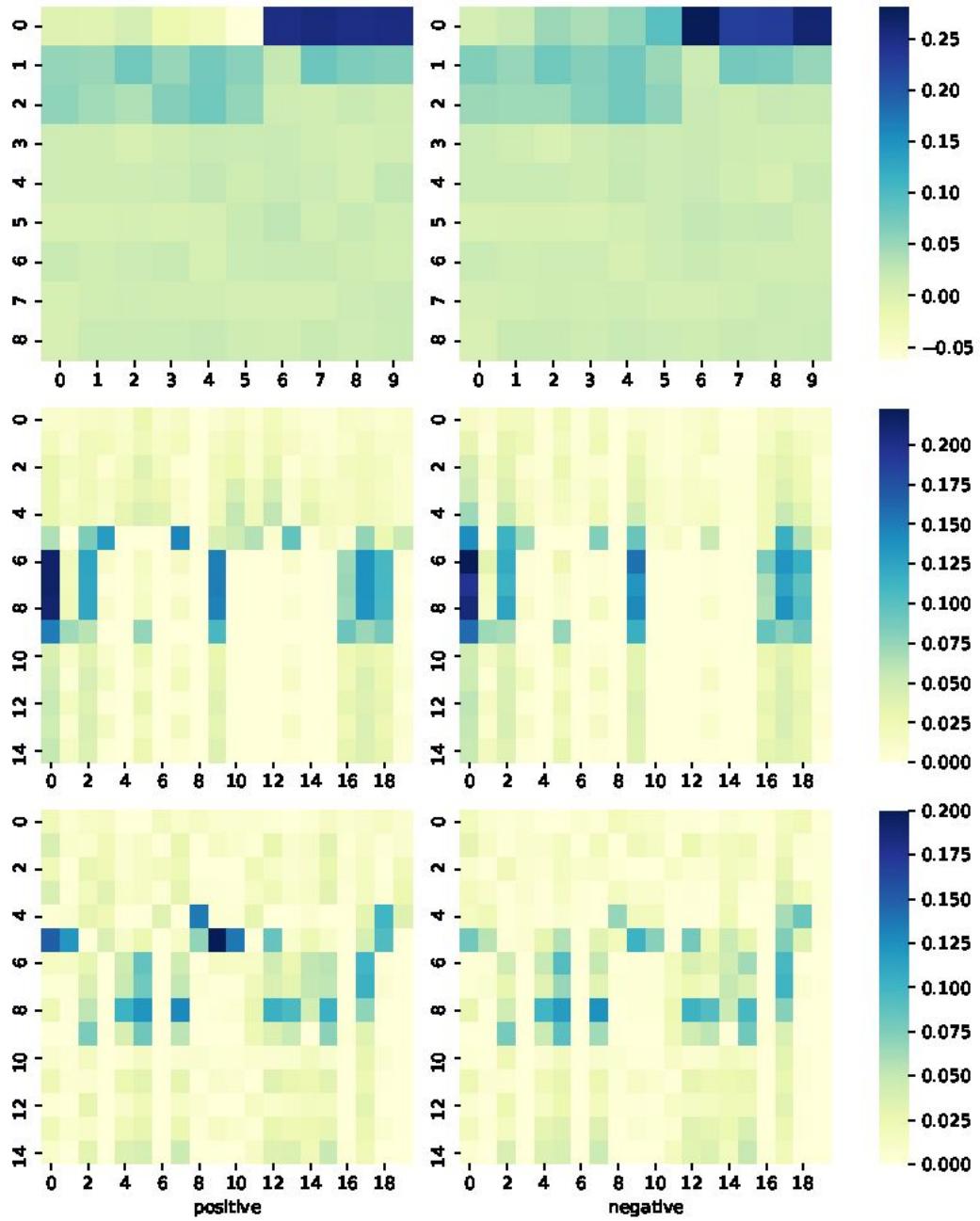


Figure 4

The average feature values in each class of triplex lncRNA prediction model. Top: original features (the 90-dim features are reshaped to 9*10), middle: features after trained with one CNN layer (x-axis: filter, y-axis: convolution values the 1st to the 15th), and bottom: features after trained with two CNN layers (x-axis: filter, y-axis: convolution values the 1st to the 15th); left: positive data, right: negative data.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [additionalfile.docx](#)