

Prediction of enhancer-promoter interactions using the cross-cell information and domain adversarial neural network

Fang Jing

Northwestern Polytechnical University School of Automation <https://orcid.org/0000-0002-1079-7207>

Shao-Wu Zhang ([✉ zhangsw@nwpu.edu.cn](mailto:zhangsw@nwpu.edu.cn))

Northwestern Polytechnical University School of Automation <https://orcid.org/0000-0003-1305-7447>

Shihua Zhang

Academy of Mathematics and Systems Science, Chinese Academy of Sciences

Methodology article

Keywords: Enhancer-promoter interactions, Cell line, Convolutional neural network, Transfer learning, Gradient reversal layer

Posted Date: November 5th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-52120/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 7th, 2020. See the published version at <https://doi.org/10.1186/s12859-020-03844-4>.

Abstract

Background: Enhancer-promoter interactions (EPIs) play key roles in transcriptional regulation and disease progression. Although several computational methods have been developed to predict such interactions, their performances are not satisfactory when training and testing data from different cell lines. Currently, it is still unclear what extent EPI prediction across cell lines can be made based on sequence-level information.

Results: In this work, we present a novel Sequence-based method (called SEPT) to predict the Enhancer-Promoter interactions in the new cell line by using the cross-cell information and Transfer learning. SEPT first learns the features of enhancer and promoter from DNA sequences with convolutional neural network (CNN), then designing the gradient reversal layer of transfer learning to reduce the cell line specific features meanwhile retaining the features associated with EPIs. When the locations of enhancers and promoters are provided in new cell line, SEPT can successfully recognize EPIs in this new cell line based on labeled data of other cell lines. The experiment results show that SEPT can effectively learn the latent import EPIs-related features between cell lines and achieves the best prediction performance in terms of AUC.

Conclusion: SEPT is an effective method for predicting EPIs in the new cell line. Domain adversarial architecture of transfer learning used in SEPT can learn the latent EPI features shared among cell lines from all other existing labeled data. It can be expected that SEPT will be of interest to researchers concerned with biological interaction prediction.

1 Background

The enhancer-promoter interactions (EPIs) play a critical role in gene regulation in eukaryotes. In genetics, a promoter is a region of DNA sequence upstream of a particular gene [1]. The length of a promoter is probably hundreds to thousands of base pairs [2]. Its function aims to initiate gene transcription of a particular gene. While an enhancer is also an important transcriptional regulatory short DNA fragments that further activate the level of transcription of its target genes by contacting close physical proximity to the promoters in the three-dimensional (3D) nuclear space [3]. Hundreds of thousands of enhancers have been estimated to be contained in the human genome. Normally, a promoter is under the control of multiple enhancers, and multiple promoters can be regulated by a single enhancer. Additionally, the distances between interacting enhancer and promoter pairs have varied widely, varying from kilobases to millions of base-pairs because of the chromatin folding in the 3D space [4-8]. Moreover, as more and more studies reported that enhancer sequence variations are associated with serious human diseases [9-12]. Thus, the importance of EPIs for gene expression is matter-of-course.

Over the past decade, many high-throughput experimental approaches, such as chromosome conformation capture-based (3C) [13] and its variants of Hi-C [14] and ChIA-PET [1], have been developed to study the chromatin interactions. Although Hi-C and ChIA-PET could measure the whole genome DNA-

DNA interactions, the genomic resolutions are often low, varying from few kilobases to tens of thousands bases [15-17]. In order to study EPIs, very high (<10kb) resolution data is needed. All these experimental approaches are technically challenging, time-consuming and have high false-negative rate. What is more, the EPIs vary across different cellular conditions and tissues [18]. While the number of 3D chromatin interaction experiments continue to increase, it is still not possible to perform chromatin interaction experiments for all types of cell and tissues. Therefore, computational approaches are urgently desired to complement experimental protocols.

Due to the limitations of experimental approaches, the number of available experimental data of EPIs is still limited. Several computational methods have been developed to predict EPIs of the genome. Depending on the type of the input data, computational methods can mainly be divided into two categories: DNA sequence-based methods and epigenomic data-based methods. For DNA sequence-based methods, PEP [19] and EP2vec [20] took advantage of natural language processing to learn the feature representation of DNA sequences, and SPEID [21] used convolutional neural network to learn the feature representation of DNA sequences. Recently, Zhuang *et al* [22] introduced a novel method to improve the prediction performance of EPIs by using the existing labeled data to pretrain a convolutional neural network (CNN), then adopting the training data from the cell line of interest to continue to train the CNN. Above these methods can accurately predict cell line-specific EPIs from genomic sequence, but they work well only when the training and testing data are from the same cell lines. For the epigenomic data-based methods, IM-PET [23], RIPPLE [24], TargetFinder [2], EpiTensor [25], and JEME [26] used many one-dimensional (1D) local chromatin states including but not limited to transcription factors (TFs) binding, histone modifications and chromatin accessibility signatures to predict EPIs. Though achieving acceptable performance, these models rely on labeled training data from the same cell line as the test data, which limits their usefulness for new cell lines. Generally speaking, high-resolution chromatin interactions experimental data are hard to get, and the cell-specific models have no acceptable generalization due to the specificity of EPIs. Therefore, how to predict the EPIs of new cell line is an urgent problem. And the very intuitive idea is that to train a generic model which learns shared features between different cell lines, and then to predict the EPIs of the new cell line. However, this idea has its drawbacks. The feature distribution learned from the general model is different from that of a particular cell line that we care about. If the distribution of features learned from multiple cell lines is as similar as possible to the distribution of features in the cell line that we care about, then we can make more accurate predictions. Therefore, it is important to develop an effective method to make feature distribution as consistent as possible for transfer knowledge from other cell lines to a specific cell line that we care about.

It is well known that transfer learning (TL), an important branch of machine learning, is widely concerned in image recognition [27] and natural language processing [28], which focuses on the application of knowledge transfer onto new problems. Inspired by the work of [29], which suggested an adversarial neural network by gradient reversal layer (GRL) to fit the feature distribution of source domain and target domain for recognizing handwritten numbers, we proposed a novel method of SEPT to predict EPIs in new cell line. There are no EPI labels information and only has the locations of enhancers and promoters

in a particular cell line. The enhancer-promoter pairs of a particular cell line with no labels are defined as target domain, and labeled enhancer-promoter pairs from other cell lines are defined as source domain. Our goal is to learn the features from the source domain, and further reduce the data distribution differences between the target domain and the source domain through adversarial learning. First, we used convolution and long short-term memory (LSTM) layers to learn the features of EP pairs from the enhancer and promoter sequences. Then, adversarial neural network with GRL was used to reduce domain-specific features. GRL could reverse the direction of the gradient by multiplying the gradient with a negative constraints value. Finally, the trained model learns the EPI-related features from the source domain to predict the EPIs of the target domain.

SEPT is of great significance for three points. i) It could be used as an alternative to the experimental methods, helping other tasks such as identifying mechanisms of SNPs from genome-wide association studies (GWAS) [30]. ii) It could reveal to what extent EPIs of one cell line could be recognized by data from other cell lines. iii) It could improve our understanding of gene regulation and disease progression. To this end, we adopted two strategies: one is to combine training data of different cell lines into one unit, and the other is to design a model of SEPT with transfer learning [31] to transfer the informative features from the combined unit to a new cell line. The experimental results show that SEPT has better performance than several other methods. Model architecture analysis shows that long short-term memory (LSTM) layer, and GRL layer are important for across cell line EPIs prediction. Convolution kernels analysis shows that SEPT can effectively capture sequence features that determine EPIs.

2 Results

2.1 Comparison with other existing state-of-the-art methods

We first compared SEPT with other state-of-the-art methods of LS-SVM [32], SPEID [21] and RIPPLE [24]. SPEID and LS-SVM are the two sequence-based methods for predicting the DNA regulatory elements. LS-SVM [32] widely used the k -mer features, and did not take into account the interactions of high-order features. Because the LS-SVM [32] can only take one DNA sequence as an input sample, so we concatenated the sequences of the enhancer and promoter to train and test LS-SVM. SPEID [21] used deep learning to capture sequence features for predicting the cross-cell-line EPIs, while it lacks the ability to contain cell-specific information. RIPPLE [24] is a supervised model based on epigenomic data from ChIP-seq and DNase-seq experiments, and it only used five cell lines data to train the model, due to HMEC cell line lack of the epigenomic data. Our SEPT method simultaneously considers the source and target domain sequence information. For each test cell line, data from the other six cell lines were merged into a training data set to train SPEID [21], LS-SVM [32] and SEPT. The average results of SEPT, RIPPLE, SPEID and LS-SVM on seven test cell lines are shown in Figure 1, from which we can see that SEPT has the highest average AUC values on seven cell lines than RIPPLE, SPEID and LS-SVM. SEPT achieves 0.72, 0.76, 0.78, 0.77, 0.78, 0.73, and 0.76 for GM12878, HMEC, HUVEC, HeLa-S3, IMR90, K562, NHEK cell line, respectively. SEPT also has the highest average AUPR and F1-score values on seven cell lines than RIPPLE, SPEID and LS-SVM (Tables S1-S2). These results demonstrate that our SEPT can effectively

predict the enhancer-promoter interactions.

2.1 Influence of different neural network architectures in feature learning phase

We designed a series of computational network architectures in feature learning phase (Table S3) to investigate the impact of network structure. The first network architecture (namely BASE) includes only one convolutional layer in feature learning phase (Figure S1(a)). The second network architecture (namely BASE+LSTM) includes one convolutional layer and one LSTM layer in feature learning phase (Figure S1(b)). The third network architecture (namely BASE+FC) include one convolutional layer and one full connection layer in feature learning phase (Figure S1(c)). SEPT includes two convolutional layers and one LSTM layer in feature learning phase (Figure S1(d)). The grid search strategy was used to optimize the hyperparameters of four models in this work.

The average results of BASE, BASE+LSTM, BASE+FC and SEPT on seven test cell lines are shown in Figure 2, from which we can see that adding the full connection layer (BASE+FC) and the LSTM layer (BASE+LSTM) in BASE model can improve the predictive performance of EPIs. Especially, the results of BASE+LSTM is better than that of BASE+FC, which indicates that the long-range dependency of DNA features such as motifs and other features can be captured by the LSTM layer. SEPT shows the best performance than BASE+LSTM, BASE+FC and BASE models, indicating that the first CNN layer learns the individual patterns in the sequences and the second CNN layer learns the high-order interactions between patterns. The high- order interactions may be commonalities among different cell lines, and deep neural network architectures can extract the high-order EPI features from DNA sequences.

2.3 Influence of domain adversarial operation

To validate the effectiveness of domain adversarial operation, we constructed the model of SEP, which has no the domain adversarial network architecture compared with SEPT. Since SEP has no the domain adversarial operation, the data of target domain cannot be utilized in SEP. That is, SEP just used the data of source domain in the training phase. Table 1 shows the average AUC results of SEPT and SEP by training model on one cell line data and test on another cell line data in running 10 times, from which we can see that AUC values of SEPT for each test cell line are higher than that of SEP. The average AUPR and accuracy values of SEPT for each test cell line are also higher than that of SEP (Table S4-S5). These results indicate that domain adversarial operation can effectively improve the predictive performance of EPIs in new cell lines, and SEPT makes use of the EPIs information in other cell lines for recognizing the EPIs in new cell lines.

We also investigated the effectiveness of domain adversarial operation by training model on the six cell lines data and test on one other cell line data. The average AUC values of SEP and SEPT in running 10 times are shown in Table 2, from which we can see that average AUC of SEPT is higher 6% ~ 9% than that of SEP on seven test cell lines. The average AUPR, F1-score and accuracy values of SEP and SEPT in running 10 times are shown in Table S6. These results show that the domain adversarial operation is also effective when using more cell lines data as the training data.

2.4 Results comparison of using different number of cell lines as the source domain data

To investigate the influence of different cell line number used in the source domain, we used different number of cell lines as the source domain data to predict the EPIs with SEPT. For selecting 1, 2, 3, 4, 5, 6 cell lines as the source domain data, each test cell line has 6, 15, 20, 15, 6, 1 results, respectively. The AUC results of using different number of cell lines as the source domain data to train SEPT are shown in Figure 3. Each boxplot in Figure 3 represents the results of one test cell line across all combinations of source domain cell lines. To make the results more reliable, we repeated 10 times, that is, the AUC value of each test cell line is the mean of 10 running results. From Figure 3, we can see that more cell lines as the source domain can achieves higher AUC than less cell lines as the source domain. Especially, when 6 cell lines are mixed as the source domain data to train the model, SEPT achieves the highest AUC values for every test cell line. With the increase of the cell lines number in source domain, AUC value is gradually larger. In addition, for different test cell line, the combination of appropriate cell lines as the source data can improve the performance of SEPT. These results show that using more existing labeled data of EPIs is helpful the prediction of EPIs in new cell line.

2.5 Motifs identified by SEPT

As to investigate the motifs, we identified sequence features for each model by comparing patterns of the convolutional kernels in the first layer with sequence motifs from the database HOCOMOCO Human v11. We reconstructed the output of the first convolutional layer for each input sample sequence, then extracted the subsequence that best match each kernel to compute the position frequency matrix (PFM) from the aligned sub sequences for each kernel. The motif comparison tool of Tomtom [33] was used to match these PFMs to known TF motifs. After obtaining the sequence motifs of the two models, we defined a formula to measure the relative importance of motifs. The formula will consider the occurrence times of the same motif in both models and the ranks of the occurrence times of the motif in SEPT. If a

motif appears more in SEPT and less in SEP then the greater the effect of the motif in SEPT, the higher its relative importance score will be. The formula is defined as: . Here n and m are the occurrence times of motif learned from SEPT and SEP, respectively; r is the rank (in descending order) of the motif in the motif set of SEPT according to the occurrence times. The closer the value is to 1, the more important the motif is in SEPT. To avoid contingency, this process was repeated five times.

We found a set of potentially important transcription factor binding motifs by the transfer model. For each test cell line, top five import motifs involved in EPIs are show in Table S7. For different test cell lines, the models learned common important TF motifs such as ZNF563, THAP1, RXRA, and SP3, which are involved in many important processes, such as the transcriptional regulation and cell-cycle regulation. Interestingly, we found many of potentially important transcription factors that are associated with the corresponding cell lines (Table 3). For instance, MAFB motif learned by SEPT in the K562 cell line is associated with regulating lineage-specific hematopoiesis. This is consistent with the fact that K562 belong to a blood related cell line. NR4A1 motif learned by SEPT in GM12878 cell line was reported to play a role in the vascular response to injury, while ZNF341 motif learned by both SEPT and SEP in GM12878 cell line was reported to involve in the regulation of immune homeostasis. It is consistent with the fact that GM12878 is lymphoblastoid related cell lines. We also provided other motifs identified by both SEPT and SPEID (Table S8). These results show that our SEPT can learn important motifs, and these motifs are relevant to enhancer-promoter loops of a novel cell line.

3 Discussions

One important factor in SEPT is how to aggregate the labeled data of different cell lines into a source domain. We investigated the aggregation of different number of cell lines, and found that aggregating all available labeled data of cell lines as the source domain for training model can yield better performance than aggregating partial cell lines data as the source domain. In addition, there is redundancy among different cell lines. Redundancy not only slows down model training, but also damages prediction performance. Thus, how to aggregate different cell lines as the source domain data is important.

How to integrate other features such as histone modifications, chromatin accessibility and DNA shape into one model for predicting EPIs is also important factors. Although we only use sequence information in this work, the results of are still better than other methods which integrate the sequence and epigenomic information. Thus, if more information related to enhancers and promoters is integrated into our SEPT, it is hope that SEPT can significantly improve the performance of EPIs prediction.

Although SEPT can predict the potential EPIs in new cell line, it needs to provide the location of the enhancers and promoters in advance. Therefore, it is need to develop new methods for identifying the EPIs in specific cell line without enhancer and promoter location information.

4 Conclusions

Although some deep learning methods have been developed to predict EPIs within the same cell, they cannot get a good performance for predicting EPIs in the unlabeled cell lines, due to lack of understanding of the interested cell lines. In this work, we proposed a transfer learning model to predict EPIs in interested (or “new”) cell lines. To better leverage the existing EPIs knowledge, we adopt the adversarial learning mechanism to learn useful features in the existing labeled cell lines and interesting unlabeled cell lines. Experiment results with the domain adversarial operation indicate that it is helpful to predict EPIs in new cell lines. We expect that the model could learn informative features cross domains and reveal some commonalities (common TFs) between source and target domains. By learning commonalities between the source and target domains, SEPT outperforms other state-of-the-art methods for predicting EPIs in new cell lines.

Although SEPT can effectively predict EPIs in specific cell lines from enhancer and promoter sequences, it can be further improved by considering the following factors. 1) SEPT just uses the sequence information of enhancers and promoters. Integrating other experimental data such as core histone modification ChIP-seq data or DNase-seq data can improve the performance of SEPT. 2) Each cell line is treated equally in the source domain, but the contribution of different cell lines should be different for the test cell line. Determining which cell lines should be used as training data is still needed to be explored, as more and more labeled data will become available. 3) Some EPIs maybe have the cell line specificity, while others are universal across many cell lines. Thus, different samples within the same cell line should have different contribution for cross-cell prediction. Assigning a proper weight to each sample can also improve the performance of SEPT. It can be expected that SEPT can be helpful in other biological interactions prediction scenarios [44-45], such as the detection lncRNA-miRNA interactions.

5 Materials And Methods

5.1 Data and preprocessing

We used the same Hi-C data as [2], and downloaded the Hi-C data of seven cell lines of K562 (mesoderm-lineage cells from a patient with leukemia), GM12878 (lymphoblastoid cells), HeLa-S3 (ectoderm-lineage cells from a patient with cervical cancer), HUVEC (umbilical vein endothelial cells), IMR90 (fetal lung fibroblasts), NHEK (epidermal keratinocytes) and HMEC (mammary epithelial cells) from Gene Expression Omnibus (GEO) GSE63525. The human reference genome hg19 was used to define the genomic locations. Promoters and activate enhancers in the first four cell lines were identified using segmentation-based annotations from both ENCODE Segway [46] and ChromHMM of Roadmap Epigenomics [47], only ChromHMM annotations were used in the other cell lines. Then, RNA-seq data from ENCODE were used to select activate promoters according to the rule of their mean FPKM >0.3 with irreproducible discovery rate <0.1 for each cell line. The genome-wide Hi-C measurements were used to annotate all enhancer-promoter pairs as interacting or non-interacting in each cell type. For each enhancer-promoter pair, the distance between promoter and enhancer of the pair is more than 10kb and less than 2Mb [2]. To exclude the effect of distance on determining EPIs, interacting enhancer-promoter pairs were assigned to one bin (the total bin number is 5) based on quantile discretization of the distance between the enhancer and

promoter. Random non-interacting pairs of active enhancers and promoters were assigned to their corresponding bin and then subsampled as the same number of positive samples within each bin. The subsampled non-interacting pairs were considered as negative samples. Table 4 gives the numbers of positive and negative pair in each cell line.

For each positive/negative sample, sequences of enhancers are extended or cut to 3kb flanking regions on location center of enhancers, and promoter are extended or cut to 2kb flanking regions on location center of promoters. One-hot coding format of enhancer and promoter sequence is used as input data of model.

We examined the overlapping number of positive EPIs between any two cell lines. For any two positive EPI pairs from any two cell lines, if the position of the two enhancers and the position of the two promoters both same, the two EPI pairs are considered to overlap. By comparison, positive samples of different cell lines have very little overlap (Table S9).

For comparison with RIPPLE based on epigenetic data, we used data sets from the Roadmap project for the six cell lines. Because we want to make prediction across cell lines, we downloaded the peak files of 14 data sets that are measured in all six cell lines. These data sets include CTCF, POLR2A, H2AZ, H3K27ac, H3K27me3, H3K36me3, H3K4me1, H3K4me2, H3K4me3, H3K79me2, H3K9ac, H3K9me3, H4K20me1 and DNase-seq. An enhancer or promoter sample is represented as a binary vector in which each dimension corresponds to one of the epigenetic data sets. The feature vectors of enhancer and promoter are concatenated to represent each EPI pair. In addition, we also used other two features (i.e., the Pearson's correlation of the 14 signals associated with the enhancer-promoter pair, and the mRNA level of the gene associated with the promoter) to represent each EPI pair.

5.2 SEPT

Domain adaptation [28] is defined as that source domain and target domain share features and categories, but feature distribution is different. The source domain samples with rich information are used to improve the performance of the model in target domain prediction. The source domain has abundant supervised labeling information, and the target domain has no or few labels. Because SEPT used the idea of domain adaptation, we therefore describe the input data in domain adaptation terms. As focusing on the task of EPIs prediction across cell lines, we assume that there is no labeled training data available in test cell line, only the locations of enhancers and promoters are provided. So, we can utilize the abundant supervised labeling information of other cell lines (called source domain) to improve the performance of EPIs prediction in new cell lines without labeled data (called target domain).

An overview of SEPT is shown in Figure 4(a). For predicting the EPIs in cell line #A, unlike the existing two methods which extract the specific features (in Figure 4(b)) or shared features (in Figure 4(c)) of cell lines, SEPT uses the rich information of cell lines #B and #C to extract the features that are relevant to the EPIs in cell line #A by using the transfer learning (TL). As shown in Figure 4(a), SEPT mainly includes feature

extractor, domain discriminator and EPI predictor. SEPT simultaneously trains two classifiers of the main label classifier and the domain discriminator. These two classifiers share feature extractor layers. It is worth mentioning that we used GRL to design the domain discriminator. GRL reverses the direction of the gradient during the back propagation, but it does nothing when forward propagation. The mixed data of labeled EP pairs of cell line #B and cell line #C are used as the source domain data, and the data of unlabeled EP pairs of cell line #A are used as the target domain data. Each training sample has a domain label, with 0 indicating that the sample belongs to the source domain, and 1 indicating that the sample belongs to the target domain. Each mini-batch training data contains an equal number of samples from both source and target domains. In each training iteration, the parameters in feature extractor layers and EPI label predictor layers are updated on the source domain data, while the parameters in feature extractor layers and domain discriminator layers are updated on both the source and target domains data. In other words, in each training iteration, the feature extractor layers learn the features related to EPI from the samples of cell line #B and #C during the first back propagation, while during the second back propagation, the features learned in the feature extractor layers cannot distinguish which cell line the samples come from due to the GRL. With the training going on, SEPT gradually learns the features which are related to EPI and not related to cell lines.

Feature extractor consists of two convolution layers, two max-pooling layers, two dropout layers, and one recurrent long short-term memory (LSTM) layer. Domain discriminator consists of the GRL, one dense layer, one dropout layer, and the output. EPIs predictor consists of one dense layer, one dropout layer, and the output. Since informative features may differ between enhancers and promoters, we use two convolution layers, rectified linear unit (ReLU) and max-pooling layers for enhancers and promoters, respectively. Thus, the inputs are two one-hot matrixes to represent enhancer and promoter sequence, respectively. Because large number of kernels can sufficiently extract the features, and motif features of DNA sequences are short than 40bp, so each convolution layer consists of 300 ‘kernels’ with length 40. Max-pooling layer reduces the output dimension with pool length 20, stride 20. The outputs of the two branches are concatenated into one tensor, which is input to the dropout layer with dropout rates of 0.25. The dropout layer randomly selects partial input data to next layer to avoid overfitting. The recurrent LSTM layer is used to extract feature combinations of the two branches, and the output dimension of LSTM is 100. For domain discriminator, the output of LSTM layer feeds into GRL, and a dense layer maps the learned distributed features to the domain label space. It contains 50 units with ReLU activations. The output feeds into a sigmoid unit to predict the domain probability after dropout layer with dropout rates of 0.5. For EPI predictor, the output of LSTM layer feeds into dense layer, which further maps the learned distributed features to the sample label space. It contains 100 units with ReLU activations. The output feeds into a sigmoid unit to predict the probability after dropout layer with dropout rates of 0.5.

Fig. 4. Deep neural network-based methods for predicting EPIs. (a) SEPT architecture. SEPT uses the feature extractor to learn EPIs-related features from the mixed cell line data, meanwhile the domain discriminator with the transfer learning mechanism is adopted to remove the cell-line-specific features and retain the features independent of cell line. EPI label predictor identifies the EPIs in new cell line based on the learned features. (b) Previous model trains a model on specific cell line data, in which the

training and test data are both from same cell line. (c) Mixed cell line data are used to train a general model for predicting EPIs in a new cell line.

5.2.1 SEPT Model training

We trained SEPT for 80 epochs with mini-batches of 64 samples by back-propagation. In the training phase, source domain data were used to train the feature extractor and the EPI predictor, and both source and target domain data were used to train the feature extractor and the domain discriminator. SEPT seeks to minimize the loss of EPIs label and domain discriminator. Binary cross-entropy loss function for both EPIs label predictor and domain discriminator is used, which is minimized by stochastic gradient descent (SGD) with initialized learning rate initialized equals 0.001. In view of the two optimization objectives, SEPT learns a discriminative representation for EPI prediction and indistinguishable representation for domain prediction. The objective function of the SEPT is defined as follows:

$$\begin{aligned}
 & E(\theta_f, \theta_y, \theta_d) \\
 &= \sum_{i=1}^N L_y(G_y(G_f(x_i; \theta_f); \theta_y), y_i) \\
 &\quad - \lambda \sum_{i=1}^M L_d(G_d(G_f(x_i; \theta_f); \theta_d), d_i) \tag{1} \\
 &= \sum_{i=1}^N L_y^i(\theta_f, \theta_y) - \lambda \sum_{i=1}^M L_d^i(\theta_f, \theta_d)
 \end{aligned}$$

Here, L_y is the loss of EPIs label predictor, L_d is the loss of domain discriminator, G_f is a mapping that maps the input x to a feature vector, G_y is a mapping that maps the feature vector to the EPIs label, G_d is a mapping that maps the feature vector to the domain label, x_i is the i -th sample, θ_f is the parameters of mapping G_f , θ_y is the parameters of mapping G_y , θ_d is the parameters of mapping G_d , y_i is the EPI label of the i -th sample, d_i is the domain label of the i -th sample, N is the number of labeled EPI training samples, M is the number of unlabeled EPI training samples but they have the domain labels, and λ is a constant that controls the tradeoff between two objectives.

It is a problem of minimax optimization, that is, we attempts to seek a saddle point of the functional $E(\theta_f, \theta_y, \theta_d)$ that is delivered by parameters $\hat{\theta}_f$, $\hat{\theta}_y$, and $\hat{\theta}_d$. At the saddle point, the loss of EPI label

predictor and domain discriminator is minimized. The loss of EPI label predictor is minimized by the feature mapping parameter θ_f , while the loss of domain discriminator is maximized by θ_f on account of the GRL that changes the sign of the gradient during the back-propagation. In the end, SEPT learns the features that are discriminative and domain-invariant. The features learned from the cell lines with label information (source domain) are effective for the new cell lines (target domain).

The training procedure of SEPT can be described as follows: i) Randomly separating the dataset of target domain into approximate equal two parts, one as the training data in which each sample has domain label but no EPI label, and the other as the testing data. ii) Mixing the data from other six cell lines and randomly shuffling the data. The mixed data are used as source domain dataset in which each sample has both domain and EPIs labels. iii) Training SEPT with the source domain data and target domain data. iv) Evaluating the performance of SEPT with the test data in target domain.

All the experiments are based on Python by using the scikit-learn machine learning library [48] and Keras framework (<https://keras.io/>) with Tensorflow as back-end [49].

5.3 Evaluation metrics

We adopted the metrics of Accuracy, Precision, Recall, F1-score, AUC and AUPR to assess the performance of SEPT. These metrics are defined respectively as the following [50-53].

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FN} + \text{FP}} \quad (2)$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (3)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (4)$$

$$\text{F1} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (5)$$

where TP is the number of correctly predicted EPIs, TN is the number of correctly predicted non-EPIs, FP is the number of incorrectly predicted EPIs and FN is the number of incorrectly predicted non-EPIs. AUC is the area under the receiver operating characteristic (ROC) curve which is the plot of the true-positive rate (i.e., sensitivity) as a function of false-positive rate (i.e., 1-specificity) based on various thresholds. AUPR is the area under the precision-recall curve which is the plot of the precision as a function of recall based on various thresholds.

Abbreviations

EPIs: Enhancer-promoter interactions; AUC: Area under the receiver operating curves; 3D: Three-dimensional; 3C: Chromosome conformation capture-based approach; CNN: Convolutional neural network; GWAS: Genome-wide association studies; TL: Transfer learning; GRL: Gradient reversal layer; LSTM: Long short-term memory; ReLU: Rectified linear unit; SGD: Stochastic gradient descent; AUPR: Area under the precision-recall curves; lncRNA: Long non-coding RNA; miRNA: Micro RNA.

Declarations

Ethics approval and consent to participate

Not applicable. All of the repository data is freely available.

Consent for publication

Not applicable.

Availability of data and materials

The related source codes and datasets are available at <https://github.com/NWPU-903PR/SEPT>.

Competing interests

The authors declare that they have no competing interests.

Funding

This work has been supported by the National Natural Science Foundation of China (No. 61873202 to SWZ; No. 11661141019 and 61621003 to SZ); the Strategic Priority Research Program of the Chinese Academy of Sciences (CAS) [No. XDB13040600], the National Ten Thousand Talent Program for Young Top-notch Talents, the Key Research Program of the Chinese Academy of Sciences, [No. KFZD-SW-219] and CAS Frontier Science Research Key Project for Top Young Scientist [No. QYZDB-SSW-SYS008]. The funding body did not play any roles in the design of the study and collection, analysis, and interpretation of data and in writing the manuscript.

Authors' contributions

FJ designed and performed the experiments, and wrote the initial manuscript. Both FJ and SZ designed the methods. SZ revised the manuscript. FJ, SWZ and SZ analysed the results. All authors participated in the definition of the process and approved the final manuscript.

Acknowledgments

The authors would like to thank our group members for their valuable suggestions.

Authors' information

Fang Jing: Email: jingnpu@foxmail.com; Address: 127 West Youyi Road, Xi'an, Shaanxi, 710072, China

Shao-Wu Zhang: Email: zhangsw@nwpu.edu.cn; Address: 127 West Youyi Road, Xi'an, Shaanxi, 710072, China

Shihua Zhang: Email: zsh@amss.ac.cn; Address: 55 Zhongguancun East Road, Beijing, 10090, China.

References

1. Fullwood MJ, Ruan Y. ChIP-based methods for the identification of long-range chromatin interactions. *J Cell Biochem*. 2009;107(1):30-9.
2. Whalen S, Truty RM, Pollard KS. Enhancer–promoter interactions are encoded by complex genomic signatures on looping chromatin. *Nat Genet*. 2016; 48(5):488-96.
3. Shlyueva D, Stampfel G, Stark A. Transcriptional enhancers: from properties to genome-wide predictions. *Nat Rev Genet*. 2014; 15(4):272-86.
4. Visel A, Rubin EM, Pennacchio LA. Genomic views of distant-acting enhancers. *Nature*. 2009; 461(7261):199-205.
5. Van Steensel B, Dekker J. Genomics tools for unraveling chromosome architecture. *Nat Biotechnol*. 2010; 28(10):1089-95.
6. Bickmore WA, van Steensel B. Genome architecture: domain organization of interphase chromosomes. *Cell*. 2013; 152(6):1270-84.
7. Dekker J, Mirny L. The 3D genome as moderator of chromosomal communication. *Cell*. 2016; 164(6):1110-21.
8. Rowley MJ, Corces VG. The three-dimensional genome: principles and roles of long-distance interactions. *Curr Opin Cell Biol*. 2016; 40:8-14.
9. Achinger-Kawecka J, Clark SJ. Disruption of the 3D cancer genome blueprint. *Epigenomics*. 2017; 9(1):47-55.
10. Franke M, Ibrahim DM, Andrey G, Schwarzer W, Heinrich V, Schöpflin R, Kraft K, Kempfer R, Jerković I, Chan W-L. Formation of new chromatin domains determines pathogenicity of genomic duplications. *Nature*. 2016; 538(7624):265-69.
11. Claussnitzer M, Dankel SN, Kim K-H, Quon G, Meuleman W, Haugen C, Glunk V, Sousa IS, Beaudry JL, Puviindran V. FTO obesity variant circuitry and adipocyte browning in humans. *New Engl J Med*. 2015; 373(10):895-907.

12. Lupiáñez DG, Kraft K, Heinrich V, Krawitz P, Brancati F, Klopocki E, Horn D, Kayserili H, Opitz JM, Laxova R. Disruptions of topological chromatin domains cause pathogenic rewiring of gene–enhancer interactions. *Cell*. 2015; 161(5):1012-25.
13. Dekker J, Rippe K, Dekker M, Kleckner N. Capturing chromosome conformation. *science*. 2002; 295(5558):1306-11.
14. Lieberman-Aiden E, Van Berkum NL, Williams L, Imakaev M, Ragoczy T, Telling A, Amit I, Lajoie BR, Sabo PJ, Dorschner MO. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *science*. 2009; 326(5950):289-93.
15. Dixon JR, Selvaraj S, Yue F, Kim A, Li Y, Shen Y, Hu M, Liu JS, Ren B. Topological domains in mammalian genomes identified by analysis of chromatin interactions. *Nature*. 2012; 485(7398):376-80.
16. De Laat W, Duboule D. Topology of mammalian developmental enhancers and their regulatory landscapes. *Nature*. 2013; 502(7472):499-506.
17. Rao SS, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, Robinson JT, Sanborn AL, Machol I, Omer AD, Lander ES. A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell*. 2014; 159(7):1665-80.
18. Andersson R, Gebhard C, Miguel-Escalada I, Hoof I, Bornholdt J, Boyd M, Chen Y, Zhao X, Schmidl C, Suzuki T. An atlas of active enhancers across human cell lines and tissues. *Nature*. 2014; 507(7493):455-61.
19. Yang Y, Zhang R, Singh S, Ma J. Exploiting sequence-based features for predicting enhancer–promoter interactions. *Bioinformatics*. 2017; 33(14):i252-i260.
20. Zeng W, Wu M, Jiang R. Prediction of enhancer–promoter interactions via natural language processing. *BMC Genomics*. 2018; 19(2):84.
21. Singh S, Yang Y, Póczos B, Ma J. Predicting enhancer–promoter interaction from genomic sequence with deep neural networks. *Quant Biol*. 2019; 7(2):122-37.
22. Zhuang Z, Shen X, Pan W. A simple convolutional neural network for prediction of enhancer–promoter interactions with DNA sequence data. *Bioinformatics*. 2019; 17(35):2899-906.
23. He B, Chen C, Teng L, Tan K. Global view of enhancer–promoter interactome in human cells. *Proceedings of the National Academy of Sciences*. 2014; 111(21):E2191-E2199.
24. Roy S, Siahpirani AF, Chasman D, Knaack S, Ay F, Stewart R, Wilson M, Sridharan R. A predictive modeling approach for cell line-specific long-range regulatory interactions. *Nucleic Acids Res*. 2015; 43(18):8694-712.
25. Zhu Y, Chen Z, Zhang K, Wang M, Medovoy D, Whitaker JW, Ding B, Li N, Zheng L, Wang W. Constructing 3D interaction maps from 1D epigenomes. *Nat Commun*. 2016; 7:10812. doi:10.1038/ncomms10812
26. Cao Q, Anyansi C, Hu X, Xu L, Xiong L, Tang W, Mok MT, Cheng C, Fan X, Gerstein M. Reconstruction of enhancer–target networks in 935 samples of human primary cells, tissues and cell lines. *Nat Genet*. 2017; 49(10):1428-36.

27. Sun B, Saenko K: Deep coral. Correlation alignment for deep domain adaptation. In: European Conference on Computer Vision. 2016; Springer: 443-50.
28. Abdelwahab M, Busso C. Domain adversarial for acoustic emotion recognition. IEEE/ACM Transactions on Audio, Speech, and Language Processing. 2018; 26(12):2423-35.
29. Ganin Y, Ustinova E, Ajakan H, Germain P, Larochelle H, Laviolette F, Marchand M, Lempitsky V. Domain-adversarial training of neural networks. The Journal of Machine Learning Research. 2016; 17(1):2096-30.
30. Guo H, Ahmed M, Zhang F, Yao CQ, Li S, Liang Y, Hua J, Soares F, Sun Y, Langstein J. Modulation of long noncoding RNAs by risk SNPs underlying genetic predispositions to prostate cancer. Nat Genet. 2016; 48(10):1142-50.
31. Pan SJ, Yang Q. A survey on transfer learning. IEEE Transactions on knowledge and data engineering. 2009; 22(10):1345-59.
32. Dongwon L. LS-GKM: a new gkm-SVM for large-scale datasets. Bioinformatics. 2016; 32(14): 2196-98.
33. Gupta, S., Stamatoyannopoulos, J.A., Bailey, T.L. et al. Quantifying similarity between motifs. Genome Biol. 2007; 8, R24. doi: <https://doi.org/10.1186/gb-2007-8-2-r24>
34. Eychène, Alain, Rocques N , Pouponnot C. A new MAFia in cancer. Nature Reviews Cancer. 2008; 8(9):683-93.
35. Fryer C, Archer T. Chromatin remodelling by the glucocorticoid receptor requires the BRG1 complex. Nature. 1998; 393(6680):88-91.
36. Papac-Milicevic N, Breuss JM, Zaujec J, et al. The interferon stimulated gene 12 inactivates vasculoprotective functions of NR4A nuclear receptors. Circulation Research. 2012; 110(8):e50-e63.
37. Stefanie FJ, Hartberger JM, Manfred F, et al. ZNF341 controls STAT3 expression and thereby immunocompetence. ence Immunology. 2018; 3(24):eaat4941. doi:10.1126/scimmunol.aat4941
38. Bowman CJ, Ayer DE, Dynlacht BD. Foxk proteins repress the initiation of starvation-induced atrophy and autophagy programs. Nature Cell Biology. 2014; 16(12):1202-14.
39. Bower KE, Fritz JM, McGuire KL. Transcriptional repression of MMP-1 by p21SNFT and reduced in vitro invasiveness of hepatocarcinoma cells. Oncogene. 2004; 23(54):8805-14.
40. Yang YJ, Baltus AE, Mathew RS, et al. Microcephaly gene links trithorax and REST/NRSF to control neural stem cell proliferation and differentiation. Cell. 2012; 151(5): 1097-112.
41. Beyer CA, Cabanelas ME, Berquist TH. Unilateral facet dislocations and fracture-dislocations of the cervical spine. The Bone & Joint Journal. 1992; 73(6):977-81.
42. Bamforth SD, Bragança J, Eloranta JJ, et al. Cardiac malformations, adrenal agenesis, neural crest defects and exencephaly in mice lacking Cited2, a new Tfap2 co-activator. Nat Genet. 2001; 29(4):469-74.
43. Takashima H, Nishio H, Wakao H, et al. Molecular cloning and characterization of a KRAB-containing zinc finger protein, ZNF317, and its isoforms. Biochemical and Biophysical Research

- Communications. 2001; 288(4): 771-9.
44. Hu P, Huang YA, Chan KCC, You ZH. Learning Multimodal Networks From Heterogeneous Data for Prediction of lncRNA–miRNA Interactions. IEEE/ACM Transactions on Computational Biology and Bioinformatics. 2020; 17(5): 1516-24.
45. Hu PW, Chan KCC, You ZH. Large-scale prediction of drug-target interactions from deep representations. 2016 International Joint Conference on Neural Networks. 2016; 1236-43.
46. Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, Aken BL, Barrell D, Zadissa A, Searle S. GENCODE: the reference human genome annotation for The ENCODE Project. Genome Res. 2012; 22(9):1760-74.
47. Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, Heravi-Moussavi A, Kheradpour P, Zhang Z, Wang J, Ziller MJ. Integrative analysis of 111 reference human epigenomes. Nature. 2015; 518(7539):317-30.
48. Swami A, Jain R. Scikit-learn: Machine learning in Python. Journal of machine learning research. 2011; 12(10):2825-30.
49. Abadi M, Barham P, Chen J, Chen Z, Davis A, Dean J, Devin M, Ghemawat S, Irving G, Isard M, et al. Tensorflow: A system for large-scale machine learning. arXiv preprint arXiv:1603.04467.
50. Jing F, Zhang S, Cao Z, Zhang S. An integrative framework for combining sequence and epigenomic data to predict transcription factor binding sites using deep learning. IEEE/ACM transactions on computational biology and bioinformatics. 2019; 1-1 doi: 10.1109/TCBB.2019.2901789
51. Fan XN, Zhang SW. LPI-BLS: Predicting lncRNA-protein interactions with a broad learning system-based stacked ensemble classifier. Neurocomputing. 2019; 370:88-93.
52. Zhang SW, Zhang XX, Fan XN, Li WN. LPI-CNNCP: Prediction of lncRNA-protein interactions by using convolutional neural network with the copy-padding trick. Analytical Biochemistry. 2020; 601: 113767.
53. Fan XX, Zhang SW, Zhang SY and Ni JJ. lncRNA_Mdeep: An Alignment-Free Predictor for Distinguishing Long Non-Coding RNAs from Protein-Coding Transcripts by Multimodal Deep Learning. Int. J. Mol. Sci. 2020; 21: 5222.

Tables

Due to technical limitations, the tables are provided in the Supplementary Files section.

Figures

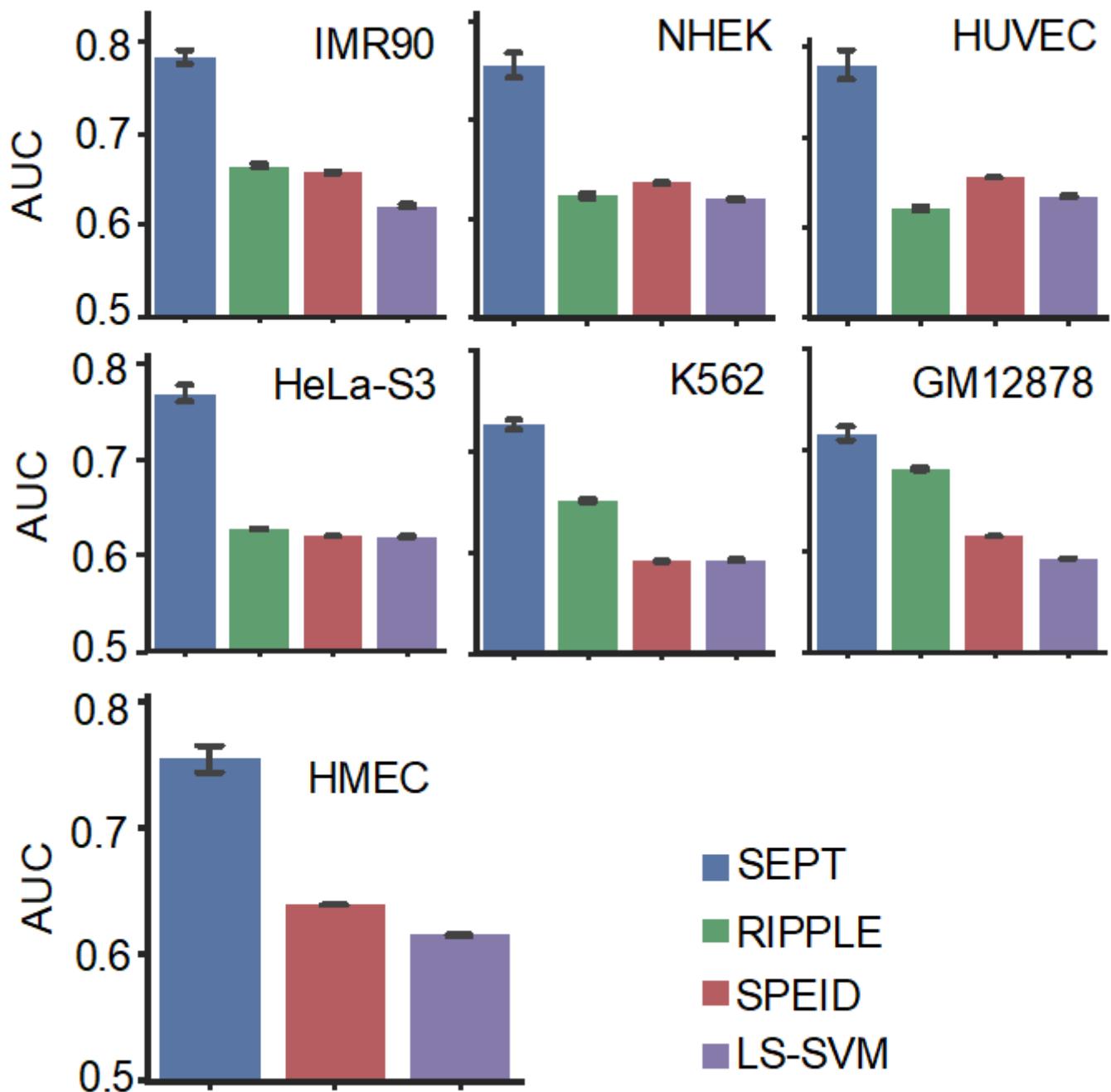


Figure 1

Average AUC of SEPT, RIPPLE, SPEID and LS-SVM on seven test cell lines.

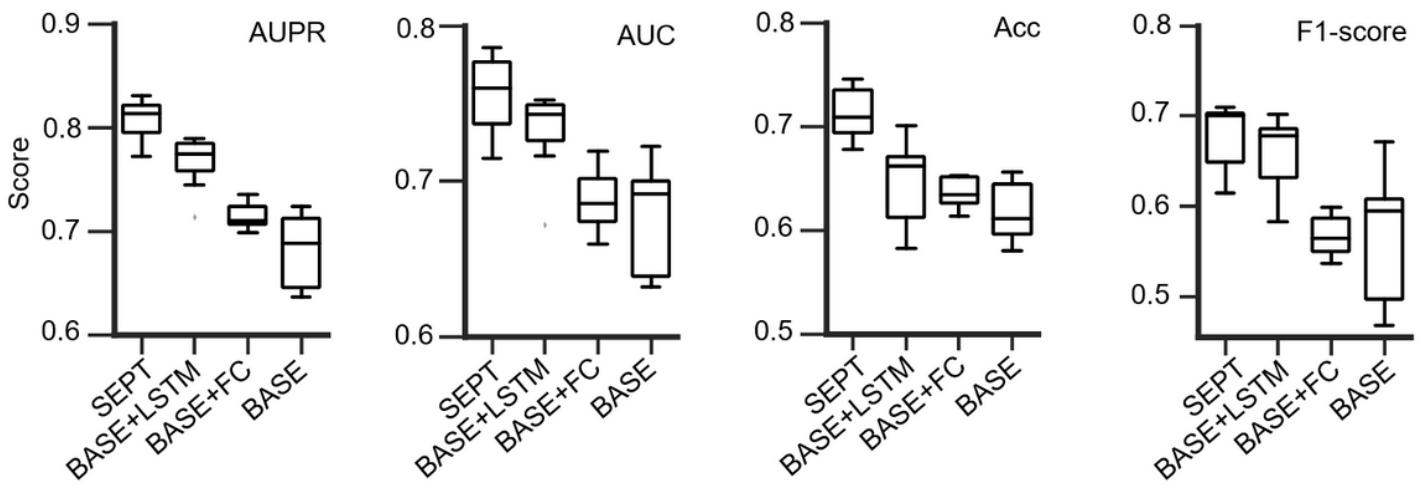


Figure 2

Average results of the BASE, BASE+LSTM, BASE+FC and SEPT on seven test cell lines with the optimal hyperparameters.

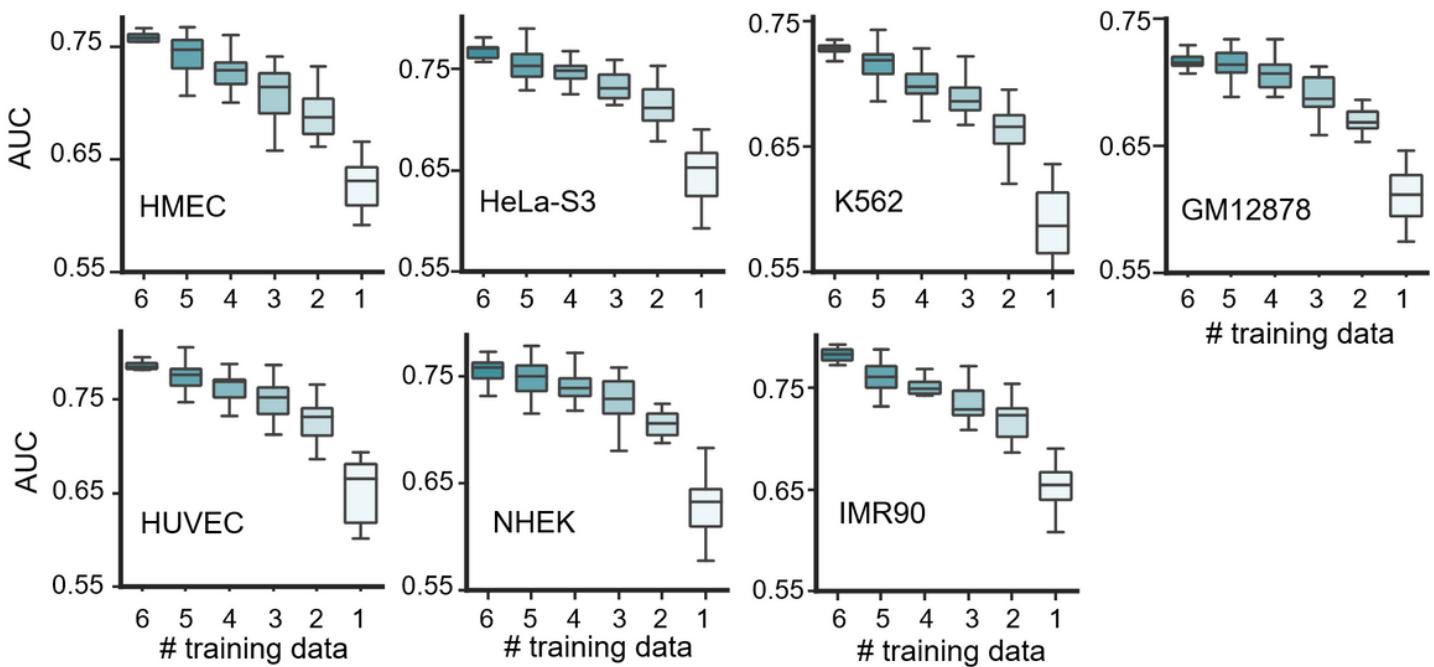


Figure 3

AUC boxplots of different number of cell lines as the source domain in SEPT. The horizontal axis represents the number of cell lines in the source domain.

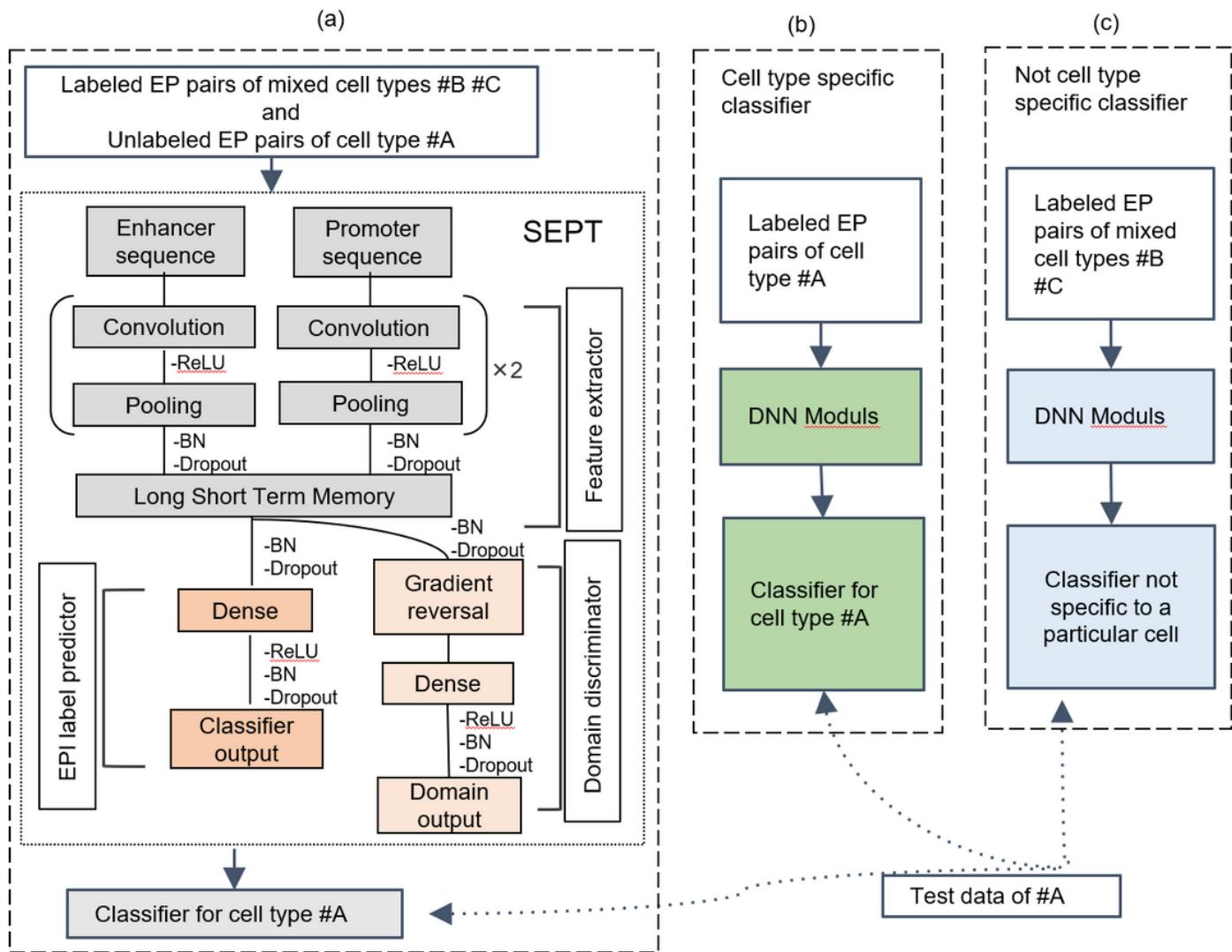


Figure 4

Deep neural network-based methods for predicting EPIs. (a) SEPT architecture. SEPT uses the feature extractor to learn EPIs-related features from the mixed cell line data, meanwhile the domain discriminator with the transfer learning mechanism is adopted to remove the cell-line-specific features and retain the features independent of cell line. EPI label predictor identifies the EPIs in new cell line based on the learned features. (b) Previous model trains a model on specific cell line data, in which the training and test data are both from same cell line. (c) Mixed cell line data are used to train a general model for predicting EPIs in a new cell line.

Supplementary Files

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

- SEPTSupplementaryMaterials201014.docx
- Tables14.docx