

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

3                                   Fang Jing<sup>1</sup>, Shao-Wu Zhang<sup>1\*</sup> and Shihua Zhang<sup>2,3,4\*</sup>

4       <sup>1</sup>Key Laboratory of Information Fusion Technology of Ministry of Education, School of  
5       Automation, Northwestern Polytechnical University, Xi'an 710072, China

6       <sup>2</sup>NCMIS, CEMS, RCSDS, Academy of Mathematics and Systems Science, Chinese Academy of  
7       Sciences, Beijing 100190, China

8       <sup>3</sup>School of Mathematical Sciences, University of Chinese Academy of Sciences, Beijing 100049,  
9       China

10      <sup>4</sup>Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences,  
11      Kunming 650223, China

12      \* Corresponding authors. Emails: [zhangsw@nwpu.edu.cn](mailto:zhangsw@nwpu.edu.cn); [zsh@amss.ac.cn](mailto:zsh@amss.ac.cn)

13      **Abstract**

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

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

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

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

37

## 38 **1 Background**

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

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

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

66 Due to the limitations of experimental approaches, the number of available  
67 experimental data of EPIs is still limited. Several computational methods have been  
68 developed to predict EPIs. Depending on the type of the input data, computational  
69 methods can mainly be divided into two categories: DNA sequence-based methods and  
70 epigenomic data-based methods. For DNA sequence-based methods, PEP [19] and  
71 EP2vec [20] took advantage of natural language processing to learn the feature  
72 representation of DNA sequences, and SPEID [21] used convolutional neural network  
73 to learn the feature representation of DNA sequences. Recently, Zhuang et al [22]  
74 introduced a novel method to improve the prediction performance of EPIs by using the  
75 existing labeled data to pretrain a convolutional neural network (CNN), then adopting  
76 the training data from the cell line of interest to continue to train the CNN. Above these  
77 methods can accurately predict cell line-specific EPIs from genomic sequence, but they  
78 work well only when the training and testing data are from the same cell lines. For the  
79 epigenomic data-based methods, IM-PET [23], RIPPLE [24], TargetFinder [2],  
80 EpiTensor [25], and JEME [26] used many one-dimensional (1D) local chromatin states  
81 including but not limited to transcription factors (TFs) binding, histone modifications  
82 and chromatin accessibility signatures to predict EPIs. Though achieving acceptable  
83 performance, these models rely on labeled training data from the same cell line as the  
84 test data, which limits their usefulness for new cell lines. Generally speaking, high-  
85 resolution chromatin interactions experimental data are hard to get, and the cell-specific  
86 models have no acceptable generalization due to the specificity of EPIs. Therefore, how

87 to predict the EPIs of new cell line is an urgent problem. And the very intuitive idea is  
88 that to train a generic model which learns shared features between different cell lines,  
89 and then to predict the EPIs of the new cell line. However, this idea has its drawbacks.  
90 The feature distribution learned from the general model is different from that of a  
91 particular cell line that we care about. If the distribution of features learned from  
92 multiple cell lines is as similar as possible to the distribution of features in the cell line  
93 that we care about, then we can make more accurate predictions. Therefore, it is  
94 important to develop an effective method to make feature distribution as consistent as  
95 possible for transfer knowledge from other cell lines to a specific cell line that we care  
96 about.

97 It is well known that transfer learning (TL), an important branch of machine  
98 learning, is widely concerned in image recognition [27] and natural language  
99 processing [28], which focuses on the application of knowledge transfer onto new  
100 problems. Inspired by the work of [29], which suggested an adversarial neural network  
101 by gradient reversal layer (GRL) to fit the feature distribution of source domain and  
102 target domain for recognizing handwritten numbers, we proposed a novel method of  
103 SEPT to predict EPIs in the new cell line. There are no EPI labels information and only  
104 has the locations of enhancers and promoters in a particular cell line. The enhancer-  
105 promoter pairs of a particular cell line with no labels are defined as target domain, and  
106 labeled enhancer-promoter pairs from other cell lines are defined as source domain. Our  
107 goal is to learn the features from the source domain, and further reduce the data  
108 distribution differences between the target domain and the source domain through  
109 adversarial learning. First, we used convolution layers and long short-term memory  
110 (LSTM) layer to learn the features of EP pairs from the enhancer and promoter  
111 sequences. Then, adversarial neural network with GRL was used to reduce domain-  
112 specific features. GRL could reverse the direction of the gradient by multiplying the  
113 gradient with a negative constraints value. Finally, the trained model learns the EPI-  
114 related features from the source domain to predict the EPIs of the target domain.

115 SEPT is of great significance for three points. i) It could be used as an alternative

116 to the experimental methods, helping other tasks such as identifying mechanisms of  
117 SNPs from genome-wide association studies (GWAS) [30]. ii) It could reveal to what  
118 extent EPIs of one cell line could be recognized with the data from other cell lines. iii)  
119 It could improve our understanding of gene regulation and disease progression. To this  
120 end, we adopted two strategies: one is to combine different cell line data as the training  
121 data, and the other is to design the SEPT model with transfer learning [31] to transfer  
122 the informative features from the data of different cell lines to a new cell line.  
123 Experimental results show that SEPT has better performance than several other state-  
124 of-the-art methods. Model architecture analysis shows that LSTM layer and GRL layer  
125 are important for the EPIs prediction of across cell lines. Convolution kernels analysis  
126 shows that SEPT can effectively capture sequence features that determine EPIs.

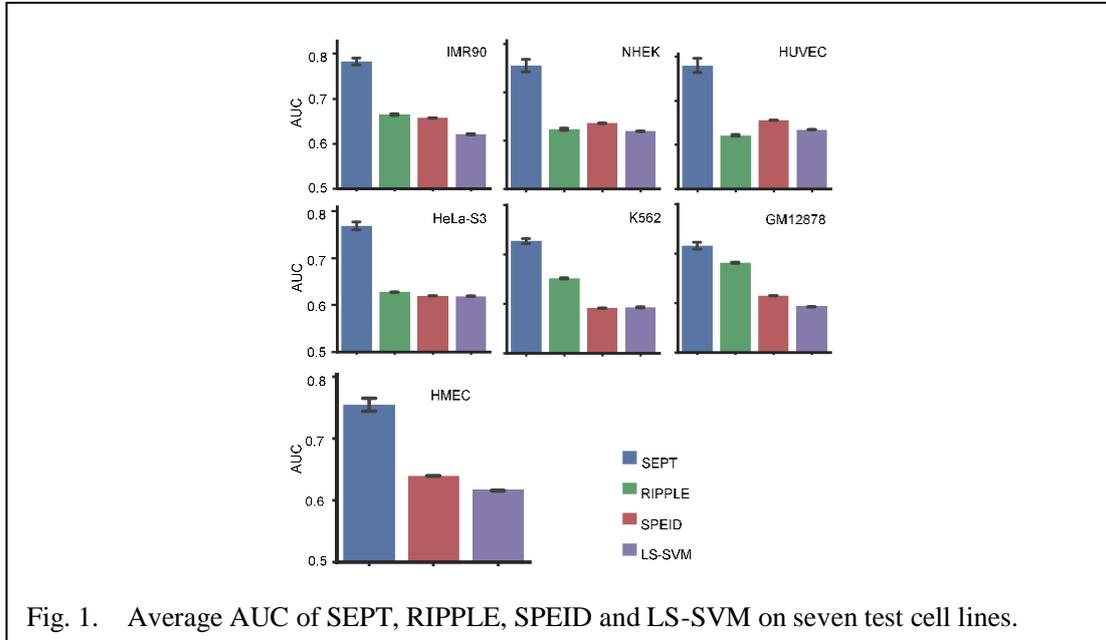
127

## 128 **2 Results**

### 129 **2.1 Comparison with other existing state-of-the-art methods**

130 We first compared SEPT with other state-of-the-art methods of LS-SVM [32], SPEID  
131 [21] and RIPPLE [24]. SPEID and LS-SVM are the two sequence-based methods for  
132 predicting the DNA regulatory elements. LS-SVM [32] widely used the  $k$ -mer features,  
133 and did not take into account the interactions of high-order features. Because the LS-  
134 SVM [32] can only take one DNA sequence as an input sample, so we concatenated the  
135 sequences of the enhancer and promoter to train and test LS-SVM. SPEID [21] used  
136 deep learning to capture sequence features for predicting the cross-cell-line EPIs, while  
137 it lacks the ability to contain cell-specific information. RIPPLE [24] is a supervised  
138 model based on epigenomic data from ChIP-seq and DNase-seq experiments, and it  
139 only used five cell lines data to train the model, due to HMEC cell line lack of the  
140 epigenomic data. Our SEPT method simultaneously considers the source and target  
141 domain sequence information. For each test cell line, data from the other six cell lines  
142 were combined the training data to train SPEID [21], LS-SVM [32] and SEPT. The  
143 average results of SEPT, RIPPLE, SPEID and LS-SVM on seven test cell lines are  
144 shown in Figure 1, from which we can see that SEPT has the highest AUC values on

145 seven cell lines than RIPPLE, SPEID and LS-SVM. SEPT achieves 0.72, 0.76, 0.78,  
146 0.77, 0.78, 0.73, and 0.76 for GM12878, HMEC, HUVEC, HeLa-S3, IMR90, K562,  
147 NHEK cell line, respectively. These results demonstrate that our SEPT can effectively  
148 predict the enhancer-promoter interactions.



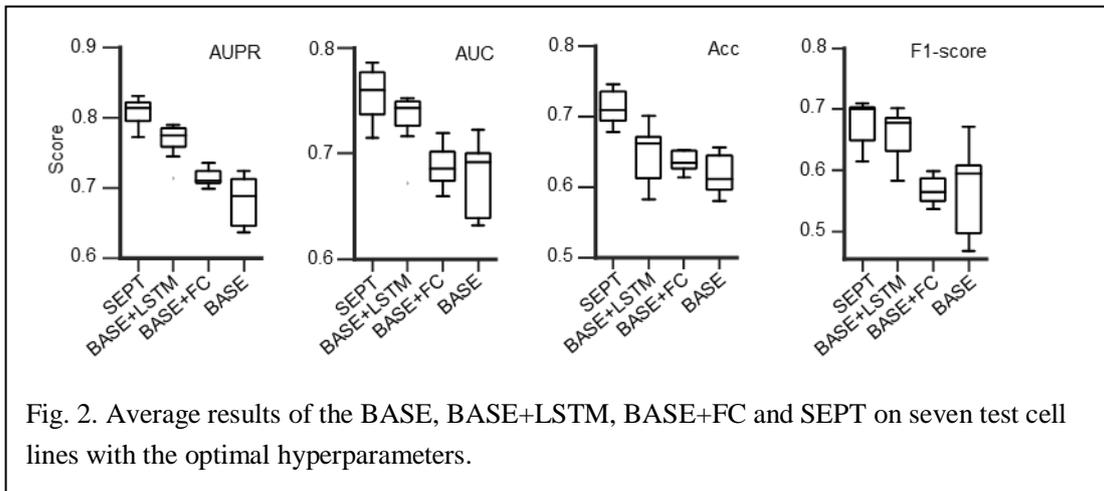
149

## 150 2.1 Influence of different neural network architectures in feature learning phase

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

161 The average results of BASE, BASE+LSTM, BASE+FC and SEPT on seven test  
162 cell lines are shown in Figure 2, from which we can see that adding the full connection

163 layer (BASE+FC) and the LSTM layer (BASE+LSTM) in BASE model can improve  
 164 the predictive performance of EPIs. Especially, the results of BASE+LSTM is better  
 165 than that of BASE+FC, which indicates that the long-range dependency of DNA  
 166 features such as motifs and other features can be captured by the LSTM layer. SEPT  
 167 shows the best performance than BASE+LSTM, BASE+FC and BASE models,  
 168 indicating that the first CNN layer learns the individual patterns in the sequences and  
 169 the second CNN layer learns the high-order interactions between patterns. The high-  
 170 order interactions may be commonalities among different cell lines, and deep neural  
 171 network architectures can extract the high-order EPI features from DNA sequences.



172

### 173 2.3 Influence of domain adversarial operation

174 To validate the effectiveness of domain adversarial operation, we constructed another  
 175 model (namely SEP), which has no the domain adversarial network architecture  
 176 compared with SEPT. Since SEP has no the domain adversarial operation, the data of  
 177 target domain cannot be used in SEP. That is, SEP just used the data of source domain  
 178 in the training phase. Table 1 shows the average AUC results of SEPT and SEP by  
 179 training model on one cell line data and test on another cell line data in running 10 times,  
 180 from which we can see that AUC values of SEPT for each test cell line are higher than  
 181 that of SEP, indicating that domain adversarial operation can effectively improve the  
 182 predictive performance of EPIs in new cell lines. SEPT makes use of the EPIs  
 183 information in other cell lines for recognizing the EPIs in new cell lines.

Table 1. Average AUC values of SEPT and SEP by training model on one cell line data and test on another cell line data in running 10 times. The row represents the labeled training cell lines and column represent the test cell lines.

Model	Test	GM12878	HMEC	HUVEC	HeLa-S3	IMR90	K562	NHEK
	Train							
SEPT	GM12878	*	0.63	0.68	0.66	0.67	0.62	0.62
	HMEC	0.59	*	0.57	0.61	0.67	0.56	0.59
	HUVEC	0.62	0.62	*	0.66	0.63	0.58	0.64
	HeLa-S3	0.61	0.62	0.69	*	0.66	0.63	0.66
	IMR90	0.58	0.65	0.62	0.59	*	0.56	0.62
	K562	0.64	0.61	0.66	0.66	0.62	*	0.64
	NHEK	0.62	0.63	0.66	0.67	0.66	0.59	*
SEP	GM12878	*	0.57	0.62	0.59	0.57	0.55	0.55
	HMEC	0.51	*	0.53	0.54	0.60	0.50	0.55
	HUVEC	0.59	0.55	*	0.58	0.57	0.56	0.55
	HeLa-S3	0.56	0.56	0.60	*	0.58	0.54	0.58
	IMR90	0.53	0.61	0.56	0.55	*	0.51	0.57
	K562	0.56	0.53	0.54	0.53	0.53	*	0.53
	NHEK	0.53	0.57	0.56	0.55	0.56	0.50	*

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

190

#### 191 **2.4 Results comparison of using different number of cell lines as the source domain** 192 **data**

193 To investigate the influence of different cell line number used in the source domain, we  
194 used different number of cell lines as the source domain data to predict the EPIs with  
195 SEPT. For selecting 1, 2, 3, 4, 5, 6 cell lines as the source domain data, each test cell  
196 line has the results of 6, 15, 20, 15, 6 and 1, respectively. The AUC results of using  
197 different number of cell lines as the source domain data to train SEPT are shown in

198 Figure 3. Each boxplot in Figure 3 represents the results of one test cell line across all  
 199 combinations of source domain cell lines. To make the results more reliable, we  
 200 repeated 10 times, that is, the AUC value of each test cell line is the mean of 10 running  
 201 results. From figure 3, we can see that more cell lines as the source domain can achieves  
 202 higher AUC than less cell lines as the source domain. Especially, when 6 cell lines are  
 203 mixed as the source domain data to train the model, SEPT achieves the highest AUC  
 204 values for every test cell line. With the increase of the cell lines number in source  
 205 domain, AUC value is gradually larger. In addition, for different test cell line, the  
 206 combination of appropriate cell lines as the source data can improve the performance  
 207 of SEPT. These results show that using more existing labeled data of EPIs is helpful  
 208 the prediction of EPIs in new cell line.  
 209

Table 2. Average AUC values of SEPT and SEP by training on the six cell lines data and test on one other cell line data in running 10 times.

Test cell line	Cell line(s) of source domain	Model	AUC
GM12878	HMEC, HUVEC, HeLa-S3, IMR90, K562, NHEK	SEP	0.65
		SEPT	0.72
HMEC	GM12878, HUVEC, HeLa-S3, IMR90, K562, NHEK	SEP	0.68
		SEPT	0.76
HUVEC	GM12878, HMEC, HeLa-S3, IMR90, K562, NHEK	SEP	0.69
		SEPT	0.78
HeLa-S3	GM12878, HMEC, HUVEC, IMR90, K562, NHEK	SEP	0.69
		SEPT	0.77
IMR90	GM12878, HMEC, HUVEC, HeLa-S3, K562, NHEK	SEP	0.72
		SEPT	0.78
K562	GM12878, HMEC, HUVEC, HeLa-S3, IMR90, NHEK	SEP	0.65
		SEPT	0.73
NHEK	GM12878, HMEC, HUVEC, HeLa-S3, IMR90, K562	SEP	0.69
		SEPT	0.76

## 210 2.5 Motifs identified by SEPT

211 As to investigate the motifs, we identified sequence features for each model by  
 212 comparing patterns of the convolutional kernels in the first layer with sequence motifs  
 213 from the database HOCOMOCO Human v11. We reconstructed the output of the first  
 214 convolutional layer for each input sample sequence, then extracted the subsequence that

215 best match each kernel to compute the position frequency matrix (PFM) from the  
 216 aligned sub sequences for each kernel. The motif comparison tool of Tomtom [33] was  
 217 used to match these PFMs to known TF motifs. After obtaining the sequence motifs of  
 218 the two models, we defined a metric of relative importance to measure the relative  
 219 importance of motifs. The metric will consider the occurrence times of the same motif  
 220 in both models and the ranks of the occurrence times of the motif in SEPT. If a motif  
 221 appears more in SEPT and less in SEP, then the greater the effect of the motif in SEPT,  
 222 the higher its relative importance score will be. The metric of relative importance is  
 223 defined as:  $\text{Relative importance} = (1 - N_{wo}/N_w)/\text{Rank}_w$ . Here  $N_w$  and  $N_{wo}$  are  
 224 the occurrence times of motif learned from SEPT and SEP, respectively;  $\text{Rank}_w$  is the  
 225 rank (in descending order) of the motif in the motif set of SEPT according to the  
 226 occurrence times. The closer the value is to 1, the more important the motif is in SEPT.  
 227 To avoid contingency, this process was repeated five times.

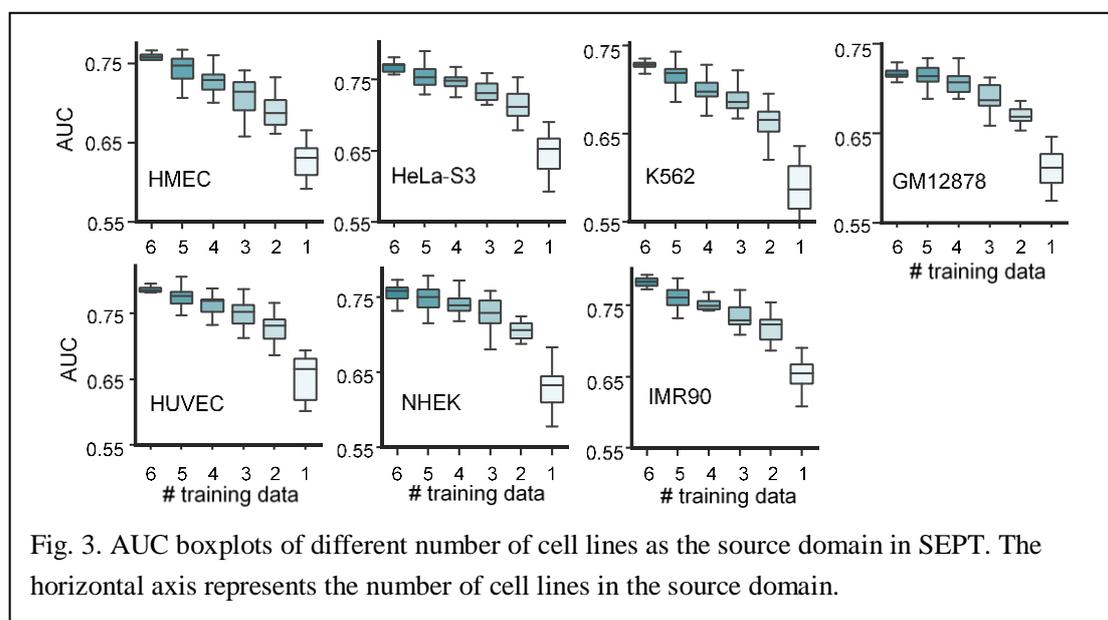


Fig. 3. AUC boxplots of different number of cell lines as the source domain in SEPT. The horizontal axis represents the number of cell lines in the source domain.

228 We found a set of potentially important transcription factor binding motifs by the  
 229 transfer model. For each test cell line, top five import motifs involved in EPIs are show  
 230 in Table S2. For different test cell lines, the models learned common important TF  
 231 motifs such as ZNF563, THAP1, RXRA, and SP3, which are involved in many  
 232 important processes, such as the transcriptional regulation and cell-cycle regulation.  
 233 Interestingly, we found many of potentially important transcription factors that are  
 234 associated with the corresponding cell lines (Table 3). For instance, MAFB motif  
 235 learned by SEPT in the K562 cell line is associated with regulating lineage-specific  
 236 hematopoiesis. This is consistent with the fact that K562 belong to a blood related cell  
 237 line. NR4A1 motif learned by SEPT in GM12878 cell line was reported to play a role  
 238 in the vascular response to injury, while ZNF341 motif learned by both SEPT and SEP  
 239 in GM12878 cell line was reported to involve in the regulation of immune homeostasis.  
 240 It is consistent with the fact that GM12878 is lymphoblastoid related cell lines. We also  
 241 provided other motifs identified by both SEPT and SPEID (Table S3). These results  
 242 show that our SEPT can learn important motifs, and these motifs are relevant to  
 243 enhancer-promoter loops of a novel cell line.

Table 3. Examples of specific motifs identified by the first convolution layer of SEPT.

Cell line	Motif	Element	Brief description	Ref.
K562	MAFB,	Promoter	It plays a pivotal role in regulating lineage-specific hematopoiesis.	[34]
K562	NR3C1,	Enhancer	It affects inflammatory responses, cellular proliferation and differentiation in target tissues, and is involved in chromatin remodeling.	[35]
GM12878	NR4A1,	Enhancer	It plays a role in the vascular response to injury.	[36]
GM12878	ZNF341,	Promoter	It is involved in the regulation of immune homeostasis.	[37]
HeLa-S3	FOXK1,	Promoter	It is involved in different processes such as glucose metabolism, aerobic glycolysis, muscle cell differentiation and autophagy.	[38]
HeLa-S3	BATF3,	Promoter	It controls the differentiation of CD8+ thymic conventional dendritic cells in the immune system.	[39]
HUVEC	ZNF335,	Enhancer	It controls the expression of genes involved in somatic development and regulates, for instance, lymphoblast proliferation.	[40]
IMR90	OLIG2,	Enhancer	It is required for oligodendrocyte and motor neuron specification in the spinal cord.	[41]
NHEK	TFAP2C,	Promoter	It is involved in important biological functions including proper eye, face, body wall, limb and neural tube development.	[42]
HMEC	ZNF317,	Enhancer	It may play an important role in erythroid maturation and lymphoid proliferation.	[43]

### 244 **3 Discussions**

245 One important factor in SEPT is how to choose the labeled data of different cell lines  
246 as the source domain data. We chose different number of cell lines as the source domain  
247 data, and found that choosing all available labeled data of cell lines as the source domain  
248 data for training model can yield better performance than choosing partial cell lines  
249 data as the source domain. In addition, there is redundancy among different cell lines.  
250 Redundancy not only slows down model training, but also damages prediction  
251 performance. Thus, how to choose the different cell lines as the source domain data is  
252 important.

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

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

### 264 **4 Conclusions**

265 Although some deep learning methods have been developed to predict EPIs within the  
266 same cells, they cannot achieve good performance for predicting EPIs in the unlabeled  
267 cell lines, due to lack of understanding of the interested cell lines. In this work, we  
268 proposed a transfer learning model to predict EPIs in new (or interested) cell lines. To  
269 better leverage the existing EPIs knowledge, we adopt the adversarial learning  
270 mechanism to learn useful features in the existing labeled cell lines and interesting  
271 unlabeled cell lines. Experiment results with the domain adversarial operation indicate

272 that it is helpful to predict EPIs in new cell lines. We expect that the model could learn  
273 informative features cross domains and reveal some commonalities (common TFs)  
274 between source and target domains. By learning commonalities between the source and  
275 target domains, SEPT outperforms other state-of-the-art methods for predicting EPIs in  
276 new cell lines.

277 Although SEPT can effectively predict EPIs in specific cell lines from enhancer  
278 and promoter sequences, it can be further improved by considering the following  
279 factors. i) SEPT just uses the sequence information of enhancers and promoters.  
280 Integrating other experimental data such as core histone modification ChIP-seq data or  
281 DNase-seq data can improve the performance of SEPT. ii) Each cell line is treated  
282 equally in the source domain, but the contribution of different cell lines should be  
283 different for the test cell line. Determining which cell lines should be used as training  
284 data is still needed to be explored, as more and more labeled data will become available.  
285 iii) Some EPIs maybe have the cell line specificity, while others are universal across  
286 many cell lines. Thus, different samples within the same cell line should have different  
287 contribution for cross-cell prediction. Assigning a proper weight to each sample can  
288 also improve the performance of SEPT.

## 289 **5 Material and Methods**

### 290 **5.1 Data and preprocessing**

291 We used the same Hi-C data as [2], and downloaded the Hi-C data of seven cell lines  
292 of K562 (mesoderm-lineage cells from a patient with leukemia), GM12878  
293 (lymphoblastoid cells), HeLa-S3 (ectoderm-lineage cells from a patient with cervical  
294 cancer), HUVEC (umbilical vein endothelial cells), IMR90 (fetal lung fibroblasts),  
295 NHEK (epidermal keratinocytes) and HMEC (mammary epithelial cells) from Gene  
296 Expression Omnibus (GEO) GSE63525. The human reference genome hg19 was used  
297 to define the genomic locations. Promoters and activate enhancers in the first four cell  
298 lines were identified using segmentation-based annotations from both ENCODE  
299 Segway [44] and ChromHMM of Roadmap Epigenomics [45], only ChromHMM  
300 annotations were used in the other cell lines. Then, RNA-seq data from ENCODE were  
301 used to select activate promoters according to the rule of their mean FPKM  $>0.3$  with

302 irreproducible discovery rate  $<0.1$  for each cell line. The genome-wide Hi-C  
 303 measurements were used to annotate all enhancer-promoter pairs as interacting or non-  
 304 interacting in each cell type. For each enhancer-promoter pair, the distance between  
 305 promoter and enhancer of the pair is more than 10kb and less than 2Mb [2]. To exclude  
 306 the effect of distance on determining EPIs, interacting enhancer-promoter pairs were  
 307 assigned to one bin (here, the total bin number sets to 5) based on quantile discretization  
 308 of the distance between the enhancer and promoter. Random non-interacting pairs of  
 309 active enhancers and promoters were assigned to their corresponding bin and then  
 310 subsampled as the same number of positive samples within each bin. The subsampled  
 311 non-interacting pairs were considered as negative samples. Table 4 gives the numbers  
 312 of positive and negative pair in each cell line.

313 For each positive/negative sample, sequences of enhancers are extended or cut to  
 314 3kb flanking regions on location center of enhancers, and promoter are extended or cut  
 315 to 2kb flanking regions on location center of promoters. One-hot coding format of  
 316 enhancer and promoter sequence is used as input data of model.

317 We examined the overlapping number of positive EPIs between any two cell lines.  
 318 For any two positive EPI pairs from any two cell lines, if the position of the two  
 319 enhancers and the position of the two promoters both same, the two EPI pairs are  
 320 considered to overlap. By removing redundancy, positive samples of different cell lines  
 321 have very little overlap (Table S4).

Table 4. Number of enhancer-promoter interactions, enhancers and promoters in each cell line.

Cell line	#true EPIs	#all EPIs	#E	#P
<b>IMR90</b>	1254	2504	108996	5253
<b>NHEK</b>	1291	2571	144302	5254
<b>HUVEC</b>	1524	3044	65358	8180
<b>HeLa-S3</b>	1740	3480	102460	7794
<b>K562</b>	1977	3952	82806	8196
<b>GM12878</b>	2113	4226	100036	8453
<b>HMEC</b>	1342	2684	155328	5267

# denotes the number.

322 To compare with RIPPLE based on epigenetic data, we used the data from the  
 323 Roadmap project for the six cell lines. Because we want to make EPI prediction across  
 324 cell lines, we downloaded the peak files of 14 datasets that are measured in all six cell

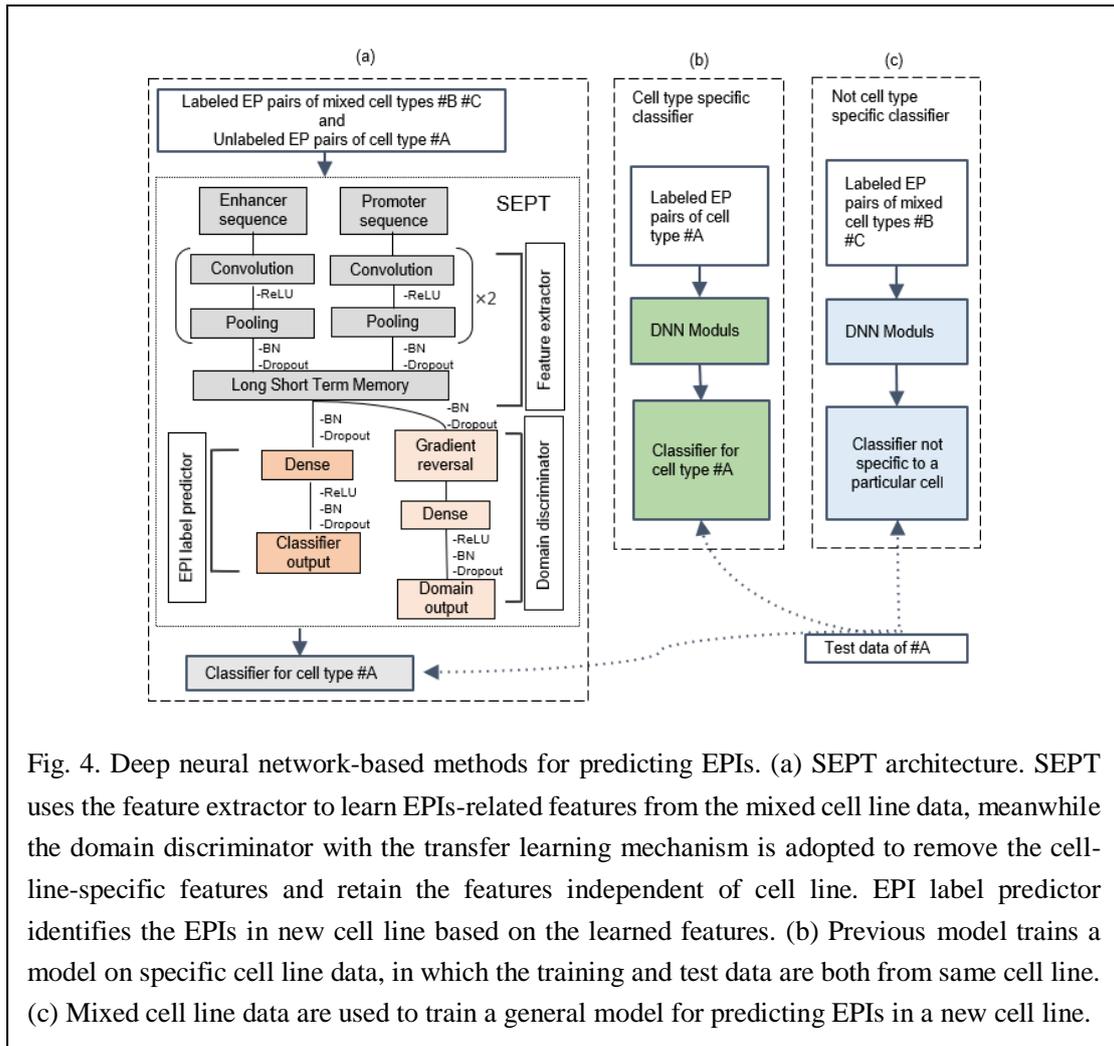
325 lines. These datasets include CTCF, POLR2A, H2AZ, H3K27ac, H3K27me3,  
326 H3K36me3, H3K4me1, H3K4me2, H3K4me3, H3K79me2, H3K9ac, H3K9me3,  
327 H4K20me1 and DNase-seq. An enhancer or promoter sample is represented as a binary  
328 vector in which each dimension corresponds to one of the epigenetic datasets. The  
329 feature vectors of enhancer and promoter are concatenated to represent each EPI pair.  
330 In addition, we also used other two features (i.e., the Pearson’s correlation of the 14  
331 signals associated with the enhancer-promoter pair, and the mRNA level of the gene  
332 associated with the promoter) to represent each EPI pair.

## 333 **5.2 SEPT**

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

345 An overview of SEPT is shown in Figure 4(a). For predicting the EPIs in cell line  
346 #A, unlike the existing two methods which extract the specific features (in Figure 4(b))  
347 or shared features (in Figure 4(c)) of cell lines, SEPT uses the rich information of cell  
348 lines #B and #C to extract the features that are relevant to the EPIs in cell line #A by  
349 using the transfer learning (TL). As shown in Figure 4(a), we used GRL to design the  
350 domain discriminator. GRL reverses the gradient direction during the back propagation,  
351 but it does nothing when forward propagation. Obviously, without GRL, SEPT will  
352 gradually learn the features related to domain in the training process. With the GRL,  
353 the model learns the features that cannot tell whether the data is from the source domain

354 (labeled data of mixed cell lines #B, #C) or the target domain (unlabeled data of cell  
 355 lines #A) in training process. SEPT simultaneously trains two classifiers of the main  
 356 label classifier and the domain discriminator. These two classifiers share feature  
 357 extractor layers. SEPT mainly includes feature extractor, domain discriminator and EPI  
 358 predictor.



359 Feature extractor consists of two convolution layers, two max-pooling layers, two  
 360 dropout layers, and one recurrent long short-term memory (LSTM) layer. Domain  
 361 discriminator consists of the GRL, one dense layer, one dropout layer, and the output.  
 362 EPIs predictor consists of one dense layer, one dropout layer, and the output. Since  
 363 informative features may differ between enhancers and promoters, we use two  
 364 convolution layers, rectified linear unit (ReLU) and max-pooling layers for enhancers

365 and promoters, respectively. Thus, the inputs are two one-hot matrixes to represent  
366 enhancer and promoter sequence, respectively. Because large number of kernels can  
367 sufficiently extract the features, and motif features of DNA sequences are short than  
368 40bp, so each convolution layer consists of 300 kernels with length 40. Max-pooling  
369 layer reduces the output dimension with pool length 20, stride 20. The outputs of the  
370 two branches are concatenated into one tensor, which is input to the dropout layer with  
371 dropout rates of 0.25. The dropout layer randomly selects partial input data to next layer  
372 to avoid overfitting. The recurrent LSTM layer is used to extract feature combinations  
373 of the two branches, and the output dimension of LSTM is 100. For domain  
374 discriminator, the output of LSTM layer feeds into GRL, and a dense layer maps the  
375 learned distributed features to the domain label space. It contains 50 units with ReLU  
376 activations. The output feeds into a sigmoid unit to predict the domain probability after  
377 dropout layer with dropout rates of 0.5. For EPI predictor, the output of LSTM layer  
378 feeds into dense layer, which further maps the learned distributed features to the sample  
379 label space. It contains 100 units with ReLU activations. The output feeds into a  
380 sigmoid unit to predict the probability after dropout layer with dropout rates of 0.5.

### 381 **SEPT Model training**

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

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

392 Here,  $L_y$  is the loss of EPIs label predictor,  $L_d$  is the loss of domain discriminator,  $G_f$   
393 is a mapping that maps the input  $x$  to a feature vector,  $G_y$  is a mapping that maps the  
394 feature vector to the EPIs label,  $G_d$  is a mapping that maps the feature vector to the  
395 domain label,  $x_i$  is the  $i$ -th sample,  $\theta_f$  is the parameters of mapping  $G_f$ ,  $\theta_y$  is the  
396 parameters of mapping  $G_y$ ,  $\theta_d$  is the parameters of mapping  $G_d$ ,  $y_i$  is the EPI label of  
397 the  $i$ -th sample,  $d_i$  is the domain label of the  $i$ -th sample,  $N$  is the number of labeled  
398 EPI training samples,  $M$  is the number of unlabeled EPI training samples but they have  
399 the domain labels, and  $\lambda$  is a constant that controls the tradeoff between two objectives.

400 It is a problem of minimax optimization, that is, we attempts to seek a saddle point  
401 of the functional  $E(\theta_f, \theta_y, \theta_d)$  that is delivered by parameters  $\widehat{\theta}_f$ ,  $\widehat{\theta}_y$ , and  $\widehat{\theta}_d$ . At  
402 the saddle point, the loss of EPI label predictor and domain discriminator is minimized.  
403 The loss of EPI label predictor is minimized by the feature mapping parameter  $\theta_f$ , while  
404 the loss of domain discriminator is maximized by  $\theta_f$  on account of the GRL that  
405 changes the sign of the gradient during the back-propagation. In the end, SEPT learns  
406 the features that are discriminative and domain-invariant. The features learned from the  
407 cell lines with label information (source domain) are effective for the new cell lines  
408 (target domain).

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

415 domain data. iv) Evaluating the performance of SEPT with the test data in target domain.

416 All the experiments are based on Python by using the scikit-learn machine learning  
417 library [46] and Keras framework (<https://keras.io/>) with Tensorflow as back-end [47].

### 418 **5.3 Evaluation metrics**

419 We adopted the metrics of Accuracy, Precision, Recall, F1 score, AUC and AUPR to  
420 assess the performance of SEPT. These metrics are defined respectively as the  
421 following [48-51].

$$422 \quad \text{Accuracy} = \frac{TP+TN}{TP+TN+FN+FP} \quad (2)$$

$$423 \quad \text{Precision} = \frac{TP}{TP+FP} \quad (3)$$

$$424 \quad \text{Recall} = \frac{TP}{TP+FN} \quad (4)$$

$$425 \quad \text{F1} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (5)$$

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

433

### 434 **Abbreviations**

435 EPIs: Enhancer-promoter interactions; AUC: Area under the receiver operating curves;  
436 3D: Three-dimensional; 3C: Chromosome conformation capture-based approach; CNN:  
437 Convolutional neural network; GWAS: Genome-wide association studies; TL: Transfer  
438 learning; GRL: Gradient reversal layer; LSTM: Long short-term memory; ReLU:

439 Rectified linear unit; SGD: Stochastic gradient descent; AUPR: Area under the  
440 precision-recall curves.

#### 441 **Ethics approval and consent to participate**

442 Not applicable.

#### 443 **Consent for publication**

444 Not applicable.

#### 445 **Availability of data and materials**

446 The related source codes and datasets are available at [https://github.com/NWPU-](https://github.com/NWPU-903PR/SEPT)  
447 903PR/SEPT.

#### 448 **Competing interests**

449 The authors declare that they have no competing interests.

#### 450 **Funding**

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

#### 459 **Authors' contributions**

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

#### 463 **Acknowledgments**

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

#### 465 **Authors' information**

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

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

470 Shihua Zhang: Email: zsh@amss.ac.cn; Address: 55 Zhongguancun East Road, Beijing,  
471 10090, China.

## 472 **References**

- 473 1. Fullwood MJ, Ruan Y. ChIP-based methods for the identification of long-range chromatin  
474 interactions. *J Cell Biochem.* 2009;107(1):30-9.
- 475 2. Whalen S, Truty RM, Pollard KS. Enhancer–promoter interactions are encoded by complex  
476 genomic signatures on looping chromatin. *Nat Genet.* 2016; 48(5):488-96.
- 477 3. Shlyueva D, Stampfel G, Stark A. Transcriptional enhancers: from properties to genome-  
478 wide predictions. *Nat Rev Genet.* 2014; 15(4):272-86.
- 479 4. Visel A, Rubin EM, Pennacchio LA. Genomic views of distant-acting enhancers. *Nature.*  
480 2009; 461(7261):199-205.
- 481 5. Van Steensel B, Dekker J. Genomics tools for unraveling chromosome architecture. *Nat*  
482 *Biotechnol.* 2010; 28(10):1089-95.
- 483 6. Bickmore WA, van Steensel B. Genome architecture: domain organization of interphase  
484 chromosomes. *Cell.* 2013; 152(6):1270-84.
- 485 7. Dekker J, Mirny L. The 3D genome as moderator of chromosomal communication. *Cell.*  
486 2016; 164(6):1110-21.
- 487 8. Rowley MJ, Corces VG. The three-dimensional genome: principles and roles of long-  
488 distance interactions. *Curr Opin Cell Biol.* 2016; 40:8-14.
- 489 9. Achinger-Kawecka J, Clark SJ. Disruption of the 3D cancer genome blueprint.  
490 *Epigenomics.* 2017; 9(1):47-55.
- 491 10. Franke M, Ibrahim DM, Andrey G, Schwarzer W, Heinrich V, Schöpflin R, Kraft K,  
492 Kempfer R, Jerković I, Chan W-L. Formation of new chromatin domains determines  
493 pathogenicity of genomic duplications. *Nature.* 2016; 538(7624):265-69.
- 494 11. Claussnitzer M, Dankel SN, Kim K-H, Quon G, Meuleman W, Haugen C, Glunk V, Sousa

- 495 IS, Beaudry JL, Puvion-Rodan V. FTO obesity variant circuitry and adipocyte browning in  
496 humans. *New Engl J Med.* 2015; 373(10):895-907.
- 497 12. Lupiáñez DG, Kraft K, Heinrich V, Krawitz P, Brancati F, Klopocki E, Horn D, Kayserili  
498 H, Opitz JM, Laxova R. Disruptions of topological chromatin domains cause pathogenic  
499 rewiring of gene-enhancer interactions. *Cell.* 2015; 161(5):1012-25.
- 500 13. Dekker J, Rippe K, Dekker M, Kleckner N. Capturing chromosome conformation. *science.*  
501 2002; 295(5558):1306-11.
- 502 14. Lieberman-Aiden E, Van Berkum NL, Williams L, Imakaev M, Ragoczy T, Telling A, Amit  
503 I, Lajoie BR, Sabo PJ, Dorschner MO. Comprehensive mapping of long-range interactions  
504 reveals folding principles of the human genome. *science.* 2009; 326(5950):289-93.
- 505 15. Dixon JR, Selvaraj S, Yue F, Kim A, Li Y, Shen Y, Hu M, Liu JS, Ren B. Topological  
506 domains in mammalian genomes identified by analysis of chromatin interactions. *Nature.*  
507 2012; 485(7398):376-80.
- 508 16. De Laat W, Duboule D. Topology of mammalian developmental enhancers and their  
509 regulatory landscapes. *Nature.* 2013; 502(7472):499-506.
- 510 17. Rao SS, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, Robinson JT, Sanborn AL,  
511 Machol I, Omer AD, Lander ES. A 3D map of the human genome at kilobase resolution  
512 reveals principles of chromatin looping. *Cell.* 2014; 159(7):1665-80.
- 513 18. Andersson R, Gebhard C, Miguel-Escalada I, Hoof I, Bornholdt J, Boyd M, Chen Y, Zhao  
514 X, Schmidl C, Suzuki T. An atlas of active enhancers across human cell lines and tissues.  
515 *Nature.* 2014; 507(7493):455-61.
- 516 19. Yang Y, Zhang R, Singh S, Ma J. Exploiting sequence-based features for predicting  
517 enhancer–promoter interactions. *Bioinformatics.* 2017; 33(14):i252-i260.
- 518 20. Zeng W, Wu M, Jiang R. Prediction of enhancer-promoter interactions via natural language  
519 processing. *BMC Genomics.* 2018; 19(2):84. doi:10.1186/s12864-018-4459-6
- 520 21. Singh S, Yang Y, Póczos B, Ma J. Predicting enhancer-promoter interaction from genomic  
521 sequence with deep neural networks. *Quant Biol.* 2019; 7(2):122-37.
- 522 22. Zhuang Z, Shen X, Pan W. A simple convolutional neural network for prediction of  
523 enhancer–promoter interactions with DNA sequence data. *Bioinformatics.* 2019;  
524 17(35):2899-906.

- 525 23. He B, Chen C, Teng L, Tan K. Global view of enhancer–promoter interactome in human  
526 cells. *Proceedings of the National Academy of Sciences*. 2014; 111(21):E2191-E2199.
- 527 24. Roy S, Siahpirani AF, Chasman D, Knaack S, Ay F, Stewart R, Wilson M, Sridharan R. A  
528 predictive modeling approach for cell line-specific long-range regulatory interactions.  
529 *Nucleic Acids Res*. 2015; 43(18):8694-712.
- 530 25. Zhu Y, Chen Z, Zhang K, Wang M, Medovoy D, Whitaker JW, Ding B, Li N, Zheng L,  
531 Wang W. Constructing 3D interaction maps from 1D epigenomes. *Nat Commun*. 2016;  
532 7:10812. doi:10.1038/ncomms10812
- 533 26. Cao Q, Anyansi C, Hu X, Xu L, Xiong L, Tang W, Mok MT, Cheng C, Fan X, Gerstein M.  
534 Reconstruction of enhancer–target networks in 935 samples of human primary cells, tissues  
535 and cell lines. *Nat Genet*. 2017; 49(10):1428-36.
- 536 27. Sun B, Saenko K: Deep coral. Correlation alignment for deep domain adaptation. In:  
537 *European Conference on Computer Vision*. 2016; Springer: 443-50.
- 538 28. Abdelwahab M, Busso C. Domain adversarial for acoustic emotion recognition.  
539 *IEEE/ACM Transactions on Audio, Speech, and Language Processing*. 2018; 26(12):2423-  
540 35.
- 541 29. Ganin Y, Ustinova E, Ajakan H, Germain P, Larochelle H, Laviolette F, Marchand M,  
542 Lempitsky V. Domain-adversarial training of neural networks. *The Journal of Machine*  
543 *Learning Research*. 2016; 17(1):2096-30.
- 544 30. Guo H, Ahmed M, Zhang F, Yao CQ, Li S, Liang Y, Hua J, Soares F, Sun Y, Langstein J.  
545 Modulation of long noncoding RNAs by risk SNPs underlying genetic predispositions to  
546 prostate cancer. *Nat Genet*. 2016; 48(10):1142-50.
- 547 31. Pan SJ, Yang Q. A survey on transfer learning. *IEEE Transactions on knowledge and data*  
548 *engineering*. 2009; 22(10):1345-59.
- 549 32. Dongwon L. LS-GKM: a new gkm-SVM for large-scale datasets. *Bioinformatics*. 2016;  
550 32(14): 2196-98.
- 551 33. Gupta, S., Stamatoyannopoulos, J.A., Bailey, T.L. et al. Quantifying similarity between  
552 motifs. *Genome Biol*. 2007; 8, R24. doi: <https://doi.org/10.1186/gb-2007-8-2-r24>
- 553 34. Eychène, Alain, Rocques N, Pouponnot C. A new MAFia in cancer. *Nature Reviews Cancer*.  
554 2008; 8(9):683-93.

- 555 35. Fryer C, Archer T. Chromatin remodelling by the glucocorticoid receptor requires the  
556 BRG1 complex. *Nature*. 1998; 393(6680):88-91.
- 557 36. Papac-Milicevic N, Breuss JM, Zaujec J, et al. The interferon stimulated gene 12 inactivates  
558 vasculoprotective functions of NR4A nuclear receptors. *Circulation Research*. 2012;  
559 110(8):e50-e63.
- 560 37. Stefanie FJ, Hartberger JM, Manfred F, et al. ZNF341 controls STAT3 expression and  
561 thereby immunocompetence. *Immunology*. 2018; 3(24):aat4941.  
562 doi:10.1126/sciimmunol.aat4941
- 563 38. Bowman CJ, Ayer DE, Dynlacht BD. Foxk proteins repress the initiation of starvation-  
564 induced atrophy and autophagy programs. *Nature Cell Biology*. 2014; 16(12):1202-14.
- 565 39. Bower KE, Fritz JM, Mcguire KL. Transcriptional repression of MMP-1 by p21SNFT and  
566 reduced in vitro invasiveness of hepatocarcinoma cells. *Oncogene*. 2004; 23(54):8805-14.
- 567 40. Yang YJ, Baltus AE, Mathew RS, et al. Microcephaly gene links trithorax and REST/NRSF  
568 to control neural stem cell proliferation and differentiation. *Cell*. 2012; 151(5): 1097-112.
- 569 41. Beyer CA, Cabanela ME, Berquist TH. Unilateral facet dislocations and fracture-  
570 dislocations of the cervical spine. *The Bone & Joint Journal*. 1992; 73(6):977-81.
- 571 42. Bamforth SD, Bragança J, Eloranta JJ, et al. Cardiac malformations, adrenal agenesis,  
572 neural crest defects and exencephaly in mice lacking Cited2, a new Tfp2 co-activator. *Nat*  
573 *Genet*. 2001; 29(4):469-74.
- 574 43. Takashima H, Nishio H, Wakao H, et al. Molecular cloning and characterization of a  
575 KRAB-containing zinc finger protein, ZNF317, and its isoforms. *Biochemical and*  
576 *Biophysical Research Communications*. 2001; 288(4): 771-9.
- 577 44. Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, Aken BL,  
578 Barrell D, Zadissa A, Searle S. GENCODE: the reference human genome annotation for  
579 The ENCODE Project. *Genome Res*. 2012; 22(9):1760-74.
- 580 45. Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, Heravi-Moussavi A, Kheradpour P,  
581 Zhang Z, Wang J, Ziller MJ. Integrative analysis of 111 reference human epigenomes.  
582 *Nature*. 2015; 518(7539):317-30.
- 583 46. Swami A, Jain R. Scikit-learn: Machine learning in Python. *Journal of machine learning*  
584 *research*. 2011; 12(10):2825-30.

- 585 47. Abadi M, Barham P, Chen J, Chen Z, Davis A, Dean J, Devin M, Ghemawat S, Irving G,  
586 Isard M, et al. Tensorflow: A system for large-scale machine learning. arXiv preprint  
587 arXiv:1603.04467.
- 588 48. Jing F, Zhang S, Cao Z, Zhang S. An integrative framework for combining sequence and  
589 epigenomic data to predict transcription factor binding sites using deep learning.  
590 IEEE/ACM transactions on computational biology and bioinformatics. 2019; 1-1. doi:  
591 10.1109/TCBB.2019.2901789
- 592 49. Fan XN, Zhang SW. LPI-BLS: Predicting lncRNA-protein interactions with a broad  
593 learning system-based stacked ensemble classifier. Neurocomputing. 2019; 370:88-93.
- 594 50. Zhang SW, Zhang XX, Fan XN, Li WN. LPI-CNNCP: Prediction of lncRNA-protein  
595 interactions by using convolutional neural network with the copy-padding trick. Analytical  
596 Biochemistry. 2020; 601: 113767.
- 597 51. Fan XX, Zhang SW, Zhang SY and Ni JJ. lncRNA\_Mdeep: An Alignment-Free Predictor  
598 for Distinguishing Long Non-Coding RNAs from Protein-Coding Transcripts by  
599 Multimodal Deep Learning. Int. J. Mol. Sci. 2020; 21: 5222.