# EPIMR: Prediction of Enhancer-Promoter Interactions by Multi-Scale ResNet on Image Representation

Qiaozhen Meng, Yinuo Lyu, Xiaoqing Peng, Junhai Xu*, Jijun Tang*, and Fei Guo*

**Abstract:** Prediction of enhancer-promoter interactions (EPIs) is key to regulating gene expression and diagnosing genetic diseases. Due to limited resolution, biological experiments perform not as well as expected while precisely identifying specific interactions, giving rise to computational biology approaches. Many EPI predictors have been developed, but their prediction accuracy still needs to be enhanced. Here, we design a new model named EPIMR to identify enhancer-promoter interactions. First, Hilbert Curve is utilized to represent sequences to images to preserve the position and spatial information. Second, a multi-scale residual neural network (ResNet) is used to learn the distinguishing features of different abstraction levels. Finally, matching heuristics are adopted to concatenate the learned features of enhancers and promoters, which pays attention to their potential interaction information. Experimental results on six cell lines indicate that EPIMR performs better than existing methods, with higher area under the precision-recall curve (AUPR) and area under the receiver operating characteristic (AUROC) results on benchmark and under-sampling datasets. Furthermore, our model is pre-trained on all cell lines, which improves not only the transferability of cross-cell line prediction, but also cell line-specific prediction ability. In conclusion, our method serves as a valuable technical tool for predicting enhancer-promoter interactions, contributing to the understanding of gene transcription mechanisms. Our code and results are available at https://github.com/guofei-tju/EPIMR.

**Key words:** enhancer-promoter interactions; Hilbert Curve; multi-scale residual neural network (ResNet)

## 1 Introduction

Promoters and enhancers are short regions of DNA that regulate gene expression in a spatiotemporal manner[1−3]. In most cases, these two elements need to make physical contact with each other to transmit transcriptional regulatory information[4−6]. One of the classical models shows that the CCCTC-binding factor (CTCF) and the cohesin complex help to stabilize the enhancer-promoter interactions (EPIs)[7]. A transcription factor (TF) binds to an enhancer to

● Qiaozhen Meng and Junhai Xu are with the School of Computer Science and Technology, College of Intelligence and Computing, Tianjin University, Tianjin 300350, China. E-mail: qiaozhenm@tju.edu.cn; jhxu@tju.edu.cn.

● Yinuo Lyu is with the Aeronautical Information Service Center of the Civil Aviation Administration of China (AISC.ATMB.CAAC), Beijing 100015, China.

● Xiaoqing Peng is with the Center for Medical Genetics & Hunan Key Laboratory of Medical Genetics, School of Life Sciences, Central South University, Changsha 410038, China.

● Jijun Tang is with the Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen 518000, China. E-mail: jj.tang@siat.ac.cn.

● Fei Guo is with the School of Computer Science and Engineering, Central South University, Changsha 410083, China. E-mail: guofei@csu.edu.cn.

∗ To whom correspondence should be addressed.

facilitate the transient recruitment of the mediator complex, which delivers RNA polymerase II (Pol II) to the target promoter to initiate transcription[8]. To date, many studies indicate that genetic mutations and disturbance of genome organization can disrupt EPIs and thus result in diseases[9−12]. Apparently, identifying enhancer-promoter interactions is of momentous practical significance for understanding mechanisms of gene regulatory networks[13−16] and helping with genetic disease diagnosis.

Recent studies have demonstrated that chromosomes form topologically associating domains (TADs) where chromatin interactions are conducted on the order of tens to hundreds of kilobases[4, 17, 18]. However, due to the presence of enhancer-promoter contacts that cross TAD boundaries, TADs do not play as significant a role in gene control as previously thought[19]. Enhancers can stimulate distal target promoters even within the same TAD by "skipping" nearby genes[20]. Moreover, one enhancer can activate more than one promoter, whereas one promoter can be regulated by several enhancers[21−23]. Overall, the multiplicity of enhancer-promoter relationships makes it difficult to identify their interactions.

Over the last few decades, chromosome conformation capture and its derivatives (3C, 4C, 5C, Hi-C, etc.)[24−28] have been developed to study genomic organization of chromatin interactions. Therefore, Hi-C, the most widely used biotechnologies, allowed kilobase-scale resolution analysis but required a significant sequencing effort with sequencing depth at a billion-read scale[17, 29]. Another limitation is the requirement for extensive input materials[30]. Different methods have been developed, for example, capture Hi-C[31] offered a high resolution to interrogate regulatory interactions but introduced additional technical biases[32]. Other alternative approaches, such as ChIA-PET[33] and HiChIP[34], reported more enhancer-promoter interactions but might require more biomaterials[29]. Hence, using these 3C-based techniques are laborious and prohibitively expensive to examine specific contacts, motivating the development of computational ways.

To circumvent the limitations of experimental approaches, machine learning algorithms have made some progress in identifying enhancer-promoter interactions. There are two types of approaches for precisely detecting EPIs on a genome-wide scale: one based on functional genomic data, and the other based on sequence information. The first approach uses data from genomic signals detected by sequencing technologies, including DNA methylation, gene expression, histone modification, transcription factors, chromatin accessibility, and so on. Early EPIs prediction models, such as RIPPLE[35] and TargetFinder[36], were based on the enrichment degree of these characteristic signals to determine whether there is an interaction. However, such data mainly come from biological experiments, which are limited and difficult to obtain. Alternatively, sequences of promoters and enhancers are directly used to extract features. One of the most intuitive strategies is to use one-hot to convert DNA sequences into binary matrices that a model can understand. For example, SPEID[37] utilizes convolutional neural network (CNN) and long short-term memory (LSTM) to identify EPIs, believing that LSTM can commendably learn the long-range dependence of sequences. In SimCNN[38] model, a simple CNN architecture combined with transfer learning is used for prediction. And EPIANN[39] incorporated an attention-based mechanism to network model to detect EPIs.

In order to extract fixed-length sequence embedding features and retain the context information, EP2vec[40] was inspired by doc2vec in natural language processing (NLP), taking DNA sequences as sentences and k-mers as words. And gradient boosted regression trees classifier (GBRT) was used to construct the prediction model. Analogously, PEP[41] used word2vec to build a PEP-Word module, together with a PEP-Motif module figuring out the occurrence frequencies of transcription factor binding site (TFBS) motifs, to elucidate sequence-based instructions.

On this basis, EPIVAN[42], EPI-DLMH[43], EPIHC[44], and EPnet[45] models used dna2vec[46] algorithm to pre-train DNA vectors on the entire human genome, so that they contained richer and more accurate sequence feature information that was more suitable for large-scale prediction models. As for classification algorithms, they adopted neural networks such as CNN and bidirectional gated recurrent unit (Bi-GRU), combined with strategies like attention mechanism and matching heuristic algorithm for further optimization. Furthermore, in the cause of preserving spatial position information between enhancers and promoters, EPIsHilbert[47] extracted features in a new way. According to Hilbert Curve[48], one-dimensional DNA sequences were mapped into

three-dimensional matrix-vector, illustrating enhancer-promoter interaction at a distance.

Taken together, these existing methods exploited different ways to ccapture more information on EPIs, and have made considerable progress on EPIs prediction. But there is still some room for improvement. First, in most of the above models, one-hot or k-mers embedding is used to transform enhancers and promoters into network inputs, which have certain limitations since they will lose some spatial information of sequences. Second, the existing model architectures are very simple, which could be detrimental to learn distinguishing features during feature extraction. Finally, most EPIs prediction methods directly concatenate the learned features of enhancers and promoters, thus neglecting potential interaction information between them.

In this paper, we propose a new deep learning model called EPIMR to identify enhancer-promoter interactions by learning information of different abstraction levels from the sequences characterized into image forms, as presented in Fig. 1. We demonstrate by experimental results that our model outperforms existing models on benchmark datasets and under-sample datasets. To be specific, we introduce a multi-scale residual neural network (ResNet) to represent different abstraction levels of enhancer or promoter features. And we utilize matching heuristic layers like concatenation, element-wise product, difference, dot, and addition to acquire interaction information. It is shown that these methods learn respective and communicative features, and improve predictive performances. Also, we adopt a pre-trained model among them, thereby improving cross-cell line and cell line-specific prediction. Finally, we verify the validity of EPIMR on the reconstructed datasets with low data dependence, indicating the practical applications of pre-trained EPIMR for predicting EPIs on cell lines that lack data.

## 2 Material and Method

### 2.1 Benchmark dataset

Following existing works, we also adopt a benchmark dataset from TargetFinder[36] in this study. Several human cell lines are used in this study, including IMR90, HUVEC, HeLa-S3, K562, GM12878, and NHEK. A total of annotated enhancers and promoters are derived from ENCODE Project[49] and Roadmap Epigenomics Project[50]. To achieve sequences with a fixed length, the enhancers and promoters are extended by adjustable flanking regions or randomly cut to 3000 bp and 2000 bp long, respectively, where bp is base pair. It is important to note that an enhancer and promoter pair that interact with each other are labeled as a positive sample. Otherwise, they are categorized as negative samples.

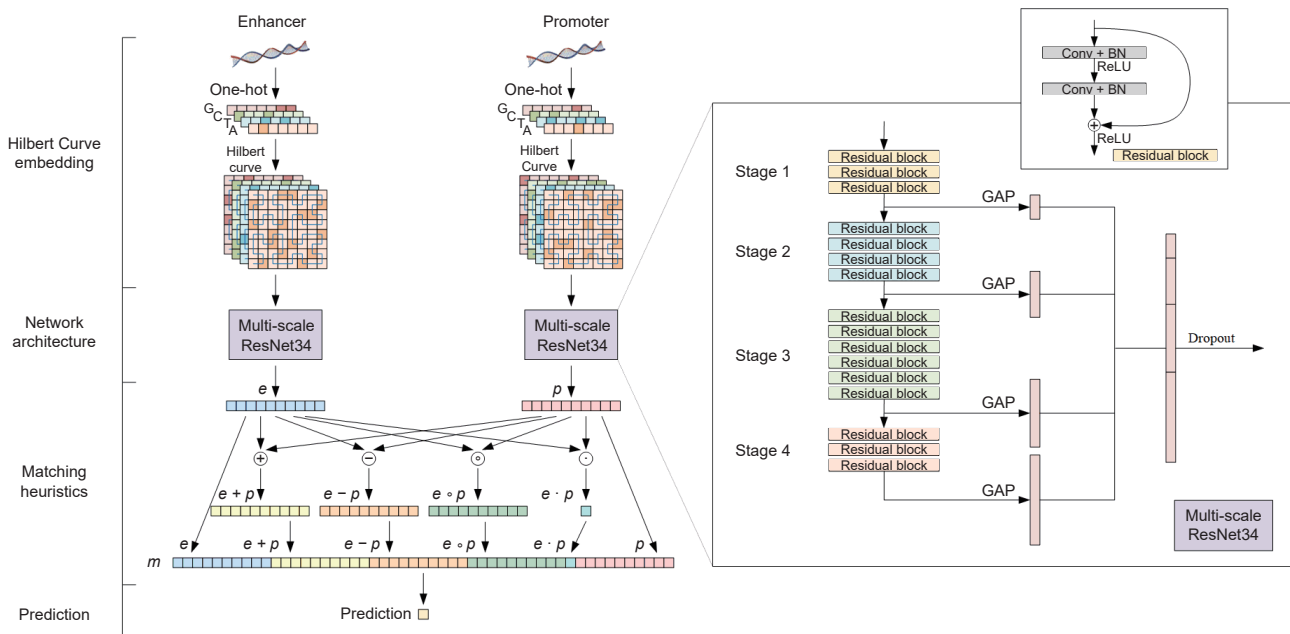As shown in Table 1, we analyze each cell line



Fig. 1 Framework of our proposed method with the detailed multi-scale representation of ResNet34.

**Table 1    EPIs benchmark datasets.**

| Cell line | Number | | |
|---|---|---|---|
| | All | Positive | Negative |
| IMR90 | 26 254 | 1254 | 25 000 |
| HUVEC | 31 924 | 1524 | 30 400 |
| HeLa-S3 | 36 540 | 1740 | 34 800 |
| K562 | 41 477 | 1977 | 39 500 |
| GM12878 | 44 313 | 2113 | 42 200 |
| NHEK | 26 891 | 1291 | 25 600 |

dataset in detail, which reveals that the ratio for EPIs and non-EPIs is 1:20, similar to the actual distribution of genes in the genome.

## 2.2    Sequence encoding

### 2.2.1    Hilbert Curve

To capture long-range interactions efficiently, we use the Hilbert Curve[48] to represent sequences. As a continuous fractal space-filling curve, each of elements is mapped to a pixel to map a 1D line to a 2D image with the Hilbert Curve.

The Hilbert Curve is constructed recursively. In the first iteration, its image is a unit square composed of 4 regions, arranged two by two. Then, each part is divided into 4 smaller regions on each iteration. Eventually, the Hilbert Curve yields an image of size $2^n \times 2^n$, where $n$ represents iteration numbers. Curve thereupon calls the rotation function to fill the entire square, as shown in Fig. 2. As discussed in Section 2.1,

the longest sequence is 3000 bp. To accommodate this length, we set the iteration number to 6 to ensure that Hilbert Curve image is enough to represent sequences with a size $2^6 \times 2^6 (= 64 \times 64)$.

### 2.2.2    Sequence representation

Here, we take inspiration from Hilbert Curve. One sequence is represented as a curve, and nucleobases are the elements filled in the sub-squares.

First, we place biological sequences in digital forms. Since DNA sequences are composed of 4 nucleobases, we use one-hot to encode enhancers and promoters, storing the nucleotides as A(1, 0, 0, 0), C(0, 0, 1, 0), G(0, 0, 0, 1), and T(0, 1, 0, 0). Thus, enhancers are mapped as (3000, 4) vectors, whereas promoters are encoded as (2000, 4) vectors. Then, we map per nucleobase to a pixel with Hilbert Curve to fill image that has a size $2^6 \times 2^6 (= 64 \times 64)$. The 4 features encoded by one-hot can be regarded as channels in an
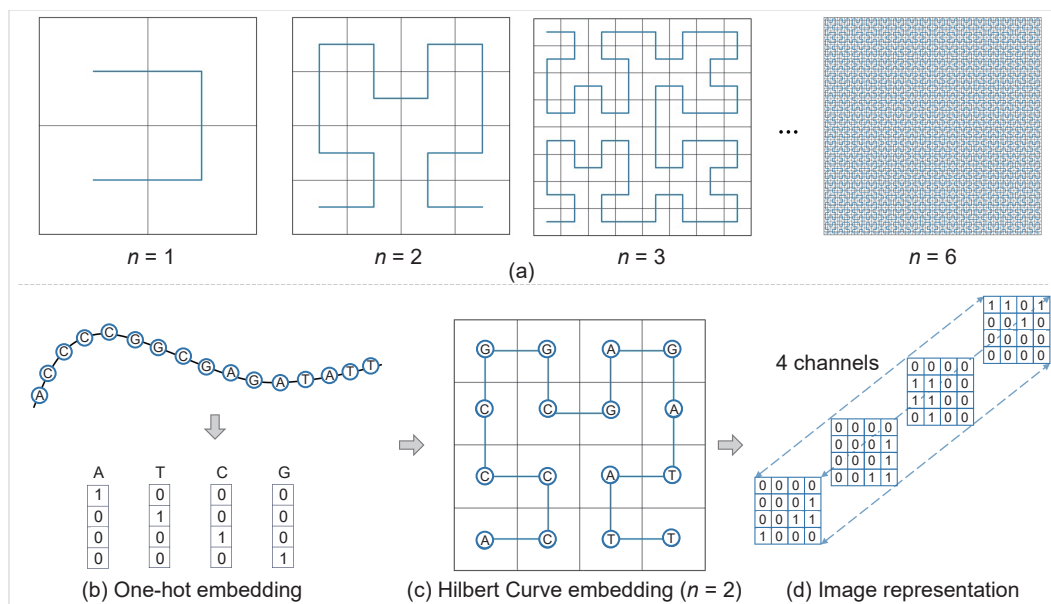


**Fig. 2    Sequence embedding of Hilbert Curve. Figure 2a descripts the different images filling with the Hilbert Curve with different iteration numbers ($n$ = 1, 2, …, 6). Figures 2b–2d show detailed steps to extract Hilbert Curve embeddings when the iteration number is set to 2 for convenience.**

image as shown in Fig. 2.

Actually, due to the chromatin folding in the genome, the proximal nuclear bases on a DNA sequence can remain adjacent to each other, while remote elements may also have very small spatial distances. The Hilbert Curve algorithm can effectively preserve the local features of the sequence, while its spatial folding can integrate the long-range characteristics, thereby obtaining more useful information. Now that the sequences are mapped to images, we can construct networks that are good at dealing with graphics.

## 2.3 Network architecture

### 2.3.1 ResNet

Deep neural networks have brought a series of breakthroughs for image classification[51−55]. So we consider that a deeper network could be used to obtain effective classification ability. However, adding too many layers to deep models may cause the problem of vanishing/exploding gradients, and thus lead to a higher training error and lower accuracy. As a way to mitigate this issue, He et al.[56] proposed ResNet, motivated by the concept of highway networks.

Shortcut connections enable ResNet to skip some layers, preserving the information from previous layers, and transferring it directly to subsequent layers.

ResNet models normally use double- or triple-layer skips containing nonlinearities (rectified linear unit (ReLU)), followed by batch normalization, that are referred to as residual blocks. A certain number of similar residual blocks are stacked in series to form four stages, where one projection shortcut is used to match dimensions between two stages, and other identity shortcuts are used to increase network depths, and finally form the ResNet. Specifically, we choose ResNet34 as the backbone of our network.

### 2.3.2 Multi-scale representation

In our model, we propose a multi-scale method to make improvements based on ResNet34 as shown in Fig. 1.

As the depth of the network increases, the features extracted from different stages of ResNet34 contain different levels of information and are highly complementary. The shallower parts of the network have smaller receptive fields and tend to focus on local details. Conversely, the deeper levels of the network have larger receptive fields and are prone to produce highly-abstracted features with more spatial

information[57, 58]. As shown in Fig. 1, we integrated four feature groups which are obtained from different deep layers to final representations[59]. These four groups, GAP1, GAP2, GAP3, and GAP4, learned information at different levels.

Global average pooling can convert feature maps to a vector by calculating the average value of these features. Accordingly, we extract the output of four scales from ResNet and concatenate them together after global average pooling layers. Then, we add a dropout layer to further decrease the phenomenon of over-fitting. In this way, we integrate the features at different levels to capture complementary information at multiple scales.

## 2.4 Matching heuristics

To capture explicit information about enhancer-promoter interactions, we use matching heuristics developed in natural language inference[60−65]. Owing to the enormous data volume, we are more inclined to adopt the approaches that are effective in capturing relationships between E-P pairs, but remain low complexity.

Now that both of the input sequences are represented in vectorial form, five heuristic matching layers are applied to the learned features of enhancers and promoters:
- Concatenation.
- Dot.
- Element-wise difference.
- Element-wise product.
- Element-wise addition.

Here, the outputs from five matching methods are concatenated together, which can be performed by the following formula:

$$m = [e, p, e - p, e \circ p, e + p, e \cdot p],$$

where "$\circ$" denotes element-wise product, and "$\cdot$" denotes dot operation. By this means, concatenation can preserve all the information of two sequences, difference and addition calculate the degree of distributional inclusion in each dimension, whereas product and dot are certain measure of the interactive communications of the two features.

At last, we feed the concatenated output $m$ for final classification.

## 2.5 Model training

The model was developed in Keras using the binary

cross-entropy loss function, a learning rate of $10^{-5}$, and an Adam optimization algorithm. It takes 50 epochs with the early-stopping method. Also, as mentioned before, we adopt a dropout layer with a probability of $p = 0.75$ to further decrease the phenomenon of over-fitting. For each cell line, using random selection, we divided the dataset into training, test, and validation sets in a ratio of 8:1:1.

As shown in Table 1, the benchmark dataset we use in this study has a drastically imbalanced pattern of negative and positive samples. The negative samples are approximately 20 times as plentiful as the positive samples, which is more in line with real EPI distribution. However, the traditional classification algorithms which emphasize total accuracy tends to pay more attention to the majority class makes learning robust models difficult. While training the classifier, positive samples, which belong to the minority class, are given greater weight in proportion to the ratio of negative to positive samples.

## 2.6　Performance evaluation

As stated, the extremely imbalanced datasets render accuracy not appropriate for performance evaluation. In light of this, we adopt the area under the precision-recall curve (AUPR)[66] and area under the receiver operating characteristic (AUROC)[67] as our evaluation metrics.

## 3　Result and Discussion

### 3.1　Ablation experiment

As shown in Tables 2 and 3, we perform ablation experiments to determine whether our model choices are effective. First, we compare our predictor with the one without a multi-scale representation of ResNet. For different cell lines, the increments of AUPR range from 4.07% to 11.28%. In four of six cell lines, AUROC increased results by about 1%, while on the other two,

**Table 2　Performance comparison with different modeling choices and various methods on six cell lines in accordance with AUPR values.**

| | Model | AUPR value | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | IMR90 | HUVEC | K562 | GM12878 | NHEK | HeLa-S3 |
| Ablation experiment | Without multi-scale ResNet | 0.601 | 0.578 | 0.661 | 0.552 | 0.770 | 0.661 |
| | Three dense layers instead of matching layers | 0.651 | 0.629 | 0.630 | 0.550 | 0.743 | 0.694 |
| | Concatenation layer only | 0.672 | 0.664 | 0.665 | 0.577 | 0.771 | 0.711 |
| | Matching layers without subtract layer | 0.686 | 0.664 | 0.684 | 0.585 | 0.777 | 0.713 |
| | Matching layers without multiply layer | 0.684 | 0.670 | 0.672 | 0.582 | 0.789 | 0.716 |
| | Matching layers without addition and dot layers | 0.692 | 0.674 | 0.711 | 0.593 | 0.809 | 0.720 |
| Existing method | SPEID | 0.313 | 0.298 | 0.314 | 0.348 | 0.394 | 0.396 |
| | SimCNN | 0.388 | 0.348 | 0.395 | 0.403 | 0.540 | 0.496 |
| | EPIsHilbert | 0.608 | 0.537 | 0.686 | 0.538 | 0.752 | 0.685 |
| Proposed method | EPIMR | **0.698** | **0.691** | **0.729** | **0.603** | **0.811** | **0.730** |

**Table 3　Performance comparison with different modeling choices and various methods on six cell lines based on AUROC values.**

| | Model | AUROC value | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | IMR90 | HUVEC | K562 | GM12878 | NHEK | HeLa-S3 |
| Ablation experiment | Without multi-scale ResNet | 0.916 | 0.898 | 0.940 | 0.896 | **0.979** | 0.911 |
| | Three dense layers instead of matching layers | 0.913 | 0.900 | 0.920 | 0.891 | 0.971 | 0.915 |
| | Concatenation layer only | 0.918 | 0.903 | 0.914 | 0.905 | 0.974 | 0.907 |
| | Matching layers without subtract layer | 0.927 | 0.913 | 0.934 | 0.905 | 0.977 | 0.923 |
| | Matching layers without multiply layer | 0.916 | 0.914 | 0.936 | 0.901 | 0.976 | 0.916 |
| | Matching layers without addition and dot layers | **0.930** | **0.916** | 0.937 | 0.906 | 0.977 | 0.924 |
| Existing method | SPEID | 0.896 | 0.885 | 0.909 | 0.869 | 0.936 | 0.897 |
| | SimCNN | 0.907 | 0.897 | 0.922 | 0.896 | 0.947 | 0.912 |
| | EPIsHilbert | 0.894 | 0.874 | 0.895 | 0.901 | 0.945 | 0.914 |
| Proposed method | EPIMR | 0.923 | **0.916** | **0.941** | **0.913** | **0.979** | **0.933** |

results increased essentially the same. In short, adding multi-scale representation in ResNet can achieve the best performances in most cases.

We also try to remove some parts of the matching heuristic modules in Section 2.4. We concatenate the enhancer and promoter features, seamlessly integrating them into one/three dense layers, which are denoted as the items "concatenation layer only" and "three dense layers instead of matching layers". We also attempt to build models without one or two layers of matching heuristics. The results show that our matching operations are different from the simple concatenation of features or the use of multiple dense layers. And four different kinds of matching layers make different contributions to the final results. But overall EPIMR with additional matching heuristic steps is effective in improving predictive performances on nearly all of the six cell lines. It can be concluded that the refinement steps employed in our EPIMR can capture more information not only for enhancer or promoter sequences but also for their interactive communications.

### 3.2 Comparison with state-of-the-art method

To verify model validity, we also compare the performances of EPIMR with that of several other existing models, including SPEID[37] and SimCNN[38], which use one-hot to represent DNA sequences, and EPIsHilbert[47]. Each predictor is trained and tested on the same original imbalanced dataset on each cell line. And all of the compared models are trained through the process described in their papers. The prediction results of AUPR and AUROC for all models on six cell lines are presented in Fig. 3, Table 2, and Table 3, respectively.

Although the existing models have reported excellent results with data augmentation and some training strategies, they did not perform well on benchmark datasets. In contrast, EPIMR achieves better results on AUPR and AUROC. For AUPR in particular, our model is at least 4.35% higher than the runner-up models on all of the cell lines. It demonstrates that EPIMR can better deal with the datasets that approximate the actual distribution of EPIs. With regard to AUROC, EPIMR also outperforms the state-of-the-art methods on six cell lines, with an over 1% increment in most cases. We are of the opinion that using Hilbert Curve to extract features can better represent sequence location information and spatial information than simply using one-hot.

### 3.3 Pre-training strategy for cross-cell line prediction

The enhancer-promoter interactions are cell line-specific, giving rise to different interaction principles among different cell lines. To confirm this, we apply models established on specific cell lines to cross-cell line prediction. Figure 4 illustrates the AUPR and AUROC results of cross-cell line prediction using cell line-specific models. Not surprisingly, our proposed model performs well when the training set and test set come from the same cell line, especially with AUROC
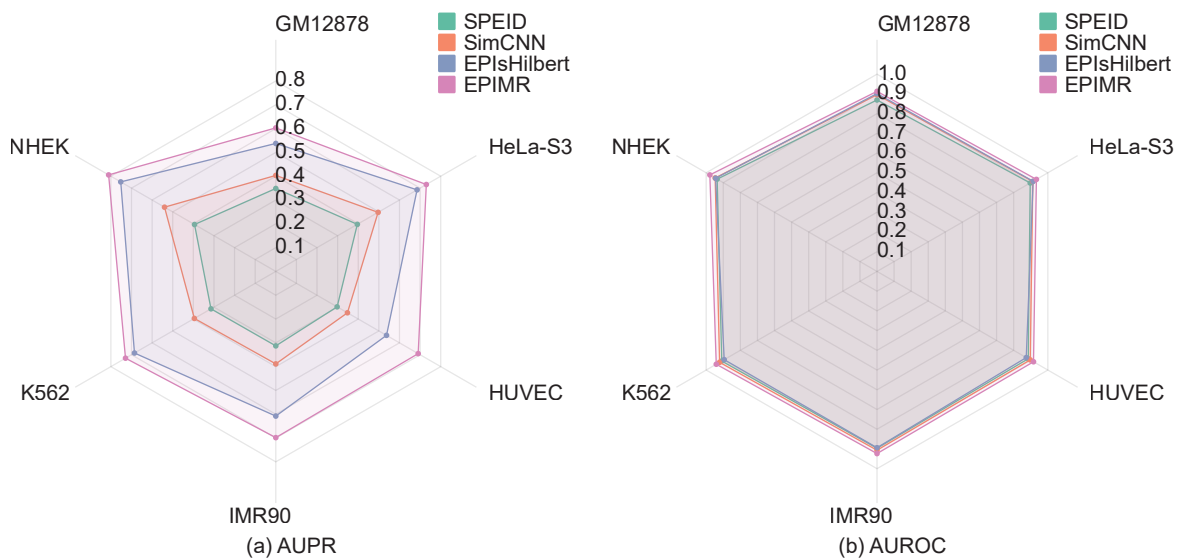


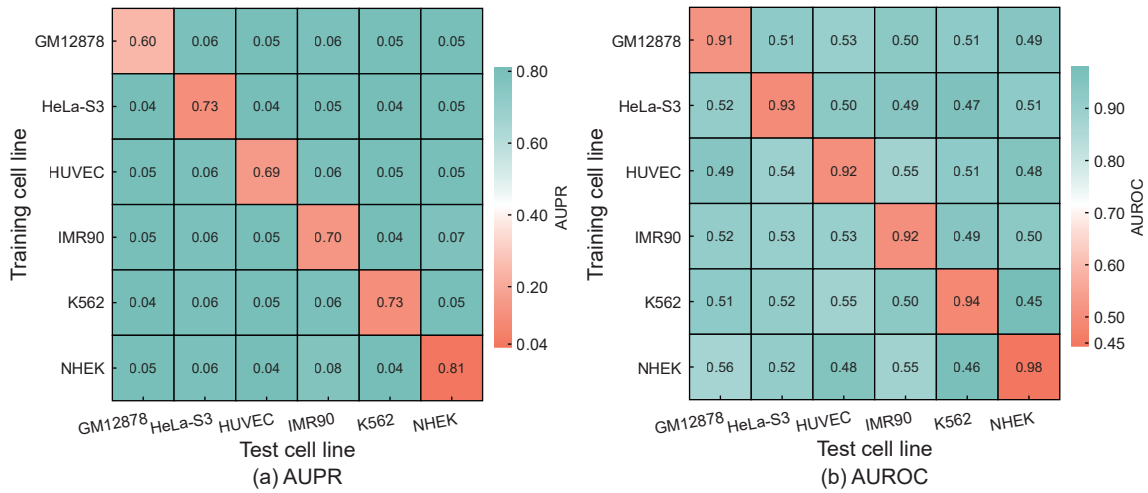**Fig. 3   Performances of different models in AUPR and AUROC metrices on six cell lines.**

**Fig. 4** **Performances of cell line-specific models for cross-cell line prediction.**

values exceeding 0.9 (the diagonal results in Fig. 4). In comparison, the prediction performances are abysmal for cross-cell line evaluation. It indicates that the model trained on a particular cell line cannot accurately predict EPIs on other cell lines, since they can only learn the interaction patterns of a single cell line and lack generalization ability across datasets of different cell lines. In other words, cell line-specific interactions exist.

It is quite clear that the corresponding models trained for each cell line can capture EPI cell line-specific features. Further, we intend to complement features that are common among all six cell lines to improve the prediction ability. Here, we adopt a pre-training strategy, and the procedure is described as follows:

● Create a new training set $D_{tr}^{all}$ by aggregating the training set of six cell lines.

● Pre-train a model for 15 epochs on the $D_{tr}^{all}$, a

training set created in the first step.

● Continue to train an additional 20 epochs using the target cell line's training set $D_{tr}^{s}$.

● Evaluate the specific test set $D_{ts}^{s}$ of this cell line.

Figure 5 reveals the AUPR and AUROC values based on the pre-trained models. Compared with the heat maps in Fig. 4, it can be perceived that the cross-cell line prediction performances are dramatically increased. Generally, AUPR scores increase by 30%−50% and AUC scores also have at least 30% increases after applying the pre-training approach. An obvious example is that when we adopt a pre-trained model with a target cell line K562 to identify EPIs of five other cell lines, we obtain AUC values of 0.864 to 0.940, a significant increase over that from the cell-specific models. For predicting particular cell lines, pre-trained models achieve better results as well, suggesting that there are many of the features extracted
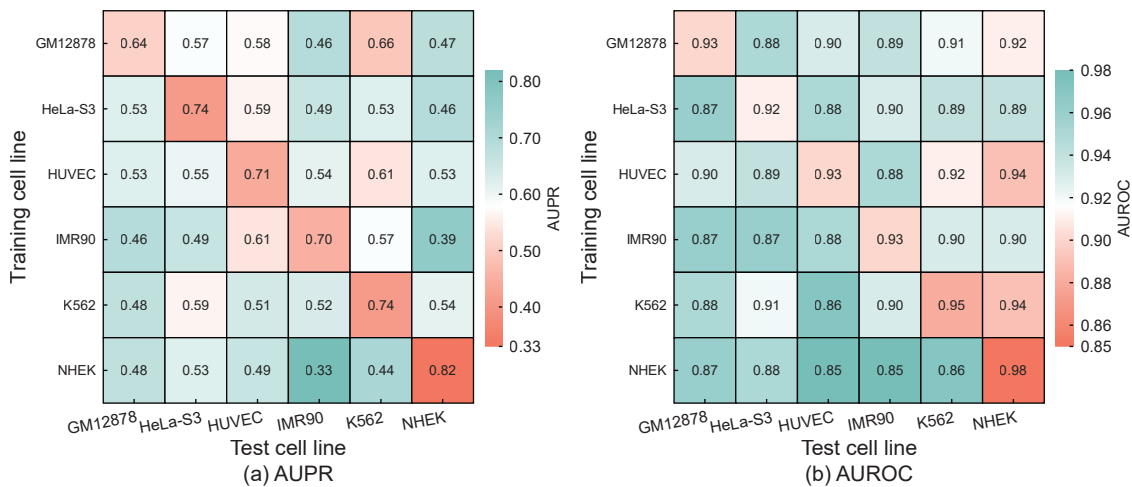


**Fig. 5** **Performances of pre-trained models for cross-cell line prediction.**

from EPIs shared in different cell lines. Therefore, our pre-trained model not only prominently outperforms the former one while recognizing EPIs in various cell lines, but also has a good ability for cell line-specific prediction.

## 3.4 Evaluation on under-sampling datasets

Since the benchmark datasets are highly imbalanced, we execute under-sampling algorithms to balance the two classes. To be specific, for each cell line, we randomly remove samples from the majority class to get a balanced dataset with equal numbers of negative samples and positive samples.

We evaluate the performances of AUPR and AUROC of our model on these under-sampling datasets, as shown in Tables 4 and 5. It is clear that

**Table 4 Performances of different models on under-sampling datasets in terms of AUPR.**

| Model | AUPR value | | | | | |
|---|---|---|---|---|---|---|
| | IMR90 | HUVEC | K562 | GM12878 | NHEK | HeLa-S3 |
| SPEID | 0.677 | 0.770 | 0.742 | 0.800 | 0.821 | 0.862 |
| SimCNN | 0.717 | 0.754 | 0.804 | 0.779 | 0.872 | 0.864 |
| EPIsHilbert | 0.880 | 0.851 | 0.880 | 0.845 | 0.865 | 0.929 |
| EPIMR | **0.899** | **0.961** | **0.932** | **0.925** | **0.961** | **0.970** |

**Table 5 Performances of different models on under-sampling datasets in terms of AUROC.**

| Model | AUROC value | | | | | |
|---|---|---|---|---|---|---|
| | IMR90 | HUVEC | K562 | GM12878 | NHEK | HeLa-S3 |
| SPEID | 0.779 | 0.847 | 0.801 | 0.836 | 0.881 | 0.901 |
| SimCNN | 0.800 | 0.834 | 0.846 | 0.827 | 0.903 | 0.892 |
| EPIsHilbert | 0.876 | 0.862 | 0.852 | 0.825 | 0.843 | 0.903 |
| EPIMR | **0.891** | **0.955** | **0.901** | **0.910** | **0.944** | **0.959** |

AUPR has made significant progress on these balanced datasets. The AUPR and AUROC values are basically above 0.9.

In addition, we compare EPIMR with other existing models on these datasets, as illustrated in Fig. 6 (full details are given in Tables 4 and 5). Even on the IMR90 cell line, EPIMR can obtain outstanding results with AUPR score 1.92% higher and AUROC score 1.51% higher than the second. The numbers leap to 11.02% and 9.23% on the HUVEC cell line. So it is evident that EPIMR also performs well on balanced datasets.

## 3.5 Verification on re-formulated datasets

Cao and Fullwood[68] raised some issues with TargetFinder datasets. They suggested that positive E-P pairs had highly overlapping windows with other positive samples in the same dataset, thus causing high similarity between training and test sets, which would exaggerate model prediction performances. In response, Whalen and Pollard[69] re-formulated datasets to predict chromatin interactions between genomic bins to reduce dependence and overcome generalization issues. Therefore, we verify the performances of EPIMR on these reconstructed sample sets.

Since all sequences are 5000 bp long, we need to set the number of iterations as 7 while constructing Hilbert Curve images, acquiring a map of $2^7 \times 2^7$ to ensure that it is enough to represent sequences. Although we chose the minimum number of iterations, most of the pixels in maps are still unused. Thus, cropping image by removing unused parts of the images can increase the proportion of pixels that are used. We end up with (64,
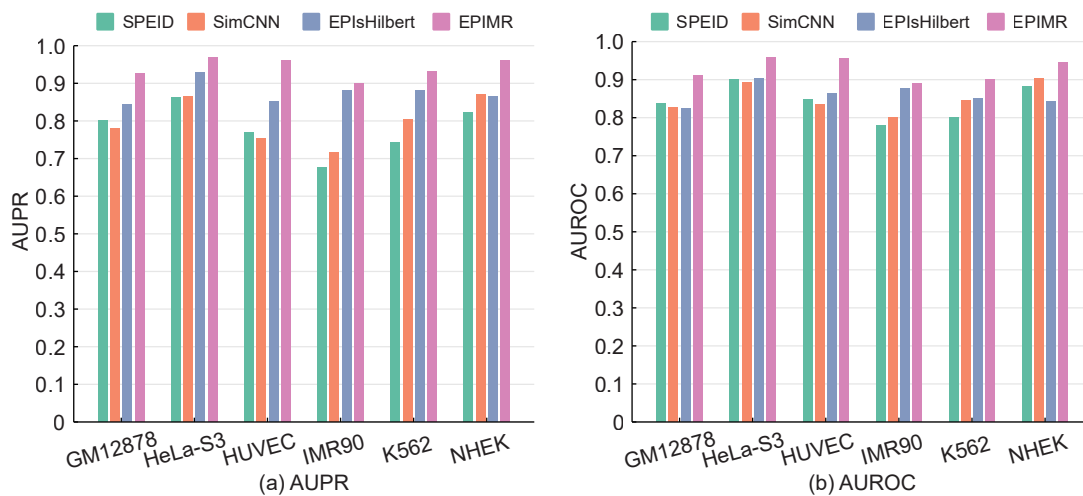


Fig. 6 Performances of different models on under-sampling datasets in terms of AUPR and AUROC.

96, 4) matrices.

We choose three cell lines from the re-formulating datasets including GM12878, IMR90, and K562. These three datasets differ greatly in size, as shown in Table 6. A 1:10 ratio of positives to negatives still exists in these new datasets.

We evaluate model performances using cell line-specific models and pre-trained models. As shown in Fig. 7, Table 7, and Table 8, we find that our model performs lower on imbalanced datasets than it does on benchmark datasets. Meanwhile, our pre-trained model still maintains a good predictive ability after pre-training on three cell lines, with AUPR and AUROC

**Table 6   Details of EPIs re-formulated datasets.**

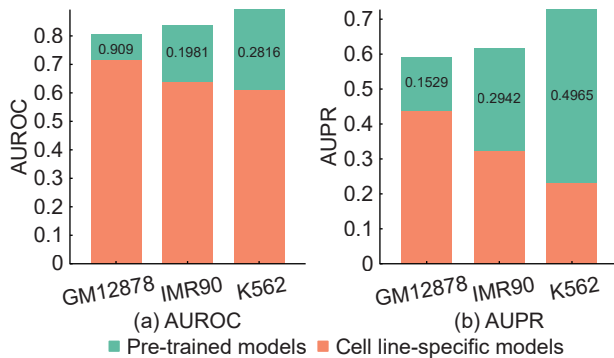| Cell line | Number | |
|---|---|---|
| | Positive | Negative |
| IMR90 | 3882 | 41 215 |
| GM12878 | 6260 | 63 676 |
| K562 | 1533 | 15 882 |



**Fig. 7   Performances of our model compared with that using pre-training strategy in terms of AUPR and AUROC on the re-formulating datasets. The values show the percentage improvement of the corresponding metric after using the pre-training strategy.**

**Table 7   Performances of our model on the re-formulating datasets.**

| Dataset | AUPR | AUROC |
|---|---|---|
| GM12878 | 0.439 | 0.714 |
| IMR90 | 0.324 | 0.639 |
| K562 | 0.231 | 0.611 |

**Table 8   Performances of our model using the pre-training strategy on the re-formulating datasets.**

| Dataset | AUPR | AUROC |
|---|---|---|
| GM12878 | 0.592 | 0.805 |
| IMR90 | 0.618 | 0.837 |
| K562 | 0.727 | 0.892 |

results rising sharply. Specifically, AUPR scores surge by about 15%−50%, whereas AUROC gains 9%−28%. Furthermore, the smaller the sample size of the cell line is, the more growth it can get through the pre-training approach. In other words, the pre-training strategy can greatly improve prediction results, especially may have practical applications for predicting EPIs in cell lines with little data.

### 3.6   Visualisation of multi-scale representation

As we all know, the shallow network can extract concrete features, whereas the deep network can extract more abstract and comprehensive features. To understand feature distribution at different scales of the model, as shown in Fig. 8, uniform manifold approximation and projection (UMAP)[70] is used to visualize the outputs of GAP1, GAP2, GAP3, and GAP4 of multi-scaling features in four stages in our model.

Based on the re-formulating datasets after pre-training, a multi-scale ResNet model was constructed for the K562 cell line. Due to the large amount of data, we randomly under-sample the dataset to ensure that there is a balance between negative and positive samples. Then we visualize the results for each GAP layer. It can be found that all these four scale features have classification abilities by observing that two categories of samples have their convergence trend and are separated. However, the classification abilities differ in performance due to their different scales of concern. For example, the features from GAP1 of enhancers that only concern local information are separated into two classes based on their labels although not as clear as GAP4 which focuses on global information. In brief, the respective features of EPIs and non-EPIs are gradually visualized into two classes, and the differences in feature distributions from GAP1 to GAP4 increase progressively. Accordingly, four scale features calculated from various depths of the network exhibit different performances and contribute distinctive contributions to classification results.

## 4   Conclusion

Here, we propose an enhanced promoter-enhancer interaction model called EPIMR, which uses only sequence information. We encode sequences into images by using the Hilbert Curve. Compared with other sequence encoding methods, it allows us to retain more sequence location information and spatial
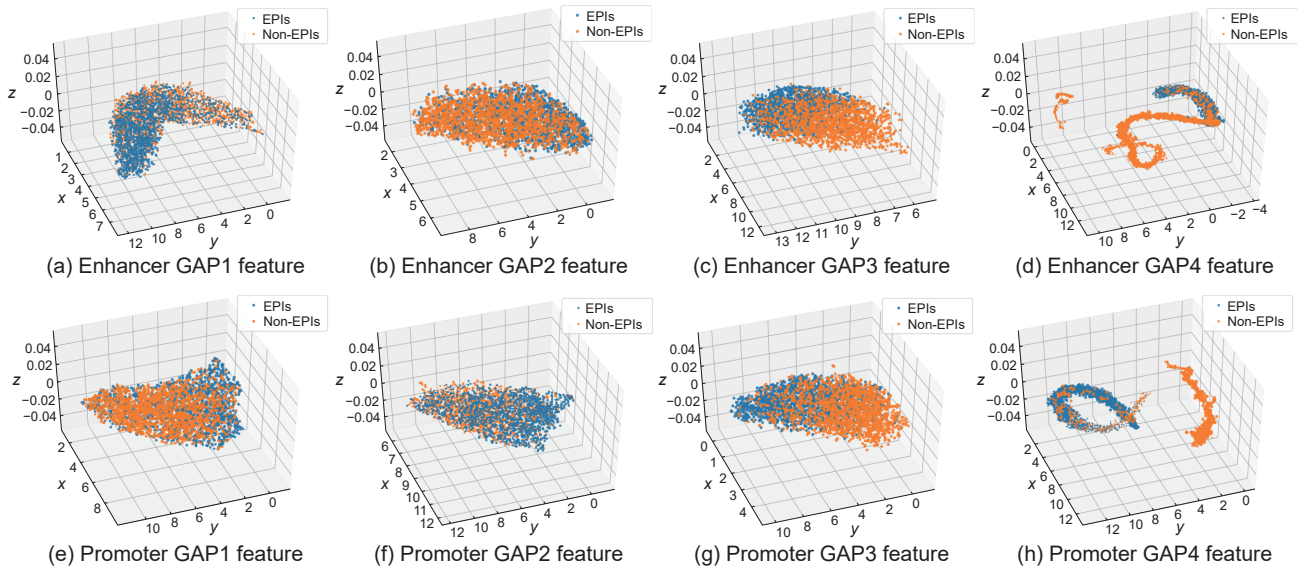
Fig. 8    Feature distributions from different stages of the network visualized by UMAP.

information. Specifically, we introduce a multi-scale ResNet that integrates local and concrete features extracted from the shallow part of the network with comprehensive and abstract features obtained from the deep network, which is not considered in previous EPI prediction models. We use UMAP for visualization and intuitively exhibit the feature distributions of different depths of the network. Finally, we use matching heuristics to match features between enhancer and promoter, focusing on the potential interaction information between them. Experiments illustrate that EPIMR produces better results on the benchmark datasets and under-sampling datasets than the most advanced methods available, and makes a great contribution to the prediction of EPIs on both the benchmark datasets and the reconstructed datasets. In addition, we focus on cross-cell line prediction and pre-train our models among all cell lines so that model can capture both specific as well as common features of cell lines, which verifies the transferability of our model. It also indicates that EPIMR has potential practical applications for identifying EPIs on cell lines that lack data.

## Acknowledgment

## References

[1]  Y. Lyu, Z. Zhang, J. Li, W. He, Y. Ding, and F. Guo, iEnhancer-KL: A novel two-layer predictor for identifying enhancers by position specific of nucleotide composition, *IEEE/ACM Trans. Comput. Biol. Bioinform.*, vol. 18, no. 6, pp. 2809–2815, 2021.

[2]  Y. Lyu, W. He, S. Li, Q. Zou, and F. Guo, iPro2L-PSTKNC: A two-layer predictor for discovering various types of promoters by position specific of nucleotide composition, *IEEE J. Biomed. Health Inform.*, vol. 25, no. 6, pp. 2329–2337, 2021.

[3]  J. Jin, Y. Yu, R. Wang, X. Zeng, C. Pang, Y. Jiang, Z. Li, Y. Dai, R. Su, Q. Zou, et al., iDNA-ABF: Multi-scale deep biological language learning model for the interpretable prediction of DNA methylations, *Genome Biol.*, vol. 23, no. 1, p. 219, 2022.

[4]  K. Monfils and T. S. Barakat, Models behind the mystery of establishing enhancer-promoter interactions, *Eur. J. Cell Biol.*, vol. 100, nos. 5&6, p. 151170, 2021.

[5]  D. T. Bergman, T. R. Jones, V. Liu, J. Ray, E. Jagoda, L. Siraj, H. Y. Kang, J. Nasser, M. Kane, A. Rios, et al., Compatibility rules of human enhancer and promoter sequences, *Nature*, vol. 607, no. 7917, pp. 176–184, 2022.

[6]  L. Liu, L. R. Zhang, F. Y. Dao, Y. C. Yang, and H. Lin, A computational framework for identifying the transcription factors involved in enhancer-promoter loop formation, *Mol. Ther. Nucleic Acids*, vol. 23, pp. 347–354, 2020.

[7]  H. Lv, F. Y. Dao, H. Zulfiqar, W. Su, H. Ding, L. Liu, and H. Lin, A sequence-based deep learning approach to

predict CTCF-mediated chromatin loop, *Brief. Bioinform.*, vol. 22, no. 5, p. bbab031, 2021.

[8] O. Kyrchanova and P. Georgiev, Mechanisms of enhancer-promoter interactions in higher eukaryotes, *Int. J. Mol. Sci.*, vol. 22, no. 2, p. 671, 2021.

[9] N. V. N. Carullo and J. J. Day, Genomic enhancers in brain health and disease, *Genes*, vol. 10, no. 1, p. 43, 2019.

[10] K. Hamamoto and T. Fukaya, Molecular architecture of enhancer-promoter interaction, *Curr. Opin. Cell Biol.*, vol. 74, pp. 62–70, 2022.

[11] C. Cao, J. Wang, D. Kwok, F. Cui, Z. Zhang, D. Zhao, M. J. Li, and Q. Zou, webTWAS: A resource for disease candidate susceptibility genes identified by transcriptome-wide association study, *Nucleic Acids Res.*, vol. 50, no. D1, pp. D1123–D1130, 2022.

[12] L. Yu, K. Yang, X. He, M. Li, L. Gao, and Y. Zha, Repositioning linifanib as a potent anti-necroptosis agent for sepsis, *Cell Death Discov.*, vol. 9, no. 1, p. 57, 2023.

[13] W. He, J. Tang, Q. Zou, and F. Guo, MMFGRN: A multi-source multi-model fusion method for gene regulatory network reconstruction, *Brief. Bioinform.*, vol. 22, no. 6, p. bbab166, 2021.

[14] M. Zhao, W. He, J. Tang, Q. Zou, and F. Guo, A comprehensive overview and critical evaluation of gene regulatory network inference technologies, *Brief. Bioinform.*, vol. 22, no. 5, p. bbab009, 2021.

[15] M. Zhao, W. He, J. Tang, Q. Zou, and F. Guo, A hybrid deep learning framework for gene regulatory network inference from single-cell transcriptomic data, *Brief. Bioinform.*, vol. 23, no. 2, p. bbab568, 2022.

[16] J. Wang, Y. Chen, and Q. Zou, Inferring gene regulatory network from single-cell transcriptomes with graph autoencoder model, *PLoS Genet.*, vol. 19, no. 9, p. e1010942, 2023.

[17] T. Borggrefe and B. D. Giaimo, *Enhancers and Promoters*. New York, NY, USA: Humana, 2021.

[18] R. Wang, Y. Jiang, J. Jin, C. Yin, H. Yu, F. Wang, J. Feng, R. Su, K. Nakai, Q. Zou, et al., DeepBIO: An automated and interpretable deep-learning platform for high-throughput biological sequence prediction, functional annotation and visualization analysis, *Nucleic Acids Res.*, vol. 51, no. 7, pp. 3017–3029, 2023.

[19] B. M. Javierre, O. S. Burren, S. P. Wilder, R. Kreuzhuber, S. M. Hill, S. Sewitz, J. Cairns, S. W. Wingett, C. Várnai, M. J. Thiecke, et al., Lineage-specific genome architecture links enhancers and non-coding disease variants to target gene promoters, *Cell*, vol. 167, no. 5, pp. 1369–1384, 2016.

[20] E. E. M. Furlong and M. Levine, Developmental enhancers and chromosome topology, *Science*, vol. 361, no. 6409, pp. 1341–1345, 2018.

[21] M. Osterwalder, I. Barozzi, V. Tissières, Y. Fukuda-Yuzawa, B. J. Mannion, S. Y. Afzal, E. A. Lee, Y. Zhu, I. Plajzer-Frick, C. S. Pickle, et al., Enhancer redundancy provides phenotypic robustness in mammalian development, *Nature*, vol. 554, no. 7691, pp. 239–243, 2018.

[22] F. Jing, S. W. Zhang, and S. Zhang, Prediction of enhancer-promoter interactions using the cross-cell type

information and domain adversarial neural network, *BMC Bioinformatics*, vol. 21, no. 1, p. 507, 2020.

[23] H. Ray-Jones and M. Spivakov, Transcriptional enhancers and their communication with gene promoters, *Cell. Mol. Life Sci.*, vol. 78, nos. 19&20, pp. 6453–6485, 2021.

[24] S. S. Rao, M. H. Huntley, N. C. Durand, E. K. Stamenova, I. D. Bochkov, J. T. Robinson, A. L. Sanborn, I. Machol, A. D. Omer, E. S. Lander, et al., A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping, *Cell*, vol. 159, no. 7, pp. 1665–1680, 2014.

[25] J. Dekker, K. Rippe, M. Dekker, and N. Kleckner, Capturing chromosome conformation, *Science*, vol. 295, no. 5558, pp. 1306–1311, 2002.

[26] J. Dostie, T. A. Richmond, R. A. Arnaout, R. R. Selzer, W. L. Lee, T. A. Honan, E. D. Rubio, A. Krumm, J. Lamb, C. Nusbaum, et al., Chromosome Conformation Capture Carbon Copy (5C): A massively parallel solution for mapping interactions between genomic elements, *Genome Res.*, vol. 16, no. 10, pp. 1299–1309, 2006.

[27] E. Lieberman-Aiden, N. L. van Berkum, L. Williams, M. Imakaev, T. Ragoczy, A. Telling, I. Amit, B. R. Lajoie, P. J. Sabo, M. O. Dorschner, et al., Comprehensive mapping of long-range interactions reveals folding principles of the human genome, *Science*, vol. 326, no. 5950, pp. 289–293, 2009.

[28] M. Simonis, P. Klous, E. Splinter, Y. Moshkin, R. Willemsen, E. de Wit, B. van Steensel, and W. de Laat, Nuclear organization of active and inactive chromatin domains uncovered by chromosome conformation capture–on-chip (4C), *Nat. Genet.*, vol. 38, pp. 1348–1354, 2006.

[29] L. Lu, X. Liu, W. K. Huang, P. Giusti-Rodríguez, J. Cui, S. Zhang, W. Xu, Z. Wen, S. Ma, J. D. Rosen, et al., Robust Hi-C maps of enhancer-promoter interactions reveal the function of non-coding genome in neural development and diseases, *Mol. Cell*, vol. 79, no. 3, pp. 521–534, 2020.

[30] J. C. Birkhoff, R. W. W. Brouwer, P. Kolovos, A. L. Korporaal, A. Bermejo-Santos, I. Boltsis, K. Nowosad, M. C. G. N. van den Hout, F. G. Grosveld, W. F. J. van IJcken, et al., Targeted chromatin conformation analysis identifies novel distal neural enhancers of ZEB2 in pluripotent stem cell differentiation, *Hum. Mol. Genet.*, vol. 29, no. 15, pp. 2535–2550, 2020.

[31] B. Mifsud, F. Tavares-Cadete, A. N. Young, R. Sugar, S. Schoenfelder, L. Ferreira, S. W. Wingett, S. Andrews, W. Grey, P. A. Ewels, et al., Mapping long-range promoter contacts in human cells with high-resolution capture Hi-C, *Nat. Genet.*, vol. 47, no. 6, pp. 598–606, 2015.

[32] Y. B. Zouari, A. M. Molitor, N. Sikorska, V. Pancaldi, and T. Sexton, ChiCMaxima: A robust and simple pipeline for detection and visualization of chromatin looping in Capture Hi-C, *Genome Biol.*, vol. 20, no. 1, p. 102, 2019.

[33] M. J. Fullwood, M. H. Liu, Y. F. Pan, J. Liu, H. Xu, Y. Bin Mohamed, Y. L. Orlov, S. Velkov, A. Ho, P. H. Mei, et al., An oestrogen-receptor-α-bound human chromatin interactome, *Nature*, vol. 462, no. 7269, pp. 58–64, 2009.

[34] M. R. Mumbach, A. J. Rubin, R. A. Flynn, C. Dai, P. A. Khavari, W. J. Greenleaf, and H. Y. Chang, HiChIP: Efficient and sensitive analysis of protein-directed genome architecture, *Nat. Meth.*, vol. 13, no. 11, pp. 919–922, 2016.

[35] S. Roy, A. F. Siahpirani, D. Chasman, S. Knaack, F. Ay, R. Stewart, M. Wilson, and R. Sridharan, A predictive modeling approach for cell line-specific long-range regulatory interactions, *Nucleic Acids Res.*, vol. 43, no. 18, pp. 8694–8712, 2015.

[36] S. Whalen, R. M. Truty, and K. S. Pollard, Enhancer-promoter interactions are encoded by complex genomic signatures on looping chromatin, *Nat. Genet.*, vol. 48, no. 5, pp. 488–496, 2016.

[37] S. Singh, Y. Yang, B. Póczos, and J. Ma, Predicting enhancer-promoter interaction from genomic sequence with deep neural networks, *Quant. Biol.*, vol. 7, no. 2, pp. 122–137, 2019.

[38] Z. Zhuang, X. Shen, and W. Pan, A simple convolutional neural network for prediction of enhancer–promoter interactions with DNA sequence data, *Bioinformatics*, vol. 35, no. 17, pp. 2899–2906, 2019.

[39] W. Mao, D. Kostka, and M. Chikina, Modeling enhancer-promoter interactions with attention-based neural networks, https://www.biorxiv.org/content/10.1101/219 667v1, 2017.

[40] W. Zeng, M. Wu, and R. Jiang, Prediction of enhancer-promoter interactions via natural language processing, *BMC Genomics*, vol. 19, no. Suppl2, p. 84, 2018.

[41] Y. Yang, R. Zhang, S. Singh, and J. Ma, Exploiting sequence-based features for predicting enhancer-promoter interactions, *Bioinformatics*, vol. 33, no. 14, pp. i252–i260, 2017.

[42] Z. Hong, X. Zeng, L. Wei, and X. Liu, Identifying enhancer-promoter interactions with neural network based on pre-trained DNA vectors and attention mechanism, *Bioinformatics*, vol. 36, no. 4, pp. 1037–1043, 2020.

[43] X. Min, C. Ye, X. Liu, and X. Zeng, Predicting enhancer-promoter interactions by deep learning and matching heuristic, *Brief. Bioinform.*, vol. 22, no. 4, p. bbaa254, 2021.

[44] S. Liu, X. Xu, Z. Yang, X. Zhao, S. Liu, and W. Zhang, EPIHC: Improving enhancer-promoter interaction prediction by using hybrid features and communicative learning, *IEEE/ACM Trans. Comput. Biol. Bioinform.*, vol. 19, no. 6, pp. 3435–3443, 2022.

[45] Z. Wang, L. Zhou, S. Jiang, and W. Huang, EPnet: A general network to predict enhancer-promoter interactions, in *Proc. 11th Int. Conf. Information Science and Technology* (*ICIST*), Chengdu, China, 2021, pp. 119–124.

[46] P. Ng, dna2vec: Consistent vector representations of variable-length k-mers, arXiv preprint arXiv: 1701.06279, 2017.

[47] M. Zhang, Y. Hu, and M. Zhu, EPIsHilbert: Prediction of enhancer-promoter interactions via Hilbert Curve encoding and transfer learning, *Genes*, vol. 12, no. 9, p. 1385, 2021.

[48] B. Yin, M. Balvert, D. Zambrano, A. Schönhuth, and S. Bohte, An image representation based convolutional network for DNA classification, arXiv preprint arXiv: 1806.04931, 2018.

[49] ENCODE Project Consortium, An integrated encyclopedia of DNA elements in the human genome, *Nature*, vol. 489, no. 7414, pp. 57–74, 2012.

[50] Roadmap Epigenomics Consortium, A. Kundaje, W. Meuleman, J. Ernst, M. Bilenky, A. Yen, A. Heravi-Moussavi, P. Kheradpour, Z. Zhang, J. Wang, et al., Integrative analysis of 111 reference human epigenomes, *Nature*, vol. 518, no. 7539, pp. 317–330, 2015.

[51] A. Krizhevsky, I. Sutskever, and G. E. Hinton, Imagenet classification with deep convolutional neural networks, *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.

[52] M. D. Zeiler and R. Fergus, Visualizing and understanding convolutional networks, in *Computer Vision—ECCV 2014*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, eds. Cham, Switzerland: Springer International Publishing, 2014, pp. 818–833.

[53] H. L. Li, Y. H. Pang, and B. Liu, BioSeq-BLM: A platform for analyzing DNA, RNA and protein sequences based on biological language models, *Nucleic Acids Res.*, vol. 49, no. 22, p. e129, 2021.

[54] K. Yan, H. Lv, Y. Guo, W. Peng, and B. Liu, sAMPpred-GAT: Prediction of antimicrobial peptide by graph attention network and predicted peptide structure, *Bioinformatics*, vol. 39, no. 1, p. btac715, 2023.

[55] Y. J. Tang, Y. H. Pang, and B. Liu, IDP-Seq2Seq: Identification of intrinsically disordered regions based on sequence to sequence learning, *Bioinformatics*, vol. 36, no. 21, pp. 5177–5186, 2021.

[56] K. He, X. Zhang, S. Ren, and J. Sun, Deep residual learning for image recognition, in *Proc. IEEE Conf. Computer Vision and Pattern Recognition* (*CVPR*), Las Vegas, NV, USA, 2016, pp. 770–778.

[57] L. Liu, Z. Qiu, G. Li, S. Liu, W. Ouyang, and L. Lin, Crowd counting with deep structured scale integration network, in *Proc. IEEE/CVF Int. Conf. Computer Vision* (*ICCV*), Seoul, Republic of Korea, 2019, pp. 1774–1783.

[58] Z. Deng, H. Sun, S. Zhou, J. Zhao, L. Lei, and H. Zou, Multi-scale object detection in remote sensing imagery with convolutional neural networks, *ISPRS J. Photogramm. Remote. Sens.*, vol. 145, pp. 3–22, 2018.

[59] E. Rumetshofer, M. Hofmarcher, C. Röhrl, S. Hochreiter, and G. Klambauer, Human-level protein localization with convolutional neural networks, presented at the Int. Conf. Learn. Represent. (ICLR), Vancouver, Canada, 2018.

[60] Y. Liu, C. Sun, L. Lin, and X. Wang, Learning natural language inference using bidirectional LSTM model and inner-attention, arXiv preprint arXiv: 1605.09090, 2016.

[61] L. Mou, R. Men, G. Li, Y. Xu, L. Zhang, R. Yan, and Z. Jin, Natural language inference by tree-based convolution and heuristic matching, arXiv preprint arXiv: 1512.08422, 2015.

[62] Y. Nie and M. Bansal, Shortcut-stacked sentence encoders for multi-domain inference, arXiv preprint arXiv: 1708.02312, 2017.

[63] C. Ao, X. Ye, T. Sakurai, Q. Zou, and L. Yu, m5U-SVM: Identification of RNA 5-methyluridine modification sites based on multi-view features of physicochemical features and distributed representation, *BMC Biol.*, vol. 21, no. 1, p. 93, 2023.

[64] H. Li and B. Liu, BioSeq-Diabolo: Biological sequence similarity analysis using Diabolo, *PLoS Comput. Biol.*, vol. 19, no. 6, p. e1011214, 2023.

[65] X. Zeng, F. Wang, Y. Luo, S. G. Kang, J. Tang, F. C. Lightstone, E. F. Fang, W. Cornell, R. Nussinov, and F. Cheng, Deep generative molecular design reshapes drug discovery, *Cell Rep. Med.*, vol. 3, no. 12, p. 100794, 2022.

[66] J. Davis and M. Goadrich, The relationship between Precision-Recall and ROC curves, in *Proc. 23rd Int. Conf. Machine learning*, Pittsburgh, PA, USA, 2006, pp. 233–240.

[67] J. A. Hanley and B. J. McNeil, The meaning and use of the area under a receiver operating characteristic (ROC) curve, *Radiology*, vol. 143, no. 1, pp. 29–36, 1982.

[68] F. Cao and M. J. Fullwood, Inflated performance measures in enhancer-promoter interaction-prediction methods, *Nat. Genet.*, vol. 51, no. 8, pp. 1196–1198, 2019.

[69] S. Whalen and K. S. Pollard, Reply to 'Inflated performance measures in enhancer-promoter interaction-prediction methods', *Nat. Genet.*, vol. 51, no. 8, pp. 1198–1200, 2019.

[70] L. McInnes, J. Healy, and J. Melville, UMAP: Uniform manifold approximation and projection for dimension reduction, arXiv preprint arXiv: 1802.03426, 2018.

**Qiaozhen Meng** received the master degree from Tianjin University, China in 2018. She currently is pursuing the PhD degree at the College of Intelligence and Computing, Tianjin University, China. Her research interests include deep learning, protein function, and bioinformatics.

**Xiaoqing Peng** received the BS and MS degrees in computer engineering from Central South University, China in 2009 and 2012, respectively, and the PhD degree in computer science from Central South University, China in 2016. She is an associate professor at the Center for Medical Genetics & Hunan Key Laboratory of Medical Genetics, School of Life Sciences, Central South University, China. Her research interests focus on epigenetics and proteomics.

**Jijun Tang** received the PhD degree from University of New Mexico, USA in 2004. He currently is a professor at the Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, China. His main research interests include computational biology.

**Yinuo Lyu** received the BEng and MEng degrees from Tianjin University, China in 2020 and 2023, respectively. She currently is working at the Aeronautical Information Service Center of the Civil Aviation Administration of China (AISC.ATMB. CAAC), Beijing, China. Her research interests include deep learning and bioinformatics.

**Junhai Xu** received the PhD degree from Shandong University, China in 2014. He currently is an associate professor at College of Intelligence and Computing, Tianjin University, China. His research interests include medical image analysis and intelligent information processing.

**Fei Guo** received the PhD degree from Shandong University, China in 2012. She is currently a professor at the School of Computer Science and Engineering, Central South University, China. Her research interests include bioinformatics and computational biology.