

1
2
3
4
5
6
7
8
9
10
11
12

Supplementary Information

EpiGePT: a Pretrained Transformer model for epigenomics

Zijing Gao^{1,#}, Qiao Liu^{2,#,*}, Wanwen Zeng², Rui Jiang^{1,*} and Wing Hung Wong^{2,3,*}

¹ Ministry of Education Key Laboratory of Bioinformatics, Bioinformatics Division at the Beijing National Research Center for Information Science and Technology, Center for Synthetic and Systems Biology, Department of Automation, Tsinghua University, Beijing 100084, China;

² Department of Statistics, Stanford University, Stanford, CA 94305, USA;

³ Department of Biomedical Data Science, Bio-X Program, Center for Personal Dynamic Regulomes, Stanford University, Stanford, CA 94305, USA;

* To whom correspondence should be addressed.

The first two authors contributed equally.

E-mail: liuqiao@stanford.edu, whwong@stanford.edu, ruijiang@tsinghua.edu.cn

13 Contents

14	Supplementary Texts	3
15	Text S1. Data splitting strategy for model training.....	3
16	Text S2. System design and implementation of the web server.....	4
17	Text S3. Case application of the EpiGePT-online.....	5
18	Text S4. Running time of the EpiGePT and baseline methods.....	6
19	Text S5. Implementation of Enformer model and Enformer+.....	7
20	Text S6. Data processing for ChromHMM annotation data.....	8
21	Supplementary Figures	9
22	Fig. S1.....	9
23	Fig. S2.....	10
24	Fig. S3.....	12
25	Fig. S4.....	14
26	Fig. S5.....	16
27	Fig. S6.....	18
28	Fig. S7.....	19
29	Fig. S8.....	20
30	Fig. S9.....	21
31	Fig. S10.....	22
32	Fig. S11.....	24
33	Fig. S12.....	25
34	Fig. S13.....	26
35	Fig. S14.....	27
36	Supplementary Tables.....	28
37	References	29

38 **Supplementary Texts**

39 **Text S1. Data splitting strategy for model training.**

40 To comprehensively validate the performance of EpiGePT in predicting chromatin accessibility,
41 we adopted three different data splitting strategies in the DNase¹ prediction experiment to
42 verify the model's prediction ability when facing new genomic regions and cell types, which
43 can meet researchers' usage needs to the maximum extent. Firstly, cross-cell type prediction
44 refers to splitting the training and testing sets according to cell types in the same genomic
45 region, where the cell types in the testing set have not appeared in the training set (Figs. S1b).
46 Secondly, cross-genomic region prediction refers to splitting the training and testing sets
47 according to genomic regions in the same cell type (Figs. S1a). Thirdly, simultaneous cross-
48 cell type and genomic region prediction, where the prediction can be performed in completely
49 novel cell types and genomic regions with the expression of transcription factors in that cell
50 type. The training set needs to subset both cell types and genomic regions (Figs. S1c). To
51 complete the latter two auxiliary predictions, we also split the data into 5 folds according to
52 both cell types and genomic regions, so that both cross-validation can be performed in one
53 round of training, but this will also reduce the amount of training and testing data.

54 **Text S2. System design and implementation of the web server.**

55 EpiGePT-online runs on a Linux-based Apache web server (<https://www.apache.org>) and
56 utilizes the Bootstrap v3.3.7 framework (<https://getbootstrap.com/docs/3.3/>) for its web-
57 frontend display. The backend of the server uses PHP v7.4.5 (<http://www.php.net>). The
58 platform is compatible with the majority of mainstream web browsers, including Google
59 Chrome, Firefox, Microsoft Edge, and Apple Safari.

60 **Text S3. Case application of the EpiGePT-online.**

61 The online prediction web service of EpiGePT enables users to predict eight types of
62 epigenomic signals using EpiGePT without the need for setting up environments, writing code,
63 or computational resources. In this section, we describe a usage scenario of EpiGePT-online
64 for epigenomic signals prediction (Fig. S13). Users are provided with the flexibility to annotate
65 either multiple genomic regions or a single locus at their discretion. Assuming an algorithmic
66 researcher is interested in determining the potential regulatory role of a specific chromatin
67 region based on its epigenetic modifications. In this case, the researcher can utilize EpiGePT-
68 online to calculate the epigenetic signals on this region, to obtain references for assessing the
69 potential regulatory role of the region. The submission prerequisites encompass two essential
70 components. 1) The expression profiles of 711 TFs, which facilitate EpiGePT in acquiring
71 precise cell type or tissue information. 2) The specific location of a locus on the genome or
72 uploading of a bed file containing the information of genomic regions. It is worth noting that
73 each line in the uploaded BED file should correspond to a 128kbp region to comply with the
74 input length requirement of EpiGePT. If users select a specific locus, we will provide the
75 predicted results for the region spanning 128kbp upstream and downstream of that locus. The
76 web server allows users to upload expression values of 711 TFs in either numpy or comma-
77 separated values (CSV) format. When predicting for N genomic regions, users can obtain a
78 downloadable matrix stored in CSV format with dimensions $(N \times 1000, 8)$. Each row denotes
79 a 128bp genomic bin, and each column denotes an epigenetic profile. The specific referents
80 of each row and column are provided in the downloadable table. This allows users to perform
81 downstream analyses, such as related analyses in the areas of gene regulation and human
82 disease.

83 **Text S4. Running time of the EpiGePT and baseline methods.**

84 To demonstrate the computational efficiency of our model, we recorded the runtime of
85 EpiGePT and baseline methods for one epoch on two sets of experiments, with different data
86 sizes and input sequence lengths. Firstly, in the DNase signal prediction experiment on 129
87 cell types, with an input sequence length of 10kbp and using the same training data, Enformer
88 requires approximately 3 hours and 4 minutes to complete one epoch, while EpiGePT only
89 takes 2 hours and 17 minutes. In contrast, ChromDragoNN², which uses a genomic bin rather
90 than a long region as the model input, requires 24 hours for pre-training and 8 hours for fine-
91 tuning. In this case, the batch size of ChromDragoNN was set to 1024, which is equivalent to
92 EpiGePT using a batch size of around 20. This modeling and computation approach presents
93 challenges in terms of computational efficiency when dealing with large amounts of data.
94 DeepCAGE³ faces similar efficiency issues using the same approach. Even with a batch size
95 of 256 on a single GPU, it still takes nearly 10 hours to complete one epoch of training.
96 Secondly, we also recorded the running time of the models under larger-scale data and longer
97 input sequences. When the number of input genomic bins increased from 50 to 1000, which
98 corresponds to an input sequence length of approximately 128k, EpiGePT took approximately
99 3 hours to complete one epoch of training on 20 cell lines and 13,300 genomic regions, while
100 Enformer required approximately 27 hours to train one epoch, as it required a longer input
101 sequence of approximately 190kbp. Furthermore, EpiGePT without TF module (EpiGePT-seq)
102 had approximately 1/4 of the parameters of Enformer and took approximately 2 hours and 40
103 minutes to train. In terms of performance, EpiGePT-seq performed similarly to Enformer on
104 this dataset. This also explains why we chose to simplify the pure sequence model rather than
105 directly adding a TF module to Enformer.

106 **Text S5. Implementation of Enformer model and Enformer+.**

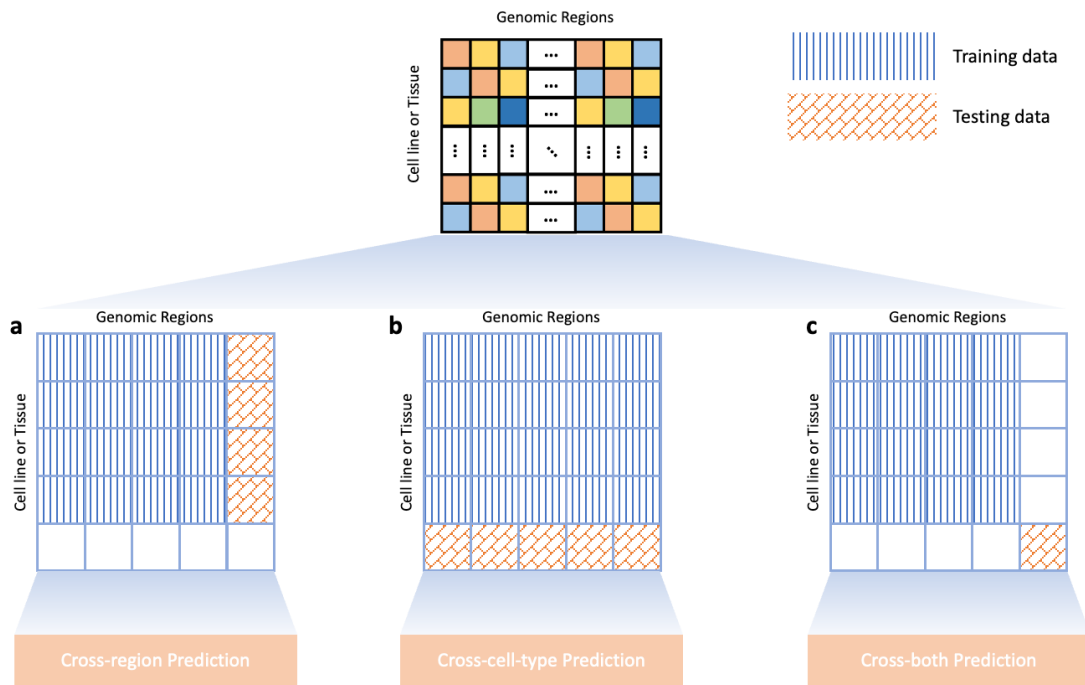
107 To ensure a fair comparison between models and prevent the possibility of information
108 leakage, we implemented the Enformer⁴ model ourselves and trained it on our own collected
109 data. Due to differences in dataset size and partitioning compared to Enformer, we reduced
110 the number of encoder layers in Enformer to prevent overfitting. Thus, we reduced the number
111 of encoder layers in Enformer to 3. Additionally, we introduced Enformer+ to enable a fair
112 comparison between EpiGePT and Enformer in bin-level prediction. As Enformer takes only
113 the DNA sequence as input, it tends to predict the same values for the same locus in different
114 cell types, resulting in a loss of locus-level prediction ability. To address this, we incorporated
115 the binding status and expression of the same transcription factors in Enformer+, and
116 compared it to EpiGePT's performance on the same tasks.

117 **Text S6. Data processing for ChromHMM annotation data.**

118 We downloaded the 15-state ChromHMM⁵ annotations across 127 epigenomes from the
119 ROADMAP project. The state of chromatin is annotated for each 200bp bin in a specific cell
120 type. RNA-seq data of TFs across 56 cell types were download and extracted from the
121 ROADMAP⁶ project (Supplementary table S10 and S11). Subsequently, we mapped the 711
122 transcription factors to the downloaded RNA-seq data, resulting in the identification of RNA-
123 seq data for 642 transcription factors. In the subsequent experiments, we utilized the
124 expression data of these 642 transcription factors. We finally calculated the normalized TPM
125 values of the 642 TFs on 56 cell types we extracted for the using in the classification model.
126 For coarse grain chromatin state prediction, we took the state 'Quies' as low signal regions
127 and other states as signal regions. For fine grain chromatin state prediction, we extracted the
128 state 'TssA', 'TssAFlnk', 'TssBiv' and 'BivFlnk' as TSS regions, state 'EnhG', 'Enh' and 'EnhBiv'
129 as enhancer regions, 'Quies' as low signal regions and other state as other regions. To balance
130 the number of different chromatin states, we downsampled the low signal regions and obtained
131 921,074 bin each cell line finally.

132 **Supplementary Figures**

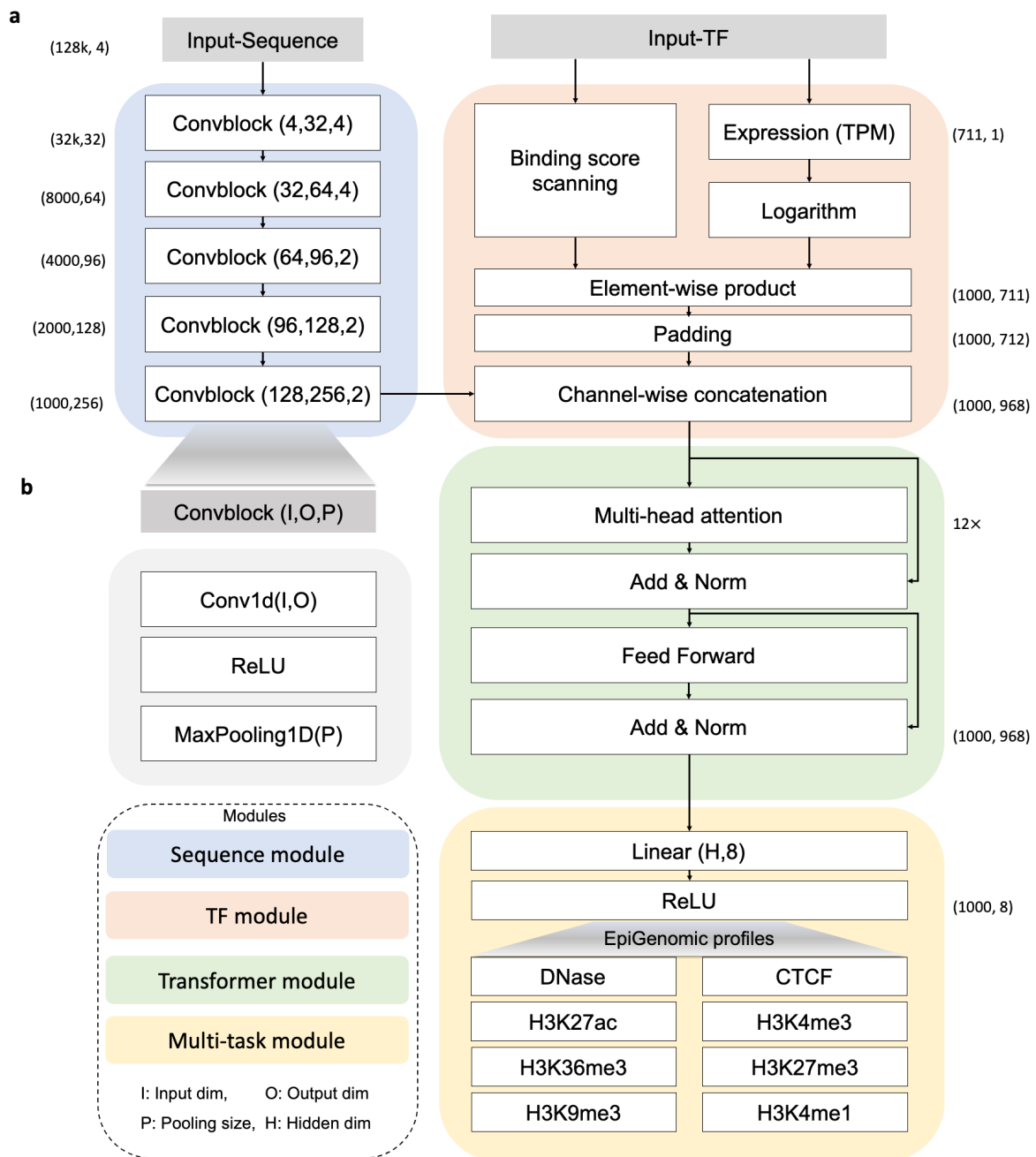
133 **Fig. S1**



134 **Fig. S1. Three data partitioning strategies for model training and testing.** **a**, Cross
135 genomic region prediction. The training and testing datasets utilized the expression profiles of
136 identical cell types, but were evaluated on novel genomic regions for prediction. **b**, Cross cell
137 type prediction. The training and testing datasets utilized the same genomic regions, but were
138 evaluated on novel cell types for prediction. **c**, Cross genomic region and cell type prediction.
139 The cell types and genomic regions used in the training and test sets were both different.

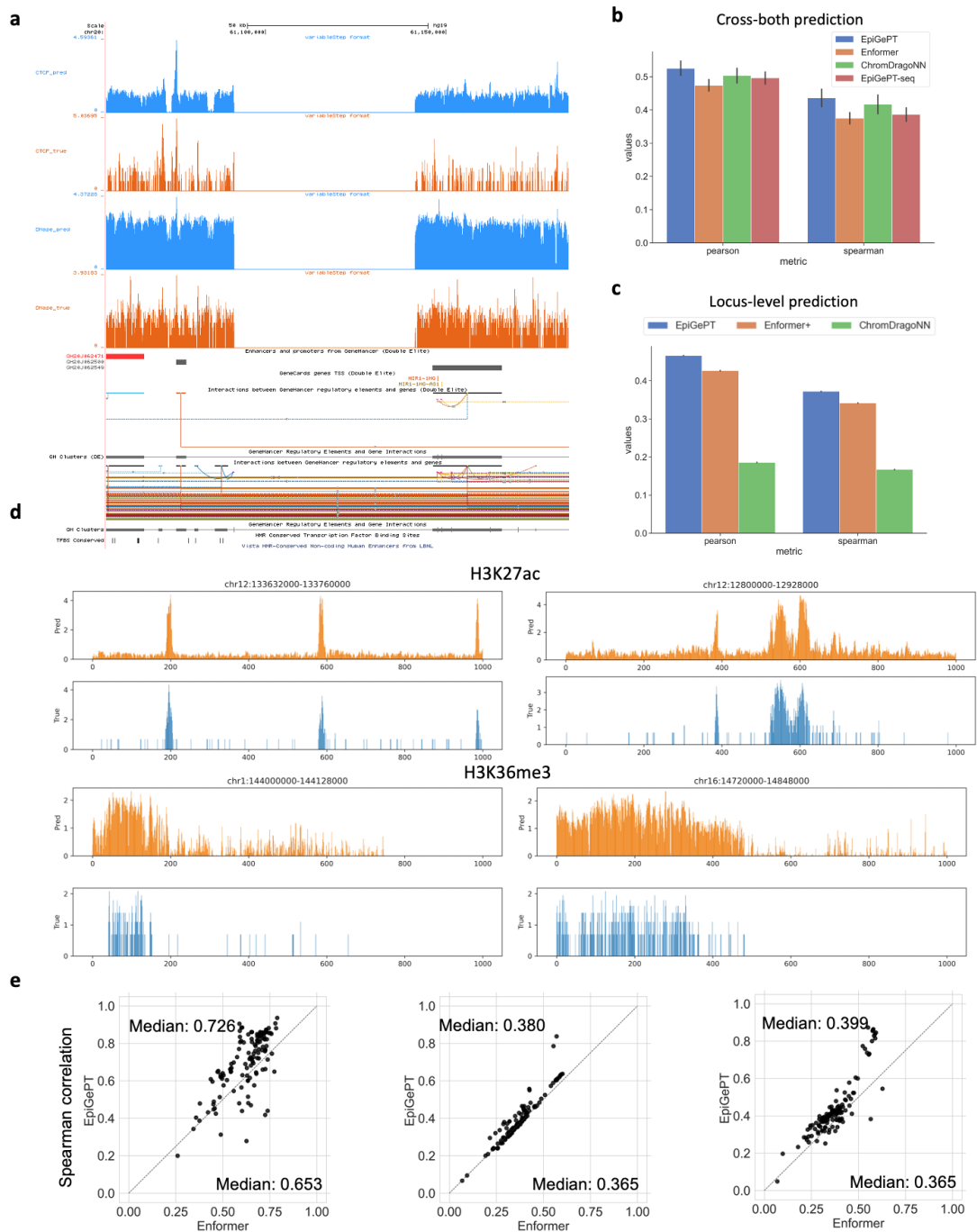
140

141 **Fig. S2**



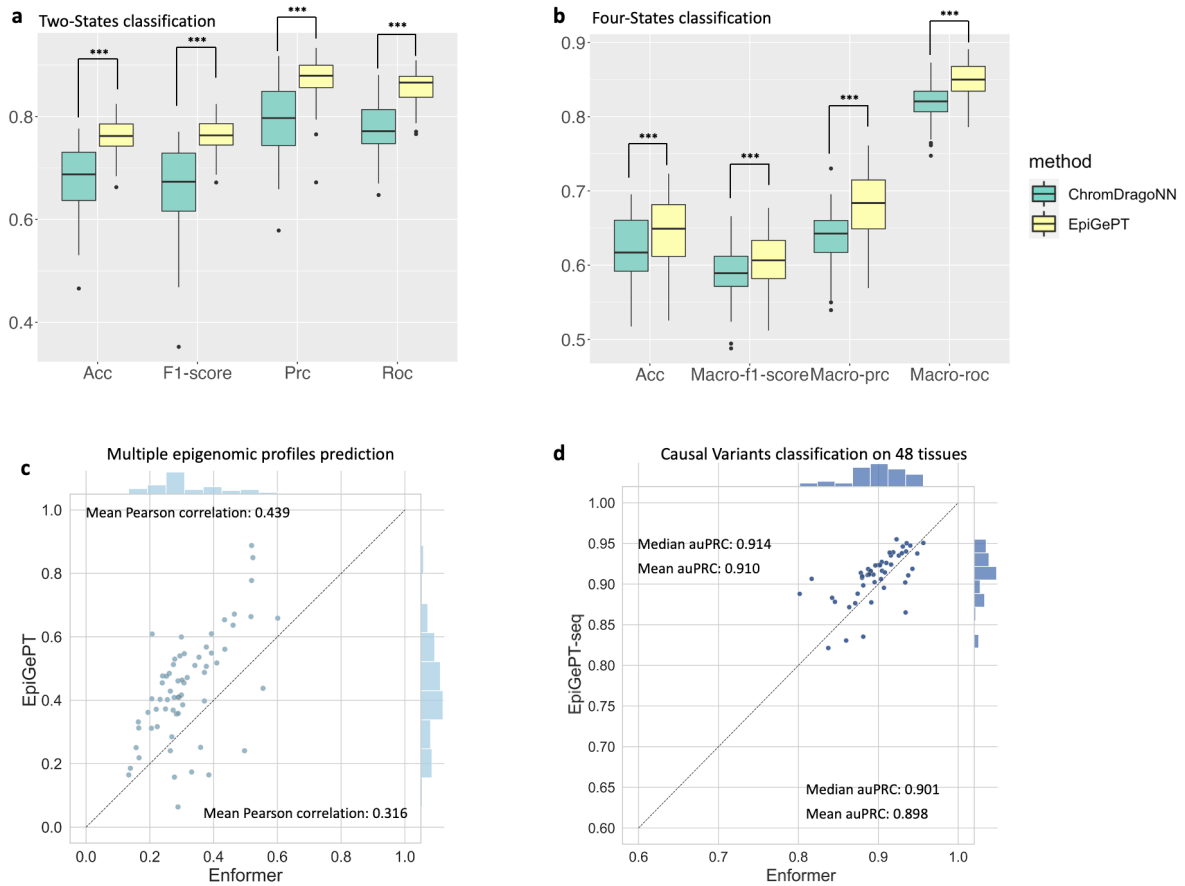
142 **Fig S2. Model architecture of EpiGePT for multiple epigenomic signals prediction. a,**
 143 **The computational process of EpiGePT. The sequence module employs a stack of five**
 144 **convolutional layers followed by pooling operations, resulting in representations that capture**
 145 **sequence patterns. The TF module integrates motif binding information and gene expression**
 146 **data to represent cell-specific information. The Transformer module takes the genomic bin**

147 sequences mentioned above as input and learns the interaction relationships between bins,
148 capturing the interactions among them. Finally, the obtained embeddings are mapped to the
149 eight types of epigenomic signals through a fully connected layer. **b**, Specific details of the
150 convolutional block involve the fusion of 1D convolution, ReLU activation function, and max
151 pooling operation to achieve changes in the feature dimension O and extract bin-level features.



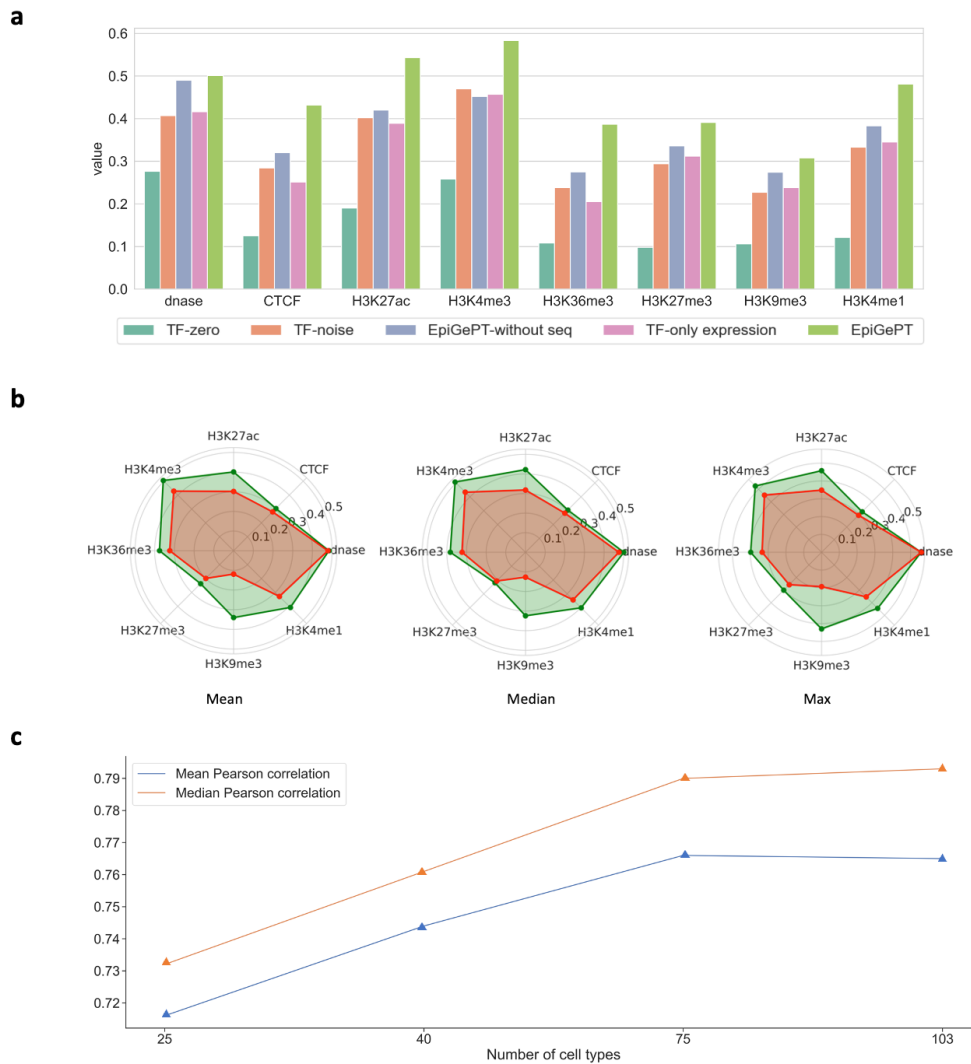
153 **Fig. S3. EpiGePT's performance in predicting DNase-seq and other epigenetic signals**
 154 **is demonstrated in a**, through visualization of predicted results for DNase and CTCF signals.
 155 EpiGePT is able to make accurate predictions for these signals, as well as for the regulatory

156 relationships within a genomic region of 20th chromosome ranging from 61,100,000 to
157 61,150,000. **b**, EpiGePT and baseline methods were compared for their performance in
158 predicting epigenetic signals in new cell types and genomic regions (cross-both prediction).
159 The left panel shows the Pearson correlation coefficient, and the right panel shows the
160 Spearman correlation coefficient. **c**, Locus level prediction of DNase signal. We predicted a
161 value for each genomic locus, and calculated the correlation coefficient between the predicted
162 values and true values for the same locus in different cell types. **d**, Visualization of predicted
163 signals, such as the comparison between predicted and true values in a 128kbp region (from
164 133,632,000 to 133,760,000) on chromosome 12, shows that the presence of a large number
165 of zeros in both the true and predicted signals can limit the correlation between the two signals.
166 **e**, Comparison of EpiGePT and Enformer performance. Each point in the scatter plot
167 represents the performance of Enformer on the data of a specific cell type (x-axis) compared
168 to the performance of EpiGePT (y-axis). The three graphs represent the prediction of
169 continuous DNase signals (spearman correlation coefficient).



171 **Fig. S4. Performance of EpiGePT and baseline methods on chromatin states**
 172 **classification, multiple epigenomic profiles prediction and causal variants**
 173 **classification. a**, Binary classification of chromatin states for distinguishing functional regions
 174 on the chromatin based on the annotation data from ChromHMM-15-states. **b**, Four-class
 175 chromatin state classification is used to distinguish functional regions on the chromatin,
 176 including TSS, potential enhancers, other functional regions, and non-functional regions
 177 based on the annotation data from ChromHMM-15-states. *** indicates that the p -value is less
 178 than $1e-3$ under one-sided Wilcoxon signed rank test. **c**, Cross-cell-type prediction of 8
 179 epigenomic signals at 8 test cell types. Each dot denotes the Pearson correlation coefficient
 180 of the predicted signals and true signals at the specific cell types on a specific epigenomic
 181 signal. **d**, The performance of EpiGePT and Enformer in discriminating causal eQTLs across

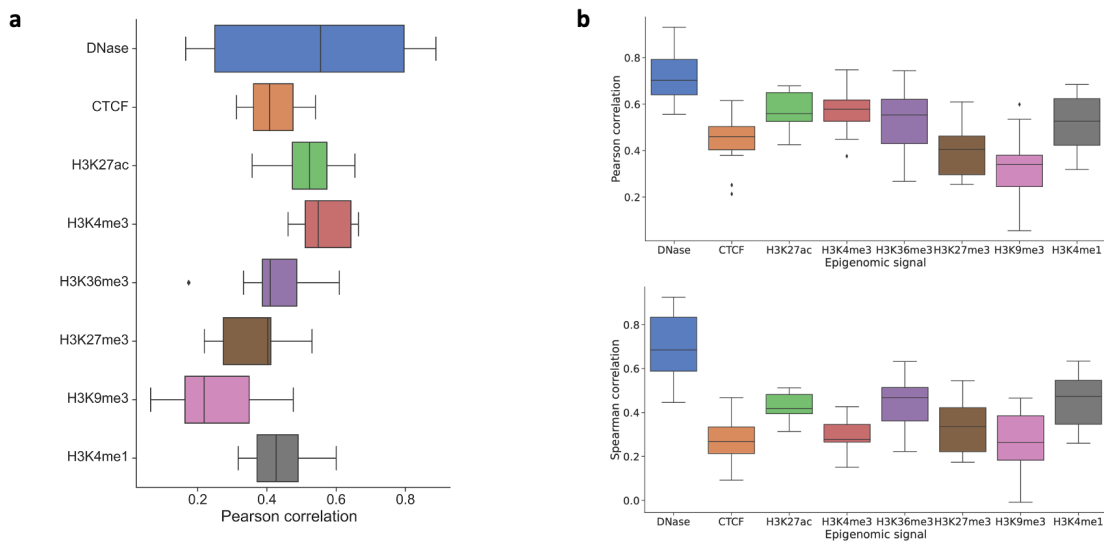
182 48 tissues, each dot representing the average auPRC obtained from 5-fold cross-validation
183 on a specific tissue.



185 **Fig. S5. Ablation analysis of the EpiGePT model.** **a**, Ablation analysis on the TF module
 186 and the Sequence module, we observed a decrease in predictive performance for each
 187 module across eight chromatin epigenetic signals, as evidenced by a reduction in Pearson
 188 correlation coefficient. **b**, Ablation analysis on the Multi-task module. The green shaded area
 189 in the figure represents the results of multi-signal cross-cell-type predictions, while the red
 190 shaded area represents the results of training and predicting on each signal individually. It can
 191 be observed that the multi-task module has a positive effect on the model performance across
 192 all signals. **c**, Ablation analysis of the number of the training cell types. When the number of

193 training cell types increases while the number of testing cell types remains constant, there is
194 an increasing trend in performance as the number of training cell types increases.

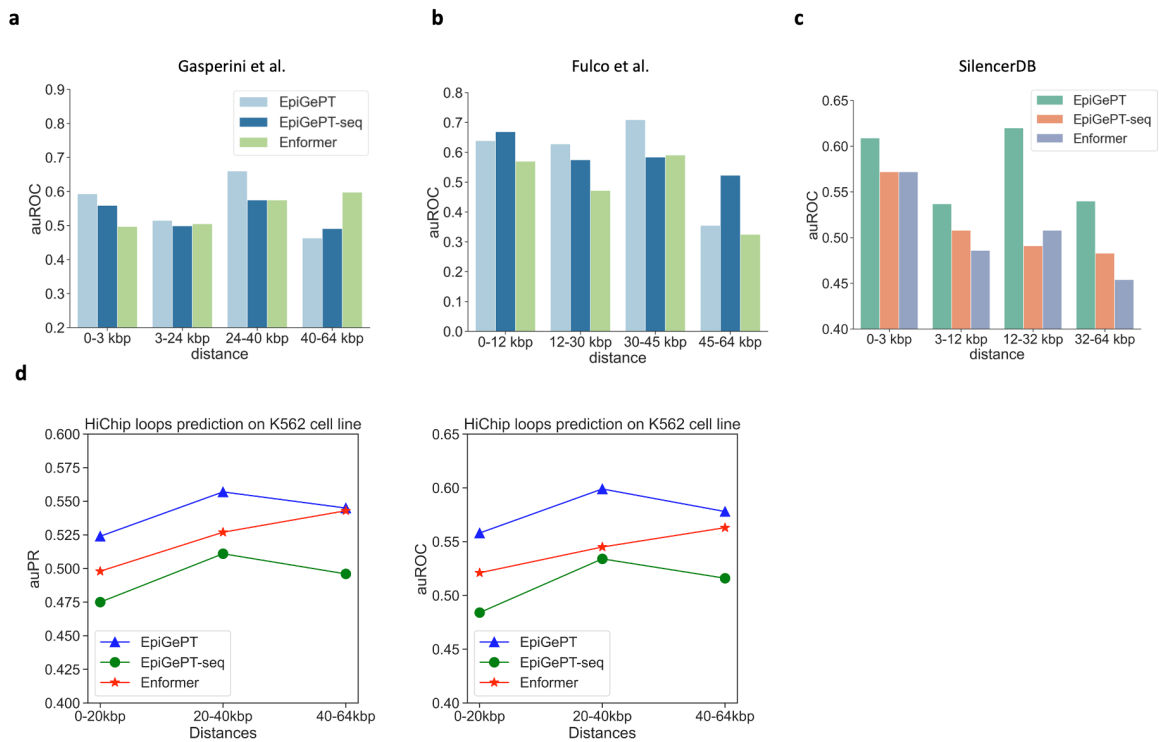
195 **Fig. S6**



196

197 **Fig. S6. Performance of EpiGePT in cross-cell-type prediction.** **a**, The predictive
 198 performance of EpiGePT on 8 unseen cell types on hg19 reference genome (pearson
 199 correlation coefficients). **b**, The predictive performance of EpiGePT on 19 new cell types on
 200 hg38 reference genome (upper: pearson correlation coefficients, lower: Spearman correlation
 201 coefficients).

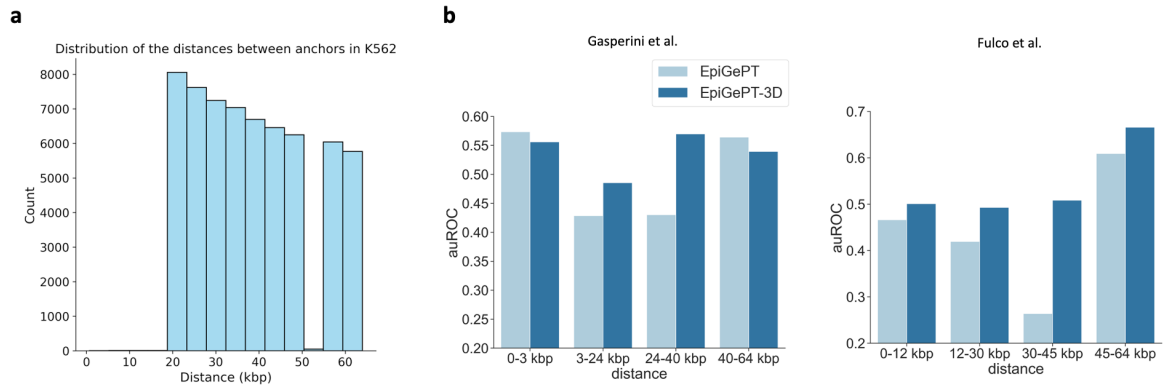
202



204 **Fig. S7. The performance (auROC) of attention score of EpiGePT in distinguishing**
 205 **regulatory element-gene pairs at different distance ranges. a,** The performance of
 206 EpiGePT in distinguishing enhancer-gene pairs at different distance ranges on the data from
 207 Gasperini et al⁷. **b,** The performance of EpiGePT in distinguishing enhancer-gene pairs at
 208 different distance ranges on the data from Fulco et al⁸. **c,** The performance of EpiGePT in
 209 distinguishing silencer-promoter pairs at different distance ranges on the data from
 210 SilencerDB⁹. **d,** The performance (auROC and auPR) of attention score of EpiGePT in
 211 distinguishing HiChIP loops of H3K27ac at different distance ranges on K562 cell line.

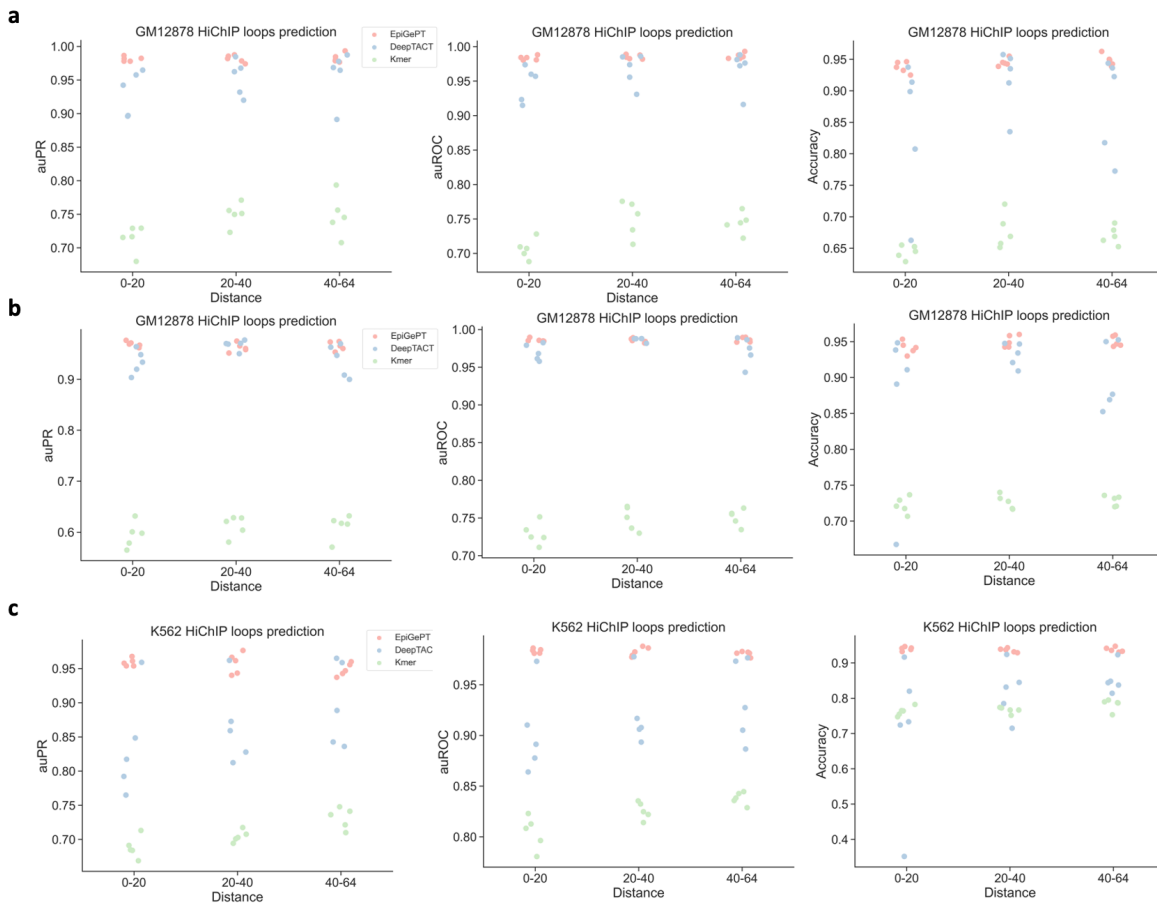
212 **Fig. S8**

213



214

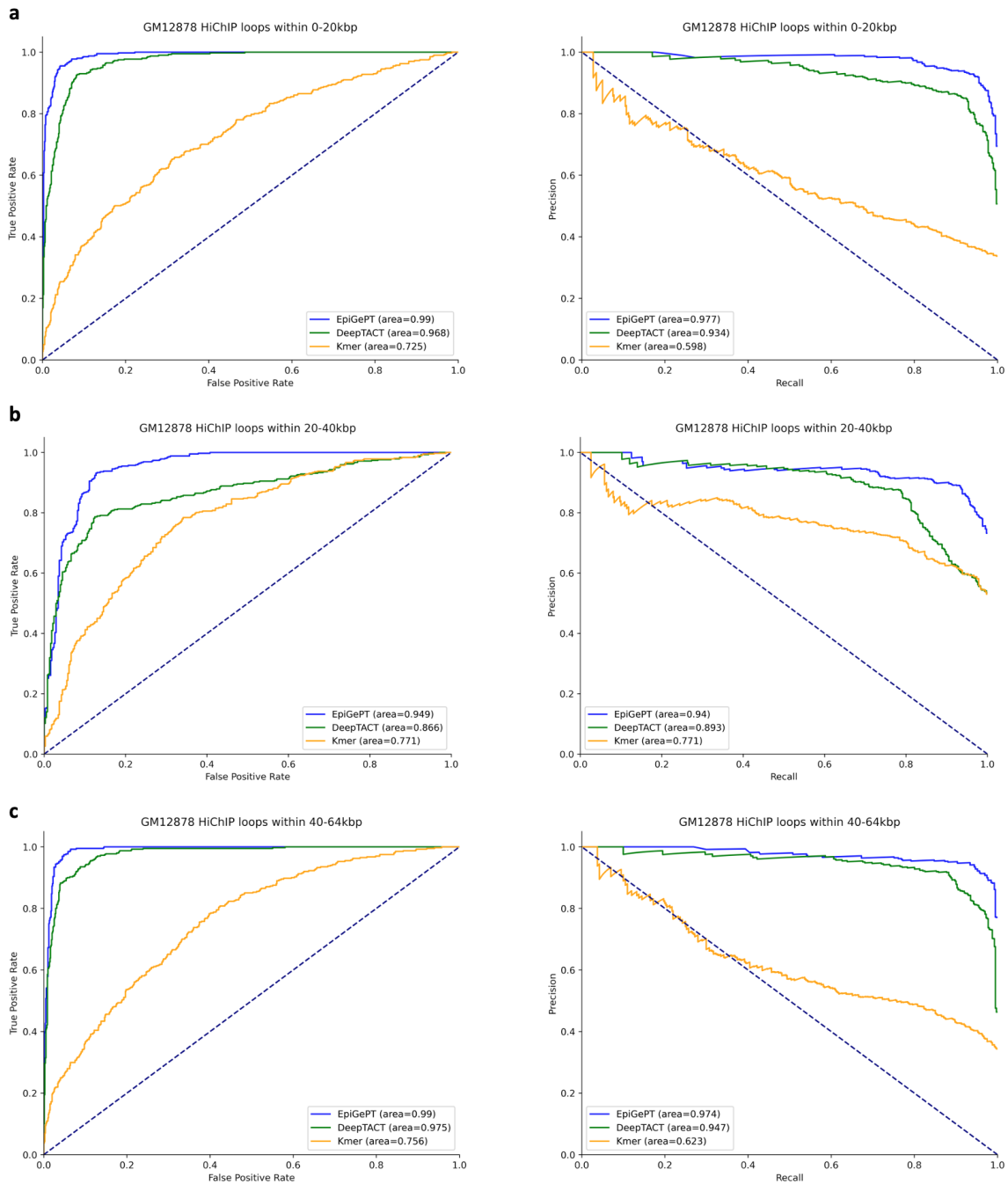
215 **Fig. S8. Incorporating 3D genomic information from HiChip data enhances the**
216 **predictive performance of EpiGePT on E-P regulatory interaction on K562 cell line. a,**
217 The distance distribution between the two anchors of the filtered loops on the K562 cell line.
218 **b,** The performance (auROC) of self-attention scores of EpiGePT and EpiGePT-3D in
219 identifying enhancer-promoter interactions across different distance ranges on the K562 cell
220 type.



222 **Fig. S9. The fine-tuning performance of the EpiGePT model on predicting potential**
 223 **enhancer-promoter regulatory networks. a,** The performance (measured by auROC and
 224 auPRC) of the fine-tuned EpiGePT model and baseline methods (DeepTACT and Kmer) on
 225 HiChIP loops data in distinguishing enhancer-gene pairs at various distance ranges (0-20 kbp,
 226 20-40 kbp and 40-64 kbp). **b,** The performance of the fine-tuned EpiGePT model and baseline
 227 methods on HiChIP loops data in distinguishing enhancer-gene pairs under 1:2 positive-
 228 negative sample ratio on GM12878 cell line. **c,** The performance of the fine-tuned EpiGePT
 229 model and baseline methods on HiChIP loops data in distinguishing enhancer-gene pairs
 230 under 1:2 positive-negative sample ratio on K562 cell line.

231 **Fig. S10**

232



233
234

Fig. S10. The ROC and PR curves of the EpiGePT model on predicting potential

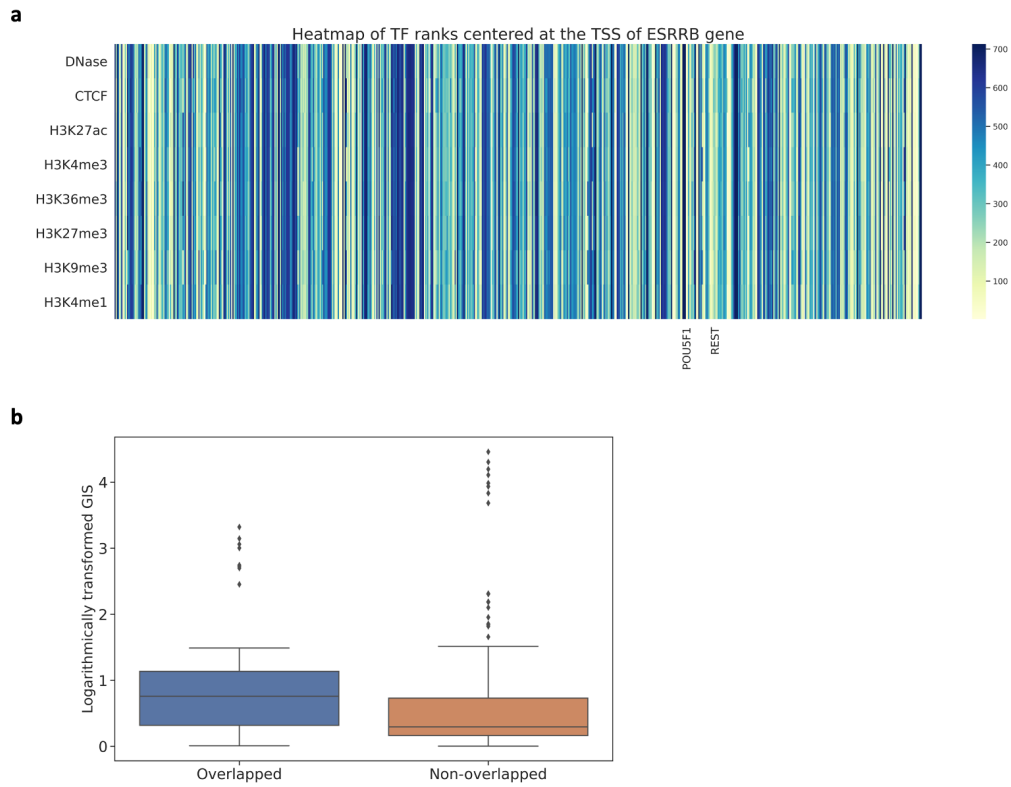
enhancer-promoter regulatory networks. a, The ROC and PR curves of EpiGePT and

baseline methods for predicting HiChIP loops from the GM12878 cell line (0-20 kbp). b, The

ROC and PR curves of EpiGePT and baseline methods for predicting HiChIP loops from the

238 GM12878 cell line (20-40 kbp). **c**, The ROC and PR curves of EpiGePT and baseline methods
239 for predicting HiChIP loops from the GM12878 cell line (40-64 kbp).

240 **Fig. S11**



241

242 **Fig. S11. The GIS of ChIP-seq overlapped bins versus non-overlapped bins of *POU5F1***

243 **centered at the TSS of *ESRRB*.** a, Heatmap of TF ranks across 128 kbp region surrounding

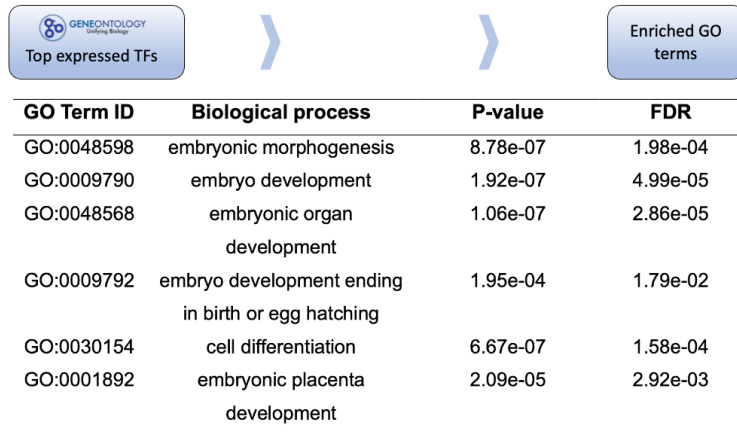
244 the TSS of *ESRRB* gene, each row denotes an epigenomic signal and each column denotes

245 a TF. b, Distribution of non-zero GIS values on overlapped and non-overlapped bins in chip-

246 seq data (ENCF696NWL).

247 **Fig. S12**

248



249

250 **Fig. S12. Gene ontology enrichment analysis based on the top 5% TFs with high**
251 **expression in ESCs.** The results showed lower significance for biological processes
252 associated with embryonic cell development compared with GO terms enriched with the top
253 5% ranked TFs.

1 **EpiGePT** ONLINE PREDICTION ABOUT

About EpiGePT

EpiGePT, a transformer-based model for **cross-cell-line** prediction of chromatin states by taking long DNA sequence and transcription factor profile as inputs. With EpiGePT, we can investigate the problem of how the trans-regulatory factors (e.g., TFs) regulate target gene by interacting with the cis-regulatory elements and further lead to the changes in chromatin states. Given the expression profile of hundreds of TFs from a cellular context, EpiGePT is able to predict the genome-wide chromatin states given the cellular context.

2 Online prediction

Input File

Region file*: C:\example\Example.bed Browse Example

TF Expression file*: C:\example\Example.csv Browse Example

Email address: Optionally provide your email address to receive notification Submit

TaskID: Input taskID to retrieve the task you submitted to our server! Example Retrieve

Multiple-regions Prediction

Chromatin	Start	End
chr1	12039	140039
chr12	100	128100
chr20	6789120	6917120
...

Single-region Prediction

Input File

TF Expression file*: C:\example\Example.csv Browse Example

Specify parameters

Chrom*: chr4 Location*: 19000

Email address: Optionally provide your email address to receive notification Submit

TaskID: Input taskID to retrieve the task you submitted to our server! Example Retrieve

3 Current status and information of the task

Task **2023050416225334** is **under calculating**, please wait for a few minutes. You can also **record the taskID** and then **retrieve** it in the analysis page.

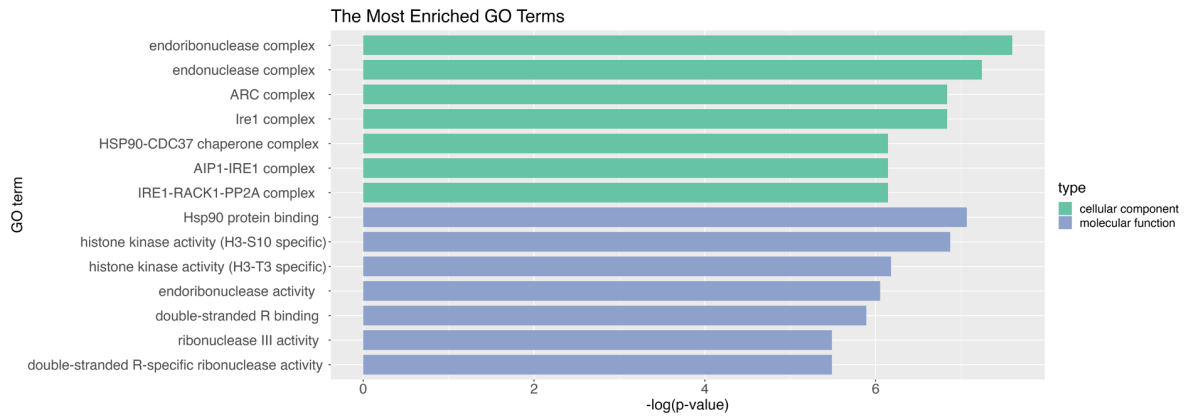
4 Current status and information of the task

Task **2023050416225334** has been **finished!** You can **download** the results directly in this page. You can also **record the taskID** and then **retrieve this task** in the analysis page.

Download

255 **Fig. S13. Case application of the EpiGePT-online.** Users can choose either single locus
 256 annotation or multi-region annotation on EpiGePT-online, and each genomic region requires
 257 a length of 128kbp. Users need to upload the TPM values of transcription factors expression
 258 simultaneously. After annotation, users can enter the result page and download the predicted
 259 files. The predictions are provided at the resolution of 128bp genomic bins, and users can
 260 obtain the predicted signals for these eight epigenomic profiles. Additionally, users have the
 261 option to download the prediction results in CSV format for further analysis and exploration.

262 **Fig. S14**



263

264 **Fig. S14 Enrichment result (Cellular component and Molecular function) of the nearest**

265 **genes of the COVID-19 associated SNPs with the low LOS.**

Supplementary Tables

Table S1. The information of DNase-seq bam file across 129 biosamples from the ENCODE¹⁰ project.

Table S2. The information of RNA-seq tab-separated values (tsv) file across 129 biosamples from the ENCODE¹⁰ project.

Table S3. The information of DNase-seq, CTCF and other six Histone markers bam file across 28 cell lines or tissues from the ENCODE¹⁰ project (hg19).

Table S4. The information of DNase-seq, CTCF and other six Histone markers bam file across 105 cell lines or tissues from the ENCODE¹⁰ project (hg38).

Table S5. The information of RNA-seq tab-separated values (tsv) file across 28 cell lines or tissues from the ENCODE¹⁰ project (hg19).

Table S6. The information of RNA-seq tab-separated values (tsv) file across 105 cell lines or tissues from the ENCODE¹⁰ project (hg38).

Table S7. The preprocessed expression data of 711 human transcription factors from the ENCODE¹⁰ project across 129 biosamples.

Table S8. The preprocessed expression data of 711 human transcription factors from the ENCODE¹⁰ project across 28 cell lines or tissues (hg19).

Table S9. The preprocessed expression data of 711 human transcription factors from the ENCODE¹⁰ project across 105 cell lines or tissues (hg38).

Table S10. The order and names of epigenomes of the expression matrices across 56 epigenomes from the ROADMAP⁶ project.

Table S11. The preprocessed expression data of 642 human transcription factors across 56 epigenomes from the ROADMAP⁶ project.

References

1. Song, L. & Crawford, G.E. DNase-seq: a high-resolution technique for mapping active gene regulatory elements across the genome from mammalian cells. *Cold Spring Harbor Protocols* **2010**, pdb. prot5384 (2010).
2. Nair, S., Kim, D.S., Perricone, J. & Kundaje, A. Integrating regulatory DNA sequence and gene expression to predict genome-wide chromatin accessibility across cellular contexts. *Bioinformatics* **35**, i108-i116 (2019).
3. Liu, Q., Hua, K., Zhang, X., Wong, W.H. & Jiang, R. DeepCAGE: incorporating transcription factors in genome-wide prediction of chromatin accessibility. *Genomics, Proteomics & Bioinformatics* **20**, 496-507 (2022).
4. Avsec, Ž. et al. Effective gene expression prediction from sequence by integrating long-range interactions. *Nature methods* **18**, 1196-1203 (2021).
5. Ernst, J. & Kellis, M. Chromatin-state discovery and genome annotation with ChromHMM. *Nature protocols* **12**, 2478-2492 (2017).
6. Bernstein, B.E. et al. The NIH roadmap epigenomics mapping consortium. *Nature biotechnology* **28**, 1045-1048 (2010).
7. Gasperini, M. et al. A genome-wide framework for mapping gene regulation via cellular genetic screens. *Cell* **176**, 377-390. e319 (2019).
8. Fulco, C.P. et al. Activity-by-contact model of enhancer–promoter regulation from thousands of CRISPR perturbations. *Nature genetics* **51**, 1664-1669 (2019).
9. Zeng, W. et al. SilencerDB: a comprehensive database of silencers. *Nucleic acids research* **49**, D221-D228 (2021).
10. Consortium, E.P. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57 (2012).