

1 **EpiGePT: a Pretrained Transformer model for epigenomics**

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

3 ¹ Ministry of Education Key Laboratory of Bioinformatics, Bioinformatics Division at the
4 Beijing National Research Center for Information Science and Technology, Center for
5 Synthetic and Systems Biology, Department of Automation, Tsinghua University, Beijing
6 100084, China;

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

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

10 * To whom correspondence should be addressed.

11 # The first two authors contributed equally.

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

13 **Abstract**

14 The inherent similarities between natural language and biological sequences have given rise to
15 great interest in adapting the transformer-based large language models (LLMs) underlying
16 recent breakthroughs in natural language processing (references), for applications in genomics.
17 However, current LLMs for genomics suffer from several limitations such as the inability to
18 include chromatin interactions in the training data, and the inability to make prediction in new
19 cellular contexts not represented in the training data. To mitigate these problems, we propose
20 EpiGePT, a transformer-based pretrained language model for predicting context-specific
21 epigenomic signals and chromatin contacts. By taking the context-specific activities of
22 transcription factors (TFs) and 3D genome interactions into consideration, EpiGePT offers
23 wider applicability and deeper biological insights than models trained on DNA sequence only.
24 In a series of experiments, EpiGePT demonstrates superior performance in a diverse set of
25 epigenomic signals prediction tasks when compared to existing methods. In particular, our
26 model enables cross-cell-type prediction of long-range interactions and offers insight on the
27 functional impact of genetic variants under different cellular contexts. These new capabilities
28 will enhance the usefulness of LLM in the study of gene regulatory mechanisms. We provide
29 free online prediction service of EpiGePT through <http://health.tsinghua.edu.cn/epigept/>.

30 **Introduction**

31 A fundamental but largely unresolved problem in genomics is to decode the information
32 residing in the non-coding part of the human genome¹. It remains incompletely understood
33 how regulatory elements govern gene expression in different contexts¹, and how noncoding
34 variants may disrupt the underlying regulatory syntax of DNA². Fortunately, recent advances
35 in epigenome sequencing^{3, 4} have resulted in the accumulation of data useful for the study of
36 these questions, including chromatin accessibility, DNA methylation, histone modifications,
37 and 3D chromatin interaction. Thus, there is great interest in performing systematic analysis of
38 these data to enhance our ability to interpret the non-coding part of the genome⁵⁻¹¹.

39 The inherent similarities between natural language and biological sequences has also stimulated
40 interest in developing large language models (LLM) for the interpretation of genome
41 sequences¹². As is well known, the development of large language model (LLM) has been the
42 main driving force behind many recent breakthroughs in artificial intelligence such as ChatGPT.
43 The architecture of the LLM is typically a multilayer transformer network, and the model is
44 trained on a very large corpus of natural language data. Such pre-trained models can be readily
45 tailored or adapted to various downstream tasks. Considering DNA sequences as the texts in
46 the genomic language, similar transformer-based approaches have been used to model DNA
47 sequences^{13, 14}. For example, the Enformer model¹⁵ takes the DNA sequence of a large genomic
48 region as input and predict thousands of epigenomic features across cellular contexts covered
49 by the training data. Although already useful in many applications, such models relying on only
50 DNA sequences as input are not capable of predicting the function of sequences in new cellular
51 contexts. Furthermore, despite the importance of 3D chromatin contacts in gene regulation, 3D

52 interaction data have not been included in the training of current genomic LLMs. Therefore,
53 there is an urgent need to further develop the core technologies of genomic LLMs to overcome
54 these limitations.

55 In this paper, we present EpiGePT, a transformer-based model for epigenomics prediction with
56 the following new capabilities. First, the inability to make predictions in novel contexts has
57 greatly limited the applicability of current methods, EpiGePT removes this limitation by
58 making both the input and output context-dependent, where the context is represented by a TF-
59 profile vector specifying the expression of key TFs in that context. This choice is motivated by
60 the fact that reference gene expression data are available for many cellular contexts that are
61 important in development and diseases, but for which few epigenomic features have been
62 measured. We note that the reference TF expression profile has been used to represent cellular
63 context in earlier works on accessibility prediction^{6, 16}, but this idea has not been explored for
64 the development of genomic LLMs. Second, a new learning algorithm is developed to enable
65 the inclusion of 3D chromatin contact data in the training data. In this way, EpiGePT can
66 predict 3D genome features such as enhancer-promoter interactions that are known to be
67 important for gene regulation but are not modeled in current genomic LLMs. By using a masked
68 training strategy, EpiGePT can be trained on a diverse set of contexts even if different sets of
69 epigenomics signals are available in different contexts. There is a profound difference in
70 training strategy between EpiGePT and current genomic LLMs. Each input genomic region
71 provides an example for training in current LLMs such as the Enformer. In contrast, each
72 combination of input region and cellular context provides an example for training in EpiGePT,
73 thus providing a much larger number of examples available for model training. As for training

74 data sets, since most cellular contexts that have epigenomic data will also have expression data,

75 we can use most available epigenomic data, such as those used by the Enformer, to train our

76 model.

77 In a series of experiments, we illustrate that our model is superior to existing methods in

78 epigenomic signals prediction, long-range chromatin interaction prediction, as well as the

79 variant effect prediction.

80

81 **Results**

82 **Overview of EpiGePT**

83 EpiGePT is a genomic language model for cross-cell-type prediction of chromatin states by
84 multi-task learning based on genome-wide pre-training on epigenomic data (Fig. 1 and Fig. S2).
85 The model is composed of four modules, including a sequence module, a TF module, a
86 transformer module, and a prediction module. The sequence module is responsible for
87 processing the long DNA sequence of interest (e.g., 128 kb) by employing a series of
88 convolutional and pooling blocks (e.g., 5) to extract a comprehensive set of sequence features.
89 The TF module is specifically designed to represent a cellular context by a TF-profile vector,
90 which specifies the state of a few hundred TFs in that context. The features computed by the
91 sequence and TF modules are then fed as input tokens to the transformer module, where each
92 token corresponds to a genomic bin (e.g., a 128 bp window) in the original DNA sequence. The
93 transformer module leverages self-attention mechanisms to learn the relationships among the
94 input bins, enabling the model to make predictions of multiple chromatin states given the
95 context information from the TF module. Importantly, by including a novel loss term that
96 involves the self-attention weights, EpiGePT is capable of learning from data on context-
97 specific chromatin interactions. Since 3D interaction is known to be a key mechanism in gene
98 regulation, the ability to learn from interaction data is an attractive feature of our approach.
99 Finally, the fourth module in EpiGePT is a predictive module which predicts epigenomic
100 signals and chromatin interactions based on the output of the transformer module.

101 **Genome-wide prediction of epigenomic signals**

102 To assess the performance on predicting epigenomic signals, we first compared EpiGePT to
103 task-specific models that are specifically designed for predicting a single epigenomic signal.
104 Taking the chromatin accessibility for instance, the performance of EpiGePT was compared
105 against existing task-specific models such as BIRD¹⁷, ChromDragoNN⁶, and DeepCAGE¹⁶.
106 The widely available public DNase-seq¹⁸ data across 129 cellular contexts on 1,175,374
107 genomic regions were collected and preprocessed from ENCODE database¹⁹ (see Methods).
108 Performance is evaluated in three prediction settings: i) “cross-region” setting where the
109 predictive model is tested on new genomic regions not seen in training, ii) “cross-cell type”
110 setting where the model is tested on new cell types, and iii) “cross-both” setting where testing
111 is done on new regions in new cell types (Fig. S1, Supplementary Text S1). In each setting, we
112 employed three evaluation metrics, namely Pearson correlation coefficient (PCC), Spearman
113 correlation coefficient (SCC) and prediction square error (PSE), to assess the similarity between
114 the predicted and true values of the DNase signals (See Methods). The results, presented in Fig.
115 2a and Fig. S3, showed that EpiGePT consistently outperformed baseline methods including
116 BIRD¹⁷, and ChromDragoNN⁶ by a relatively large margin under the above settings. For
117 example, EpiGePT achieved a cross-cell type prediction PCC of 0.787, demonstrated a 6.9%
118 higher performance than the best baseline method, ChromDragoNN. In addition, we also
119 evaluated the prediction of binary chromatin accessibility status i.e. predicting whether a peak
120 exists within the corresponding genomic bin (>50% overlap). For binary prediction, EpiGePT
121 again achieved a superior performance with an average auPRC (area under the precision-recall
122 curve) of 0.767 compared to 0.623 of DeepCAGE¹⁶ and 0.476 of ChromDragoNN⁶ (Fig. 2c).
123 Finally, we compared EpiGePT and ChromDragoNN⁶ in the binary classification of functional

124 regions versus nonfunctional regions, using the functional chromatin status derived from
125 ChromHMM²⁰ annotations as ground truth (Supplementary Text S6). EpiGePT achieved an
126 average 8.1% higher auROC (area under the receiver operating characteristic curve) than
127 ChromDragoNN⁶, and an average 2.3% higher macro-auROC than ChromDragoNN⁶ (p -value
128 < 0.001 under one-sided Wilcoxon signed rank test) in a finer-grained classification for different
129 types of regulatory elements (Fig. S4). These results demonstrate that EpiGePT provides better
130 predictions than task-specific models.

131 Next, we compared EpiGePT with a state-of-the-art genomic LLM, Enformer¹⁵, in two different
132 ways. First, we trained an Enformer model from scratch with only the aforementioned DNase-
133 seq data (Supplementary Text S5). EpiGePT demonstrates a 3.3% to 5.2% higher performance
134 than Enformer in terms of the median Pearson correlation coefficient under the three prediction
135 settings (Fig. 2b). Second, we compared EpiGePT directly to the pretrained Enformer model
136 provided by the original paper. To do this, we collected eight different epigenomic signals from
137 104 different cellular contexts (Supplementary Table S4, S6 and S9). We first left out 13 of
138 these contexts where HiChIP data are also available for downstream chromatin interactions
139 validation. Then, EpiGePT model was trained across 72 training cellular contexts (without
140 using HiChIP-based chromatin contacts data in the training) and subsequently compared
141 against pre-trained Enformer on the remaining 19 test cellular contexts, on 15,870 training
142 genomic regions with 128kbp length. Since most of the cellular contexts have missing
143 epigenomic signals, we designed a masked training strategy to handle this issue (See Methods).
144 Under the test cellular contexts, EpiGePT exhibited superior performance with higher PCC than
145 Enformer in 60 out of 78 matched epigenomic signals across 19 test cellular contexts by

146 achieving an average PCC of 0.510, compared to 0.440 of Enformer (Fig. 2d and Fig. S6b). For
147 DNase-seq specifically, the average PCC of EpiGePT reached 0.710 and the average SCC
148 reached 0.664 across 7 cell types, compared to the average PCC of 0.455 and the average SCC
149 of 0.488 of Enformer. In the above comparison, we are in fact comparing out-sample prediction
150 by EpiGePT to in-sample prediction by Enformer. The favorable results achieved by EpiGePT
151 in this experimental setting suggests that our model enables prediction in novel contexts without
152 sacrificing performance. To illustrate the prediction performance further, several tracks of
153 predicted chromatin states and the corresponding ground truth chromatin states were displayed
154 in Fig.2e.

155 **EpiGePT enables the prediction of chromatin interactions**

156 We examined the capacity of EpiGePT for predicting long-range chromatin interactions, which
157 is important for understanding chromatin architecture and relations between regulatory
158 elements and target genes. We employed several experimental settings to examine the ability
159 of EpiGePT in capturing long-range chromatin interactions. In setting (A), we directly utilized
160 the self-attention weights extracted from the pretrained EpiGePT model (without including
161 HiChIP data in the training) to predict enhancer-promoter (E-P) interactions and silencer-
162 promoter (S-P) interactions. In setting (B), we integrated HiChIP-derived 3D chromatin
163 contacts into the training of the model and then use the model to predict E-P interactions in
164 novel contexts not seen in the training. In setting (C), we designed a pretrain-finetune strategy
165 for EpiGePT model to predict E-P interactions. The results under each setting are discussed
166 below.

167 Setting (A): prediction by EpiGePT not trained with 3D contact data. In this setting, we use the
168 cell-type specific self-attention scores to predict chromatin interactions, including E-P and S-P
169 interactions (see Methods). Two sets of interactions containing 664 and 5,091 candidate
170 element-gene interactions, obtained by CRISPRi²¹ experiments on K562 cell line, were
171 collected and further filtered and divided into positive and negative samples, for use as ground
172 truths to evaluate E-P prediction performance. In the Gasperini et al²². dataset, EpiGePT
173 consistently outperformed Enformer by achieving the highest auPRC in most cases (Fig. 3a).
174 For instance, EpiGePT achieved auPRC of 0.647 to 0.887 for identifying enhancer-gene
175 transcription start site (TSS) pairs in different distance groups (Fig. 3a and Fig. S7). In the Fulco
176 et al.²³ dataset, EpiGePT also outperformed other competing methods. For example, EpiGePT
177 achieves an auPRC of 0.504, compared to 0.307 of Enformer in the 30-45kbp group (Fig. 3a).
178 Next, to assess performance on S-P interactions., we downloaded putative silencers from the
179 SilencerDB²⁴ and used the TSS of annotated nearest gene as the potential target. We selected
180 the same number of negative pairs randomly while conserving the distance distribution. The
181 results show that EpiGePT achieved a better performance in distinguishing positive S-P pairs
182 from negative pairs than Enformer. For instance, EpiGePT achieves an auROC of 0.575 in long-
183 range S-P interactions (32-64kbp) compared to 0.483 of Enformer (Fig. 3b). Finally, to assess
184 performance in predicting chromatin interactions, we collected HiChIP²⁵ loops on K562 and
185 GM12878 cell lines from the HiChIPdb²⁶. EpiGePT achieves a superior performance by
186 discerning HiChIP loops from randomly selected loops with the same distance distribution. For
187 instance, EpiGePT achieves an auPRC of 0.520 for long range loops (40-64kbp) prediction in
188 GM12878 cell line, surpassing that of Enformer (0.484) by a large margin (Fig. 3g). These

189 results clearly demonstrated the utility of EpiGePT attention scores in capturing functional
190 chromatin interactions.

191 To better understand the self-attention mechanism of EpiGePT, we showed the attention
192 weights (averaged across heads) for the bin containing the TSS of the gene *CHD4*. The attention
193 weights were computed based on the pretrained EpiGePT model with K562 cell line as the
194 context of interest. We also display chromatin interactions detected under K562 as well as
195 regulatory elements annotations from the GeneHancer²⁷, It is seen that both the interaction data
196 and the regulatory element annotations are consistent with the attention weights learned by
197 EpiGePT (Fig. 3c and Fig. 3f).

198 Setting (B): Prediction by EpiGePT-3D, which include Hi-C data in its training. The above
199 results suggest that in a good transformer-based genomic language model, the attention weight
200 given by one bin to another bin (within the input region) should be consistent with the strength
201 of 3D interaction between them. Thus, when experimental data on 3D interaction are available,
202 we can leverage this data to improve the learning of the parameters of our genomic language
203 model, by penalizing parameter values that resulted in poor correlation between the attention
204 weights and the interaction data (see Methods). To obtain such training data, we collected
205 4,107,687 H3K27ac-based HiChIP loops across 13 cell lines or tissues from HiChIPdb²⁶, which
206 denote potential E-P interactions. Setting aside loops from K562 cell line as test data, other
207 HiChIP loops are incorporated into the training. The resulting model is denoted as EpiGePT-
208 3D. We found that adding 3D interaction data in the training can lead to a noticeable
209 improvement for cross-cell-type prediction (3.3% higher in PCC) (Fig. 3e). Moreover,
210 EpiGePT-3D demonstrated improved predictive performance on E-P interactions identified by

211 HiChIP loops in new cellular contexts. For instance, the auPRC increased from 0.652 to 0.695
212 for Gasperini et al.'s dataset, which is on a context not covered by the Hi-C data in the training,
213 in 24-40kbp group when incorporating 3D genome data.

214 Setting (C): Prediction by fine-tuning pretrained EpiGePT. Fine-tuning is an strategy that
215 transfers the knowledge of a pretrained model to new tasks, which is particularly prevalent in
216 language models such as GPT²⁸ and BERT²⁹. Here, we explore the performance of fine-tuning
217 given a pretrained EpiGePT model on downstream tasks, such as predicting 3D genome
218 interaction. Specifically, we fixed the weights of a pretrained EpiGePT model and added an
219 additional finetune network for predicting E-P interactions. We compared EpiGePT with
220 finetuning strategy (EpiGePT-finetune) to two baselines, DeepTACT³⁰ and a k-mer frequency
221 based method²⁹ with HiChIP H3K27ac loops from K562 and GM12878 cell lines (see Methods).
222 The results illustrate that EpiGePT-finetune exhibited a superior classification performance
223 across diverse distance ranges compared to baselines. For example, EpiGePT-finetune achieved
224 an auROC of 0.949, surpassing 0.866 of DeepTACT³⁰ and 0.771 of Kmer by a large margin in
225 the GM12878 cell line within the 20-40kbp distance range (Fig. 3h, Fig. S9 and Fig. S10). This
226 significant improvement demonstrates the power of fine-tuning a base pretrained genomic
227 language model on a downstream task with limited data.

228 **EpiGePT unveils the regulatory relationships between TFs and target genes**

229 In this section, we further explored the TF module to see whether EpiGePT is able to learn the
230 regulatory relationships between TFs and target genes (TGs). We defined gradient importance
231 scores (GIS) based on the absolute gradient values of predicted epigenomic signals with respect

232 to the expression of a TF in the input TF profile, to rank the TFs for their potential to regulate
233 a given TG (see Methods). Particularly, we use the TF profile of embryonic stem cell (ESC) to
234 specify the context in the EpiGePT model. We selected the important ESC regulator *POU5F1*
235 as the target gene and calculated the GIS for identifying TF-TF interactions (see Methods).
236 Multiple potential regulators for *POU5F1* identified by EpiGePT in ESC context are consistent
237 with literatures, such as *ESRRB-POU5F1*³¹ (rank 2nd), and *ETV5-POU5F1*³² (rank 5th). Next,
238 we focus on *ESRRB* which plays essential role for balancing pluripotency of ESCs³³. Treating
239 *ESRRB* as the target gene, our GIS-based ranking identified several key TFs, such as *POU5F1*
240 and *REST*, that have significantly higher ranks than other TFs (Fig. 4a). By using ChIP-seq data
241 of *POU5F1* for validation, we observed significantly higher GIS in bins overlapping with the
242 ChIP-seq data (Fig. S11, p -value < 0.00018 under one-sided Mann-Whitney U test). Next, we
243 visualized the TF ranks obtained from eight epigenomic profiles across 1000 bins surrounding
244 the TSS of *ESRRB*. By averaging ranks across these signals and bins among all the 711 TFs,
245 the important ESC regulator *POU5F1* ranks 3 out of 711 (Fig. 4b). We further collected the top
246 5% of TFs for each bin and conducted gene ontology (GO) enrichment analysis based on these
247 TF coding genes. Interestingly, the GO terms enriched also included biological processes of
248 embryonic cell differentiation and development. However, using the top 5% of TFs with high
249 expression in ESCs resulted in lower significance for biological processes associated with
250 embryonic cell development (Fig. 4c and Fig. S12), which again demonstrates the effectiveness
251 of the GIS-based ranking. Furthermore, we use TF-TG relationships from either ChIP-seq data
252 or external databases as ground truth to validate the TF-TG relationships inferred by EpiGePT.
253 We defined potential TF-target gene pairs based on TF ChIP-seq data specific to certain cell

254 types among all human genes (see Methods). The results demonstrated a significant higher rank
255 of TF-target gene pairs, compared to TF-non-target gene pairs based on the integrated GIS-
256 based ranking (Fig. 4d, p -value < 0.001 under one-sided Mann-Whitney U test). Second, we
257 collected TF-TG regulatory network data from two publicly available databases. We obtained
258 a total of 1,066 TF-TG pairs from the GRNdb³⁴ database based on liver-specific GTEx data,
259 and 2,705 TF-TG pairs from the TRRUST³⁵ database after filtering. Then we calculated the
260 rank of each TF based on either integrated GIS or the TF expression value by using the liver
261 expression as the TF reference profile. Interestingly, we found that the median ranking
262 percentile of TFs from TRRUST was 3.1%, significantly higher than the percentile of 20.4%
263 based on expression values (Fig. 4e, p -value $< 1e-5$ under one-sided Wilcoxon signed rank test).
264 with a similar result was obtained using another database GRNdb, where EpiGePT is seen to
265 achieve a median ranking percentile of 6.3%, compared to 36.0% by gene expression value.
266 For instance, *TMEM55B*, which plays a significant role in lysosome movement, and is regulated
267 by sterol response element binding factor 2 (*SREBF2*)³⁶. Consistently, GIS ranking identified
268 *SREBF2* as the top-ranked TF associated with *TMEM55B*. Overall, the validation results from
269 both ChIP-seq datasets and external databases support the effectiveness of GIS in identifying
270 context-specific TF-TG relationships.

271 **EpiGePT improves variant effect prediction**

272 Context-specific prediction of the functional impact of genetic variants is important for genetic
273 studies. To test the utility of EpiGePT in this task, we first collected an eQTLs dataset³⁷ that
274 contains 20,913 causal and non-causal variant-gene pairs across 49 different tissues from the
275 supplementary data of Wang et al³⁷. EpiGePT, EpiGePT-seq (i.e. EpiGePT without the TF

276 module) and Enformer were then applied to estimate the context-specific log-ratio scores (LOS)
277 between the alternative DNA sequence and the reference DNA sequence, (see Methods, Fig.
278 5a). Finally, a random forest classifier is trained based on these LOS's to distinguish causal
279 variant-gene pairs from non-causal pairs. The experimental results show that better prediction
280 performance can be achieved when the LOS is based on EpiGePT than when the LOS is based
281 on Enformer. For example, in the lung tissue, EpiGePT achieved an auPRC of 0.922, compared
282 to 0.873 of Enformer, for the classification of casual SNPs vs non-causal SNPs. To verify the
283 effectiveness of TF module, we replace the TF reference profile of lung with a less relevant cell
284 type, stomach, and the auPRC decreases from 0.922 to 0.892 (Fig. 5b). Similar results were
285 seen for other tissue contexts—across 48 tissues, EpiGePT-seq achieved an average auPRC of
286 0.910, compared to 0.898 of Enformer (Fig. S4d). The above experiments demonstrated the
287 usefulness of EpiGePT in assessing variant effects.

288 To further evaluate the performance of EpiGePT in predicting disease-associated variants, we
289 extracted 52, 876 pathogenic SNPs from the ClinVar³⁸ database and 418, 863 benign SNPs from
290 the ClinVar database, also with 84, 095 benign SNPs from the ExAC database³⁹ as positive and
291 negative sets, respectively. We defined a 128kbp region surrounding each pathogenic SNP as
292 the risk region. We extracted all benign or likely benign SNPs that fall within the risk region as
293 the positive samples. As the relevant tissue or cell type information is not available, we
294 concatenated the LOS of the eight epigenomic signals and also the self-attention scores, across
295 multiple cellular contexts, and then evaluated whether the constructed features are beneficial in
296 distinguishing pathogenic SNPs from benign ones in a classification analysis. To achieve this,
297 we augmented the popular CADD-derived features (CADD⁴⁰ scores) by concatenating them

298 to the EpiGePT-derived features discussed in the above, to obtain a comprehensive feature
299 vector (see Methods). Subsequently, we compared the performance of the multi-layer
300 perceptron (MLP) classifier based on the comprehensive feature vector to that based on CADD-
301 derived features alone. The results demonstrated that incorporating EpiGePT-derived features
302 significantly enhance the performance in predicting pathogenic SNPs. Specifically, when the
303 positive-to-negative sample ratio was set to be 1:1, the average auROC increased from 0.772
304 to 0.806, and the average accuracy increased from 0.690 to 0.723 (Fig. 5c). This observation
305 indicates that features extracted by EpiGePT provide a valuable complement to CADD scores,
306 enabling a more comprehensive interpretation of disease-associated variants.

307 **EpiGePT prioritizes potential SNPs associated with comorbidities of COVID-19**

308 We investigated whether using EpiGePT to predict variant effects could help in the discovery
309 of key SNPs related to COVID-19. COVID-19 is an infectious disease caused by the SARS-
310 CoV-2 virus, which emerged in late 2019 and quickly spread around the world, causing a global
311 pandemic⁴¹. In order to validate the ability of EpiGePT in identifying key SNPs, we collected
312 GWAS data from a COVID-19 genetic study⁴², including 9,484 variants derived from 4,933
313 patients with confirmed severe respiratory symptoms and 1,398,672 control individuals without
314 COVID-19 symptoms. To validate the ability of the model to identify COVID-19-associated
315 SNPs, we firstly defined a risk region around the selected COVID-19-associated SNPs and
316 computed the rank of the variant score of pathogenic SNPs within the surrounding benign SNPs
317 from the ClinVar database. Note that the expected percentile rank for random guessing (uniform
318 distribution) is 0.5 (see Methods). Previous studies^{43, 44} suggested that COVID-19 infection
319 could potentially impair the function of the heart or the lungs, leading to congestive heart failure

320 or decreased lung function. Interestingly, we found that the average rank of COVID-19-
321 associated SNPs was 0.250 when lung expression data was employed for the TF reference
322 profile and a 6-kbp risk region was examined (Fig. 5d, p -value < 0.05 under one-sided Binomial
323 exact test). However, when we employed the expression data from less relevant contexts, such
324 as K562 cells or Testis cells, the median rank is close to random guessing (i.e. 0.5), indicating
325 its ineffectiveness in discerning SNPs pertinent to COVID-19. These results suggest that
326 EpiGePT model is able to prioritize the COVID-19-associated SNPs thus shedding lights on
327 finding the potential disease-associated variants and the relevant tissue contexts with our
328 pretrained large language model.

329 Next, we examine whether the genes close to max-LOS SNPs are likely be associated with
330 biological processes and functions relevant to COVID-19, when compared with genes close to
331 low scores SNPs or not closed to associated SNPs. Since the genetic pathology of COVID-19
332 is not yet clear and the earliest lesion is in the lungs, we ranked all 9,484 possible SNPs using
333 lung expression data as the TF reference profile. We then identified the SNPs with the highest
334 ranks and performed GO enrichment analysis on nearest genes of the top-30 scored SNPs (Fig.
335 5e). The enrichment results revealed potential biological processes that are relevant to COVID-
336 19, such as the regulation of glucokinase activity which is associated with the homeostasis of
337 human blood glucose⁴⁵. Notably, diabetes mellitus, a condition closely associated with
338 hyperglycemia, is a typical comorbidity of COVID-19⁴⁶. However, GO enrichment analysis
339 based on the nearest genes of the lowest-scored 30 SNPs resulted in enrichment outcomes that
340 were less relevant to COVID-19 or its complications (Supplementary Fig. S14). Among the
341 potential genes around the top-10 scored SNPs, we identified that the *TBC1D4* gene, which

342 regulates glucose homeostasis, is potentially associated with COVID-19 comorbidities. Our
343 findings are consistent with previous research by Pellegrina et al.⁴⁷ and highlights the potential
344 of our EpiGePT approach in discovering new genetic markers that may be implicated in the
345 pathogenesis of COVID-19. Overall, our EpiGePT model provides new perspectives for
346 understanding how the genetic variants could contribute to the COVID-19 susceptibility and
347 severity.

348 **Model ablation analysis**

349 To verify the roles of the main modules in the model, we conducted ablation experiments on
350 the model architecture (Fig. S5). For TF module ablation, the results compared to EpiGePT
351 without TF module (EpiGePT-seq) and the inclusion of the TF module led to improvement in
352 cross-cell-type prediction of DNase signals, with a median PCC of 0.787 of EpiGePT,
353 compared to 0.74 for EpiGePT-seq. We additionally examined the impact of the TF module by
354 employing three methods, namely replacing TF scores with zero, replacing TF scores with
355 random noise, and removing motif binding scores. The results again confirmed the positive
356 impact of the TF module (Fig. S5a). For sequence module ablation, we trained a TF-only model
357 without the sequence module. The results indicated that removing the sequence module resulted
358 in an average decrease of 0.084 in the PCCs of the epigenomic signals on a cell-type wise basis
359 (Fig. S5a). For multi-task module ablation, we trained eight separate predictive models for each
360 of the eight epigenomic signals. In the case of the H3K4me1 signal prediction, the performance
361 of the single-task prediction model exhibited an average PCC decrease from 0.408 to 0.329
362 compared to the multi-task prediction model. Similarly, the overall prediction performance for
363 the eight signals declined by 0.074 (Fig. S5b). This decrease may be attributed to the intricate

364 nature of gene regulation that multiple epigenomic signals can synergize with each other,
365 allowing their joint modeling to gain deeper biological insights.

366 **Online prediction tool for EpiGePT**

367 In order to facilitate the utilization of EpiGePT for the prediction of multiple chromatin states
368 of any cellular context and genomic regions, we have developed a user-friendly web server,
369 named EpiGePT-online (<http://health.tsinghua.edu.cn/epigept/>) (Supplementary Text S2). The
370 web server was developed using PHP, JavaScript and HTML, which provides an interactive
371 web interface for efficiently online prediction of epigenomic profiles (Fig. 6). The web server
372 includes a built-in kernel that encompasses the framework for data preprocessing, TF motif
373 binding scores calculating, and prediction of epigenomic signals for both hg19 and hg38 human
374 reference genome. Users can obtain the predicted signals for multiple genomic regions by
375 submitting a region file and a TF expression file in Numpy or CSV formats (Supplementary
376 Table S5), or predicted signals for a specific region by submitting a TF expression file (Fig.
377 S13). We provided TF expression profile across more than 100 cellular contexts from ENCODE
378 on the download page. Users can download the results in csv format for further applications
379 such as genetics analysis. Furthermore, we provide a case application of the EpiGePT-online to
380 enable users to quickly learn how to use our website (Supplementary Text S3). We anticipate
381 that this web server will assist researchers in deepening their understanding of gene regulatory
382 mechanisms.

383 **Discussion**

384 In this paper, we introduced a pretrained transformer-based language model for epigenomics.

385 Compared with the existing task-specific models and sequence-based language model,
386 EpiGePT has the added capability to make predictions on novel contexts. Furthermore,
387 EpiGePT is able to incorporate a new type of data (3D genome interaction data) during model
388 training, which enables the identifying functional regulatory interactions such as enhancer-
389 promoter interactions. EpiGePT demonstrates state-of-art performance in diverse experimental
390 settings compared to existing methods. Based on the predicted epigenomic features and 3D
391 interactions from EpiGePT, we performed two investigations on how information is encoded
392 in the human genome sequence: First we identify the interactions of cis-regulatory elements
393 and their target genes with the help of self-attention mechanism in EpiGePT. Through direct
394 utilization of self-attention scores, model fine-tuning, and leveraging 3D genome interactions,
395 we validated the capacity of EpiGePT to capture regulatory interactions. Second, to assist the
396 identification and interpretation of human disease-associated SNPs, we estimate the effect of a
397 variant on the epigenomic features around the variant, based on the LOS computed by the
398 outputs of EpiGePT. Such variant effect prediction by EpiGePT establishes a foundation for
399 understanding the underlying relationship between genetic variations and disease mechanisms.

400 There exist several extensions and refinements that can be applied to further improve the
401 EpiGePT model. Firstly, the incorporation of chromatin regulators (CRs) as trans-acting factors
402 into the TF module could enhance the modeling of regulated transcription processes, thereby
403 increasing the accuracy of the predictions. Second, the integration of DNA methylation
404 information⁴⁸ while modeling DNA sequences allows for a more comprehensive and accurate
405 decoding of the epigenomic language, providing a more comprehensive model of gene
406 regulation states compared to the analysis solely based on DNA sequences. Third, the rapid

407 advancements in sequencing technologies have enabled the accumulation of vast amounts of
408 multi-omics data, encompassing different scales from biomolecules to single cells, tissues, and
409 organs⁴⁹. The integration of multiscale and multi-omics information is a trend and a major
410 challenge in deciphering gene regulatory landscapes. Integrating single-cell level data into
411 EpiGePT is an important direction for future improvement. For example, utilizing clustered
412 single-cell multi-omics data as pseudo-bulk data can further expand the training context of
413 EpiGePT. The application of EpiGePT to single-cell epigenomics could enable the profiling of
414 chromatin signals at single-cell resolution, facilitating a holistic understanding of regulatory
415 heterogeneity in different cell subpopulations.

416 Based on EpiGePT, users are able to predict multiple chromatin profiles in different cell lines
417 or tissues, which could provide a foundation for biological discovery, decoding transcriptional
418 regulation mechanisms, and investigating disease mechanisms. We anticipate EpiGePT will
419 furnish researchers with valuable insights into understanding regulatory mechanisms.

420 **Methods**

421 **Data processing**

422 **Chromatin accessibility data and Expression data** We used three different datasets in the
423 experiments. For chromatin accessible data, we downloaded DNase bam files and narrow peaks
424 across 129 human biosamples from ENCODE¹⁹ project (Supplementary table S1 and S2). We
425 divided the human hg19 genome into 200bp non-overlapping bins, and we assigned the label
426 for each bin in each cell type. For the regression design, we pooled the bam files of multiple
427 replicates for a cell type (Supplementary table S1 and S2), and obtain the raw read count n_{lk}
428 for bin l in cell type k . We normalized the raw read count in order to eliminate the effect of
429 sequencing depths, in the form of $\tilde{n}_{lk} = Nn_{lk}/N_k$, where N_k denotes the total number of
430 pooled reads for cell type k and $N = \min_k N_k$ denotes the minimal number of pooled reads
431 across all cell types. The normalized read counts are further log transformed with pseudo count
432 1, which represent the continuous level of chromatin accessibility. For binary classification
433 design, we assigned a binary label y_{lk} to 1 if the number of raw read counts of the bin l in the
434 cell type k greater than 30, which represent the bin is an accessible region in this cell type,
435 resulting in the identification of regions as accessible in 13% on average and 8% at median in
436 the screened genomic regions across 129 cell types. The proportion of open regions varies
437 among different cell types, and the average openness level mentioned above is generally
438 consistent with that maintained in ChromDragoNN⁶.

439 RNA-seq data of the 711 human transcription factors were downloaded and extracted from the
440 ENCODE project (Supplementary table S5 and S6). We perform log transformation with

441 pseudo count 1 and quantile normalization based on TPM values. The normalized TPM values
442 were averaged across replicates and mean expression profile after normalization of each cell
443 type was finally used to calculated of the transcription feature.

444 **Multiple chromatin signals data** For the human reference genome hg19 (GRCh37), DNase-
445 seq, RNA-seq and ChIP-seq data were also downloaded from ENCODE project
446 (Supplementary table S3, S4 and S6). We applied the same process to these data as above, and
447 finally we obtained the 8 epigenomic signals of 13,300,000 bins of 128bp in 28 cell types. The
448 continuous level of chromatin signals we extracted were 'DNase', 'CTCF', 'H3K27ac',
449 'H3K4me3', 'H3K36me3', 'H3K27me3', 'H3K9me3' and 'H3K4me1', which includes crucial
450 epigenetic modifications and markers for gene regulation and transcription.

451 For the collected the data of human reference genome hg38 (GRCh38), we adopted a data
452 collection strategy that includes missing data. Specifically, within a particular tissue or cell type,
453 we ensured the presence of at least one ChIP-seq signal. Then, epigenomic profiles of 8 signals
454 for 15,870,000 bins of 128bp across 104 cell types were obtained.

455 **Model architecture**

456 **Sequence module** As shown in Fig. 1 and Fig. S2a, the sequence module receives a one-hot
457 matrix ($A = [0,0,0,1]$, $C = [0,1,0,0]$, $G = [0,0,1,0]$, $T = [0,0,0,1]$) of size (128000,4) as input,
458 representing a sequence of 128 kilobase pairs (kbps) and contains five 1-dimensional
459 convolutional blocks to extract DNA sequence features. Each block includes a convolutional
460 layer and a maxpooling layer (Fig. S2b). The first convolutional layer considers the input
461 channels as 4 and performs convolution along the sequence direction. The input sequence

462 features are one-hot embeddings of size $L \times 4$, where L denotes the length of the input long
463 range DNA sequence. After 5 maxpooling layers, the output size of sequence feature is
464 $L/N \times C$, where C denotes the hyper-parameter for sequence embedding and N denotes the
465 length of locus to predict. We set C to 256 in the pre-training stage of chromatin accessibility
466 prediction experiments. Rectified linear units (ReLU) are used after each convolution operation
467 for keeping positive activations and setting negative activation values to zeros. By reducing the
468 input length by 128 times through pooling operations, this module effectively compresses the
469 input information while retaining essential features. Sequence features were then concatenated
470 with TF expression features, and we finally obtained a vector of size $L/N \times (C + n_{TF})$, where
471 n_{TF} denotes the dimension of the transcription factors features after padding. In our model,
472 after adding padding to the 711 TFs, the n_{TF} is set to 712. Therefore, the input token number
473 for the transformer module is 1000, and each token embedding has a dimensionality of 968.

474 **Transformer module** We utilize the transformer module to integrate information from both
475 the sequence and transcription factors (TFs), enabling the capturing of long-range interactions
476 between genomic bins. We applied N_t layers of Transformer encoder with n_h different
477 attention heads to the token embedding sequence. The input word embedding (X) of the
478 transformer encoder is a genomic bin sequence with dimensions
479 (*Sequence length, embedding dim*). Specifically, this dimension is (1000, 968) in EpiGePT,
480 indicating that input genomic bin sequence has a length of 1000, and each genomic bin has an
481 embedded representation that combines the sequence information with cell-type-specific
482 features with dimension of 968. For position embedding, we employed absolute position
483 embedding to represent the positional information of the 1000 genomic bins in the input 128kbp

484 DNA sequence, with dimensions of (1000, 968). Each Transformer encoder includes a multi-
485 head self-attention mechanism and a feed-forward neural network. For self-attention in each
486 head, the calculation is based on the matrix operation.

$$487 \quad \text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

488 For multi-head attention, Transformer encoder learns parameter matrices $W_i^Q \in$
489 $\mathbb{R}^{d_{model} \times d_K}$, $W_i^K \in \mathbb{R}^{d_{model} \times d_K}$ and $W_i^V \in \mathbb{R}^{d_{model} \times d_V}$ for the i_{th} head and concatenate the
490 multiple heads to do the projection, then learns parameter matrices $W^O \in \mathbb{R}^{n_h d_v \times d_{model}}$ to
491 obtain the output of multi-head attention layer.

$$492 \quad Q_i = XW_i^Q, K_i = XW_i^K, V_i = XW_i^V$$

$$493 \quad A_i = \text{softmax}\left(\frac{Q_i K_i^T}{\sqrt{d_k}}\right)$$

$$494 \quad \text{head}_i = \text{Attention}(Q_i, K_i, V_i) = A_i V_i$$

$$495 \quad \text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_{n_h})W^O$$

496 where d_{model} denotes the dimension of token embedding X , which is 968 in EpiGePT X
497 denotes the embeddings from the sequence module for the first attention layer or the output of
498 previous attention layer. n_h denotes the number of head in Transformer encoder, which is 8 in
499 EpiGePT, and $d_K = d_V = d_{model}/n_h = 121$. The matrix A_i is called the self-attention matrix
500 for head i . The outputs of n_h heads are then concatenated, and a mapping function represented
501 by W^O is applied to obtain the output of the multi-head attention. After passing through an add
502 & norm layer, the multi-head attention output is used as input to the feed-forward layer, where
503 more comprehensive features of the input sequence are extracted. The above describes the

504 computational workflow of a single Transformer encoder layer. We set N_t to 16 for the
505 chromatin accessible prediction experiments, N_t to 12 for the chromatin state classification and
506 multiple chromatin signals prediction experiments.

507 **Prediction Module** For regression model, the output layer uses a linear transformation and
508 use mean square error (MSE) as the loss function. For classification model, the output layer
509 uses a linear transformation combined with a sigmoid function, and use the cross-entropy loss
510 for classification experiments.

511 **TF module** For binding status, we scanned the input bins for potential binding sites for a set
512 of 711 human transcription factors from HOCOMOCO database⁵⁰ with the tool Homer⁵¹ (Table
513 S5). We then selected the maximum score of reported binding status for each transcription
514 factor to obtain a vector of 711 dimensions as the motif feature for each DNA bin. For gene
515 expression, we focused on log-transformed TPM values of the 711 transcription factors and
516 obtained a vector of 711 dimensions after quantile normalization as the expression feature. With
517 these data, we combined the two vectors of motif and expression features by taking the element-
518 wise product, and we concatenated the result to the output of sequence module.

519 **Model evaluation**

520 To evaluate our model, we applied five-fold cross-validation in the different experiments on
521 cell-type level. For chromatin accessible experiments, the 129 cell lines are partitioned into a
522 training set and a testing set randomly.

523 Cell-type-wise metrics are defined to evaluate our method in different experiments, which were
524 calculate with the data within a test cell type across all genomic locus. For binary classification

525 design, we used cell-type-wise auPRC and auROC to evaluate our EpiGePT. Let $Y_{L \times K}$ and $\hat{Y}_{L \times K}$
526 be the true and predicted matrix, where L denotes the number of locus and K denotes the number
527 of test cell types. We calculated the auPRC and auROC for each $(y_{1i}, y_{2i}, \dots, y_{Li})$ and
528 $(\hat{y}_{1i}, \hat{y}_{2i}, \dots, \hat{y}_{Li})$. For multiple classification, we use macro average of the auPRC and auROC
529 to evaluate the classification performance, which compute the metric independently for each
530 class and then take the average hence treating all classes equally. For regression design, we
531 used two metrics for model evaluation, which are cell-type-wise Pearson correlation coefficient
532 and prediction squared error. Prediction square error (PSR) is calculated as $PSR = 1 -$
533 $\sum_k \sum_l (y_{lk} - \hat{y}_{lk})^2 / (y_{lk} - \bar{y}_{*k})^2$, where $\bar{y}_{*k} = \sum_l y_{lk} / L$ denotes the mean of the true level of
534 the response in the cell type k.

535 To compare the performance of our method with other baseline methods, we conducted
536 hypothesis testing on the metrics based on cell types. Since the metrics on a given cell type
537 across different methods are paired data and the statistical distribution is unknown, we
538 employed both Binomial and Wilcoxon tests, with the alternative hypothesis being that
539 EpiGePT outperforms the other methods. If we reject the null hypothesis, it provides
540 compelling evidence to support the claim that EpiGePT performs better than the other methods.

541 To evaluate the computational efficiency, we recorded the running time of a single epoch of
542 EpiGePT and baseline methods (Supplementary Text S4). Compared to traditional CNN models
543 such as DeepCAGE¹⁶ and ChromDragoNN⁶, as well as larger sequence models like Enformer,
544 EpiGePT demonstrates a balance between high computational efficiency and performance.

545 **Model training strategy**

546 As our proposed model is designed for cross-cell-type prediction of epigenomic signals by
547 multi-task learning. Some of the target epigenomic signals are missing in the existing ENCODE
548 database. For instance, there are 104 cellular contexts with both gene expression and at least
549 one of the epigenomic data. However, this number will decrease from 104 to 28 if we consider
550 eight epigenomic signals simultaneously. The proposed model takes each cellular context and
551 genomic region pair as a training instance, which ensures the availability of a very large number
552 of training instances. To utilize the data from the cellular contexts where some signals are not
553 available (missing data), we will use a new training strategy to handle the missing data where
554 the loss function is designed as

$$555 \quad L = \frac{1}{J} \sum_{j=1}^J \frac{1}{|B_i|} \sum_{k=1}^K \|y_{i,j,k} - \hat{y}_{i,j,k}\|_2^2 \cdot I(k \in B_i)$$

556 where $y_{i,j,k}$ and $\hat{y}_{i,j,k}$ denote the k^{th} true and predicted signal from the j^{th} genomic bin in the
557 i^{th} context, and $I(\cdot)$ is an indicator function and B_i denotes the index set that contains all
558 available signals in the i^{th} context. We update the parameters in the model through stochastic
559 gradient descent based on minibatches. We utilized the Adam optimizer with a batch size of 10
560 and a learning rate set to 5×10^{-5} . This training strategy provide us with a significantly larger
561 training sample size and allows us to utilize much more available data from the public databases,
562 and we enable EpiGePT to learn broader patterns of epigenetic states across diverse cell types.

563 **Incorporation of 3D chromatin interaction data**

564 With the emergence of methodologies like Hi-C and HiChIP for genome-wide chromatin
565 interaction measurement, a substantial volume of 3D chromatin interaction data has been
566 produced across various cellular contexts. Clearly, this data can provide highly valuable

567 information for identifying functional elements in the genome and for understanding gene
568 regulation, but this information has not been captured by current genomic LLMs such as the
569 Enformer¹⁵ or earlier CNN-based genomic models^{6, 7, 16, 52}.

570 We propose here to exploit the self-attention weights of the transformer model to design a
571 learning strategy that would allow EpiGePT to capture interaction information from Hi-C or
572 HiChIP data. Specifically, we propose to use the ground truth 3D genome interaction to guide
573 the self-attention matrices in the transformer module during the training process. First, we
574 obtained loop information at 5k resolution from the HiChIPdb database²⁶. Given potential noise
575 within HiChIP data, we selectively filtered potential H3K27ac-based HiChIP loops using a
576 stringent q-value threshold of 0.001. This curation aimed to utilization of highly confident loops,
577 safeguarding the model's ability to capture regulatory information without interference from
578 noise. In this way, we acquired corresponding HiChIP loop data for 13 out of 104 cell types.
579 Next, we mapped these loops onto the genomic bins used for pre-training. Specifically, we
580 employed the normalized count as a metric to gauge the likelihood score for each loop. During
581 the mapping process, we aggregated all loops based on this score to each specific genomic bin,
582 and then we obtained the HiChIP interaction matrix H_i . Based on the self-attention matrix
583 $A_{p,q}^i \in R^{J \times J}$ and the HiChIP interaction matrix H_i from the i^{th} cell type/tissue where p, q are
584 indexes for transformer layer and multi-heads, we apply a row-wise normalization to H_i (row
585 sum to 1) to obtain \tilde{H}_i and average the self-attention matrices across the heads in the last
586 transformer layers to obtain \tilde{A}^i . Since elevated attention weights are expected between regions
587 that interacts in 3D, we will compute a new loss term CSL, which is defined as cosine similarity
588 loss between the rows of \tilde{H}_i and \tilde{A}^i . Through the guidance of 3D genome interaction data, our

589 approach can learn a more comprehensive model for gene regulation. For example, it will
590 enable prediction of cell-type specific enhancer-promoter interaction, which is a task beyond
591 current models such as the Enformer. Note that the CSL term does not alter the architecture of
592 the model. It simply put some soft constraints on the attention weights according to the
593 experimental data on chromatin interactions, so that the optimized model will give predictions
594 that are more consistent with the context-specific interaction data. During training, the weight
595 α for 3d genome loss was chosen as 2.

596 **Fine-tuning for predicting E-P interaction**

597 For the fine-tuning process, we kept the parameters of the pre-trained model fixed without
598 making any updates. For the specific fine-tuning task of chromatin interaction prediction based
599 on HiChIP data, the multi-task prediction module was replaced with a two-layer MLP network,
600 containing 256 hidden nodes for each layer. During the training process, only the weights in the
601 MLP network in the prediction module were updated. Notably, when utilizing HiChIP data at a
602 resolution of 5k, both the enhancer and promoter anchors spanned 5kbp. Then we use a region
603 extending 128kbp from the center of the anchor of the neighboring gene, as input region for
604 EpiGePT. Consequently, a 968-dimensional feature vector for each genomic bin was derived
605 from the output of the last transformer encoder layer. These feature vectors from all bins within
606 the two anchors were concatenated, resulting in a high-dimensional vector of size 76,472. To
607 ensure the fairness of validating EpiGePT-finetune in capturing E-P interaction relationships,
608 we fine-tuned the model separately on the HiChIP data of each cell line during the fine-tuning
609 process. The test cell lines K562 and GM12878 were excluded from the pretrained EpiGePT
610 training cell types.

611 **Baseline methods**

612 Four baselines were introduced for epigenetic signals prediction. BIRD¹⁷ is a multiple linear
613 regression model that only takes gene expression data as input and makes predictions on a fixed
614 locus. ChromDragoNN⁶ is a deep neural network that takes gene expression of 1630 TFs and
615 DNA sequence as input. Specifically, ChromDragoNN⁶ uses a ResNet⁵³ to extract sequence
616 features and use linear transformation to combine the TF gene expression feature and sequence
617 feature to make the final prediction. DeepCAGE¹⁶ is a deep densely connected convolutional
618 network for predicting chromatin accessibility. Enformer¹⁵ is a deep neural network that
619 integrates convolutional neural network and transformer, and only takes DNA sequence as input.
620 Enformer takes DNA sequence of length 196kbp as input to predict 5,313 genomic tracks of
621 human and 1,643 tracks of mouse genome simultaneously. Enformer can only model and
622 predict cell types in the training data and cannot be applied to new cell types. In order to ensure
623 fairness in some of the benchmark experiment, we retrained the Enformer model with the same
624 input and output data as EpiGePT with Pytorch-lightning and made modifications on the
625 number of encoder layers when reproduce the Enformer model (Supplementary Text S5).
626 Besides, comparison with the pretrained Enformer model was also provided in Fig.2d where
627 we strictly used the ENCODE experiment ID to obtain the matched experiments for comparison.

628 Two baseline methods were introduced for predicting HiChIP interaction. DeepTACT³⁰ is a
629 deep learning method for predicting 3D chromatin contacts using both DNA sequence and
630 chromatin accessibility. We adopted the structure of DeepTACT³⁰ and kept the anchor length at
631 5k. The input to the model consists of two anchor sequences represented as one-hot matrices
632 and the two openness scores of the two anchors on the corresponding cell type extracted from

633 OpenAnnotate⁵⁴. Regarding the Kmer features⁵⁵, K is chosen as 5 to extract sequence features.
634 For each anchor, a vector of dimension $4^5 = 1024$ was obtained. Further training was
635 performed using an MLP with a hidden layer dimension of 256.

636 **Prediction of 3D genome interaction**

637 We collected cis-regulatory elements-gene pairs in K562 cells from other studies and public
638 database to demonstrate the interpretability of self-attention mechanisms in the EpiGePT.
639 Enhancers and silencers are typical *cis*-regulatory elements known play important roles in
640 transcriptional control during normal development and disease. For enhancers, we downloaded
641 enhancer-gene pairs from two studies: Gasperini et al.²² and Fulco et al.²³, both of which were
642 tested using a CRISPRi²¹ assay perturbation. Two datasets contain 664 and 5,091 element-gene
643 interactions. For silencers, we obtained and random sampled 831 validated silencers-gene pairs
644 with distance within 64kbp in K562 cells curated from high-throughput experiments from
645 SilencerDB²⁴. As there are no experimentally validated interaction relationships between these
646 silencers and genes, we generated silencer-gene pairs by associating the nearest neighbor genes
647 for classification purposes. Similarly, negative samples were generated by constructing DNase-
648 seq, ATAC-seq and nearest genes using the same approach. Ultimately, we obtained a dataset
649 comprising 1,662 silencer-gene pairs, encompassing both positive and negative instances.

650 To obtain scores for regulatory element-gene pairs, we first used the region extending 128kbp
651 from the center of the enhancer as input and extracted the token where the interacting genes
652 reside, so that we could filter out regulatory element-gene pairs that were located further than
653 64kbp apart. Subsequently, we stratified the remaining pairs based on their distance. Since the

654 positive and negative sample ratios varied across datasets, we adopted different stratification
655 strategies for different distance ranges (Fig. 3). Next, we averaged the attention matrices of the
656 Transformer encoder across all layers and heads. The summed attention scores from other
657 tokens to the key token containing the gene TSS were used as the attention score of this element-
658 gene pair. This score represents the attention value that the enhancer-centered region receives
659 for the TSS of the gene. We also calculated the attention score from the bin containing the center
660 of the regulatory element to the bin containing the TSS, which only slightly affects the
661 experimental results of regulatory element prioritization.

662 We collected 5k resolution data from the HiChIPdb (<http://health.tsinghua.edu.cn/hichipdb/>)
663 database, specifically from K562 and GM12878 cell lines. We filtered the data to include only
664 loops where at least one anchor falls within a gene region. We stratified the loops based on
665 distance into three categories: 0-20kbp, 20-40kbp, and 40-64kbp. For each distance category,
666 we selected 2000 positive pairs with most significant q-value. To ensure consistency in the
667 distance distribution, we selected negative pairs by fixing a gene and choosing anchors at
668 equidistant locations in the opposite direction. These are then used to as test data to evaluate
669 the prediction methods.

670 **Gradient importance scores**

671 EpiGePT possesses the capability to assign priority rankings to transcription factors by utilizing
672 gradient importance scores (GIS), taking into account specific cell types and chromatin regions.
673 The GIS were employed to identify potential functional relationships between specific TFs and
674 target genes. Specifically, for a given TF-target gene pair, the TSS of genes were used as central

675 loci, and the regions spanning 128 kbp upstream and downstream of the TSS were selected as
676 input. Next, we selected bins with motif binding scores indicating potential binding for the
677 given TF. For these selected bins, we calculated the GIS for the predictions of eight epigenomic
678 signals, for each of 711 core TFs.

$$679 \quad GIS_{ijk} = \frac{1}{|\zeta|} \sum_{l \in \zeta} \left| \frac{\partial \hat{y}_{ljk}}{\partial tf_{ij}} \right|$$

680 Where, i denotes the i th TF in the set of core TFs, j denotes the j th cell type, k denotes the k th
681 predicted epigenomic signal, and ζ denotes the set of genomic bins that have binding for the
682 given TF. In the calculation of the gradient, \hat{y}_{ljk} denotes the predicted value of the k th
683 epigenomic signal by the model using the expression in the j th cell type at the l th bin. On the
684 other hand, tf_{ij} denotes the product of the expression of i th TF in the j th cell type and the
685 corresponding TF binding score.

686 If we consider the GIS for the prediction of all 8 epigenomic signals simultaneously, we can
687 prioritize the TFs by calculating their ranks based on each signal separately. Then, we can
688 calculate an integrated gradient importance score (IGIS) for each TF by averaging the ranks
689 from all 8 signals.

$$690 \quad IGIS_{ij} = \frac{1}{8} \sum_k rank(GIS_{ijk})$$

691 Both the GIS and the IGIS are capable of capturing the significance of a transcription factor
692 (TF) in regulating a specific gene within the context of a specific cell type. Consequently, these
693 scores hold potential value in the discovery of TFs that play crucial roles in the regulation of
694 specific genes, thereby contributing to our understanding of essential regulatory mechanisms.

695 In the context of validating TF-TG pairs in the GRNdb and TRRUST databases, we opted to
696 utilize liver expression data as a representative example due to the unavailability of cell type
697 information for TRRUST. Furthermore, in this experimental setup, the tf_{ij} denotes the
698 expression of i th TF in the j th cell type and ζ denotes the set of genomic bins that have binding
699 for the TF of the given TF-target gene pair.

700 **Potential TF-target gene pairs from ChIP-seq data**

701 In this study, we utilized three distinct cell types to conduct a comprehensive screening of TF-
702 target gene pairs and non-target gene pairs across the human genome. Initially, we obtained the
703 narrow peak files (ENCFF388AJH, ENCFF717IXP, and ENCFF885KLR) from ChIP-seq
704 experiments across three cell types from the ENCODE project. Subsequently, we examined the
705 number of peaks within a 128kbp region both upstream and downstream of the TSS for each
706 gene. Different thresholds were applied to the ChIP-seq data of various TFs. Genes lacking any
707 peaks within the defined region were classified as non-target genes, while genes surpassing the
708 threshold in terms of peak counts were designated as target genes. Specifically, for the
709 aforementioned three cell types, threshold values of 10, 15, and 6 were respectively employed.
710 Finally, the IGIS approach was employed to determine the corresponding ranks of TFs in the
711 TF-target gene pairs.

712 **Pathogenic SNPs prioritization**

713 We collected single nucleotide polymorphisms (SNPs) data from the ClinVar and ExAC
714 databases, which include both potentially pathogenic and benign SNPs. To evaluate the ability
715 of EpiGePT to predict variant effects, we computed the log-ratio scores (LOS) for multiple

716 chromatin signals using EpiGePT on these SNPs. Subsequently, we utilized these scores to
717 distinguish between pathogenic and benign SNPs. The LOS for each chromatin signal was
718 defined by computing a forward pass through the model using the reference and alternative
719 alleles.

$$720 \quad \Delta O_{signal} = \log \left(\frac{output(I_{alt})}{output(I_{ref})} \right)$$

721 Where I_{ref} denotes the input DNA sequence based on the reference genome, and I_{alt} denotes
722 the input DNA sequence containing variants. Each chromatin epigenomic profile in each cell
723 line or tissue predicted by EpiGePT can be used to compute a specific variant score. We did not
724 take the absolute value in this calculation, so the resulting LOS indicates the direction of change
725 in the model output after the appearance of the variant. In addition to the predicted chromatin
726 signals output by the eight models, attention score changes based on self-attention are also
727 noteworthy. We computed the log-ratio scores for attention by summing the attention scores of
728 the 10 bins upstream and downstream of the locus of the SNP, to evaluate the effect of the
729 variant.

$$730 \quad \Delta O_{attention} = \sum_{i=-5}^5 \left| \log \left(\frac{attn(bin_i)_{I(alt)}}{attn(bin_i)_{I(ref)}} \right) \right|$$

731 Where i represents the index of the neighboring bins relative to the locus of the SNP. To avoid
732 the variant effects of different bins from cancelling each other out during the summation process,
733 we computed the absolute value of the change in attention scores for each bin and then summed
734 the scores of the 10 adjacent bins centered at the SNP position. For the classification of
735 pathogenic SNPs, we calculated these nine LOS for attention separately for each of the 28
736 tissues or cell lines in training data. As a result, we obtained a feature vector of 252 dimensions

737 for each SNP. Then a classifier with 252 features computed by EpiGePT and 52 annotations
738 from CADD score as inputs are used to predict pathogenic SNPs against benign or likely benign
739 SNPs. Here, we employed MLP as classifier to validate the effectiveness of the features we
740 obtained. A five-fold cross-validation experiment is employed for validation, and we utilize two
741 different positive-to-negative sample ratios, namely 1:1 and 1:2. For each sample ratio, we
742 randomly sample 32,000 positive samples. The effectiveness of the variant score in identifying
743 pathogenic SNPs is evaluated using the area under the auROC and the auPRC. Additionally,
744 we also utilized the logistic regression (LR) as the classifier, consistent with the LR classifier
745 used in CADD, and found a similar improvement when predicting pathogenic SNPs.

746 **COVID-19-associated SNPS prioritization.** We applied the same method to calculate the
747 LOS of the 8 epigenomic signals for the COVID-19 GWAS data. The absolute values of the
748 scores were summed as the overall score for each SNP. Then, we use the absolute sum as the
749 effect score of the SNP and prioritize the COVID-19-associated SNPs based on this score. For
750 each significant SNP associated with COVID-19 severity obtained from the GWAS data, we
751 selected normal SNPs within a 128kb region around the SNP as background to calculate the
752 rank of the LOS for the COVID-19 associated SNP in this region. Furthermore, we calculated
753 the LOS for all 9,484 COVID-19 associated SNPs and ranked them accordingly. The top 10
754 SNPs with the highest LOS were selected, which are considered to have potential genetic
755 associations with COVID-19 severity and complications.

756 **GTEx classification**

757 We collected eQTL data from the supplementary materials of Wang et al³⁷. In their study, the

758 authors identified causal eQTLs through statistical fine-mapping, using a posterior inclusion
759 probability (PIP) threshold of >0.9 for putative causal variants based on expression modifier
760 score (EMS), and a PIP threshold of <0.9 for putative non-causal variants. To validate the ability
761 of EpiGePT to distinguish potential causal variants, we perform a classification task on these
762 variants. For each variation, 128kbp sequence regions near it were selected as the input of the
763 model, and a score of variation was given by EpiGePT model. For each variant under each
764 tissue, we can obtain an 8-dimensional vector of genomic features including DNase, CTCF and
765 other ChIP-seq signals. Based on the LOS, separate random forest classifiers consisting of 10
766 decision trees are trained for each tissue in order to distinguish between causal and non-causal
767 variants. The models are evaluated using 5-fold validation on each tissue, with area under the
768 auPRC and auROC as metrics for assessing their ability to distinguish between causal and non-
769 causal variants.

770 **Code availability**

771 All components of EpiGePT are freely available at <https://github.com/ZjGaothu/EpiGePT>.
772 Here, users can access the code for reproducing EpiGePT, as well as the data collection and
773 preprocessing pipelines used for model training in benchmark experiments.

774 **Data availability**

775 Information and processed data of multiple chromatin signals of whole genome, motif binding
776 status and expression data of TFs in the corresponding cell lines/tissues, which are used in
777 EpiGePT are available at Supplementary Tables. The information about the cell lines/tissues
778 used and the 711 filtered transcription factors is available in the supplementary table. The High

779 throughput validated silencers of K562 cell line are download from SilencerDB
780 (<http://health.tsinghua.edu.cn/silencerdb>) database. The HiChIP data of K562 cell line and
781 GM12878 cell line are downloaded from HiChIPdb (<http://health.tsinghua.edu.cn/hichipdb/>)
782 database. The DNase-seq peak and ATAC-seq peak data are obtained from the ENCODE
783 project. Enhancer-gene pairs of CRISPRi²³ experiments are obtained from the supplementary
784 information of Gasperini et al. and Fulco et al. The regulatory network data for transcription
785 factors and target genes were obtained from the TRRUST³⁵ database
786 (<https://www.grnpedia.org/trrust/>) and the GRNdb³⁴ database (<http://www.grndb.com>). The
787 annotated chromatin states for whole genome are downloaded from the ROADMAP
788 epigenomics project (https://egg2.wustl.edu/roadmap/web_portal/chr_state_learning.html).
789 The RNA-seq read counts matrix for protein coding genes used for the prediction of the
790 chromatin 15-states annotated by ChromHMM are downloaded from the ROADMAP project
791 (<https://egg2.wustl.edu/roadmap/data/byDataType/rna/expression/57epigenomes.N.pc.gz>).
792 The GWAS data of COVID-19 are download from the COVID-19 Host Genetics Initiative
793 (<https://www.covid19hg.org/>).

794 **Ethics declarations**

795 **Competing interests**

796 The authors have declared no competing interests.

797 **Acknowledgments**

798 Z.G and R.J. was supported by the National Key Research and Development Program of China
799 [2021YFF1200902] and [2023YFF1204802], and the National Natural Science Foundation of

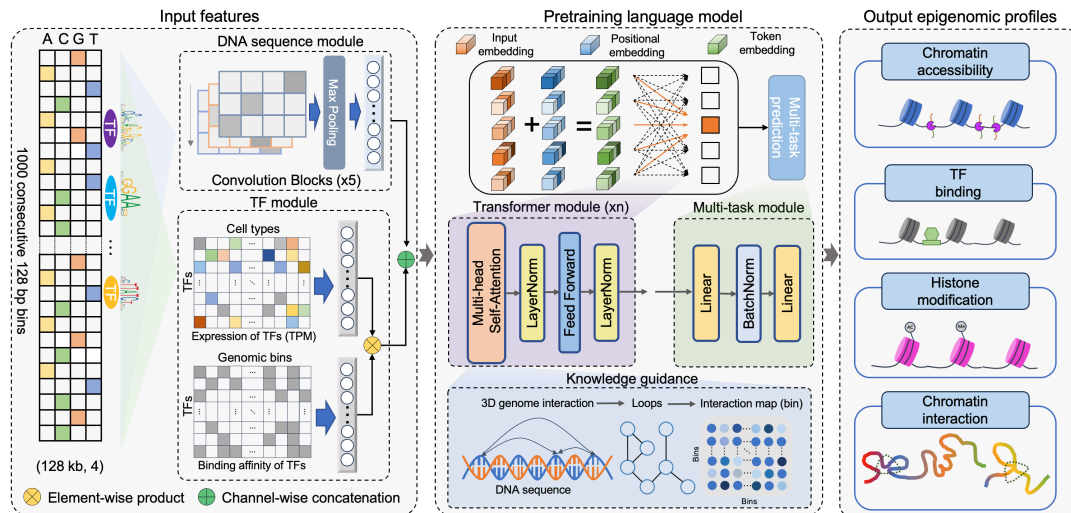
800 China [62203236 and 62273194]. Q.L., W.Z and W.H.W were supported by NIH grants R01

801 HG010359, P50 HG007735 and NSF DMS 1952386.

802

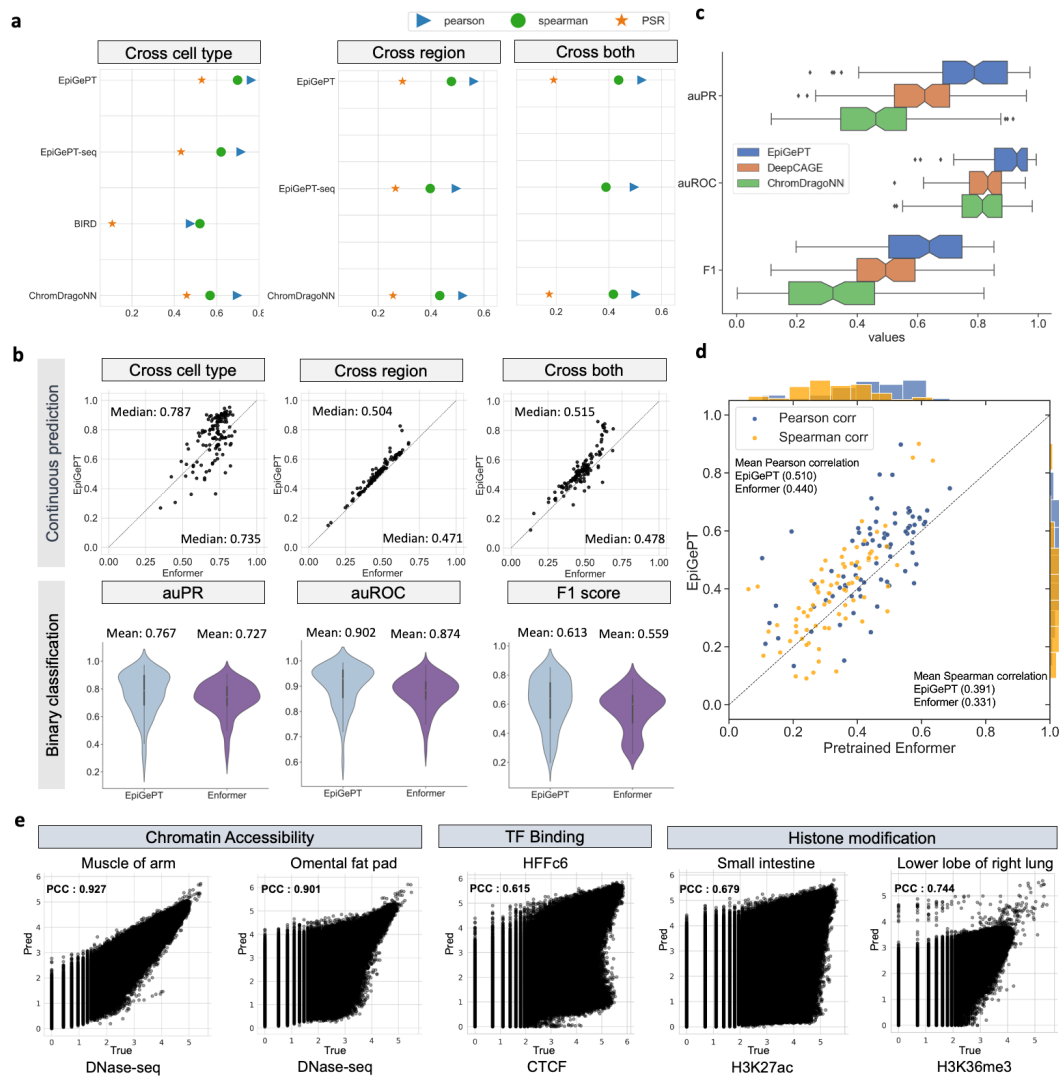
803 **Figures**

804 **Figure 1**



805 **Fig. 1 Overview of the EpiGePT model for multiple epigenomic signals prediction.** The
806 EpiGePT model consists of four modules, namely the Sequence module, the TF module, the
807 Transformer module, and the Multi-task prediction module. The sequence module comprises
808 multiple layers of convolution applied to the one-hot encoded DNA sequence input. The input
809 sequence length consists of 1000 genomic bins of 128bp for the prediction of multiple signals
810 and 50 bins of 200bp for the prediction of DNase signal alone. The TF module encompasses
811 the binding status and expression of 711 transcription factors. The Transformer module consists
812 of a series of consecutive transformer encoders, while the multi-task module is composed of a
813 fully connected layer. Additionally, the EpiGePT framework integrates an optional knowledge
814 guidance module that enhances the interpretability of the model by incorporating three-
815 dimensional chromatin interaction data into the attention layer, thus improving its
816 understanding of regulatory mechanisms.

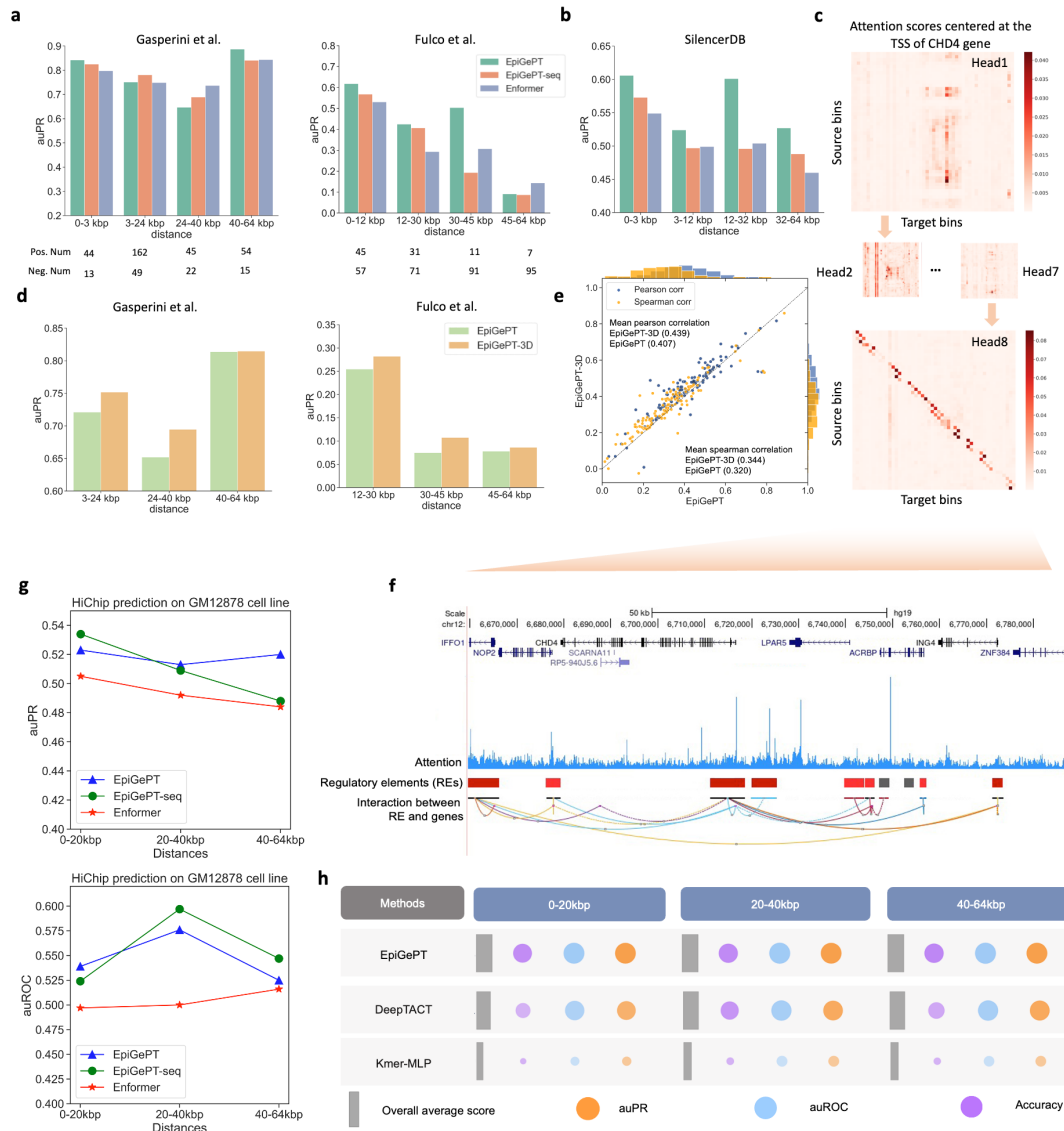
Figure 2



817 **Fig. 2 Performance of EpiGePT and baseline methods on the benchmark experiment. a,**
 818 EpiGePT and baseline methods were compared in terms of their regression performance for
 819 DNase signal regression across cell types, genomic regions, and combined cell type and
 820 genomic regions. **b,** Comparison of EpiGePT and Enformer performance. Each point in the
 821 scatter plot represents the performance of Enformer on the data of a specific cell type (x-axis)
 822 compared to the performance of EpiGePT (y-axis). The top three graphs represent the
 823 prediction of continuous DNase signals (pearson correlation coefficient), while the bottom three
 824 graphs represent the binary classification of chromatin accessibility regions. **c,** EpiGePT and

825 baseline methods' performance on binary prediction of DNase-seq signals. **d**, EpiGePT
826 demonstrates more excellent performance in predicting diverse epigenetic signals across
827 various cell types, compared with the pre-trained Enformer on 78 genomic tracks across 19
828 unseen cell types. The orange points represent Spearman correlation coefficient, and the blue
829 points represent Pearson correlation coefficient. **e**, EpiGePT cross-cell-type predictions
830 compared to experimental signals visualized for a representative example. The predictions
831 specific to DNase are based on the hg19 reference genome, while predictions for multiple
832 epigenomic profiles are conducted using the hg38 reference genome.

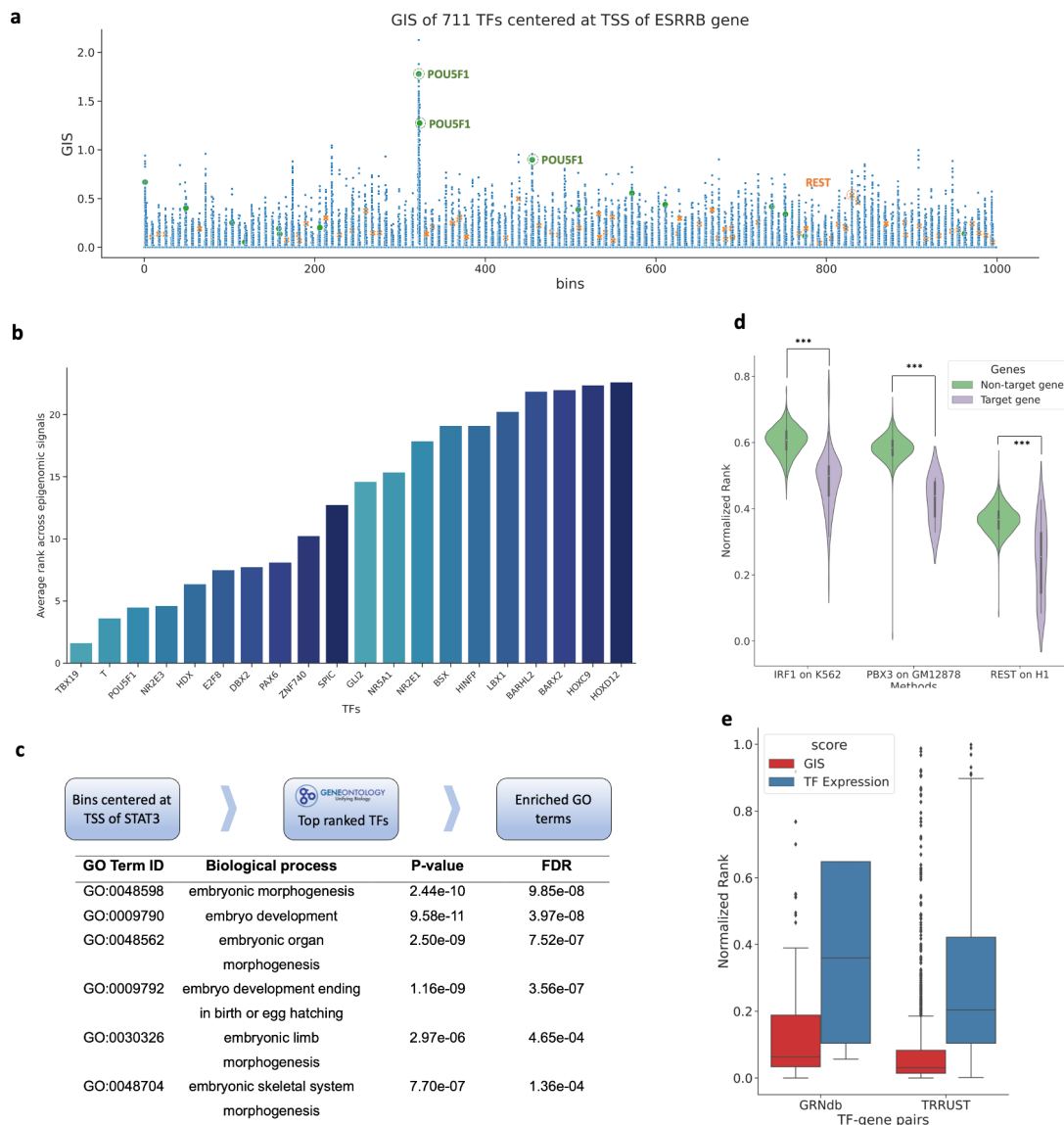
Figure 3



833 **Fig. 3 Application of self-attention mechanism in EpiGePT for long-range chromatin**
 834 **interaction identification.** **a**, The performance (auPRC) of attention score of EpiGePT in
 835 distinguishing enhancer-gene pairs at different distance ranges on two different datasets. **b**, The
 836 performance (auPRC) of attention score of EpiGePT in distinguishing silencer-gene pairs at
 837 different distance ranges based on the data from SilencerDB²⁴. **c**, Heatmap of the self-attention
 838 matrix of each attention head centered at the TSS of the *CHD4* gene, the (i, j) element in the

839 matrix denotes the average attention score between the i th genomic bin and the j th genomic bin
840 across all layers. **d**, The performance (auPR) of self-attention scores of EpiGePT and EpiGePT-
841 3D in identifying enhancer-promoter interactions across different distance ranges on the K562
842 cell type. **e**, The predictive performance (blue points denote pearson correlation coefficients
843 and orange points denote spearman correlation coefficients) of EpiGePT with knowledge
844 guidance across 19 cell types and 15,870 long sequences (128kbp). **f**, Attention scores centered
845 at the TSS of the *CHD4* gene, and putative enhancer regions in its vicinity. **g**, The performance
846 (auROC and auPR) of attention score of EpiGePT in distinguishing HiChIP loops of H3K27ac
847 at different distance ranges on GM12878 cell line. **h**, The performance (auROC and auPRC) of
848 the fine-tuned EpiGePT model and baseline methods (DeepTACT and Kmer) in distinguishing
849 enhancer-gene pairs at various distance ranges (0-20 kbp, 20-40 kbp and 40-64 kbp) on K562
850 cell line under a 5-fold cross validation setting. The size of the bubbles in the plot represents
851 the magnitude of the metric values, while the width of the gray rectangles along the x-axis
852 signifies the overall average values of the three metrics.

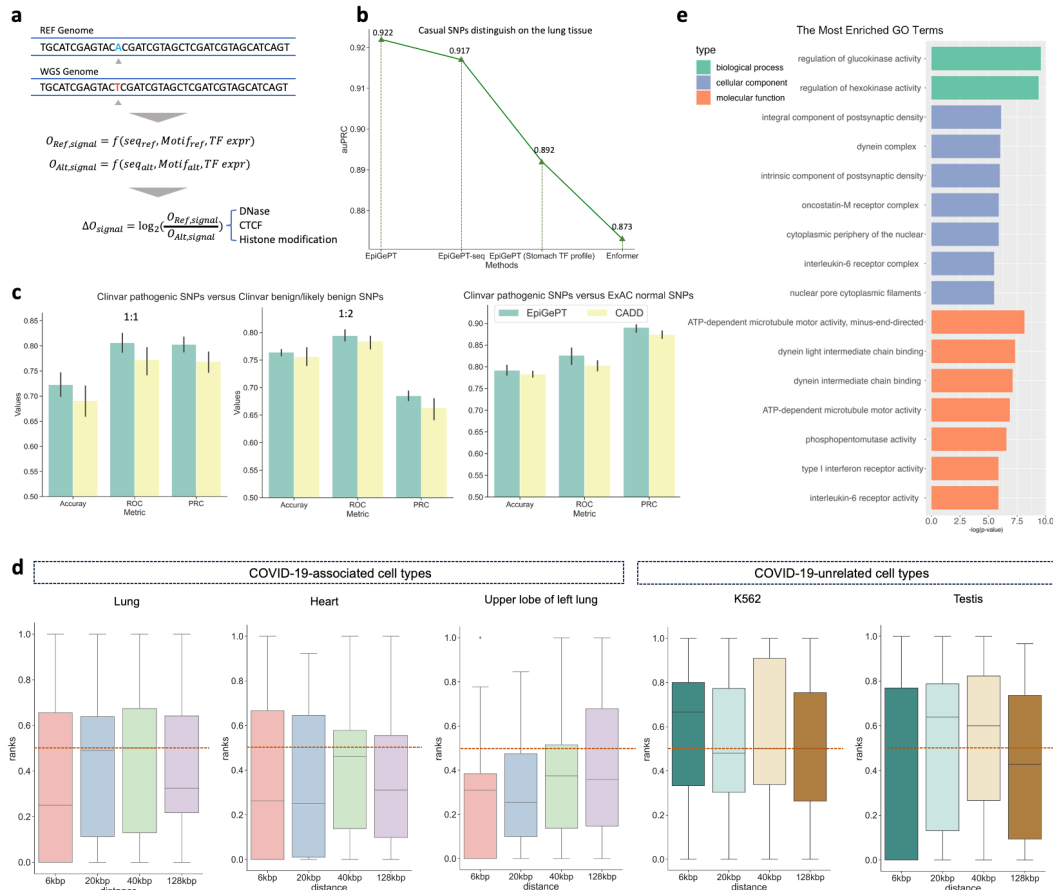
853 **Figure 4**



854 **Fig. 4 Gradient importance scores (GIS) uncover regulatory transcription factors. a,**
 855 Genomic regions around TSS of the *ESRRB* gene and TF expression data on ESC were used in
 856 EpiGePT. The scatter plot represents the GIS scores of 711 TFs on each genomic bin. Each dot
 857 represents the GIS score of a core TF on a specific genomic bin. Two important ESC regulators
 858 *REST* and *POU5F1* are highlighted. **b,** Bar plot of the top 5% ranked TFs, based on the average
 859 ranks from the GIS of eight epigenomic signals across bins (below). **c,** Based on the top 5%
 860 ranked TFs in 128kbp region centered at TSS of the *ESRRB* gene, gene ontology enrichment

861 analysis revealed significant enrichment in biological processes related to embryonic
862 development and cellular differentiation. **d**, Based on TF ChIP-seq data, all 23,635 human
863 genes were classified into target genes and non-target genes. The results revealed that TFs
864 exhibited significantly higher ranks on potential target genes compared to non-target genes. **e**,
865 The distribution of the rank of TFs in the GIS and expression value among the 2,705 TF-gene
866 pairs from the TRRUST database and 1,066 TF-gene pairs derived from genotype-tissue
867 expression (GTEx) data of the liver sourced from the GRNdb database.
868

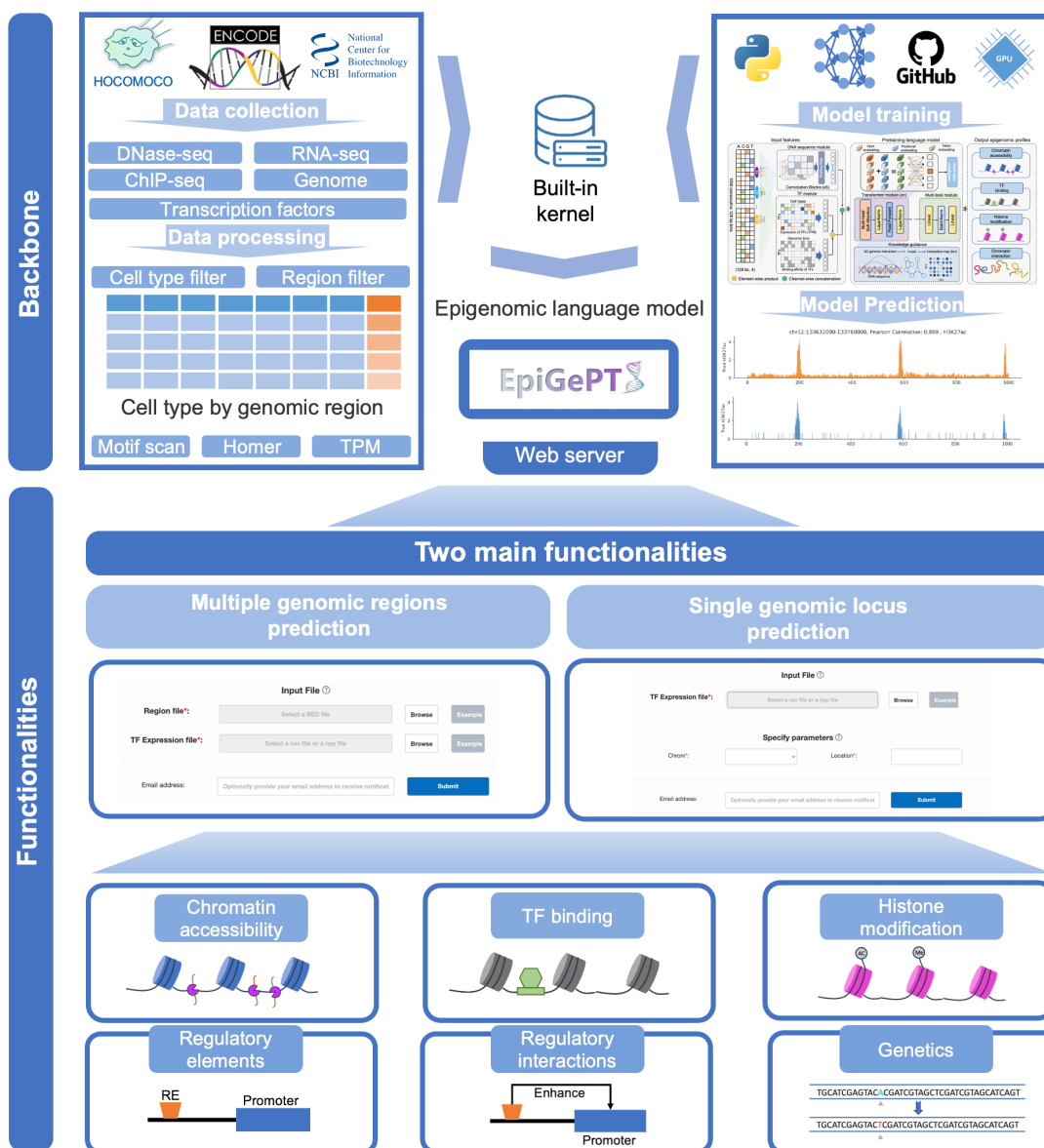
869 **Figure 5**



870 **Fig. 5 Variant effect prediction of EpiGePT.** **a**, The LOS for each epigenomic signal is
 871 calculated by the log change fold of the predicted epigenomic signal for reference genome and
 872 WGS genome. **b**, The performance of EpiGePT and Enformer in discriminating causal SNPs
 873 on the Lung tissue. **c**, The three subplots from left to right respectively depict the classification
 874 results for disease-related SNPs and benign SNPs down-sampled sourced from the ClinVar
 875 database, with balanced positive and negative samples (1:1 and 1:2 ratio), as well as normal
 876 SNPs sourced from the ExAC database with a MLP classifier. **d**, The ranked position of
 877 COVID-19 related GWAS data among surrounding benign SNPs based on their LOS, as

878 determined using different tissue or cell-type expression data. The results were stratified based
879 on the distance range of the risk region. The resulting mean and median ranks were both below
880 0.5. **e**, Enrichment result (Biological process, Cellular component and Molecular function) of
881 the nearest genes of the COVID-19 associated SNPs with the max LOS.

882 **Figure 6**



883 **Fig. 6 Overview of the online prediction web server of EpiGePT.** We collected eight types
 884 of epigenetic genome modification signals and corresponding expression data of transcription
 885 factors in different cell types or tissues from the ENCODE project. Based on these data, we
 886 trained the EpiGePT model and deployed it as a built-in kernel on an Apache server. Users
 887 without much coding experience can also access the web server in two ways to obtain the eight

888 types of epigenetic genome modification signals for specified cell types and genomic regions

889 without programming or installation.

890

891 Reference

- 892 1. Preissl, S., Gaulton, K.J. & Ren, B. Characterizing cis-regulatory elements using single-cell
893 epigenomics. *Nature Reviews Genetics* **24**, 21-43 (2023).
- 894 2. O'Malley, R.C. et al. Cistrome and epicistrome features shape the regulatory DNA landscape.
895 *Cell* **165**, 1280-1292 (2016).
- 896 3. Vandereyken, K., Sifrim, A., Thienpont, B. & Voet, T. Methods and applications for single-cell
897 and spatial multi-omics. *Nature Reviews Genetics*, 1-22 (2023).
- 898 4. Wang, K.C. & Chang, H.Y. Epigenomics: technologies and applications. *Circulation research*
899 **122**, 1191-1199 (2018).
- 900 5. Kelley, D.R., Snoek, J. & Rinn, J.L. Basset: learning the regulatory code of the accessible
901 genome with deep convolutional neural networks. *Genome research* **26**, 990-999 (2016).
- 902 6. Nair, S., Kim, D.S., Perricone, J. & Kundaje, A. Integrating regulatory DNA sequence and gene
903 expression to predict genome-wide chromatin accessibility across cellular contexts.
904 *Bioinformatics* **35**, i108-i116 (2019).
- 905 7. Zhou, J. & Troyanskaya, O.G. Predicting effects of noncoding variants with deep learning-
906 based sequence model. *Nature methods* **12**, 931-934 (2015).
- 907 8. Avsec, Ž. et al. Base-resolution models of transcription-factor binding reveal soft motif syntax.
908 *Nature Genetics* **53**, 354-366 (2021).
- 909 9. Fudenberg, G., Kelley, D.R. & Pollard, K.S. Predicting 3D genome folding from DNA sequence
910 with Akita. *Nature methods* **17**, 1111-1117 (2020).
- 911 10. Zhou, J. Sequence-based modeling of three-dimensional genome architecture from kilobase to
912 chromosome scale. *Nature genetics* **54**, 725-734 (2022).
- 913 11. Chen, K.M., Wong, A.K., Troyanskaya, O.G. & Zhou, J. A sequence-based global map of
914 regulatory activity for deciphering human genetics. *Nature genetics* **54**, 940-949 (2022).
- 915 12. Zhang, S. et al. Applications of transformer-based language models in bioinformatics: a survey.
916 *Bioinformatics Advances* **3**, vbad001 (2023).
- 917 13. Ji, Y., Zhou, Z., Liu, H. & Davuluri, R.V. DNABERT: pre-trained Bidirectional Encoder
918 Representations from Transformers model for DNA-language in genome. *Bioinformatics* **37**,
919 2112-2120 (2021).
- 920 14. Zhou, Z. et al. Dnabert-2: Efficient foundation model and benchmark for multi-species genome.
921 *arXiv preprint arXiv:2306.15006* (2023).
- 922 15. Avsec, Ž. et al. Effective gene expression prediction from sequence by integrating long-range
923 interactions. *Nature methods* **18**, 1196-1203 (2021).
- 924 16. Liu, Q., Hua, K., Zhang, X., Wong, W.H. & Jiang, R. DeepCAGE: incorporating transcription
925 factors in genome-wide prediction of chromatin accessibility. *Genomics, Proteomics &*
926 *Bioinformatics* **20**, 496-507 (2022).

- 927 17. Zhou, W. et al. Genome-wide prediction of DNase I hypersensitivity using gene expression.
928 *Nature communications* **8**, 1-17 (2017).
- 929 18. Song, L. & Crawford, G.E. DNase-seq: a high-resolution technique for mapping active gene
930 regulatory elements across the genome from mammalian cells. *Cold Spring Harbor Protocols*
931 **2010**, pdb. prot5384 (2010).
- 932 19. Consortium, E.P. An integrated encyclopedia of DNA elements in the human genome. *Nature*
933 **489**, 57 (2012).
- 934 20. Ernst, J. & Kellis, M. Chromatin-state discovery and genome annotation with ChromHMM.
935 *Nature protocols* **12**, 2478-2492 (2017).
- 936 21. Larson, M.H. et al. CRISPR interference (CRISPRi) for sequence-specific control of gene
937 expression. *Nature protocols* **8**, 2180-2196 (2013).
- 938 22. Gasperini, M. et al. A genome-wide framework for mapping gene regulation via cellular genetic
939 screens. *Cell* **176**, 377-390. e319 (2019).
- 940 23. Fulco, C.P. et al. Activity-by-contact model of enhancer–promoter regulation from thousands
941 of CRISPR perturbations. *Nature genetics* **51**, 1664-1669 (2019).
- 942 24. Zeng, W. et al. SilencerDB: a comprehensive database of silencers. *Nucleic acids research* **49**,
943 D221-D228 (2021).
- 944 25. Mumbach, M.R. et al. HiChIP: efficient and sensitive analysis of protein-directed genome
945 architecture. *Nature methods* **13**, 919-922 (2016).
- 946 26. Zeng, W., Liu, Q., Yin, Q., Jiang, R. & Wong, W.H. HiChIPdb: a comprehensive database of
947 HiChIP regulatory interactions. *Nucleic Acids Research* **51**, D159-D166 (2023).
- 948 27. Fishilevich, S. et al. GeneHancer: genome-wide integration of enhancers and target genes in
949 GeneCards. *Database* **2017**, bax028 (2017).
- 950 28. Radford, A., Narasimhan, K., Salimans, T. & Sutskever, I. Improving language understanding
951 by generative pre-training. (2018).
- 952 29. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. Bert: Pre-training of deep bidirectional
953 transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- 954 30. Li, W., Wong, W.H. & Jiang, R. DeepTACT: predicting 3D chromatin contacts via bootstrapping
955 deep learning. *Nucleic acids research* **47**, e60-e60 (2019).
- 956 31. van den Berg, D.L. et al. An Oct4-centered protein interaction network in embryonic stem cells.
957 *Cell stem cell* **6**, 369-381 (2010).
- 958 32. Zhang, J. et al. The oncogene Etv5 promotes MET in somatic reprogramming and orchestrates
959 epiblast/primitive endoderm specification during mESCs differentiation. *Cell death & disease*
960 **9**, 224 (2018).
- 961 33. Levy, S.H. et al. Esrrb is a cell-cycle-dependent associated factor balancing pluripotency and
962 XEN differentiation. *Stem Cell Reports* **17**, 1334-1350 (2022).
- 963 34. Fang, L. et al. GRNdb: decoding the gene regulatory networks in diverse human and mouse

- 964 conditions. *Nucleic acids research* **49**, D97-D103 (2021).
- 965 35. Han, H. et al. TRRUST v2: an expanded reference database of human and mouse transcriptional
966 regulatory interactions. *Nucleic acids research* **46**, D380-D386 (2018).
- 967 36. Willett, R. et al. TFEB regulates lysosomal positioning by modulating TMEM55B expression
968 and JIP4 recruitment to lysosomes. *Nature communications* **8**, 1580 (2017).
- 969 37. Wang, Q.S. et al. Leveraging supervised learning for functionally informed fine-mapping of cis-
970 eQTLs identifies an additional 20,913 putative causal eQTLs. *Nature Communications* **12**, 3394
971 (2021).
- 972 38. Landrum, M.J. et al. ClinVar: public archive of interpretations of clinically relevant variants.
973 *Nucleic acids research* **44**, D862-D868 (2016).
- 974 39. Karczewski, K.J. et al. The ExAC browser: displaying reference data information from over 60
975 000 exomes. *Nucleic acids research* **45**, D840-D845 (2017).
- 976 40. Rentzsch, P., Witten, D., Cooper, G.M., Shendure, J. & Kircher, M. CADD: predicting the
977 deleteriousness of variants throughout the human genome. *Nucleic acids research* **47**, D886-
978 D894 (2019).
- 979 41. Li, J., Lai, S., Gao, G.F. & Shi, W. The emergence, genomic diversity and global spread of
980 SARS-CoV-2. *Nature* **600**, 408-418 (2021).
- 981 42. org, C.-H.G.I.a.b. The COVID-19 host genetics initiative, a global initiative to elucidate the role
982 of host genetic factors in susceptibility and severity of the SARS-CoV-2 virus pandemic.
983 *European Journal of Human Genetics* **28**, 715-718 (2020).
- 984 43. Wang, W., Wang, C.-Y., Wang, S.-I. & Wei, J.C.-C. Long-term cardiovascular outcomes in
985 COVID-19 survivors among non-vaccinated population: a retrospective cohort study from the
986 TriNetX US collaborative networks. *EClinicalMedicine* **53** (2022).
- 987 44. Huang, C. et al. Clinical features of patients infected with 2019 novel coronavirus in Wuhan,
988 China. *The lancet* **395**, 497-506 (2020).
- 989 45. Agius, L. Targeting hepatic glucokinase in type 2 diabetes: weighing the benefits and risks.
990 *Diabetes* **58**, 18-20 (2009).
- 991 46. Singh, A.K., Gupta, R., Ghosh, A. & Misra, A. Diabetes in COVID-19: Prevalence,
992 pathophysiology, prognosis and practical considerations. *Diabetes & Metabolic Syndrome:
993 Clinical Research & Reviews* **14**, 303-310 (2020).
- 994 47. Pellegrina, D., Bahcheli, A.T., Krassowski, M. & Reimand, J. Human phospho-signaling
995 networks of SARS-CoV-2 infection are rewired by population genetic variants. *Molecular
996 Systems Biology* **18**, e10823 (2022).
- 997 48. Loyfer, N. et al. A DNA methylation atlas of normal human cell types. *Nature* **613**, 355-364
998 (2023).
- 999 49. Gao, Z. et al. scEpiTools: a database to comprehensively interrogate analytic tools for single-
1000 cell epigenomic data. *Journal of Genetics and Genomics* (2023).

- 1001 50. Kulakovskiy, I.V. et al. HOCOMOCO: towards a complete collection of transcription factor
1002 binding models for human and mouse via large-scale ChIP-Seq analysis. *Nucleic acids research*
1003 **46**, D252-D259 (2018).
- 1004 51. Heinz, S. et al. Simple combinations of lineage-determining transcription factors prime cis-
1005 regulatory elements required for macrophage and B cell identities. *Molecular cell* **38**, 576-589
1006 (2010).
- 1007 52. Kelley, D.R. et al. Sequential regulatory activity prediction across chromosomes with
1008 convolutional neural networks. *Genome research* **28**, 739-750 (2018).
- 1009 53. He, K., Zhang, X., Ren, S. & Sun, J. in Proceedings of the IEEE conference on computer vision
1010 and pattern recognition 770-778 (2016).
- 1011 54. Chen, S. et al. OpenAnnotate: a web server to annotate the chromatin accessibility of genomic
1012 regions. *Nucleic Acids Research* **49**, W483-W490 (2021).
- 1013 55. Chor, B., Horn, D., Goldman, N., Levy, Y. & Massingham, T. Genomic DNA k-mer spectra:
1014 models and modalities. *Genome biology* **10**, 1-10 (2009).

1015 **Supplementary Materials**

1016 Text S1. Data splitting strategy for model training.

1017 Text S2. System design and implementation of the web server.

1018 Text S3. Case application of the EpiGePT-online.

1019 Text S4. Running time of the EpiGePT and baseline methods.

1020 Text S5. Implementation of Enformer model and Enformer+.

1021 Text S6. Data processing for ChromHMM annotation data.

1022 Fig. S1. Three data partitioning strategies for model training and testing.

1023 Fig. S2. Model architecture of EpiGePT for multiple epigenomic signals prediction.

1024 Fig. S3. EpiGePT's performance in predicting DNase-seq and other epigenetic signals.

1025 Fig. S4. Performance of EpiGePT and baseline methods on chromatin states classification,
1026 multiple epigenomic profiles prediction and causal variants classification.

1027 Fig. S5. Ablation analysis of the EpiGePT model.

1028 Fig. S6. Performance of EpiGePT in cross-cell-type prediction.

1029 Fig. S7. The performance (auROC) of attention score of EpiGePT in distinguishing regulatory
1030 element-gene pairs at different distance ranges.

1031 Fig. S8. Incorporating 3D genomic information from HiChip data enhances the predictive
1032 performance of EpiGePT on E-P regulatory interaction on K562 cell line.

1033 Fig. S9. The fine-tuning performance of the EpiGePT model on predicting potential enhancer-
1034 promoter regulatory networks.

1035 Fig. S10. The ROC and PR curves of the EpiGePT model on predicting potential enhancer-
1036 promoter regulatory networks.

1037 Fig. S11. The GIS of ChIP-seq overlapped bins versus non-overlapped bins of POU5F1
1038 centered at the TSS of ESRRB.

1039 Fig. S12. Gene ontology enrichment analysis based on the top 5% TFs with high expression in
1040 ESCs.

1041 Fig. S13. Case application of the EpiGePT-online.

1042 Fig. S14. Enrichment result (Cellular component and Molecular function) of the nearest genes
1043 of the COVID-19 associated SNPs with the low LOS.

1044 Table S1. The information of DNase-seq bam file across 129 biosamples from the ENCODE7
1045 project.

1046 Table S2. The information of RNA-seq tab-separated values (tsv) file across 129 biosamples
1047 from the ENCODE project.

1048 Table S3. The information of DNase-seq, CTCF and other six Histone markers bam file across
1049 28 cell lines or tissues from the ENCODE project (hg19).

1050 Table S4. The information of DNase-seq, CTCF and other six Histone markers bam file across
1051 104 cell lines or tissues from the ENCODE project (hg38).

1052 Table S5. The information of RNA-seq tab-separated values (tsv) file across 28 cell lines or
1053 tissues from the ENCODE project (hg19).

1054 Table S6. The information of RNA-seq tab-separated values (tsv) file across 104 cell lines or
1055 tissues from the ENCODE project (hg38).

1056 Table S7. The preprocessed expression data of 711 human transcription factors from the
1057 ENCODE project across 129 biosamples.

1058 Table S8. The preprocessed expression data of 711 human transcription factors from the
1059 ENCODE project across 28 cell lines or tissues (hg19).

1060 Table S9. The preprocessed expression data of 711 human transcription factors from the
1061 ENCODE project across 104 cell lines or tissues (hg38).

1062 Table S10. The order and names of epigenomes of the expression matrices across 56
1063 epigenomes from the ROADMAP project.

1064 Table S11. The preprocessed expression data of 642 human transcription factors across 56
1065 epigenomes from the ROADMAP project.