# EpiGePT: a Pretrained Transformer model for epigenomics

Zijing Gao[1,#], Qiao Liu[2,#,*], Wanwen Zeng[2], Wing Hung Wong[2,3,*] and Rui Jiang[1,*]

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

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

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

[*] To whom correspondence should be addressed.

[#] The first two authors contributed equally.

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

# Abstract

The transformer-based models, such as GPT-3[1] and DALL-E[2], have achieved unprecedented breakthroughs in the field of natural language processing and computer vision. The inherent similarities between natural language and biological sequences have prompted a new wave of inferring the grammatical rules underneath the biological sequences. In genomic study, it is worth noting that DNA sequences alone cannot explain all the gene activities due to epigenetic mechanism. To investigate this problem, we propose EpiGePT, a new transformer-based language pretrained model in epigenomics, for predicting genome-wide epigenomic signals by considering the mechanistic modeling of transcriptional regulation. Specifically, EpiGePT takes the context-specific activities of transcription factors (TFs) into consideration, which could offer deeper biological insights comparing to models trained on DNA sequence only. In a series of experiments, EpiGePT demonstrates state-of-the-art performance in a diverse epigenomic signals prediction tasks as well as new prediction tasks by fine-tuning. Furthermore, EpiGePT is capable of learning the cell-type-specific long-range interactions through the self-attention mechanism and interpreting the genetic variants that associated with human diseases. We expect that the advances of EpiGePT can shed light on understanding the complex regulatory mechanisms in gene regulation. We provide free online prediction service of EpiGePT through https://health.tsinghua.edu.cn/epigept/.

## Introduction

One of the fundamental problems in genomic study is how to decode and interpret the human genome sequences in a complex manner. Progress toward this goal is largely hindered by the vast majority of non-coding regions[3]. For example, it remains unclear how the genomic variants in the noncoding regions lead to malfunctions of regulatory elements by disrupting the underlying regulatory syntax of DNA[4]. Inspired from the field of natural language processing, there exists a natural analogy between human language and DNA sequence where texts are made of words and DNA sequence can be characterized by nucleotides or k-mers. The inherent similarities between natural language and biological sequences provide new perspectives towards better understanding the complex DNA language.

Recently, generative pre-trained transformer (GPT) models have achieved unprecedented success in various domains, including computer vision and natural language processing (NLP)[1, 5]. Such pre-trained models can be readily tailored or adapted to various downstream tasks. To date, the application of generative pre-trained models in genomic study remains largely unexplored. It is noticeable that a number of machine learning-based approaches have been proposed for predicting various genomic and epigenomic signals, such as chromatin accessibility[6, 7], histone modification[8] or chromatin interactions[9, 10]. However, these methods are rather scattered, with specific models designed for specific prediction tasks. It is in an urgent need to develop a foundation model to facilitate multiple genomic and epigenomic prediction tasks and unveil the universal gene regulation rules.

To design a genomic foundation model, it is worth noting that the existing large language

3

52    models (LLMs) purely rely on the language context consisting of words and sentences while

53    the DNA sequences cannot explain all the heritable and stable changes in gene activity due to

54    epigenetic mechanisms. In other words, the genomic foundation model based on pure DNA

55    sequence may largely ignore the context-specific information, thus lacking mechanistic

56    interpretation of context-specific gene regulation. For example, using transformer-based

57    language model to decode genome sequence has been attempted by a recent work Enformer[11].

58    However, Enformer is not capable of predicting the function of sequences in new cellular

59    contexts, which largely limits its generalization power.

60    To overcome the above limitation, we proposed EpiGePT, a new transformer-based deep

61    learning framework, to predict genome-wide epigenomic signals by taking the mechanistic

62    modeling of transcriptional regulation into consideration. With EpiGePT, we are able to

63    investigate how to utilize the power of transformer-based language model to help researchers

64    uncover how trans-regulatory factors (e.g., TFs) regulate target genes by interacting with cis-

65    regulatory elements and further lead to changes in different chromatin states. After pretraining

66    on a diverse panel of cell line and tissue level data from the Encode database[12], EpiGePT is able

67    to directly predict the genome-wide chromatin states in any new cellular context given the

68    expression profile of a few hundreds of TFs or facilitate new prediction tasks (e.g., 3D genome

69    interaction) with finetuning.

70    To the best of our knowledge, EpiGePT is the first pretrained Transformer model for

71    epigenomics with mechanistic modeling of transcriptional regulation. EpiGePT differs from

72    existing methods in the following three aspects. First, unlike the methods that take pure DNA

73    sequence as input, EpiGePT additionally takes the context-specific information (e.g., TF

74    activities) as input, thus enabling genome-wide prediction power in any new cellular context.

75    Second, instead of using task-specific model to predict a single genomic and epigenomic signal,

76    EpiGePT is designed for simultaneously predicting multiple epigenomic signals of the same

77    genomic region through multi-task learning, thus improving learning efficiency and prediction

78    accuracy compared to the task-specific models. Third, many methods typically take short DNA

79    sequence (e.g., a few hundred or a thousand base pair) as input, which may not be adequate to

80    capture the complex syntax of DNA due to truncation. The long input DNA sequence (e.g.,

81    128kb) for EpiGePT greatly enhances the ability for the model to capture the long-range

82    interaction in the genome, which are crucial for understanding the gene regulation mechanism.

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

84    various tasks of chromatin states prediction, as well as the variant effect prediction. We also

85    show that the self-attention mechanism greatly helps unveil the complex code in the

86    conformation of long-range chromatin interactions, such as promoter-enhancer interactions and

87    promoter-silencer interactions. EpiGePT is an example of how transformer-based language

88    model and large-scale pretrain can be used in genomics research to provide biological insights.

89    With the help of EpiGePT model, it is expected that researchers can dissect the comprehensive

90    genomic regulatory code given the cellular context information and accelerate research findings

91    in genomic study.

# Results

## Overview of EpiGePT model

We developed a novel Transformer[13]-based language model named EpiGePT to predict multiple chromatin states across different cell types. EpiGePT is a language model for cross-cell-type prediction of chromatin states by multi-task learning based on genome-wide pre-training on epigenomic data (**Fig. 1 and** Fig. S6). EpiGePT is composed of four modules, including a sequence module, a TF module, a transformer module, and a prediction module. The sequence module is responsible for processing the long DNA sequence of interest (e.g., 128 kb) by employing a series of convolutional and pooling blocks (e.g., 5) to extract a comprehensive set of sequence features. By reducing the input length by $2^5=128$ times through pooling operations, this module effectively compresses the input information while retaining essential features. The TF module is specifically designed to extract cell-type-specific features by taking the expression of transcription factors in the given context, as well as their corresponding motif score into account. This module helps capture the unique characteristics of each cell type by considering the binding status of TFs involved in gene regulation. In the transformer module, each token corresponds to a genomic bin in the original DNA sequence and has hybrid features derived from both sequence and TFs. The module leverages self-attention mechanisms to learn the comprehensive relationships among the input bins, enabling the model to make predictions of multiple chromatin states under the given context cellular. By taking advantage of this approach, EpiGePT provides a powerful tool for predicting multiple chromatin states and enables researchers to gain insights into the underlying regulatory mechanisms of the genome.

6

## EpiGePT enables genome-wide prediction of chromatin states

To assess the predictive performance for epigenomic signals of EpiGePT, specifically in predicting chromatin accessibility, a comprehensive evaluation was conducted. We first applied EpiGePT to predict the chromatin accessibility based on the widely available public DNase-seq[14] data across diverse cell types or tissues. In brief, DNase-seq data across 129 cell types were collected from the ENCODE[12] project. After data preprocessing and normalization (see Methods), 1,175,374 genomic regions were extracted where each pair of cell type and genomic regions constitutes a training instance. We meticulously devised comprehensive experimental settings by partitioning the training and test sets based on either genomic regions or cell types. In detail, we employed the following three data partitioning settings for a comprehensive evaluation (Fig. S1, Text S1). For "cross-cell-type" prediction, we partitioned the data into training and testing sets based on cell types. For "cross-region" prediction, we partitioned the data into training and testing sets based on the genomic regions. For "cross-both" prediction, we conducted rigorous data split to ensure that both the cell types and genomic regions in the test stage are unseen during the training process. We employed three evaluation metrics, namely Pearson correlation coefficient, Spearman correlation coefficient and prediction square error, to assess the similarity between the predicted and true values of DNase signals (See Methods). It is shown that EpiGePT consistently outperforms other competing methods, including Enformer[11], BIRD[15], and ChromDragoNN[16] by a relatively large margin under the above experimental settings (**Fig. 2A** and Fig. S2). EpiGePT achieves 5.0%, 8.9%, and 5.2 % higher performance than Enformer, the best baseline method, in terms of the mean Pearson correlation coefficient under three data-partition settings, respectively (**Fig. 2B**). Besides the chromatin

135    accessibility regression task, we also designed binary chromatin accessibility status prediction

136    task by assessing whether a peak exists within the corresponding genomic bin (>50% overlap).

137    We made slight adjustments to the regression model by modifying the activation function and

138    loss function to accommodate the binary classification task (See methods). The results show

139    that EpiGePT achieves an average auPRC (area under the precision-recall curve) of 0.767

140    compared to 0.727 of Enformer[11], 0.623 of DeepCAGE[17] and 0.476 of ChromDragoNN[16] (**Fig.**

141    **2C**).

142    Next, we extended the chromatin state prediction task from a single target to multiple targets

143    by predicting multiple chromatin states, including chromatin accessibility, CTCF[18], ChIP-seq,

144    and six types of different histone modifications[19] (See Methods). When considering eights

145    different chromatin states, only 28 cell types have the corresponding available data

146    simultaneously. After preprocessing, 13,300 genomic regions each with a length of 128 kbp

147    were extracted, which cover 56.7% of the whole genome. However, compared to the data in

148    the DNase-only prediction experiment, the correlation coefficient was reduced due to the

149    prevalence of a substantial number of zero signals in the genomic regions being predicted.

150    Using a similar data split strategy as the single target for cross-cell-type prediction, EpiGePT

151    demonstrated a mean Pearson correlation coefficient between 0.259 to 0.566 of different

152    chromatin states in the test cell types (**Fig. 2D**). Specifically, EpiGePT achieves remarkably

153    high performance in predicting chromatin state signals for certain cell types, such as the colon

154    tissue, with a Pearson correlation coefficient of 0.888. Furthermore, as it shown in Fig. S3C, it

155    significantly outperformed Enformer in terms of performance across these tested cell types and

156    different signals (one-side $p$-value < 2.79e-10 under binomial hypothesis test). To make the

8

157    chromatin states prediction task more illustrative, several tracks of predicted chromatin states

158    and the corresponding ground truth chromatin states were displayed. For instance, at the

159    position of from 61,056,000 to 61,184,000 on chromosome 20, we used the UCSC genome

160    browser[20] to show the predicted values and true values of CTCF (Pearson correlation coefficient

161    of 0.518) and DNase signals (Pearson correlation coefficient of 0.869), as well as the regulatory

162    relationships within this region. (Fig. S2A). In addition, we also compared EpiGePT with

163    ChromDragoNN[16] on binary and quaternary classification tasks based on ChromHMM[21]

164    annotations (Fig. S3). EpiGePT achieved an average auROC (area under the receiver operating

165    characteristic curve) of 0.855 in binary classification, significantly higher than that of

166    ChromDragoNN[16] (0.774). In quaternary classification, EpiGePT achieved a macro-auROC of

167    0.879, also significantly higher than ChromDragoNN[16] (0.856, one-side $p$-value < 0.001).

168    These results demonstrate the effectiveness and accuracy of EpiGePT in predicting multiple

169    chromatin states, leveraging the four modules. The effectiveness and prediction power achieved,

170    in conjunction with the self-attention mechanism, lays the foundation for deciphering

171    regulatory relationships.

172    To further verify the roles of the main modules in the model, we conducted the ablation

173    experiments on the model architecture. For TF module ablation, the above experimental results

174    compared to EpiGePT without TF module (EpiGePT-seq) and EpiGePT have demonstrated that

175    EpiGePT outperforms EpiGePT-seq in cross-cell-type prediction of DNase signals, with an

176    average Pearson correlation coefficient of 0.714 and a median of 0.74 for EpiGePT-seq, while

177    EpiGePT achieves 0.756 (average) and 0.787 (median). In addition, the inclusion of the TF

178    module enables EpiGePT to predict chromatin states at the locus level for different cell types.

179    However, like Enformer, EpiGePT-seq predicts the same values for different cell types at the

180    same locus, resulting in a zero correlation for cross-cell-type prediction. We examined the

181    impact of the TF module on multi-task prediction by employing three methods, namely

182    replacing TF scores with zero, adding random noise to TF, and removing motif binding scores.

183    The results indicated that when the expression of TFs was set to zero, the prediction of H3K27ac

184    yielded a Pearson correlation coefficient of 0.190. However, incorporating the TF module

185    significantly improved the coefficient to 0.543, demonstrating a beneficial impact of the TF

186    module.

187    For sequence module ablation, we randomly subsampled 10,000 genomic bins and 20 cell

188    types to train a TF-only model. The results indicated that removing the sequence

189    module resulted in an average decrease of 0.084 in the Pearson correlation coefficients

190    of the eight signals on a cell-type wise basis, and with a particularly significant decrease

191    of 0.13 in predicting H3K4me3 signals (Fig. S4A).

192    For multi-task module, we predicted the eight chromatin states involved in training

193    using eight individual models for single-task prediction. The results were evaluated on

194    a cross-cell type prediction manner. In the case of predicting the signal of H3K4me1,

195    the average Pearson correlation decreased from 0.408 to 0.329. When predicting the

196    H3K4me1 signal, the average Pearson correlation coefficient decreased from 0.408 to

197    0.329. Similarly, the overall prediction performance for the eight signals declined by

198    0.074 (**Fig. 4B**). This decrease may be attributed to the intricate nature of gene

199    regulation. The distinct chromatin states can complement and synergize with each other

200    through multi-task learning, allowing the model to gain deeper biological insights

201    compared to a single-target prediction model.

202    Furthermore, we performed additional experiments to investigate the effect of the

203    number of cell lines on the prediction performance. Specifically, we focused on DNase

204    predictions and randomly downsampled the training cell types from 103 to 75, 50, 25

205    for each of the five folds in the cross-validation experiment. The results demonstrated

206    a strong positive correlation between the number of cell lines and the prediction

207    performance. Under five-fold cross-validation, the median Pearson correlation

208    coefficient on the test set across 129 cell types decreased from 0.793 to 0.790, 0.761,

209    and 0.732, respectively. These findings suggest that our current model has potential

210    room for improvement and additional training with more cell lines will lead to even

211    better predictive performance, thereby offering more comprehensive insights into the

212    regulatory mechanisms for researchers (Fig. S4C). In summary, EpiGePT demonstrated

213    superior performance in predicting both single and multiple epigenomic signals over

214    existing methods, providing a robust foundation for decoding the complex landscape

215    of gene regulation.

216    **EpiGePT facilitates long-range chromatin interaction identification**

217    We examined the capacity of EpiGePT for predicting long-range chromatin interactions, which

218    play a critical role in preserving chromatin architecture and elucidating 3D contacts between

219    distal regulatory elements and target genes. Traditional methods typically take short DNA

220    sequence (e.g., 1kbp) as input, thus cannot take the long-range chromatin interactions into

221    consideration. During the training of EpiGePT model, the self-attention mechanism in the

11

222  transformer module plays an important role in capturing the potential interactions between

223  different DNA bins. We utilized the cell-type specific self-attention scores to predict chromatin

224  interactions, including enhancer-promoter and silencer-promoter interactions (see Methods).

225  We initially investigated whether EpiGePT can differentiate experimentally validated enhancer-

226  promoter interactions from other interactions. Two datasets containing 664 and 5,091 candidate

227  enhancer-promoter interactions or element-TSS interactions obtained by CRISPRi[22]

228  experiments were used and further filtered and stratified by the distance. In the Gasperini et

229  al[23]. dataset, EpiGePT consistently outperform EpiGePT-seq and Enformer by achieving the

230  highest auPRC in all groups. For instance, EpiGePT achieved auPRC of 0.949, 0.726, and 0.810

231  for identifying enhancer-promoter pairs in the 0-3 kbp, 3-20 kbp, and 20-64 kbp ranges,

232  respectively (**Fig. 3A and** Fig. S8). In the Fulco et al.[24] dataset, EpiGePT obtains better

233  performance than EpiGePT-seq in most groups, which illustrates the positive benefit of the cell-

234  type-specificity brought in the TF module. EpiGePT consistently outperforms Enformer across

235  different groups by a relatively large margin. As shown in Fig. 3A, EpiGePT achieves an auPRC

236  of 0.618, compared to 0.531 of Enformer, and 0.568 of EpiGePT-seq in 0-12kbp group.

237  Next, we explored whether EpiGePT is also capable of predicting the promoter-silencer

238  interactions. Since there is very limited experiment-validated silencer-promoter interactions,

239  we downloaded putative silencers from the SilencerDB[25] and used the promoter of annotated

240  nearest gene as the potential target. For negative silencer-promoter pairs, we selected the same

241  promoter and equidistant genomic regions in the opposite direction to ensure the consistency

242  of distance distributions between positive and negative sample pairs at different distance levels.

243  As a result, EpiGePT achieves a better performance in discerning positive silencer-promoter

244 pairs from negative pairs than Enformer by a relatively large margin. For instance, EpiGePT

245 displays an auROC of 0.575 in long-range interactions (32-64kbp) with positive-to-negative

246 ratio 1:1 setting, compared to 0.547 of EpiGePT-seq, and 0.483 of Enformer (**Fig. 3B**).

247 According to these results, the self-attention mechanism significantly enhances the ability to

248 identify potential chromatin interactions and increases the interpretability of the model.

249 The HiChIP[26] sequencing technology provides unprecedented opportunities to uncover 3D

250 genomic interactions. We aim to investigate the predictive performance of EpiGePT on 3D

251 genome interaction based on HiChIP data. Here, we employ the same strategy as described

252 above to calculate attention scores for the regulatory element-promoter pairs and collected

253 HiChIP loops on K562 and GM12878 cell lines from the HiChIPdb[27]. The results demonstrate

254 that incorporating TF expression data into EpiGePT leads to enhanced predictive performance

255 for HiChIP loops compared to the pure sequence models Enformer and EpiGePT-seq, across

256 diverse distance ranges and in two distinct cell lines. Specifically, within the 20-40 kbp distance

257 range on K562 cell line, EpiGePT achieves an auROC of 0.599 for the 1:1 positive-to-negative

258 ratio, surpassing Enformer's performance of 0.545 (**Fig. 3E**). These findings suggest that, even

259 without any fine-tuning, EpiGePT's attention scores encompass more accurate and

260 comprehensive biological information, underscoring its potential for capturing intricate

261 genomic interactions.

262 To better understand the self-attention mechanism of EpiGePT and bridge the gap between the

263 model and its interpretability, we visualized the attention matrices after normalization (**Fig. 3C**).

264 The visualization shows prominent scores between certain genomic bins, indicating the

265 potential presence of interactions. We centered on the transcription start site (TSS) of the CHD4

266    gene and calculated the self-attention scores between the genomic bins within its upstream and

267    downstream 128kbp. The attention scores exhibited peaks near the regulatory elements in the

268    vicinity of the TSS, which further validates the feasibility and accuracy of our prediction of

269    enhancer-promoter interactions (**Fig. 3D**).

270    **EpiGePT improves variant effect prediction**

271    One of the most essential tasks for EpiGePT is to dissect the effect of genetic variants that occur

272    in different genomic regions. As most of the variants identified by the GWAS studies lie in the

273    non-coding regions of the genome, which makes it difficult to interpret the variant effect, most

274    sequence-based computational models directly take the alleles sequence as input and compare

275    the difference in the predicted regulatory activity. The advantage of EpiGePT model comes

276    from the TF module where variant effect can be estimated under any given cellular context.

277    This is extremely helpful when predicting the effect of the disease- or phenotype-associated

278    SNPs. To test the ability of EpiGePT in variant effect prediction, we first collected an eQTLs

279    dataset[28] that contains 20,913 causal and non-causal variant-gene pairs in total across 49

280    different tissues from the supplementary data of Wang et al[28]. EpiGePT and EpiGePT-seq are

281    then applied to estimate the log-odds scores (LOS) given both the reference and alterative DNA

282    sequence and the corresponding relevant TF profile (see Methods, **Fig. 4A**). Finally, a random

283    forest classifier is trained based on the LOS scores across different chromatin states. The

284    experimental results show that in the lung tissue, EpiGePT demonstrates and auPRC of 0.922,

285    compared to 0.873 of Enformer in distinguishing casual SNPs. To verify the effectiveness of

286    TF module, we replace the TF profile of lung with stomach, which is much less relevant to the

287    lung tissue. The auPRC decreases from 0.922 to 0.892 (**Fig. 4B**). Similarly, EpiGePT-seq can

288    achieve an average auPRC of 0.910, compared to 0.898 of Enformer using 5-fold cross-

289    validation for predicting causal variants on 48 extracted tissues (Fig. S3D). In the adrenal gland

290    tissue, EpiGePT-seq demonstrates an average auPRC of 0.883, compared to only 0.842 of

291    Enformer. The above experiments show the predictive power of EpiGePT in estimating the

292    variant effect.

293    To further evaluate the performance of EpiGePT in predicting disease-associated variants, we

294    extracted 52, 876 pathogenic SNPs from the ClinVar[29] database and 418, 863 benign SNPs from

295    the ClinVar database, also with 84, 095 benign SNPs from the ExAC database[30] as positive and

296    negative sets, respectively. We defined a 64kbp region surrounding each pathogenic SNP as the

297    risk region. We screened all benign and likely benign SNPs that fall within the risk region from

298    the negative sets for classification. As the relevant tissue or cell type information is not available,

299    we concatenated the LOS of the eight epigenomic signals and self-attention scores across 28

300    cell types into a single 252-dimensional vector and then train a classifier for predicting whether

301    the given SNP is pathogenic (see Methods). To assess the whether the 252-dimensional features

302    are beneficial in predicting pathogenic SNPs, we concatenated it with 52 annotations from

303    CADD[31], resulting in a comprehensive feature vector. Subsequently, we compared the

304    performance of this combined feature vector with that of the individual features derived

305    exclusively from CADD. We then utilized these two sets of features to train multi-layer

306    perceptron (MLP) classifiers separately. The results demonstrate that incorporating EpiGePT's

307    variant effect features from multiple cell types significantly enhances the performance of the

308    classifier in predicting pathogenic SNPs. Specifically, when the positive-to-negative sample

309    ratio was set to 1:1, the average auROC increased from 0.772 to 0.806, and the average accuracy

310   increased from 0.690 to 0.723 (**Fig. 4C**). This observation indicates that features extracted by

311   EpiGePT provide a valuable complement to CADD annotation, enabling a more comprehensive

312   depiction of variant characteristics, and thereby facilitating the discovery of disease-associated

313   variants.

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

315   We investigated whether the ability of EpiGePT to predict variant effects could help in the

316   discovery of key SNPs related to COVID-19. COVID-19 is an infectious disease caused by the

317   SARS-CoV-2 virus, which emerged in late 2019 and quickly spread around the world, causing

318   a global pandemic[32]. In order to validate the ability of EpiGePT in identifying key SNPs, we

319   collected GWAS data from the COVID-19 host genetics[33], including 9,484 variants. These

320   variants were derived from 4,933 patients with confirmed severe respiratory symptoms and

321   1,398,672 control individuals without COVID-19 symptoms. To validate the ability of the

322   model to identify COVID-19-related SNPs, we firstly defined a risk region around the selected

323   COVID-19-related SNPs and computed the rank of the variant score of pathogenic SNPs within

324   the surrounding benign SNPs from the ClinVar database. The expected rank for random

325   guessing (uniform distribution) is 0.5. Interestingly, we found that the average rank of COVID-

326   19-related SNPs was significantly lower than 0.5 across several tissues or cell types (**Fig. 4D**).

327   For instance, when lung expression data was employed and a 6-kbp risk region was examined,

328   the median rank was 0.250, and when expression data of esophagus squamous epithelium was

329   used the median rank was 0.333, significantly lower than 0.5 (one-side $p$-value of 0.013 under

330   one-sided Binomial Test). However, when we employed the expression data from smooth

331   muscle cells, which are a more widespread cell type with lower relevance to COVID-19, the

16

332  median rank exhibited a notable decrease to 0.381. Notably, when focusing on the 40-kbp risk

333  region, the median rank further declined to 0.850, higher than 0.5. These findings suggest that

334  EpiGePT model is able to prioritize the COVID-19-related SNPs thus shedding lights on

335  finding the potential disease-associated variants with our pretrained large language model.

336  Next, based on the aforementioned findings, we aimed to use EpiGePT to identify genes that

337  are highly related to COVID-19. Since the genetic pathology of COVID-19 is not yet clear and

338  the earliest lesion is in the lungs, we ranked all 9,484 possible SNPs using lung expression data.

339  We then identified the SNPs with the highest ranks and performed gene ontology enrichment

340  analysis on nearest genes of the 100 top ranked SNPs (**Fig. 4E**). The enrichment results revealed

341  potential biological processes that are relevant to COVID-19, such as the regulation of

342  glucokinase activity which is associated with the homeostasis of human blood glucose[34].

343  Notably, diabetes mellitus, a condition closely associated with hyperglycemia, is a typical

344  comorbidity of COVID-19[35]. Besides, among the top 10 potential genes that scored the highest,

345  we identified that the TBC1D4 gene, which regulates glucose homeostasis, is potentially

346  associated with COVID-19 comorbidities. Our findings are consistent with previous research

347  by Pellegrina et al.[36] and highlights the potential of our EpiGePT approach in discovering new

348  genetic markers that may be implicated in the pathogenesis of COVID-19. Overall, our

349  EpiGePT model provides new perspectives for understanding how the genetic variants could

350  contribute to the COVID-19 susceptibility and severity.

351  **Fine-tuning on EpiGePT enables accurate prediction of regulatory interactions**

352  Fine-tuning is an strategy that transfers the knowledge of a pretrained model to new tasks,

353    which is particularly prevalent in language models such as GPT[37] and BERT[38]. Here, we

354    finetuned a pretrained EpiGePT model on a new task for predicting the 3D genome interaction.

355    Given the HiChIP H3K27ac data from K562 and GM12878 cell lines, the features of two

356    anchors were extracted from a pretrained EpiGePT model and then fed to a finetune network

357    to predict whether it is a HiChIP loop. We compared EpiGePT with pretrain and finetuning

358    strategy to two baselines, DeepTACT[39] and a k-mer frequency[40] based method (see Methods).

359    The results illustrate that EpiGePT exhibits a superior classification performance across diverse

360    distance ranges compared to baselines. For example, in the GM12878 cell line within the 20-

361    40kbp distance range, EpiGePT demonstrates a significantly improved predictive performance

362    with an auROC of 0.949, surpassing 0.866 of DeepTACT [39] and 0.771 of Kmer (**Fig. 4F** and

363    Fig. S7). This significant improvement achieved through fine-tuning EpiGePT on a limited

364    dataset aligns well with the concept of few-shot learners[1], highlighting the power of the

365    pretrained EpiGePT model.

366    **EpiGePT encompasses the regulatory relationships between TFs and target genes.**

367    One of the key differentiating factors of EpiGePT compared to other sequence models lies in

368    its integration of TF binding status and TF expression. This unique feature empowers EpiGePT

369    to capture potential regulatory relationships embedded within the genomic sequence. In this

370    study, we specifically aimed to validate whether EpiGePT learns the regulatory relationships

371    between TFs and target genes (TGs). We defined gradient importance scores (GIS) based on

372    the absolute gradient values of predicted epigenomic signals w.r.t. TF profile to rank TFs give

373    a TG (see Methods). Particularly, we collected 15 TFs that play critical regulatory roles in

374    embryonic stem cells (ESC)[41, 42] and validated their interaction relationships using computed

375  GIS. As an example, for a target gene STAT3 that plays essential role for ESC pluripotency[43],

376  we computed GIS for each core TF across 1000 genomic bins and find other key TFs in ESC

377  ranked 1st (REST) and 2nd (POU5F1) at specific bins (**Fig. 5A-B**). Interestingly, the GO terms

378  enriched by the top 10% prioritized TF coding genes also included biological processes of

379  embryonic cell differentiation and development when we focus on the genomic bin that

380  POU5F1 ranked 2nd (**Fig. 5C**). We selected a TF as the target gene and calculated the integrated

381  GIS (IGIS) score for another key TF across eight epigenomic signals. Multiple TF-gene pairs

382  identified by IGIS that showed significant associations with POU5F1, such as ESRRB-

383  POU5F1[44] (rank 2nd), and ETV5-POU5F1[45] (rank 5th). Furthermore, we use TF-TG

384  relationships from either ChIP-seq data or external databases as ground truth to validate

385  whether IGIS is effective in prioritizing the TFs given a TG.  First, we defined potential TF-

386  target gene pairs based on TF ChIP-seq data specific to certain cell types among all human

387  genes (see Methods). The results demonstrated a significant difference in rank between TF-

388  target gene pairs and TF-non-target gene pairs based on the IGIS score (**Fig. 5D**), with the

389  former exhibiting considerably higher ranks (one-side $p$-value < 0.001). Second, we collected

390  TF-target regulatory network data from two publicly available databases. We obtained a total

391  of 1,066 TF-gene pairs from the GRNdb[46] database based on liver-specific GTEx data, and

392  2,705 TF-gene pairs from the TRRUST[47] database after filtering. Then we calculated the rank

393  of each TF based on GIS of the TF expression and a genomic bin-level mask for each pair.

394  Interestingly, when using liver expression data, we found that the average rank of TFs from

395  TRRUST was 7.9%, significantly lower than the rank based on expression values (one-side $p$-

396  value < 1e-5). Similarly, based on the GRNdb data, the majority of TF-gene pairs obtained had

397    TF ranks within the range of 20%, and the mean of this distribution was significantly lower

398    than the rank based on expression values with one-side *p*-value < 1e-5 (**Fig. 5E**). For instance,

399    TMEM55B plays a significant role in regulating lysosome movement, and is regulated by sterol

400    response element binding factor 2 (SREBF2)[48], while GIS enable the identification of SREBF2

401    as the top-ranked TF associated with TMEM55B, further validating the role of GIS in

402    prioritizing functional TFs. The comprehensive validation from both ChIP-seq datasets and

403    external databases further support the effectiveness of GIS in identifying context-specific TF-

404    TG relationships.

### Online prediction tool for EpiGePT

406    In order to facilitate the utilization of EpiGePT for the prediction of multiple chromatin

407    states of any cellular context and any genomic regions, especially for research

408    personnel who lack coding expertise, we have developed a user-friendly online web

409    server, named EpiGePT-online (https://health.tsinghua.edu.cn/epigept/) (Text S2). The

410    online web server was developed using PHP, JavaScript and HTML, which provides an

411    interactive web interface for online prediction of 8 chromatin states of specific genomic regions

412    (**Fig. 6**). Users can obtain the predicted signals of multiple genomic regions by submitting a

413    region file and a TF expression file of 711 selected TFs (Supplementary Table S3), or obtain

414    predicted signals of specific regions directly by selecting genomic locus. As to the two types of

415    input files, we provide example files to demonstrate their formats, and accept expression files

416    stored in either numpy or csv formats, to increase the universality of the web server (Fig. S5).

417    The web server outputs the results in a web summary html, which saves significant amount of

418    time for installation and implementation. Furthermore, we provide a detailed tutorial to enable

419    users to quickly learn how to use our website. We anticipate that this web server will assist

420    researchers in predicting chromatin states of specific cell types and further deepening

421    their understanding of gene regulatory mechanisms.

## Discussion

423    In this paper, we introduced EpiGePT, a transformer-based large language model, for predicting

424    the chromatin states given any cellular context. Compared with the existing machine learning

425    based computational methods, EpiGePT takes transcription factor profile and DNA sequence

426    information as inputs by multitask learning and self-attention mechanism within a unified

427    model. With these two types of input information and four modules of network architecture,

428    EpiGePT overcomes the limitation of the existing models and demonstrates state-of-art

429    performance in prediction of multiple chromatin signals in diverse experimental settings. With

430    the superior predictive performance of EpiGePT, we are able to investigate one of the

431    fundamental questions in functional genomics: how transcription factors and cis-regulatory

432    elements regulate gene activity. In this work, we investigated this question from two aspects:

433    1) identifying the interactions of cis-regulatory elements and their target genes with the help of

434    self-attention mechanism in EpiGePT; 2) estimating the variant effect based on the LOS scores

435    computed by the outputs of EpiGePT to assist in discovering human disease-associated SNPs.

436    First, the self-attention scores between tokens can provides us with an intuitive and quantitative

437    measure of the interaction level between different genomic regions, which offers new

438    opportunity to discover the target gene of cis-regulatory elements and find the interpretability

439    of EpiGePT. Second, the LOS scores of the multiple chromatin signals from different tissue or

440    cell lines are complementing each other, which provides us with a more comprehensive

21

441 characterization of the variant and enables accurate prediction of the variant impact. Such

442 variant effect prediction by EpiGePT establishes a foundation for understanding the underlying

443 relationship between genetic variations and disease mechanisms.

444 There exist several extensions and refinements that can be applied to further improve the

445 EpiGePT model. Firstly, the incorporation of chromatin regulators (CRs) as trans-acting factors

446 into the TF module could enhance the modeling of regulated transcription processes, thereby

447 increasing the accuracy of the predictions. Secondly, the inclusion of high-order interactions

448 between TFs in the framework could provide a more comprehensive representation of the

449 regulatory relationships, and potentially enhance the predictive performance. Third, the

450 application of EpiGePT to single-cell genomics could enable the profiling of chromatin signals

451 at single-cell resolution, facilitating a holistic understanding of regulatory heterogeneity in

452 different cell subpopulations for researchers.

453 Based on EpiGePT, users are able to predict multiple chromatin profiles in different cell lines

454 or tissues, which could provide a foundation for biological discovery, decoding transcriptional

455 regulation mechanisms, and investigating disease mechanisms. We anticipate that EpiGePT can

456 provide valuable insights to researchers in understanding regulatory mechanisms.

# Methods

**Data processing**

**Chromatin accessibility data and Expression data** We used three different datasets in the experiments. For chromatin accessible data, we downloaded DNase bam files and narrow peaks across 129 human biosamples from ENCODE[12] project (Supplementary table S1 and S2). We divided the human hg19 genome into 200bp non-overlapping locus (use bin instead), and we assigned the label for each locus in each cell type. For the regression design, we pooled the bam files of multiple replicates for a cell type (Supplementary table S1 and S2), and obtain the raw read count $n_{lk}$ for locus $l$ in cell type $k$. We normalized the raw read count in order to eliminate the effect of sequencing depths, in the form of $\tilde{n}_{lk} = N n_{lk}/N_k$, where $N_k$ denotes the total number of pooled reads for cell type k and $N = \min_k N_k$ denotes the minimal number of pooled reads across all cell types. The normalized read counts are further log transformed with pseudo count 1, which represent the continuous level of chromatin accessibility. For binary classification design, we assigned a binary label $y_{lk}$ to 1 if the number of raw read counts of the locus $l$ in the cell type $k$ greater than 30, which represent the locus is an accessible region in this cell type, resulting in the identification of regions as accessible in 13% on average and 8% at median in the screened genomic regions across 129 cell types. The proportion of open regions varies among different cell types, and the average openness level mentioned above is generally consistent with that maintained in ChromDragoNN[16].

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

478     pseudo count 1 and quantile normalization based on TPM values. The normalized TPM values

479     were averaged across replicates and mean expression profile of each cell type was finally used

480     to calculated of the transcription feature.

481     **Multiple chromatin signals data** DNase-seq, RNA-seq and ChIP-seq data were also

482     downloaded from ENCODE project (Supplementary table S3, S4 and S6). We applied the same

483     process to these data as above, and finally we obtained the 8 chromatin signals of 13,300,000

484     bins of 128bp in 28 cell types. The continuous level of chromatin signals we extracted were

485     'DNase', 'CTCF', 'H3K27ac', 'H3K4me3', 'H3K36me3', 'H3K27me3', 'H3K9me3' and

486     'H3K4me1', which includes crucial epigenetic modifications and markers for gene regulation

487     and transcription.

488     **Chromatin states data** We downloaded the 15-state ChromHMM[21] annotations across 127

489     epigenomes from the ROADMAP project. The state of chromatin is annotated for each 200bp

490     bin in a specific cell type. RNA-seq data across 56 cell types of TFs was download and extracted

491     from the ROADMAP[49] project (Supplementary table S7 and S8). Subsequently, we mapped

492     these 711 transcription factors to the downloaded RNA-seq data, resulting in the identification

493     of RNA-seq data for 642 transcription factors. In the subsequent experiments, we utilized the

494     expression data of these 642 transcription factors. We finally calculated the normalized TPM

495     values of the 642 TFs on 56 cell types we extracted for the using in the classification model.

496     For coarse grain chromatin state prediction, we took the state 'Quies' as low signal regions and

497     other states as signal regions. For fine grain chromatin state prediction, we extracted the state

498     'TssA', 'TssAFInk', 'TssBiv' and 'BivFInk' as TSS regions, state 'EnhG', 'Enh' and 'EnhBiv' as

499     enhancer regions, 'Quies' as low signal regions and other state as other regions. To balance the

500 number of different chromatin states, we downsampled the low signal regions and obtained

501 921,074 locus each cell line finally.

**Model architecture**

**Sequence module and Transformer module**

504 As shown in Figure 1 and Fig. S6A, the sequence module receives a one-hot matrix (A =

505 [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, representing a

506 sequence of 128 kilobase pairs (kbps) and contains five 1-dimentional convolutional blocks to

507 extract DNA sequence features. Each block includes a convolutional layer and a maxpooling

508 layer (Fig. S6B). The first convolutional layer considers the input channels as 4 and performs

509 convolution along the sequence direction. The input sequence features are one-hot embeddings

510 of size $L \times 4$, where $L$ denotes the length of the input long range DNA sequence. After 5

511 maxpooling layers, the output size of sequence feature is $L/N \times C$, where $C$ denotes the hyper-

512 parameter for sequence embedding and N denotes the length of locus to predict. We set C to

513 256 in the pre-training stage of chromatin accessibility prediction experiments. Rectified linear

514 units (ReLU) are used after each convolution operation for keeping positive activations and

515 setting negative activation values to zeros. Sequence features were then concatenated with

516 transcriptomics features, and we finally obtained a vector of size $L/N \times (C + n_{TF})$, where $n_{TF}$

517 denotes the dimension of the transcription factors features after padding. In our model, after

518 adding padding to the 711 TFs, the $n_{TF}$ is set to 712. Therefore, the input token number for the

519 transformer module is 1000, and each token embedding has a dimensionality of 968.

520 We utilize the transformer module to integrate information from both the sequence and

25

521      transcription factors (TFs), enabling the capturing of long-range interactions between genomic

522      bins. We applied $N_t$ layers of Transformer encoder with $n_{head}$ different attention heads to the

523      token embedding sequence. The input $X$ of the transformer encoder is a genomic bin sequence

524      with dimensions ($Sequence\ length, embedding\ dim$). Specifically, this dimension is (1000,

525      968) in EpiGePT, indicating that input genomic bin sequence has a length of 1000, and each

526      genomic bin has an embedded representation that combines the sequence information with cell-

527      type-specific features with dimension of 968. Each Transformer encoder includes a multi-head

528      self-attention mechanism and a feed-forward neural network. For self-attention in each head,

529      the calculation is based on the matrix operation.

530

$$Attention(Q, K, V) = softmax(\frac{QK^T}{d_k})V$$

531      For multi-head attention, Transformer encoder learns parameter matrices $W_i^Q \in$

532      $\mathbb{R}^{d_{model} \times d_Q}, W_i^K \in \mathbb{R}^{d_{model} \times d_K}$ and $W_i^V \in \mathbb{R}^{d_{model} \times d_V}$ for the $i_{th}$ head and concatenate the

533      multiple heads to do the projection.

534

$$Q_i = X \times W_i^Q, K_i = X \times W_i^K, Q_i = X \times W_i^V$$

535      Where $d_{model}$ denotes the dimension of token in the input sequence $X$, which is 968 in

536      EpiGePT and $d_Q = d_K = d_V = 512$. The matrices $Q$, $K$, and $V$ are obtained by the application

537      of mapping functions represented by $W_i^Q, W_i^K$ and $W_i^V$, followed by concatenating of $Q_i$ to $Q$,

538      $K_i$ to $K$, and $V$ to $V$. These mapping functions serve to transform the concatenated embeddings

539      into the resulting matrices. We set $N_t$ to 16 for the chromatin accessible prediction experiments,

540      $N_t$ to 12 for the chromatin state classification and multiple chromatin signals prediction

541      experiments, and set $n_{head}$ to 8 for all experiments.

542 The regression model, the output layer uses a linear transformation and use mean square error

543 (MSE) as the loss function. For classification model, the output layer uses a linear

544 transformation combined with a sigmoid function, and use the cross-entropy loss for

545 classification experiments.

**TF module**

547 For binding status, we scanned the input bins for potential binding sites for a set of 711 human

548 transcription factors from HOCOMOCO database[50] with the tool Homer[51] (Table S5). We then

549 selected the maximum score of reported binding status for each transcription factor to obtain a

550 vector of 711 dimensions as the motif feature for each DNA bin. For gene expression, we

551 focused on log-transformed TPM values of the 711 transcription factors and obtained a vector

552 of 711 dimensions after quantile normalization as the expression feature. With these data, we

553 combined the two vectors of motif and expression features by taking the element-wise product,

554 and we concatenated the result to the output of sequence module.

**Model evaluation**

556 To evaluate our model, we applied five-fold cross-validation in the different experiments on

557 cell-type level. For chromatin accessible experiments, the 129 cell lines are partitioned into a

558 training set and a testing set randomly.

559 Cell-type-wise metrics are defined to evaluate our method in different experiments, which were

560 calculate with the data within a test cell type across all genomic locus. For binary classification

561 design, we used cell-type-wise auPRC and auROC to evaluate our EpiGePT. Let $Y_{L \times K}$ and $\hat{Y}_{L \times K}$

562 be the true and predicted matrix, where L denotes the number of locus and K denotes the number

563    of test cell types. We calculated the auPRC and auROC for each $(y_{1i}, y_{2i}, \cdots, y_{Li})$ and

564    $(\hat{y}_{1i}, \hat{y}_{2i}, \cdots, \hat{y}_{Li})$. For multiple classification, we use macro average of the auPRC and auROC

565    to evaluate the classification performance, which compute the metric independently for each

566    class and then take the average hence treating all classes equally. For regression design, we

567    used two metrics for model evaluation, which are cell-type-wise Pearson correlation coefficient

568    and prediction squared error. Prediction square error (PSR) is calculated as $PSR = 1 -$

569    $\sum_k \sum_l (y_{lk} - \hat{y}_{lk})^2 / (y_{lk} - \bar{y}_{*k})^2$, where $\bar{y}_{*k} = \sum_l y_{lk} / L$ denotes the mean of the true level of

570    the response in the cell type k.

571    To compare the performance of our method with other baseline methods, we conducted

572    hypothesis testing on the metrics based on cell types. Since the metrics on a given cell type

573    across different methods are paired data and the statistical distribution is unknown, we

574    employed both Binomial and Wilcoxon tests, with the alternative hypothesis being that

575    EpiGePT outperforms the other methods. If we reject the null hypothesis, it provides

576    compelling evidence to support the claim that EpiGePT performs better than the other methods.

577    To evaluate the computational efficiency, we recorded the running time of a single epoch of

578    EpiGePT and baseline methods (Supplementary Text S3). Compared to traditional CNN models

579    such as DeepCAGE[17] and ChromDragoNN[16], as well as larger sequence models like Enformer,

580    EpiGePT demonstrates a balance between high computational efficiency and performance.

581    **Model fine-tunning**

582    For the fine-tuning process, we kept the parameters of the pre-trained model fixed without

583    making any updates. For the specific fine-tuning task of chromatin interaction prediction based

584   on HiChIP data, the multi-task module was replaced with a two-layer MLP network, containing

585   256 hidden nodes for each layer. During the training process, only the weights in the MLP

586   network were updated. Notably, when utilizing HiChIP data at a resolution of 5k, both the

587   enhancer and promoter anchors spanned 5kbp. Then we input a region extending 128kbp from

588   the center of the anchor of the neighboring gene into the EpiGePT. Consequently, a 968-

589   dimensional feature vector for each genomic bin was derived from the output of the last

590   transformer encoder layer. These feature vectors from all bins within the two anchors were

591   concatenated, resulting in a high-dimensional vector of size 76,472.

592   **Baseline methods**

593   Four baselines were introduced for epigenetic signals prediction. BIRD[15] is a multiple linear

594   regression model that only takes gene expression data as input and makes predictions on a fixed

595   locus. ChromDragoNN[16] is a deep neural network that takes gene expression of 1630 TFs and

596   DNA sequence as input. Specifically, ChromDragoNN[16] uses a ResNet[52] to extract sequence

597   features and use linear transformation to combine the TF gene expression feature and sequence

598   feature to make the final prediction. DeepCAGE[17] Integrating regulatory DNA sequence is a

599   deep densely connected convolutional network for predicting chromatin accessibility. The

600   dense-connected neural network architecture used by DeepCAGE[17] may struggle to capture the

601   complex interactions between genomic regions. Enformer[11] is a deep neural network that

602   integrates convolutional neural network and transformer, and only takes DNA sequence as input.

603   Enformer takes DNA sequence of length 196kbp as input to predict 5,313 genomic tracks of

604   human and 1,643 tracks of mouse genome simultaneously. However, one of the limitations of

605   Enformer is that it can only model and predict cell types in the training data and cannot be

606    applied to new cell types. In order to ensure the fairness of the benchmark experiment, we

607    retrained the Enformer model with the same input and output data as EpiGePT when reproduce

608    the Enformer model (Text S4).

609    Two baseline methods were introduced for predicting HiChIP interaction. DeepTACT[39] is a

610    deep learning method for predicting 3D chromatin contacts using both DNA sequence and

611    chromatin accessibility. We adopted the structure of DeepTACT[39] and kept the anchor length at

612    5k. The input to the model consists of two anchor sequences represented as one-hot matrices

613    and the two openness scores of the two anchors on the corresponding cell type extracted from

614    OpenAnnotate[53]. Regarding the Kmer features[40], K is chosen as 5 to extract sequence features.

615    For each anchor, a vector of dimension $4^5 = 1024$ was obtained. Further training was

616    performed using an MLP with a hidden layer dimension of 256.

## Enhancer, Silencer and HiChIP loop prioritization

618    We collected cis-regulatory elements-gene pairs in K562 cells from other studies and public

619    database to demonstrate the interpretability of self-attention mechanisms in the EpiGePT.

620    Enhancers and silencers are typical *cis*-regulatory elements known play important roles in

621    transcriptional control during normal development and disease. For enhancers, we downloaded

622    enhancer-gene pairs from two studies: Gasperini et al.[23] and Fulco et al.[24], both of which were

623    tested using a CRISPRi[22] assay perturbation. Two datasets contain 664 and 5,091 enhancer-

624    promoter interactions or element-TSS interactions. For silencers, we obtained and random

625    sampled 831 validated silencers-gene pairs with distance within 64kbp in K562 cells curated

626    from high-throughput experiments from SilencerDB. As there are no experimentally validated

627    interaction relationships between these silencers and genes, we generated silencer-gene pairs

628    by associating the nearest neighbor genes for classification purposes. Similarly, negative

629    samples were generated by constructing DNase-seq, ATAC-seq and nearest genes using the

630    same approach. Ultimately, we obtained a dataset comprising 1,662 silencer-gene pairs,

631    encompassing both positive and negative instances.

632    To obtain scores for regulatory element-gene pairs, we first used the region extending 128kbp

633    from the center of the enhancer as input and extracted the token where the interacting genes

634    reside, so that we could filter out regulatory element-gene pairs that were located further than

635    64kbp apart. Subsequently, we stratified the remaining pairs based on their distance. Since the

636    positive and negative sample ratios varied across datasets, we adopted different stratification

637    strategies for different distance ranges (**Fig. 3**). Next, we averaged the attention matrices of the

638    Transformer encoder across all layers and heads. The summed attention scores from other

639    tokens to the key token containing the gene TSS were used as the attention score of this element-

640    gene pair. This score represents the attention value that the enhancer-centered region receives

641    for the transcription start site (TSS) of the gene. We also calculated the attention score from the

642    bin containing the center of the regulatory element to the bin containing the TSS, which only

643    slightly affects the experimental results of regulatory element prioritization.

644    We collected 5k resolution data from the HiChIPdb (http://health.tsinghua.edu.cn/hichipdb/)

645    database, specifically from K562 and GM12878 cell lines. We filtered the data to include only

646    loops where at least one anchor falls within a gene region. We stratified the loops based on

647    distance into three categories: 0-20kbp, 20-40kbp, and 40-64kbp. For each distance category,

648    we selected 2000 positive pairs with most significant q-value. To ensure consistency in the

649 distance distribution, we selected negative pairs by fixing a gene and choosing anchors at

650 equidistant locations in the opposite direction.

## Gradient importance scores

652 EpiGePT possesses the capability to assign priority rankings to transcription factors by utilizing

653 gradient importance scores (GIS), taking into account specific cell types and chromatin regions.

654 The GIS were employed to identify potential functional relationships between specific

655 transcription factors (TFs) and target genes. Specifically, for a given TF-target gene pair, the

656 transcription start sites (TSS) of genes were used as central loci, and the regions spanning 128

657 kbp upstream and downstream of the TSS were selected as input. Next, we filtered out bins

658 with motif binding scores indicating potential binding for the given TF. For these selected bins,

659 we calculated the GIS for the predictions of eight epigenomic signals across the 711 core TFs.

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

661 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

662 predicted epigenomic signal, and $\zeta$ denotes the set of genomic bins that have binding for the

663 given TF. In the calculation of the gradient, $\hat{y}_{ljk}$ denotes the predicted value of the $k$th

664 epigenomic signal by the model using the expression in the $j$th cell type at the $l$th bin. On the

665 other hand, $tf_{ij}$ denotes the product of the expression of $i$th TF in the $j$th cell type and the

666 corresponding TF binding score.

667 If we consider the GIS for the prediction of all 8 epigenomic signals simultaneously, we can

668 prioritize the TFs by calculating their ranks based on each signal separately. Then, we can

669    calculate an integrated gradient importance score (IGIS) for each TF by aggregating the ranks

670    from all 8 signals.

671
$$IGIS_{ij} = \frac{1}{8} \sum_k rank(GIS_{ijk})$$

672    Both the GIS and the IGIS are capable of capturing the significance of a transcription factor

673    (TF) in regulating a specific gene within the context of a specific cell type. Consequently, these

674    scores hold potential value in the discovery of TFs that play crucial roles in the regulation of

675    specific genes, thereby contributing to our understanding of essential regulatory mechanisms.

676    In the context of validating TF-TG pairs in the GRNdb and TRRUST databases, we opted to

677    utilize liver expression data as a representative example due to the unavailability of cell type

678    information for TRRUST. Furthermore, in this experimental setup, the $tf_{ij}$ denotes the

679    expression of $i$th TF in the $j$th cell type and $\zeta$ denotes the set of genomic bins that have binding

680    for the TF of the given TF-target gene pair.

## Potential TF-target gene pairs from ChIP-seq data

682    In this study, we utilized three distinct cell types to conduct a comprehensive screening of TF-

683    target gene pairs and non-target gene pairs across the human genome. Initially, we obtained the

684    narrow peak files (ENCFF388AJH, ENCFF717IXP, and ENCFF885KLR) from ChIP-seq

685    experiments across three cell types from the ENCODE project. Subsequently, we meticulously

686    examined the number of peaks within a 128kbp region both upstream and downstream of the

687    transcription start site (TSS) for each gene. Different thresholds were applied to the ChIP-seq

688    data of various TFs. Genes lacking any peaks within the defined region were classified as non-

689   target genes, while genes surpassing the threshold in terms of peak counts were designated as

690   target genes. Specifically, for the aforementioned three cell types, threshold values of 10, 15,

691   and 6 were respectively employed. Finally, the IGIS approach was employed to determine the

692   corresponding ranks of TFs in the TF-target gene pairs.

## Pathogenic SNPs prioritization

694   We collected single nucleotide polymorphisms (SNPs) data from the ClinVar and ExAC

695   databases, which include both potentially pathogenic and benign SNPs. To evaluate the ability

696   of EpiGePT to predict variant effects, we computed the log-odds scores (LOS) for multiple

697   chromatin signals using EpiGePT on these SNPs. Subsequently, we utilized these scores to

698   distinguish between pathogenic and benign SNPs. The LOS score for each chromatin signal

699   was defined by computing a forward pass through the model using the reference and alternative

700   alleles.

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

702   Each chromatin epigenomic profile in each cell line or tissue predicted by EpiGePT can be used

703   to compute a specific variant score. We did not take the absolute value in this calculation, so

704   the resulting LOS score indicates the direction of change in the model output after the

705   appearance of the variant. In addition to the predicted chromatin signals output by the eight

706   models, attention score changes based on self-attention are also noteworthy. We computed the

707   log-odds scores for attention by summing the attention scores of the 10 bins upstream and

708   downstream of the locus of the SNP, to evaluate the effect of the variant.

709
$$\Delta O_{attention} = \sum_{i=-5}^{5} |log\left(\frac{attn(bin_i)_{I(alt)}}{attn(bin_i)_{I(ref)}}\right)|$$

710 Where $i$ represents the index of the neighboring bins relative to the locus of the SNP. To avoid

711 the variant effects of different bins from cancelling each other out during the summation process,

712 we computed the absolute value of the change in attention scores for each bin and then summed

713 the scores of the 10 adjacent bins centered at the SNP position. For the classification of

714 pathogenic SNPs, we calculated these nine LOS for attention separately for each of the 28

715 tissues or cell lines in training data. As a result, we obtained a feature vector of 252 dimensions

716 for each SNP. Then a classifier with 252 features computed by EpiGePT and 52 annotations

717 from CADD score as inputs are used to predict pathogenic SNPs against benign or likely benign

718 SNPs. Here, we employed MLP as classifier to validate the effectiveness of the features we

719 obtained. A five-fold cross-validation experiment is employed for validation, and we utilize two

720 different positive-to-negative sample ratios, namely 1:1 and 1:2. For each sample ratio, we

721 randomly sample 32,000 positive samples. The effectiveness of the variant score in identifying

722 pathogenic SNPs is evaluated using the area under the auROC and the auPRC. Additionally,

723 we also utilized the logistic regression (LR) as the classifier, consistent with the LR classifier

724 used in CADD, and found a similar improvement when predicting pathogenic SNPs.

725 We applied the same method to calculate the LOS scores of the 8 predicted chromatin signals

726 for the COVID-19 GWAS data. The absolute values of the scores were summed as the overall

727 score for each SNP. For each significant SNP associated with COVID-19 severity obtained

728 from the GWAS data, we selected normal SNPs within a 64kb region around the SNP as

729 background to calculate the rank of the LOS score for the COVID-19 associated SNP in this

35

730  region. Furthermore, we calculated the LOS scores for all 9, 484 COVID-19 associated SNPs

731  and ranked them accordingly. The top 10 SNPs with the highest LOS scores were selected,

732  which are considered to have potential genetic associations with COVID-19 severity and

733  complications.

## GTEx classification

735  We collected eQTL data from the supplementary materials of Wang et al[28]. In their study, the

736  authors identified causal eQTLs through statistical fine-mapping, using a posterior inclusion

737  probability (PIP) threshold of >0.9 for putative causal variants based on expression modifier

738  score (EMS), and a PIP threshold of <0.9 for putative non-causal variants. To validate the ability

739  of EpiGePT to distinguish potential causal variants, we perform a classification task on these

740  variants. For each variation, 128kbp sequence regions near it were selected as the input of the

741  model, and a score of variation was given by EpiGePT model. For each variant under each

742  tissue, we can obtain an 8-dimensional vector of genomic features including DNase, CTCF and

743  other ChIP-seq signals. Based on the LOS score, separate random forest classifiers consisting

744  of 10 decision trees are trained for each tissue in order to distinguish between causal and non-

745  causal variants. The models are evaluated using 5-fold validation on each tissue, with area under

746  the auPRC and auROC as metrics for assessing their ability to distinguish between causal and

747  non-causal variants.

## Code availability

749  All components of EpiGePT are freely available at https://github.com/ZjGaothu/EpiGePT.

36

750    Here, users can access the code for reproducing EpiGePT, as well as the data collection and

751    preprocessing pipelines used for model training in benchmark experiments.

## Data availability

753    Information and processed data of multiple chromatin signals of whole genome, motif binding

754    status and expression data of TFs in the corresponding cell lines/tissues, which are used in

755    EpiGePT are available at Supplementary Tables. The information about the cell lines/tissues

756    used and the 711 filtered transcription factors is available in the supplementary table. The High

757    throughput    validated    silencers    of    K562    cell    line    are    download    from    SilencerDB

758    (http://health.tsinghua.edu.cn/silencerdb) database. The HiChIP data of K562 cell line and

759    GM12878 cell line are downloaded from HiChIPdb (http://health.tsinghua.edu.cn/hichipdb/)

760    database. The DNase-seq peak and ATAC-seq peak data are obtained from the ENCODE

761    project. Enhancer-gene pairs of CRISPRi[24] experiments are obtained from the supplementary

762    information of Gasperini et al. and Fulco et al. The regulatory network data for transcription

763    factors    and    target    genes    were    obtained    from    the    TRRUST[47]    database

764    (https://www.grnpedia.org/trrust/) and the GRNdb[46] database (http://www.grndb.com). The

765    annotated    chromatin    states    for    whole    genome    are    downloaded    from    the    ROADMAP

766    epigenomics    project    (https://egg2.wustl.edu/roadmap/web_portal/chr_state_learning.html).

767    The RNA-seq read counts matrix for protein coding genes used for the prediction of the

768    chromatin 15-states annotated by ChromHMM are downloaded from the ROADMAP project

769    (https://egg2.wustl.edu/roadmap/data/byDataType/rna/expression/57epigenomes.N.pc.gz).

770    The GWAS data of COVID-19 are download from the COVID-19 Host Genetics Initiative

771    (https://www.covid19hg.org/).

## Competing interests
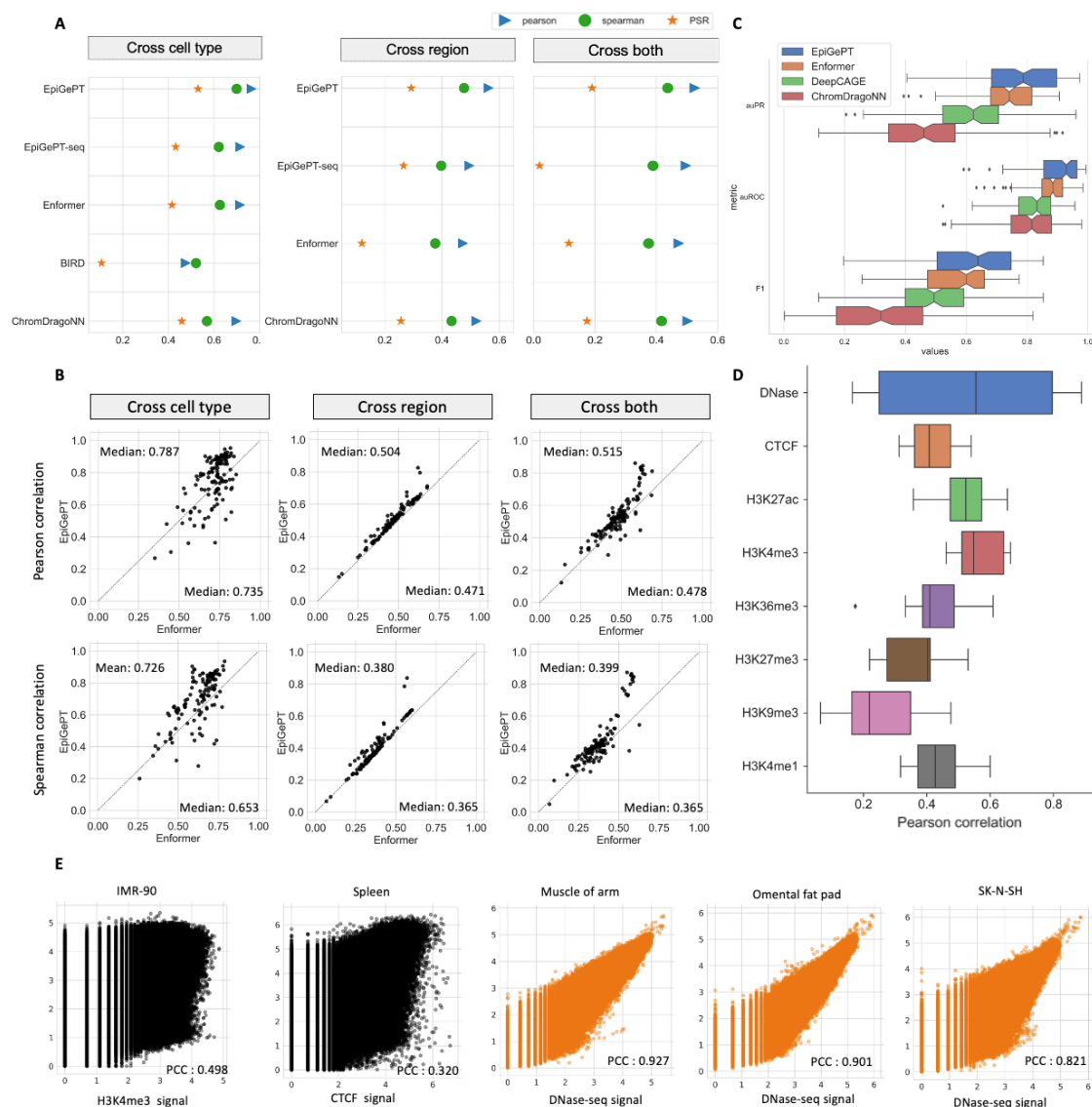
The authors have declared no competing interests.

## Acknowledgments

## Figures

### Figure 1



**Fig. 1 Overview of the EpiGePT model for multiple epigenomic signals prediction**. The model consists of four modules, namely the Sequence module, the TF module, the Transformer module, and the Multi-task module. The sequence module comprises multiple layers of convolution applied to the one-hot encoded DNA sequence input. The input sequence length consists of 1000 genomic bins of 128bp for the prediction of multiple signals and 50 bins of 200bp for the prediction of DNase signal alone. The TF module encompasses the binding status and expression of 711 transcription factors. The Transformer module consists of a series of consecutive transformer encoders, while the multi-task module is composed of a fully connected layer.

**Figure 2**



**Fig. 2 Performance of EpiGePT and baseline methods on the benchmark experiment.** (A)

EpiGePT and baseline methods were compared in terms of their regression performance for

DNase signal regression across cell types, genomic regions, and combined cell type and

genomic regions. (B) Comparison of EpiGePT and Enformer performance. Each point in the

scatter plot represents the performance of Enformer on the data of a specific cell type (x-axis)

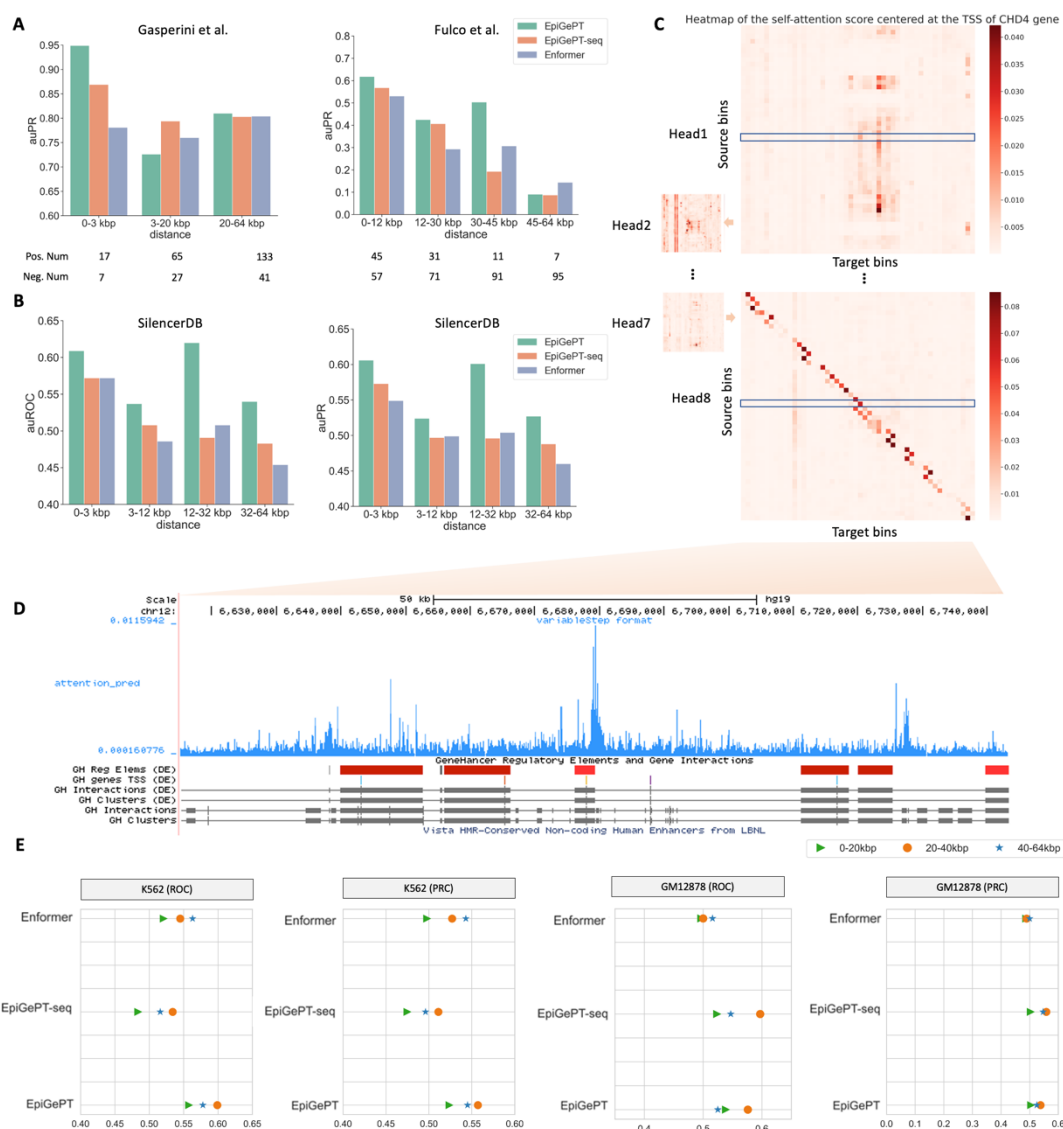compared to the performance of EpiGePT (y-axis). (C) EpiGePT and baseline methods'

performance on binary prediction of DNase-seq signals. (D) EpiGePT demonstrates excellent

799    performance in predicting diverse epigenetic signals across various cell types, including

800    DNase-seq, CTCF, and histone modifications. (E) EpiGePT predictions compared to

801    experimental signals visualized for a representative example. Genome-wide multi-signal

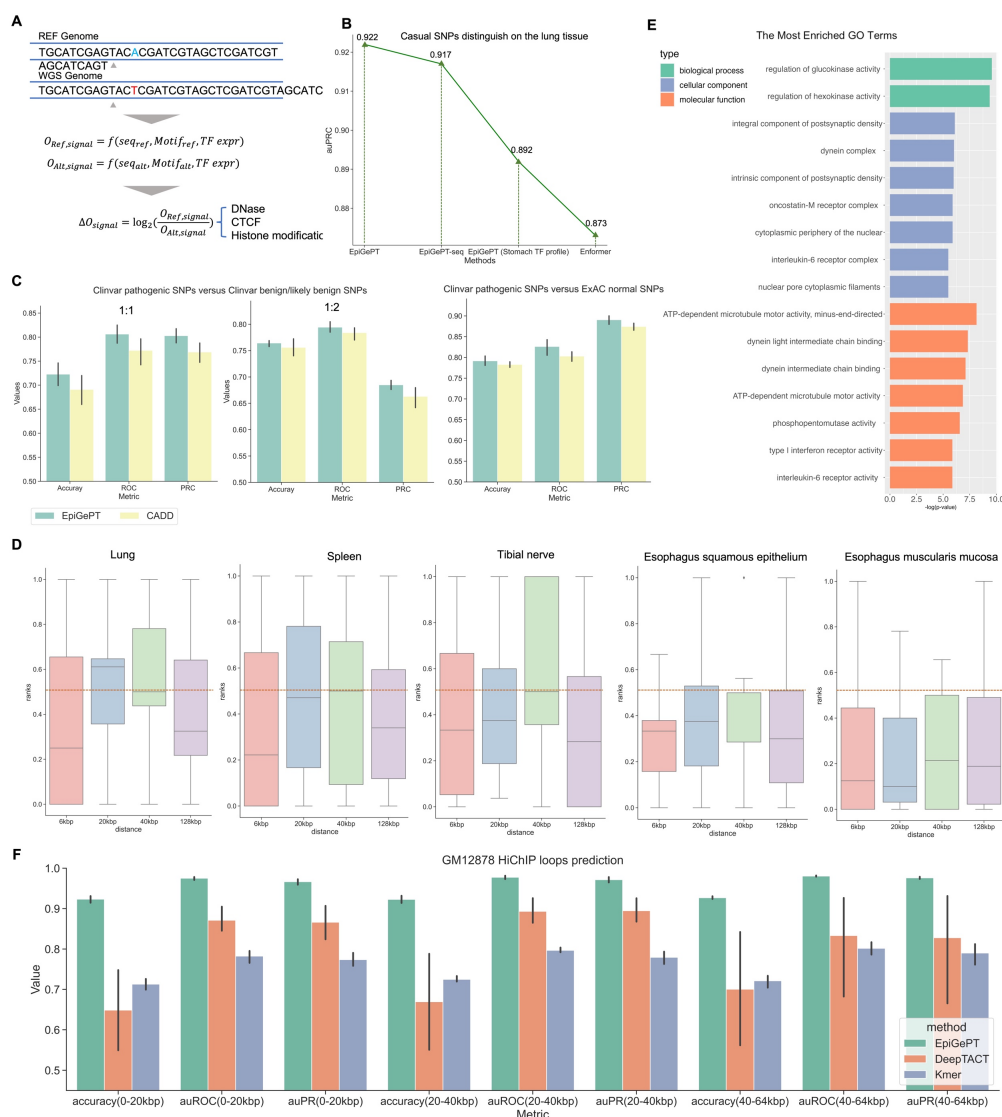802    predictions (black) and DNase-specific predictions (orange).

41

## Figure 3



**Fig. 3 Application of self-attention mechanism in EpiGePT for long-range chromatin interaction identification.** (A) The performance (auPRC) of attention score of EpiGePT in distinguishing enhancer-gene pairs at different distance ranges on two different datasets. (B) The performance (auROC and auPRC) of attention score of EpiGePT in distinguishing silencer-gene pairs at different distance ranges based on the data from SilencerDB[25]. (C) Heatmap of the self-attention matrix of each attention head centered at the TSS of the CHD4 gene, the ($i$, $j$) element in the matrix denotes the average attention score between the $i$th genomic bin and the

42

810    *j*th genomic bin across all layers. (D) Attention scores centered at the TSS of the CHD4 gene,

811    and putative enhancer regions in its vicinity. (E) The performance (auROC and auPRC) of

812    attention score of EpiGePT in distinguishing HiChIP loops of H3K27ac at different distance

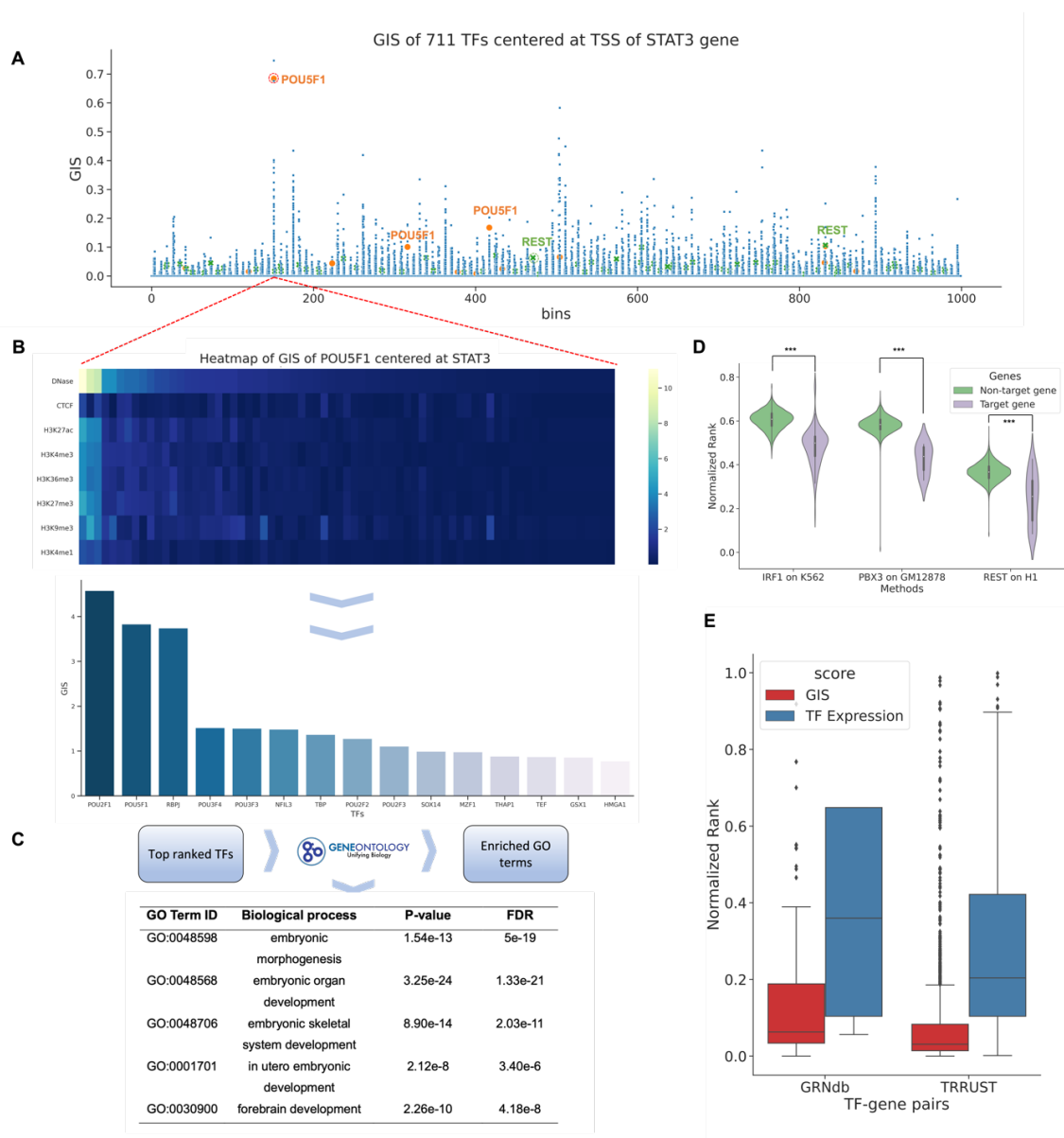813    ranges on K562 cell line and GM12878 cell line.

814 **Figure 4**



815 **Fig. 4 Variant effect prediction of EpiGePT**. (A) The LOS score for each epigenomic signal

816 is calculated by the log change fold of the predicted epigenomic signal for reference genome

817 and WGS genome. (B) The performance of EpiGePT and Enformer in discriminating causal

818 SNPs on the Lung tissue. (C) The three subplots from left to right respectively depict the

819 classification results for disease-related SNPs and benign SNPs down-sampled sourced from

820 the ClinVar database, with balanced positive and negative samples (1:1 and 1:2 ratio), as well

821    as normal SNPs sourced from the ExAC database with a MLP classifier. (D) The ranked

822    position of COVID-19 related GWAS data among surrounding benign SNPs based on their

823    LOS scores, as determined using different tissue or cell-type expression data. The results were

824    stratified based on the distance range of the risk region. The resulting mean and median ranks

825    were both below 0.5. (E) Enrichment result (Biological process, Cellular component and

826    Molecular function) of the nearest genes of the COVID-19 associated SNPs with the max LOS

827    scores. (F) The performance (auROC and auPRC) of the fine-tuned EpiGePT model and

828    baseline methods (DeepTACT and Kmer) in distinguishing enhancer-gene pairs at various

829    distance ranges (0-20 kbp, 20-40 kbp and 40-64 kbp) on K562 cell line.

830 **Figure 5**



831 **Fig. 5 Gradient importance scores (GIS) uncover regulatory transcription factors. (A)**

832 Genomic regions around TSS of the STAT3 gene and expression data on ESC were fed into

833 EpiGePT. The scatter plot represents the GIS scores of core TFs on each genomic bin. Each dot

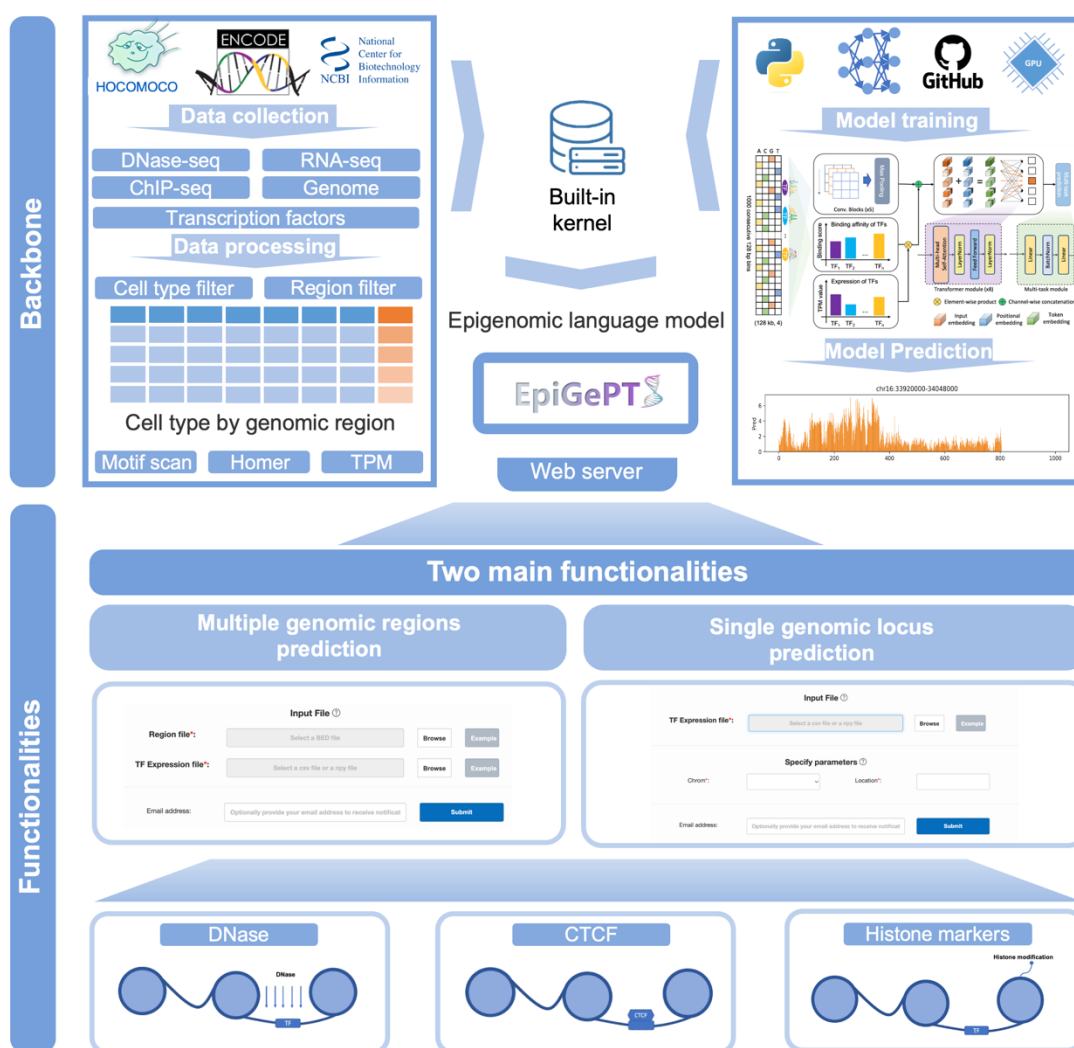834 represents the GIS score of a core TF on a specific genomic bin. In specific bins, key TFs in

835 ESC, such as NR5A2 and POU1F5 highlighted in the figure, exhibit high ranks in the GIS

836 scores. (B) Heatmap of GIS for 10% important TFs surrounding STAT3 gene, specifically

837 focusing on bins with POU1F5 ranked 3rd. Each row represents a predicted epigenomic signal,

46

838    and the TFs are sorted based on their GIS on DNase signals. Bar plot of the top 15 TFs with the

839    highest GIS scores. (C) Based on the top 10% ranked TFs mentioned above, gene ontology

840    enrichment analysis revealed significant enrichment in biological processes related to

841    embryonic development and cellular differentiation. (D) Based on TF ChIP-seq data, all 23,635

842    human genes were classified into target genes and non-target genes. The results revealed that

843    TFs exhibited significantly higher ranks on potential target genes compared to non-target genes.

844    (E) The distribution of the rank of TFs in the GIS and expression value among the 2,705 TF-

845    gene pairs from the TRRUST database and 1,066 TF-gene pairs derived from genotype-tissue

846    expression (GTEx) data of the liver sourced from the GRNdb database. The analysis reveals

847    that the median rank of TFs from the TRRUST database is significantly lower than 0.06 (one-

848    side $p$-value < 3.12 e-18) and the median rank of TFs obtained from the GRNdb database is

849    significantly lower than 0.5 (one-side $p$-value < 2.50 e-140).

850 **Figure 6**



851 **Fig. 6 Overview of the online prediction web server of EpiGePT.** We collected eight types

852 of epigenetic genome modification signals and corresponding expression data of transcription

853 factors in different cell types or tissues from the ENCODE project. Based on these data, we

854 trained the EpiGePT model and deployed it as a built-in kernel on an Apache server. Users

855 without much coding experience can also access the web server in two ways to obtain the eight

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

857 without programming or installation.

858

48

## Reference

1.  Brown, T. et al. Language models are few-shot learners. *Advances in neural information processing systems* **33**, 1877-1901 (2020).

2.  Ramesh, A. et al. in International Conference on Machine Learning 8821-8831 (PMLR, 2021).

3.  Alexander, R.P., Fang, G., Rozowsky, J., Snyder, M. & Gerstein, M.B. Annotating non-coding regions of the genome. *Nature Reviews Genetics* **11**, 559-571 (2010).

4.  O'Malley, R.C. et al. Cistrome and epicistrome features shape the regulatory DNA landscape. *Cell* **165**, 1280-1292 (2016).

5.  Yenduri, G. et al. Generative Pre-trained Transformer: A Comprehensive Review on Enabling Technologies, Potential Applications, Emerging Challenges, and Future Directions. *arXiv preprint arXiv:2305.10435* (2023).

6.  Liu, Q., Xia, F., Yin, Q. & Jiang, R. Chromatin accessibility prediction via a hybrid deep convolutional neural network. *Bioinformatics* **34**, 732-738 (2018).

7.  Min, X., Zeng, W., Chen, N., Chen, T. & Jiang, R. Chromatin accessibility prediction via convolutional long short-term memory networks with k-mer embedding. *Bioinformatics* **33**, i92-i101 (2017).

8.  Yin, Q., Wu, M., Liu, Q., Lv, H. & Jiang, R. DeepHistone: a deep learning approach to predicting histone modifications. *BMC genomics* **20**, 11-23 (2019).

9.  Liu, Q., Lv, H. & Jiang, R. hicGAN infers super resolution Hi-C data with generative adversarial networks. *Bioinformatics* **35**, i99-i107 (2019).

10. Zeng, W., Wu, M. & Jiang, R. Prediction of enhancer-promoter interactions via natural language processing. *BMC Genomics* **19**, 84 (2018).

11. Avsec, Ž. et al. Effective gene expression prediction from sequence by integrating long-range interactions. *Nature methods* **18**, 1196-1203 (2021).

12. Consortium, E.P. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57 (2012).

13. Vaswani, A. et al. in Advances in neural information processing systems 5998-6008 (2017).

14. Song, L. & Crawford, G.E. DNase-seq: a high-resolution technique for mapping active gene regulatory elements across the genome from mammalian cells. *Cold Spring Harbor Protocols*

888         **2010**, pdb. prot5384 (2010).

889    15.   Zhou, W. et al. Genome-wide prediction of DNase I hypersensitivity using gene expression.

890         *Nature communications* **8**, 1-17 (2017).

891    16.   Nair, S., Kim, D.S., Perricone, J. & Kundaje, A. Integrating regulatory DNA sequence and gene

892         expression to predict genome-wide chromatin accessibility across cellular contexts.

893         *Bioinformatics* **35**, i108-i116 (2019).

894    17.   Liu, Q., Hua, K., Zhang, X., Wong, W.H. & Jiang, R. DeepCAGE: incorporating transcription

895         factors in genome-wide prediction of chromatin accessibility. *Genomics, Proteomics &*

896         *Bioinformatics* **20**, 496-507 (2022).

897    18.   Holwerda, S.J.B. & de Laat, W. CTCF: the protein, the binding partners, the binding sites and

898         their chromatin loops. *Philosophical Transactions of the Royal Society B: Biological Sciences*

899         **368**, 20120369 (2013).

900    19.   Bannister, A.J. & Kouzarides, T. Regulation of chromatin by histone modifications. *Cell*

901         *research* **21**, 381-395 (2011).

902    20.   Karolchik, D. et al. The UCSC genome browser database. *Nucleic acids research* **31**, 51-54

903         (2003).

904    21.   Ernst, J. & Kellis, M. Chromatin-state discovery and genome annotation with ChromHMM.

905         *Nature protocols* **12**, 2478-2492 (2017).

906    22.   Larson, M.H. et al. CRISPR interference (CRISPRi) for sequence-specific control of gene

907         expression. *Nature protocols* **8**, 2180-2196 (2013).

908    23.   Gasperini, M. et al. A genome-wide framework for mapping gene regulation via cellular genetic

909         screens. *Cell* **176**, 377-390. e319 (2019).

910    24.   Fulco, C.P. et al. Activity-by-contact model of enhancer–promoter regulation from thousands

911         of CRISPR perturbations. *Nature genetics* **51**, 1664-1669 (2019).

912    25.   Zeng, W. et al. SilencerDB: a comprehensive database of silencers. *Nucleic acids research* **49**,

913         D221-D228 (2021).

914    26.   Mumbach, M.R. et al. HiChIP: efficient and sensitive analysis of protein-directed genome

915         architecture. *Nature methods* **13**, 919-922 (2016).

916    27.   Zeng, W., Liu, Q., Yin, Q., Jiang, R. & Wong, W.H. HiChIPdb: a comprehensive database of

917         HiChIP regulatory interactions. *Nucleic Acids Research* **51**, D159-D166 (2023).

918   28.   Wang, Q.S. et al. Leveraging supervised learning for functionally informed fine-mapping of cis-
919         eQTLs identifies an additional 20,913 putative causal eQTLs. *Nature Communications* **12**, 3394
920         (2021).

921   29.   Landrum, M.J. et al. ClinVar: public archive of interpretations of clinically relevant variants.
922         *Nucleic acids research* **44**, D862-D868 (2016).

923   30.   Karczewski, K.J. et al. The ExAC browser: displaying reference data information from over 60
924         000 exomes. *Nucleic acids research* **45**, D840-D845 (2017).

925   31.   Rentzsch, P., Witten, D., Cooper, G.M., Shendure, J. & Kircher, M. CADD: predicting the
926         deleteriousness of variants throughout the human genome. *Nucleic acids research* **47**, D886-
927         D894 (2019).

928   32.   Li, J., Lai, S., Gao, G.F. & Shi, W. The emergence, genomic diversity and global spread of
929         SARS-CoV-2. *Nature* **600**, 408-418 (2021).

930   33.   org, C.-H.G.I.a.b. The COVID-19 host genetics initiative, a global initiative to elucidate the role
931         of host genetic factors in susceptibility and severity of the SARS-CoV-2 virus pandemic.
932         *European Journal of Human Genetics* **28**, 715-718 (2020).

933   34.   Agius, L. Targeting hepatic glucokinase in type 2 diabetes: weighing the benefits and risks.
934         *Diabetes* **58**, 18-20 (2009).

935   35.   Singh, A.K., Gupta, R., Ghosh, A. & Misra, A. Diabetes in COVID-19: Prevalence,
936         pathophysiology, prognosis and practical considerations. *Diabetes & Metabolic Syndrome:
937         Clinical Research & Reviews* **14**, 303-310 (2020).

938   36.   Pellegrina, D., Bahcheli, A.T., Krassowski, M. & Reimand, J. Human phospho-signaling
939         networks of SARS-CoV-2 infection are rewired by population genetic variants. *Molecular
940         Systems Biology* **18**, e10823 (2022).

941   37.   Radford, A., Narasimhan, K., Salimans, T. & Sutskever, I. Improving language understanding
942         by generative pre-training.  (2018).

943   38.   Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. Bert: Pre-training of deep bidirectional
944         transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

945   39.   Li, W., Wong, W.H. & Jiang, R. DeepTACT: predicting 3D chromatin contacts via bootstrapping
946         deep learning. *Nucleic acids research* **47**, e60-e60 (2019).

947   40.   Chor, B., Horn, D., Goldman, N., Levy, Y. & Massingham, T. Genomic DNA k-mer spectra:

948        models and modalities. *Genome biology* **10**, 1-10 (2009).

949   41.   Zhou, Q., Chipperfield, H., Melton, D.A. & Wong, W.H. A gene regulatory network in mouse

950        embryonic stem cells. *Proceedings of the National Academy of Sciences* **104**, 16438-16443

951        (2007).

952   42.   Sharov, A.A. et al. Identification of Pou5f1, Sox2, and Nanog downstream target genes with

953        statistical confidence by applying a novel algorithm to time course microarray and genome-

954        wide chromatin immunoprecipitation data. *BMC genomics* **9**, 1-19 (2008).

955   43.   Raz, R., Lee, C.-K., Cannizzaro, L.A., d'Eustachio, P. & Levy, D.E. Essential role of STAT3 for

956        embryonic stem cell pluripotency. *Proceedings of the National Academy of Sciences* **96**, 2846-

957        2851 (1999).

958   44.   van den Berg, D.L. et al. An Oct4-centered protein interaction network in embryonic stem cells.

959        *Cell stem cell* **6**, 369-381 (2010).

960   45.   Zhang, J. et al. The oncogene Etv5 promotes MET in somatic reprogramming and orchestrates

961        epiblast/primitive endoderm specification during mESCs differentiation. *Cell death & disease*

962        **9**, 224 (2018).

963   46.   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   47.   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   48.   Willett, R. et al. TFEB regulates lysosomal positioning by modulating TMEM55B expression

968        and JIP4 recruitment to lysosomes. *Nature communications* **8**, 1580 (2017).

969   49.   Bernstein, B.E. et al. The NIH roadmap epigenomics mapping consortium. *Nature*

970        *biotechnology* **28**, 1045-1048 (2010).

971   50.   Kulakovskiy, I.V. et al. HOCOMOCO: towards a complete collection of transcription factor

972        binding models for human and mouse via large-scale ChIP-Seq analysis. *Nucleic acids research*

973        **46**, D252-D259 (2018).

974   51.   Heinz, S. et al. Simple combinations of lineage-determining transcription factors prime cis-

975        regulatory elements required for macrophage and B cell identities. *Molecular cell* **38**, 576-589

976        (2010).

977   52.   He, K., Zhang, X., Ren, S. & Sun, J. in Proceedings of the IEEE conference on computer vision

978        and pattern recognition 770-778 (2016).

979    53.     Chen, S. et al. OpenAnnotate: a web server to annotate the chromatin accessibility of genomic

980        regions. *Nucleic Acids Research* **49**, W483-W490 (2021).

# Supplementary Materials

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

Fig. S2. EpiGePT's performance in predicting DNase-seq and other epigenetic signals.

Fig. S3. Performance of EpiGePT and baseline methods on chromatin states classification, multiple epigenomic profiles prediction and causal variants classification.

Fig. S4. Ablation analysis of the EpiGePT model.

Fig. S5. Case application of the EpiGePT-online.

Fig S6. Model architecture of EpiGePT for multiple epigenomic signals prediction.

Fig. S7. The fine-tuning performance of the EpiGePT model on predicting potential enhancer-promoter regulatory networks.

Fig. S8. The performance (auROC) of attention score of EpiGePT in distinguishing regulatory element-gene pairs at different distance ranges.

Text S1. Data splitting strategy for model training.

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

Text S3. Running time of the EpiGePT and baseline methods.

Text S4. Implementation of Enformer model and Enformer+.

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

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

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

Table S4. The information of RNA-seq tab-separated values (tsv) file across 28 cell lines or tissues from the ENCODE project.

Table S5. The preprocessed expression data of 711 human transcription factors from the ENCODE project across 129 biosamples.

Table S6. The preprocessed expression data of 711 human transcription factors from the ENCODE project across 28 cell lines or tissues.

Table S7. The order and names of epigenomes of the expression matrices across 56 epigenomes

1010      from the ROADMAP project.

1011      Table S8. The preprocessed expression data of 642 human transcription factors across 56

1012      epigenomes from the ROADMAP project.