

1 **CREATE: cell-type-specific cis-regulatory elements identification via** 2 **discrete embedding**

3 Xuejian Cui¹, Qijin Yin¹, Zijing Gao¹, Zhen Li¹, Xiaoyang Chen¹, Shengquan Chen², Qiao Liu³,
4 Wanwen Zeng^{3,*}, and Rui Jiang^{1,*}

5 ¹ Ministry of Education Key Laboratory of Bioinformatics, Bioinformatics Division at the
6 Beijing National Research Center for Information Science and Technology, Center for
7 Synthetic and Systems Biology, Department of Automation, Tsinghua University, Beijing
8 100084, China

9 ² School of Mathematical Sciences and LPMC, Nankai University, Tianjin 300071, China

10 ³ Department of Statistics, Stanford University, Stanford, CA 94305, USA

11 * Corresponding author: wanwen@stanford.edu, ruijiang@tsinghua.edu.cn

12

13 **Identifying cis-regulatory elements (CREs) within non-coding genomic regions—such as**
14 **enhancers, silencers, promoters, and insulators—is pivotal for elucidating the intricate**
15 **gene regulatory mechanisms underlying complex biological traits. The current prevalent**
16 **sequence-based methods often focus on singular CRE types, limiting insights into cell-**
17 **type-specific biological implications. Here, we introduce CREATE, a multimodal deep**
18 **learning model based on the Vector Quantized Variational AutoEncoder framework,**
19 **designed to extract discrete CRE embeddings and classify multiple CRE classes using**
20 **genomic sequences, chromatin accessibility, and chromatin interaction data. CREATE**
21 **excels in accurate CRE identification and exhibits strong effectiveness and robustness.**
22 **We showcase CREATE's capability in generating comprehensive CRE-specific feature**
23 **spectrum, offering quantitative and interpretable insights into CRE specificity. By**
24 **enabling large-scale prediction of CREs in specific cell types, CREATE facilitates the**
25 **recognition of disease- or phenotype-related biological variabilities of CREs, thereby**
26 **expanding our understanding of gene regulation landscapes.**

27

28 **Introduction**

29 Gene regulation is a fundamental biological process that orchestrates gene expression through
30 a sophisticated network of interactions among biomolecules, including regulatory factors and
31 cis-regulatory elements (CREs) located in non-coding genomic regions^{1, 2}. CREs, such as
32 silencers, enhancers, promoters, and insulators^{3, 4}, typically locate in the chromatin accessible
33 areas^{5, 6} and play crucial roles in modulating gene expression by interacting with target genes
34 through chromatin loops^{7, 8} and other regulatory mechanisms. These characteristics are
35 essential for controlling cell-type-specific gene expression patterns, which contribute to
36 cellular diversity, tissue homeostasis, and the development of complex biological traits⁹⁻¹¹.
37 Consequently, identifying and characterizing cell-type-specific CREs is vital for advancing our
38 understanding of gene regulation in normal physiology and disease states.

39 Silencers, enhancers, promoters, and insulators each have distinct roles in gene regulation.
40 Silencers suppress gene transcription, enhancers boost transcriptional activity, promoters
41 initiate transcription, and insulators act as boundary elements to regulate gene expression^{3, 4}.
42 Due to the restricted understanding of CRE-specific genetic signatures, identifying and
43 validating CREs through biological experiments is cumbersome, time- and resource-
44 consuming^{9, 12}. Massive genomic and epigenomic data benefited from the rapid advancement
45 of high-throughput sequencing technologies¹³⁻¹⁶, has provided valuable opportunities for
46 identifying cell-type-specific CREs using computational methods. For example, DeepSEA is a
47 convolutional neural network (CNN) model based on genomic sequences, which can
48 simultaneously predict chromatin-profiling data such as transcriptional factors (TFs) binding
49 sites, histone modification sites, and chromatin accessible regions¹⁷. DanQ is a hybrid deep
50 neural network that merges convolutional and recurrent architectures, aimed at quantifying the
51 non-coding function of DNA sequences¹⁸. Enhancer-Silencer transition (ES-transition) is a
52 deep learning model based on CNN for identifying enhancers and silencers specific to cell
53 types in the human genome, and has been utilized to uncover the unique phenomenon of
54 enhancer-silencer transitions¹⁹. DeepICSH integrates DNA sequences with various epigenetic
55 features including histone modifications, chromatin accessibility and TF binding to predict
56 cell-type-specific silencers²⁰.

57 However, recently innovated computational methods encounter many limitations and
58 challenges. First, most current efforts are dedicated to the identification of a single type of
59 CRE²⁰⁻²³, particularly enhancers. In the past few decades, there has been extensive research on
60 enhancers²⁴⁻³⁰, while silencers, which generally share many properties with enhancers³¹, have
61 received little attention. Numerous undiscovered CREs and uncharacterized chromatin regions
62 suggest an urgent need for a comprehensive and scalable method of multi-class CRE
63 identification. Second, mainstream methods prevalently extract information from DNA
64 sequences to distinguish CREs¹⁷⁻¹⁹, overlooking the cell type specificity of CREs.
65 Incorporating multi-omics data, including chromatin accessibility and chromatin interaction,
66 for characterizing cell-type-specific CREs can provide valuable insights into gene regulatory
67 mechanisms and cell heterogeneity. Third, deriving interpretable biological implications from
68 conventional deep learning models remains challenging^{32, 33}, hindering the meaningful large-
69 scale identification of CREs and the understanding of model-related biological variabilities of
70 CREs.

71 To bridge these gaps, we propose CREATE (Cis-Regulatory Elements identificAtion via
72 discreTe Embedding), a novel CNN-based supervised learning model that leverages the Vector
73 Quantized Variational AutoEncoder (VQ-VAE)³⁴⁻³⁶ framework. CREATE integrates genomic
74 sequences with epigenetic features to offer a comprehensive approach for the identification and
75 classification of multi-class CREs. The VQ-VAE framework is particularly suited for this task
76 because it can distill genomic and epigenomic data into discrete CRE embeddings, capturing
77 the nuanced differences between various CRE types. This discrete representation facilitates the
78 generation of a CRE-specific feature spectrum, providing both quantitative and intuitive
79 insights into CRE specificity. CREATE's ability to integrate multi-omics data enables it to
80 overcome the limitations of previous methods by offering a more holistic view of CRE
81 functionality and their cell-type-specific roles. Furthermore, CREATE demonstrates superior
82 performance in accurately identifying CREs and exhibits robustness across diverse input types
83 and hyperparameters. Its capability to perform large-scale predictions of CREs in specific cell
84 types and to uncover disease- or phenotype-related biological variabilities in CREs underscores
85 its potential as a powerful tool for constructing a comprehensive CRE atlas. In summary,
86 CREATE represents a significant advancement in the computational identification of CREs,

87 and provides a robust foundation for future research in gene regulation and its implications for
88 human health and disease.

89

90 **Results**

91 **Overview of CREATE.** CREATE is an advanced CNN model based on the VQ-VAE
92 framework^{34,35} to predict and classify multi-class CREs from multi-omics data. Taking as input
93 the one-hot encoded genomic sequence, the vector representing the chromatin accessibility
94 scores for that sequence, and the vector representing the chromatin interaction scores for that
95 sequence, CREATE is specifically crafted to capture discrete CRE embeddings, providing a
96 comprehensive and interpretable characterization of CREs (Fig. 1a and Methods).

97 The architecture of CREATE includes (Fig. 1b and Methods): 1) Encoder Module: Each
98 type of input data is initially processed by dedicated omics-specific encoders that transform the
99 raw data into feature representations suitable for integration. Following this, the processed
100 features are concatenated and passed through the integration encoder module, which
101 synthesizes information from all input modalities to create a unified representation of the
102 genomic context. 2) Vector Quantization Module: In this module, the output embeddings of
103 encoder module are substituted with the closest counterpart in the discrete embedding space
104 called “codebook”. In brief, the features in the codebook are concatenated to form the final
105 CRE embeddings. Unlike traditional VAE-based models³² with fixed prior distributions,
106 CREATE’s codebook is dynamic and updated during training. This flexibility allows the model
107 to refine the discrete embeddings to better represent the underlying biological data. 3) Decoder
108 Module: The decoder reconstructs the original multi-omics input data from the discrete
109 embeddings. It consists of two stages: the integration decoder reconstructs the integrated
110 feature representation from the discrete embeddings. The omics-specific decoders transform
111 the integrated representation back into the respective omics data types, ensuring that the
112 reconstructed data aligns with the original input features. 4) Classifier: To enhance the model’s
113 ability to distinguish between different CRE types, CREATE includes a classifier that enforces
114 the separation of CREs into distinct vectors in the codebook. CREs of the same type are
115 encouraged to map to similar vectors, while those of different types are spread out across

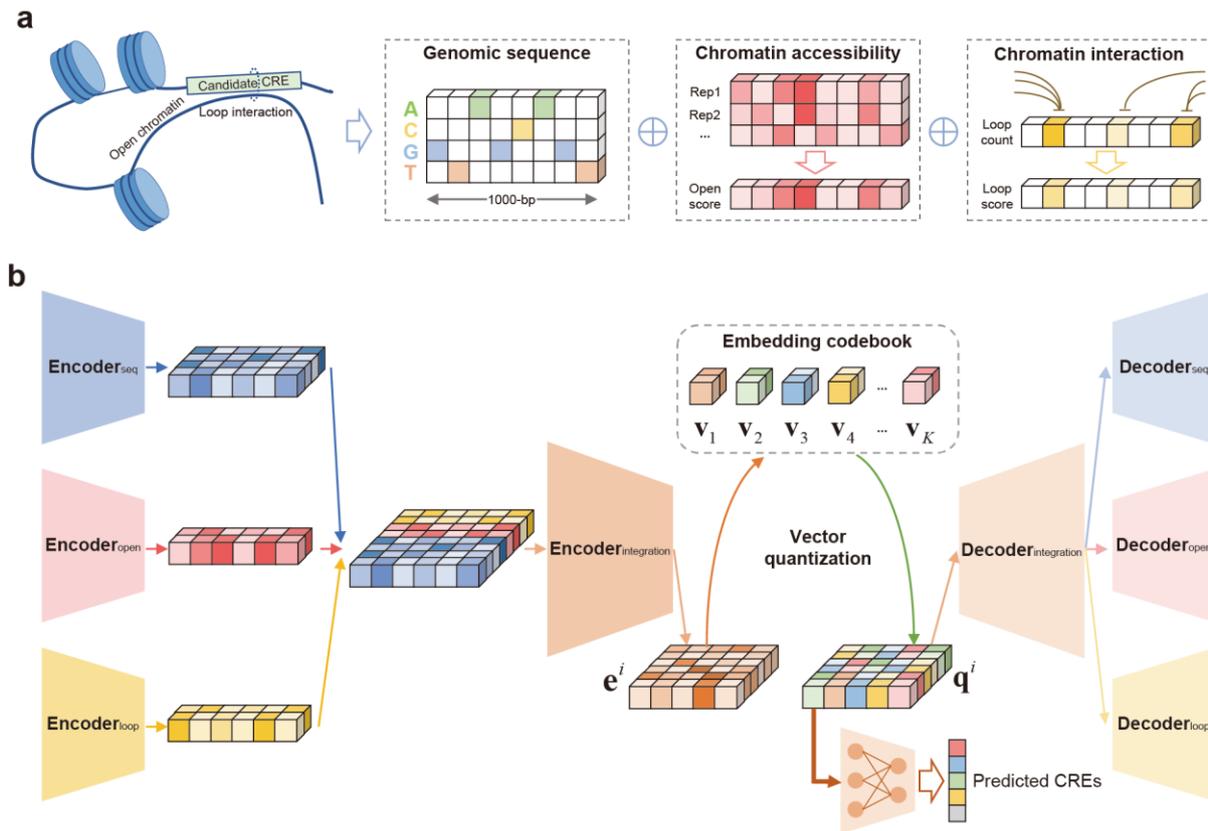


Fig. 1 | Overview of CREATE. **a**, The input of CREATE model. CREATE takes as input the genomic sequence, chromatin accessibility score and chromatin interaction score. **b**, The architecture of CREATE model. CREATE consists of encoders, a vector quantization module and decoders. The encoder module of CREATE combines encoders for multiple input-specific learning and an encoder for multiple input integration. For the i -th CRE, the encoder outputs the latent embedding \mathbf{e}^i of dimension $L \times D'$. By adopting split quantization, the latent embedding will be split into $L \times M$ vectors $\mathbf{e}_{l,j}^i$ of dimension D and then quantized to $\mathbf{q}_{l,j}^i$ for the i -th CRE using embedding codebook with the size of K .

116

117 different vectors. This helps in achieving accurate and interpretable classifications.

118 CREATE offers several key advantages compared to existing methods: 1) Comprehensive
 119 Data Integration: By incorporating multiple omics data types, CREATE captures a more
 120 complete picture of the genomic context and CRE functionality. 2) Dynamic Codebook: The
 121 updateable codebook allows for flexible and accurate representation of CREs, overcoming
 122 limitations of fixed latent spaces in traditional VAE models. 3) Interpretable Embeddings: The
 123 discrete embeddings and their organization in the codebook provide clear and interpretable
 124 insights into CRE specificity and classification.

125 Overall, CREATE represents a significant advancement in computational CRE
 126 identification. Its ability to integrate multi-omics data, produce discrete embeddings, and offer
 127 interpretable results makes it a powerful tool for understanding gene regulation and its

128 implications in complex biological processes.

129

130 **Cis-regulatory elements identification with CREATE.** We comprehensively evaluated the
131 performance of CREATE in identifying cell-type-specific CREs, including silencers, enhancers,
132 promoters, insulators, and background regions, on the K562 and HepG2 cell types (Methods).
133 To assess CREATE's effectiveness, we conducted 10-fold cross-validation experiments and
134 compared its performance with four baseline methods, including DeepSEA¹⁷, DanQ¹⁸, ES-
135 transition¹⁹ and DeepICSH²⁰. The primary evaluation metrics were area under the Receiver
136 Operating Characteristic Curve (auROC), the area under the Precision-Recall Curve (auPRC)
137 and the F1-score (Methods).

138 CREATE significantly surpasses the baseline methods by achieving the best classification
139 performance on both K562 and HepG2 cell types (one-sided paired Wilcoxon signed-rank tests
140 P -values $< 1e-3$), whereas the performance of baseline methods fluctuates across different
141 cross-validation experiments (Fig. 2a-b and Supplementary Fig. 2a-b). For the K562 cell type,
142 CREATE achieves a 10-fold macro-averaged auROC of 0.964 ± 0.002 (mean \pm s.d.),
143 outperforming the second-best method, ES-transition (0.928 ± 0.002) (Fig. 2c). Similarly,
144 CREATE acquires a 10-fold macro-averaged auPRC of 0.848 ± 0.004 , reflecting a substantial
145 improvement of 10.5% compared to the second-best method, DeepICSH (0.743 ± 0.003) (Fig.
146 2d). A comparable performance enhancement is observed for the HepG2 cell type, with
147 CREATE overtaking the baseline methods by a noticeable margin (Supplementary Fig. 2c-d).

148 Among the various CRE types, silencers and enhancers present unique challenges due to
149 their similar epigenetic signatures³¹. Despite this, CREATE demonstrated a clear distinction
150 between these difficult-to-differentiate elements. For the K562 cell type, CREATE achieves
151 notable improvements in identifying silencers, with a mean auPRC that was 13.9% higher than
152 the second-best method, which shows a greater advantage than the macro-averaged results of
153 all CREs (Fig. 2f and Supplementary Fig. 2e). Similarly, CREATE attains a remarkable 22.1%
154 improvement in mean auPRC for enhancers compared with the second-best method
155 (Supplementary Fig. 3a-b). For other CRE types—promoters, insulators, and background
156 regions—CREATE also demonstrates optimal classification performance, although baseline
157 methods provide competitive results (Supplementary Fig. 3c-h). The performance trends for

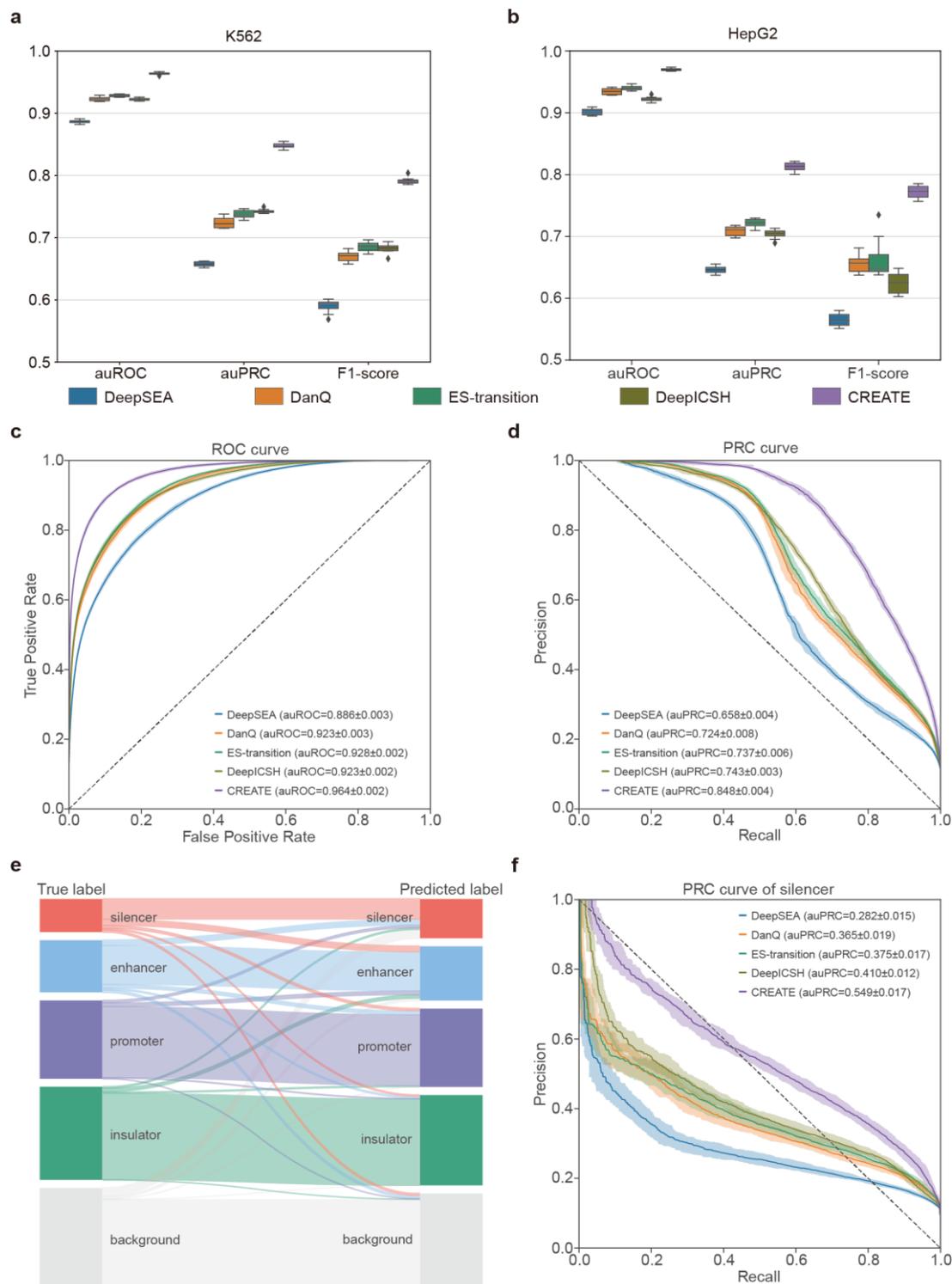


Fig. 2 | Evaluation of CREATE compared with the baseline methods. **a-b**, Boxplot of classification performance evaluated by auROC, auPRC and F1-score on the K562 cell type (**a**) and HepG2 cell type (**b**). Each box plot ranges from the upper to lower quartiles with the median as the horizontal line, whiskers extend to 1.5 times the interquartile range, and points represent outliers. **c-d**, Receiver Operating Characteristic curve (**c**) and Precision-Recall curve (**d**) comparing CREATE and baseline methods on the K562 cell type. **e**, The mapping between the true CRE labels and CREATE-predicted CRE labels on the testing data in one of the 10-fold cross-validation experiments of K562 cell type. **f**, Precision-Recall curve comparing CREATE and baseline methods for silencers in the K562 cell type. The mean and standard error of auROC or auPRC are reported in the legend. The confidence band shows ± 1 s.d. for the averaged curve.

159 the HepG2 cell type mirrored those observed for the K562 cell type, further validating
160 CREATE's robustness across different cell types (Supplementary Fig. 4).

161 The results highlight CREATE's exceptional capability in accurately identifying and
162 distinguishing between various CRE types, particularly those less studied or less abundant,
163 such as silencers and enhancers. CREATE's superior performance in capturing CRE variability
164 and cell type specificity underscores its potential as a powerful tool for advancing our
165 understanding of gene regulation mechanisms.

166

167 **Robustness and effectiveness of CREATE.** CREATE integrates genomic sequences,
168 chromatin accessibility, and chromatin interactions to deliver a thorough characterization of
169 gene regulatory processes. To assess the contributions of these different inputs, we conducted
170 extensive ablation experiments. We referred to the models employing a single type of omics
171 data as CREATE(seq), CREATE(open) and CREATE(loop), and those incorporating two
172 different types of omics data as CREATE(seq+open), CREATE(seq+loop) and
173 CREATE(open+loop), respectively. Among the seven models evaluated, CREATE consistently
174 demonstrates the highest classification performance, confirming the importance of
175 incorporating chromatin accessibility and chromatin interactions for superior CRE
176 identification (Fig. 3a and Supplementary Fig. 5a-b). Specifically, CREATE shows substantial
177 improvements in identifying challenging CRE types such as silencers and enhancers
178 (Supplementary Fig. 5c-d). Notably, CREATE(seq), which relies solely on genomic sequences,
179 achieves a 10-fold macro-averaged auPRC of 0.800 ± 0.004 , surpassing baseline methods by
180 5.7% in mean auPRC (Supplementary Fig. 5b). This underscores CREATE's robust
181 performance even when using genomic sequences alone. Incorporating additional omics data,
182 such as chromatin accessibility or chromatin interactions, further enhances performance,
183 though the inclusion of only these inputs without genomic sequences results in relatively poorer
184 outcomes (Fig. 3a and Supplementary Fig. 5a). This indicates that while chromatin accessibility
185 and interactions are valuable, genomic sequences are indispensable for optimal CRE
186 identification, particularly contributing to better identification of silencers and enhancers.

187 Based on the VQ-VAE framework^{34, 35}, discrete embedding allows the latent space of
188 CREATE to be a learnable discrete distribution, as opposed to the fixed Gaussian distribution

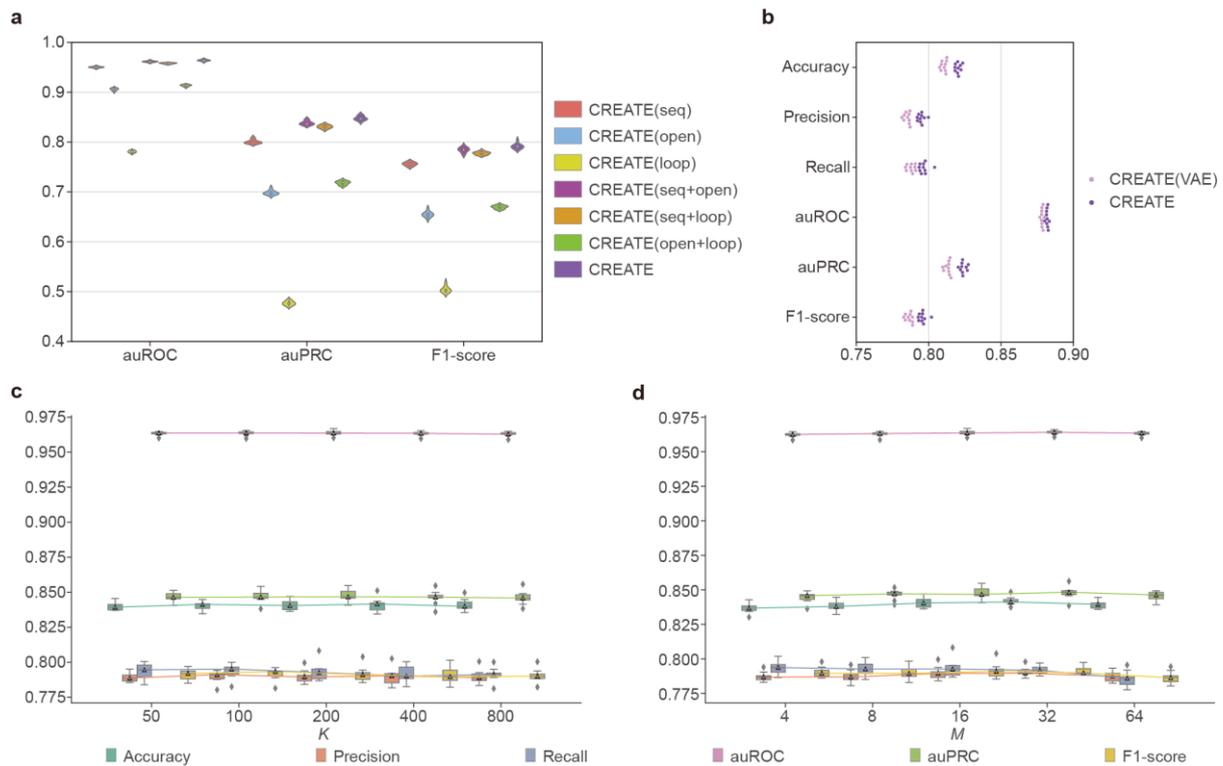


Fig. 3 | Robustness analysis of CREATE. **a**, Violin plot of classification performance evaluated by auROC, auPRC and F1-score for model ablation of CREATE on the K562 cell type. **b**, Swarm plot of classification performance evaluated by accuracy, precision, recall, auROC, auPRC and F1-score for CREATE compared with CREATE (VAE) on the K562 cell type. **c**, Classification performance of CREATE under different values of K (size of codebook) on the K562 cell type. **d**, Classification performance of CREATE under different values of M (time of split quantization) on the K562 cell type. Each box plot ranges from the upper to lower quartiles with the median as the horizontal line, whiskers extend to 1.5 times the interquartile range, and points represent outliers.

189

190 as in VAE models³². To verify the efficiency of discrete embedding in CREATE, we compared
 191 CREATE with a variant using VAE latent space (CREATE(VAE)) while keeping other modules
 192 and training strategies unchanged. CREATE significantly outperformed CREATE(VAE) across
 193 all evaluation metrics (one-sided paired Wilcoxon signed-rank tests P -values $< 1e-3$) (Fig. 3b).
 194 This result highlights the effectiveness of discrete embeddings in capturing complex CRE
 195 features.

196 To validate the stability and effectiveness of CREATE, we designed comprehensive
 197 robustness analyses for the hyperparameters in CREATE, including K denoting the size of
 198 codebook, M denoting the time of split quantization^{32, 36}, α denoting the weight of $L_{encoder}$, μ
 199 denoting the update ratio of codebook. First, to evaluate the robustness of CREATE to the size
 200 of codebook, we trained CREATE with different values of K (50, 100, 200, 400 and 800) on
 201 the K562 cell type. CREATE exhibited consistent classification performance across these

202 values, demonstrating its insensitivity to codebook size variations (Fig. 3c). Taking into
203 account the balance of CRE specificity preservation and codebook utilization, we set the
204 default value of K to 200. Second, to evaluate the stability of CREATE to the time of split
205 quantization, we trained CREATE with different values of M (4, 8, 16, 32 and 64) on the K562
206 cell type. The results show that CREATE attains highly stable classification performance across
207 different values of M (Fig. 3d). Evidently, the lower the time of split quantization, the higher
208 the dimension of codebook features. With the consideration that it is obviously challenging to
209 look up the nearest neighbors for high-dimensional vectors, we set the default value of M to
210 16. Third, following the original studies of VQ-VAE, we aimed for the codebook to have less
211 impact on the output of encoder so that we set the default value of α , the weight of $L_{encoder}$, to
212 0.25. To validate the robustness of CREATE with different weights of $L_{encoder}$, we trained
213 CREATE with different values of α (0.05, 0.1, 0.25, 0.5 and 1.0) on the K562 cell type. The
214 results demonstrate that CREATE consistently obtains stable classification performance under
215 different values of α (Supplementary Fig. 5e). Fourth, similar to the original studies of VQ-
216 VAE, we set the default value of μ , the update ratio of codebook, to 0.01. To assess the stability
217 of CREATE with different update ratios, we trained CREATE under a series of μ , 0.001, 0.005,
218 0.01, 0.05 and 0.1, on the K562 cell type. The results demonstrate the stability of the
219 classification performance under different values of μ (Supplementary Fig. 5f). To summarize,
220 the effective integration of multiple omics inputs, stable hyperparameters, and efficient discrete
221 embedding all demonstrate the robustness and usability of CREATE.

222

223 **Feature spectrum for unveiling CRE specificity.** Discrete latent embedding of CREATE can
224 reveal biological insights in an interpretable and intuitive manner. Using the latent embeddings
225 of CREs, we built a uniform manifold approximation and projection (UMAP)³⁷ plot (Fig. 4a).
226 Clearly, promoters, insulators and background regions are effectively separated, while there is
227 some degree of overlap between silencers and enhancers, which is consistent with the
228 classification results. To further validate the capability of CREATE in quantitatively
229 articulating CRE specificity, we obtained specific feature spectrum for each type of CRE
230 (Supplementary Fig. 6a and Methods). Briefly, each element in the CRE-specific feature
231 spectrum represents the probability of a codebook feature occurring in that particular CRE

232 embeddings. We can always discover a set of particular features that are uniquely associated
233 with a specific CRE and have the highest probability scores on that CRE, and we refer to these
234 features as CRE-specific features. Concretely, for the K562 cell type, there are specific features
235 uniquely enriched in the feature spectrum of each CRE (Fig. 4b). For example, we definitely
236 observe different sets of specific features corresponding to promoters, insulators and
237 background regions, which are clearly separated in the UMAP visualization (Fig. 4a) and
238 Sankey diagram (Fig. 2e) as well. For the most difficult-to-distinguish two types of CREs, the
239 feature spectrum of silencers contains a set of features (to the left of the blue dashed line) with
240 notably higher probability scores compared to their scores in the feature spectrum of enhancers.
241 Similarly, there is a set of features (between the blue dashed line and the purple dashed line)
242 with notably higher probability scores in the feature spectrum of enhancers than those of
243 silencers. In short, there is a relatively clear difference between the feature spectra of silencers
244 and enhancers while they are connected together in the UMAP visualization. A similar result
245 also occurred on the HepG2 cell type (Supplementary Fig. 6b-c). The CRE-specific feature
246 spectrum, derived from discrete latent embedding of CREATE, has the potential to depict the
247 general and comprehensive patterns of a type of CRE, further unveiling the CRE specificity
248 quantitatively and interpretably.

249 To demonstrate the potential of codebook features in the CRE-specific feature spectrum for
250 capturing key biological patterns, we identified the codebook feature with the highest
251 probability score in the CRE-specific feature spectrum of K562 cell type as the major feature
252 of that CRE, and we then zeroed it out before passing the CRE embeddings through the decoder
253 again to generate the reconstructed output. To better understand the relationship between multi-
254 omics input and the major feature of silencers, we designed comparative experiments between
255 the original and reconstructed genomic sequences, chromatin accessibility scores and
256 chromatin interaction scores. First, we conducted motif enrichment analysis for the original
257 and reconstructed silencers (Methods). It is worth noting that the motif enrichment significance
258 ($-\log_{10}P$ -value) of MAFA, LHX6 and PAX8, which were reported as repressors in the previous
259 literature³⁸⁻⁴⁰, is obviously higher in the original sequences compared to the reconstructed
260 sequences (one-sided Wilcoxon rank-sum tests P -values $< 7e-53$) (Fig. 4c and Supplementary
261 Fig. 6d-e), whereas similar comparison results were not observed for known activators, such

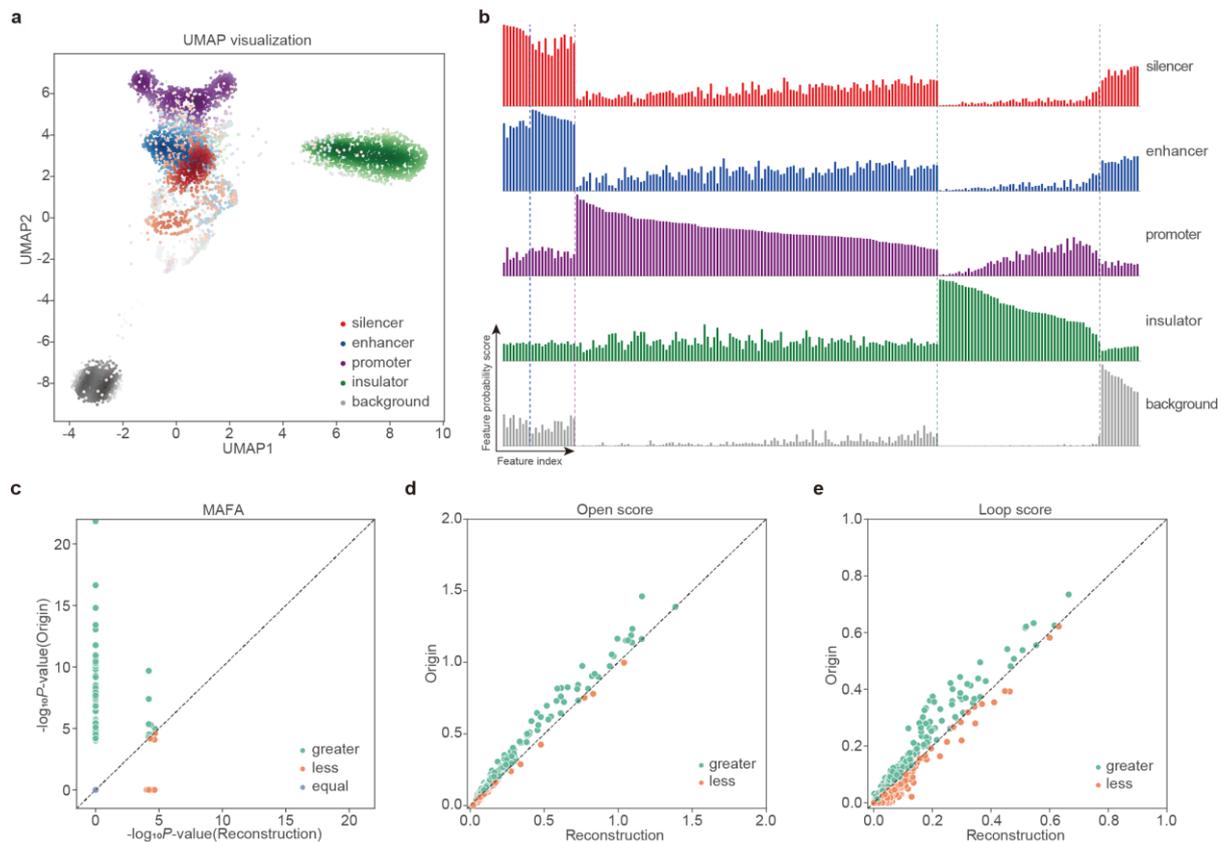


Fig. 4 | Generation and interpretation of CRE-specific feature spectrum. **a**, UMAP visualization of the CRE embeddings from CREATE on the testing data in one of the 10-fold cross-validation experiments of K562 cell type. **b**, CRE-specific feature spectrum. There are a distinct set of specific features that are enriched or depleted in the feature spectrum of each CRE on the K562 cell type. **c**, Comparison of MAFA motif enrichment significance ($-\log_{10}P$ -value) between original input and reconstructed output when information derived from the major feature in the silencer-specific feature spectrum of K562 cell type is removed by zeroing it out before passing the CRE embeddings again through the decoder. **d**, Comparison of open scores between original input and reconstructed output when information derived from the major feature in the silencer-specific feature spectrum of K562 cell type is removed. **e**, Comparison of loop scores between original input and reconstructed output when information derived from the major feature in the silencer-specific feature spectrum of K562 cell type is removed.

262

263 as POU6F1⁴¹ and MYC⁴² (Supplementary Fig. 6f-g). This also demonstrates that the identified
 264 major feature of silencer-specific feature spectrum plays a crucial role in distinguishing
 265 between silencers and enhancers, as it indeed captures the motif information of some repressors,
 266 aligning with the repressive function of silencers. Simultaneously, TFs with the most
 267 significant difference between the original and reconstructed sequences, such as PRDM4,
 268 ZNF582 and SCRT2 (one-sided Wilcoxon rank-sum tests P -values $< 6e-79$) (Supplementary
 269 Fig. 7a-c), are considered to be novel silencer-related TFs. PRDM4 has been linked with
 270 recruiting chromatin modifiers, suggesting its involvement in establishing repressive chromatin
 271 states⁴³. ZNF582 has been implicated in DNA methylation processes, which are crucial for

272 maintaining silencer function⁴⁴. SCRT2 is less characterized, but its differential binding
273 indicates a possible regulatory role in silencing mechanisms⁴⁵. Similarly, the CRE-specific
274 motif information is also harbored in the major feature of enhancers, promoters and insulators
275 (Supplementary Fig. 7d-i), demonstrating that these features catch CRE-specific sequence
276 patterns. Additionally, the unique motif patterns associated with silencers compared to
277 enhancers, promoters, and insulators provide further evidence that these elements are distinct
278 regulatory modules with specific TF associations. This distinction underscores the importance
279 of considering a broader scope of CREs, including dual-function regulatory elements that might
280 act as silencers under certain conditions and enhancers under others. Furthermore, the zeroing
281 operation led to a reduction in the reconstructed chromatin accessibility scores and chromatin
282 interaction scores (Fig. 4d-e), indicating that the major feature also captures silencer-specific
283 epigenomic characteristics. Through the extensive comparative experiments that we designed,
284 the CRE-specific feature spectrum generated by CREATE interpretably reveals the CRE
285 specificity and is potentially involved in the gene regulation process in specific cell types. In
286 conclusion, CREATE not only identifies known regulatory elements but also sheds light on less
287 understood elements like silencers, filling a critical gap in the current landscape of gene
288 regulation studies.

289

290 **Large-scale prediction of cis-regulatory elements.** The emergence of extensive epigenomic
291 sequencing data across various cell types has enabled us to leverage a wealth of information
292 for identifying cell-type-specific CREs on a large scale and establishing regulatory elements
293 maps. CREATE proves to be a powerful tool for a comprehensive characterization of gene
294 regulatory processes, revealing CREs with high accuracy and interpretability. In our study, we
295 collected 270,259 candidate CREs on the K562 cell type and 232,456 candidate CREs on the
296 HepG2 cell type for large-scale prediction (Supplementary Table 1 and Methods). For each
297 cross-validation experiment, based on the trained CREATE model, we calculated a cutoff score
298 for each type of CRE according to the validation set with a false positive rate (FPR) not
299 exceeding 0.01. Candidate regions in Supplementary Table 1 exceeding the silencer cutoff
300 score are marked as predicted silencers, and other CREs are labeled similarly. This approach
301 led to the identification of 26,012 predicted silencers, 29,423 predicted enhancers, 2,057

302 predicted promoters, and 10,558 predicted insulators in the K562 cell type using the models
303 trained on the K562 cell type, and the remaining sequences were classified as background
304 regions. Similarly, in the HepG2 cell type, we identified 16,000 predicted silencers, 49,145
305 predicted enhancers, 4,422 predicted promoters and 13,122 predicted insulators using the
306 models trained on the HepG2 cell type. The predicted CREs showed strong correlations with
307 known epigenomic markers. For example, H3K27me₃, a key histone modification associated
308 with silencers⁴⁶⁻⁴⁸, was found at higher proportions in predicted silencers (9.0% in K562 and
309 16.0% in HepG2) compared to other CREs (average 2.2% in K562 and average 4.1% in HepG2)
310 (highlighted with red dashed box; Fig. 5a and Supplementary Fig. 8a). Similarly, histone
311 modifications associated with enhancers, such as H3K9ac^{49, 50}, H3K27ac⁵⁰⁻⁵², H3K4me1^{50, 53,}
312 ⁵⁴, H3K4me2⁵⁴ and H3K4me3^{54, 55}, were more prevalent in predicted enhancers compared to
313 other CRE types (average 18.5% in K562 and average 36.6% in HepG2) (highlighted with blue
314 dashed box; Fig. 5a and Supplementary Fig. 8a).

315 To further validate the epigenomic characteristics of predicted CREs on a large scale, we
316 conducted extensive comparative analyses. First, TFs play a crucial regulatory role in gene
317 transcription by binding to CREs, and sequence-specific TF motifs can be considered key
318 factors in identifying CREs⁵⁶⁻⁵⁸. We performed motif enrichment analysis on both true CREs
319 (experimentally validated CREs) and predicted CREs (Methods). Compared to other true CREs
320 and background regions, true silencers are enriched with the binding motifs of repressive TFs
321 previously reported in the literature (Fig. 5b), such as FOXD1⁵⁹, TFAP2A⁶⁰, MAFA³⁸, MAFB⁶¹,
322 LHX6³⁹, PAX8⁴⁰, NFIA⁶² and PRDM6⁶³, which are also enriched in the predicted silencers
323 (Fig. 5c). Motifs belonging to active TFs, including POU6F1⁴¹, MYC⁴², ZFH3⁶⁴ and SOX8⁶⁵,
324 are enriched consistently across the true and predicted enhancers (Supplementary Fig. 8b-c).
325 Notably, silencers and enhancers, the two most similar types of CREs, are enriched with the
326 same TFs, such as MYC and ZFH3. These TFs have been validated to act as either activators
327 or repressors^{42, 64, 66, 67}, which aligns with the potential conversion between silencers and
328 enhancers under different conditions^{19, 68}. Similar motif enrichment results, consistent between
329 true and predicted CREs, have also been observed in promoters (Supplementary Fig. 8d-e).
330 These results indicate that CREATE effectively captured the CRE-specific sequence
331 characteristics.

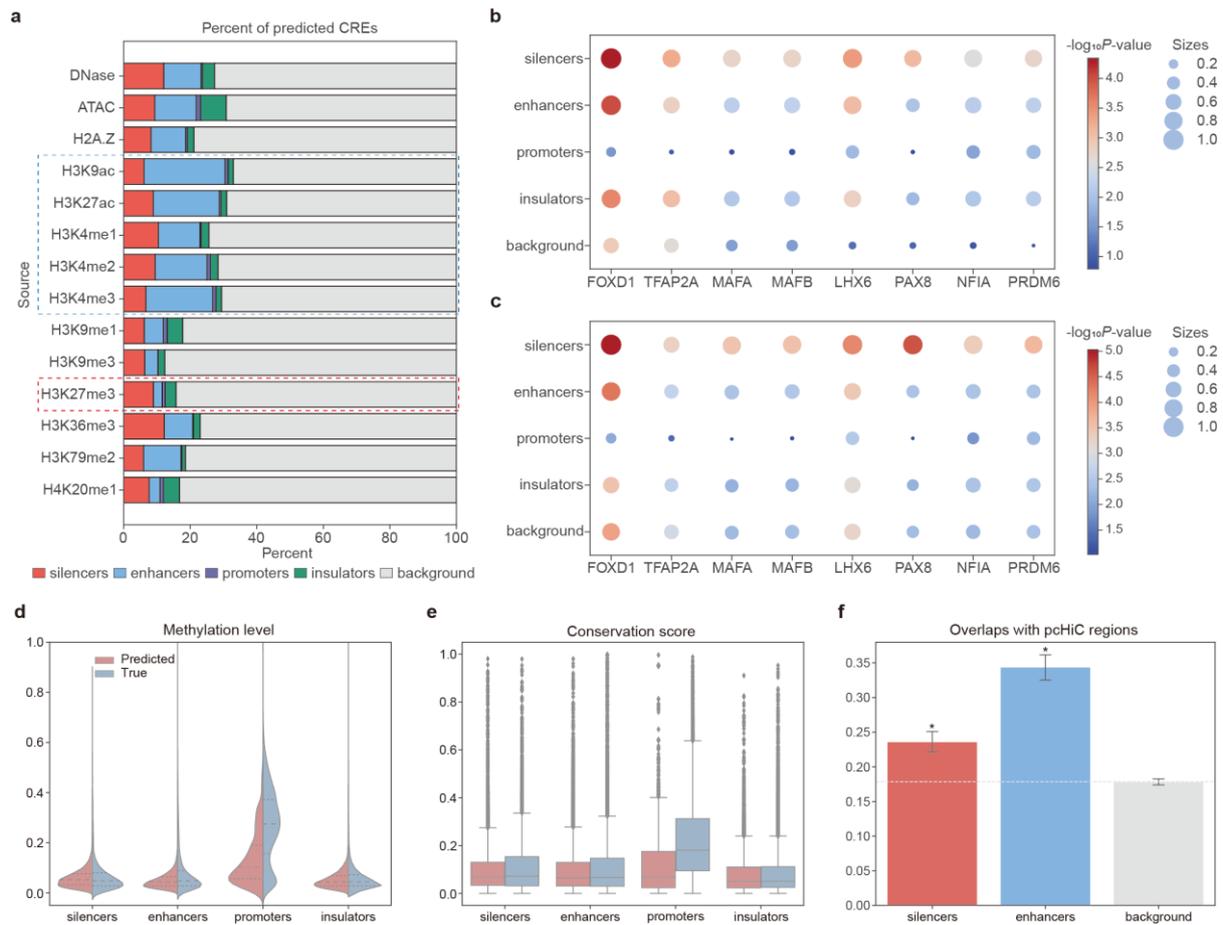


Fig. 5 | Characteristics of predicted CREs by CREATE. **a**, Percentage of predicted CREs and background regions from different candidate sources in the K562 cell type. **b-c**, Bubble plot of motif enrichment significance ($-\log_{10}P$ -value) of repressive TFs (silencer-related TFs) at true CREs (**b**) and predicted CREs (**c**) on the K562 cell type. The size of bubbles represents the proportion of CREs with P -value < 0.01 . **d**, Violin plot of methylation levels at true CREs and predicted CREs on the K562 cell type. Each violin plot contains three horizontal dashed lines denoting the median, the upper quartile, and the lower quartile. **e**, Box plot of conservation scores at true CREs and predicted CREs on the K562 cell type. Each box plot ranges from the upper to lower quartiles with the median as the horizontal line, whiskers extend to 1.5 times the interquartile range, and points represent outliers. **f**, Bar plot of overlaps between the pChIC regions and predicted silencers, enhancers or background regions on the K562 cell type. The asterisks above the bars indicate the significant enrichments compared with the background regions. (*) P -value $< 2e-6$. The error bars denote the 95% confidence interval, and the centers of error bars denote the average value.

332

333

334

335

336

337

338

Second, DNA methylation is an important epigenetic modification involved in gene regulation, particularly gene silencing^{69, 70}. We calculated the methylation levels for both true CREs and predicted CREs (Methods), and observed the consistency between them except for promoters (Fig. 5d), which may be due to the complete collection of experimentally validated promoters and the limited number of predicted promoters. Specifically, 65.1% of the predicted promoters are adjacent to the promoters of non-coding genes, which is much higher than the

339 9.6% of randomly sampled genomic regions. In addition, the methylation levels of predicted
340 CREs are significantly higher than those of predicted background regions (one-sided Wilcoxon
341 rank-sum tests P -values $< 2e-6$) (Supplementary Fig. 8f).

342 Third, CREs are usually conserved in the evolutionary process of vertebrates, and the
343 conserved regions are essential for deciphering the landscapes of gene regulation⁷¹⁻⁷³. We
344 computed the phastCons scores for the true CREs and the predicted CREs (Methods), and
345 noticed that the conservation scores exhibit strong consistency between them except for
346 promoters (Fig. 5e). Compared with the predicted background regions, the predicted CREs
347 except insulators are significantly more conserved (one-sided Wilcoxon rank-sum tests P -
348 values $< 1e-13$) (Supplementary Fig. 8g).

349 Fourth, CREs frequently regulate gene expression by connecting promoters through
350 chromatin loops, which can be identified by promoter-capture HiC (pcHiC)⁷⁴. We counted the
351 number of overlaps between the pcHiC regions and the true CREs or the predicted CREs
352 (Methods), and perceived the predicted silencers and enhancers harbor significantly more
353 overlaps with pcHiC regions than the predicted background regions (one-sided Wilcoxon rank-
354 sum tests P -values $< 2e-6$) (Fig. 5f), which aligns with the functional roles of these CREs in
355 influencing target genes.

356 We further tested CREATE's cross-cell type prediction capabilities by using models trained
357 on the K562 cell type to predict CREs in the HepG2 cell type. The results (auROC of $0.964 \pm$
358 0.002 and auPRC of 0.792 ± 0.005) demonstrate that CREATE maintains high performance
359 across different cell types, confirming its robustness and generalizability (Supplementary Fig.
360 2c-d). Collectively, CREATE precisely extracts the CRE-specific epigenomic characteristics,
361 enabling the construction of a comprehensive CRE atlas.

362

363 **Characterization of dual-function regulatory elements.** Dual-function regulatory elements
364 (DFREs) are regions that exhibit dual roles as either silencers or enhancers depending on the
365 cellular context^{9, 19, 68, 75}. Understanding DFREs is crucial for unraveling the complexity of
366 gene regulation, as these elements can significantly impact gene expression by switching
367 functions based on the cellular environment. A K562 silencer overlapping a HepG2 enhancer

368 by more than 600 bp was considered a DFRE, resulting in 2,409 DFREs (9.3% of all predicted
369 silencers) and 23,603 normal silencers. Conversely, we identified 36,448 HepG2 enhancers
370 that do not function as enhancers or silencers in the K562 cell type, categorizing them as normal
371 enhancers. Reasonably, the ability of CREATE to effectively characterize DFREs is
372 demonstrated by its capacity to assign higher CREATE enhancer scores to DFREs than normal
373 silencers (one-sided Wilcoxon rank-sum test P -value $< 6e-47$) (Fig. 6a), and higher CREATE
374 silencer scores than normal enhancers (one-sided Wilcoxon rank-sum test P -value $< 2e-5$) (Fig.
375 6b). This differentiation underscores CREATE's effectiveness in distinguishing between
376 multifunctional and context-specific regulatory elements.

377 To further explore the biological significance of DFREs, we conducted several comparative
378 analyses. First, DFREs exhibit the highest conservation scores (one-sided Wilcoxon rank-sum
379 tests P -values $< 4e-3$) and the background regions have the lowest conservation scores (one-
380 sided Wilcoxon rank-sum tests P -values $< 2e-31$) (Fig. 6c), suggesting that these elements are
381 evolutionarily preserved due to their critical roles in gene regulation. This high conservation
382 underscores their functional importance across species and reinforces the value of identifying
383 these elements for understanding gene regulatory mechanisms. Second, DFREs possess higher
384 methylation levels in the K562 cell type compared to normal silencers (one-sided Wilcoxon
385 rank-sum test P -value $< 7e-4$) (Fig. 6d). This observation highlights the unique epigenomic
386 signatures of DFREs, suggesting that their dual functionality is associated with distinct
387 methylation patterns, which may influence their regulatory roles. Third, DFREs show a strong
388 preference with more overlaps with pcHiC regions of K562 cell type than normal silencers (Fig.
389 6e). This indicates that DFREs are actively involved in chromatin looping interactions, which
390 are critical for mediating gene expression and regulatory network organization. Fourth, we
391 computed the number of overlaps with expression quantitative trait loci (eQTLs) from whole-
392 blood or liver tissues from GTEx^{76, 77} for DFREs, normal silencers, normal enhancers and the
393 background regions (Methods). DFREs are significantly enriched with more whole-blood
394 eQTLs than normal silencers (one-sided Wilcoxon rank-sum test P -value $< 3e-5$) (Fig. 6f), and
395 more liver eQTLs than normal enhancers (one-sided Wilcoxon rank-sum test P -value $< 2e-3$)
396 (Fig. 6g). This enrichment demonstrates the tissue-specific regulatory potential of DFREs,
397 highlighting their role in fine-tuning gene expression across different biological contexts.

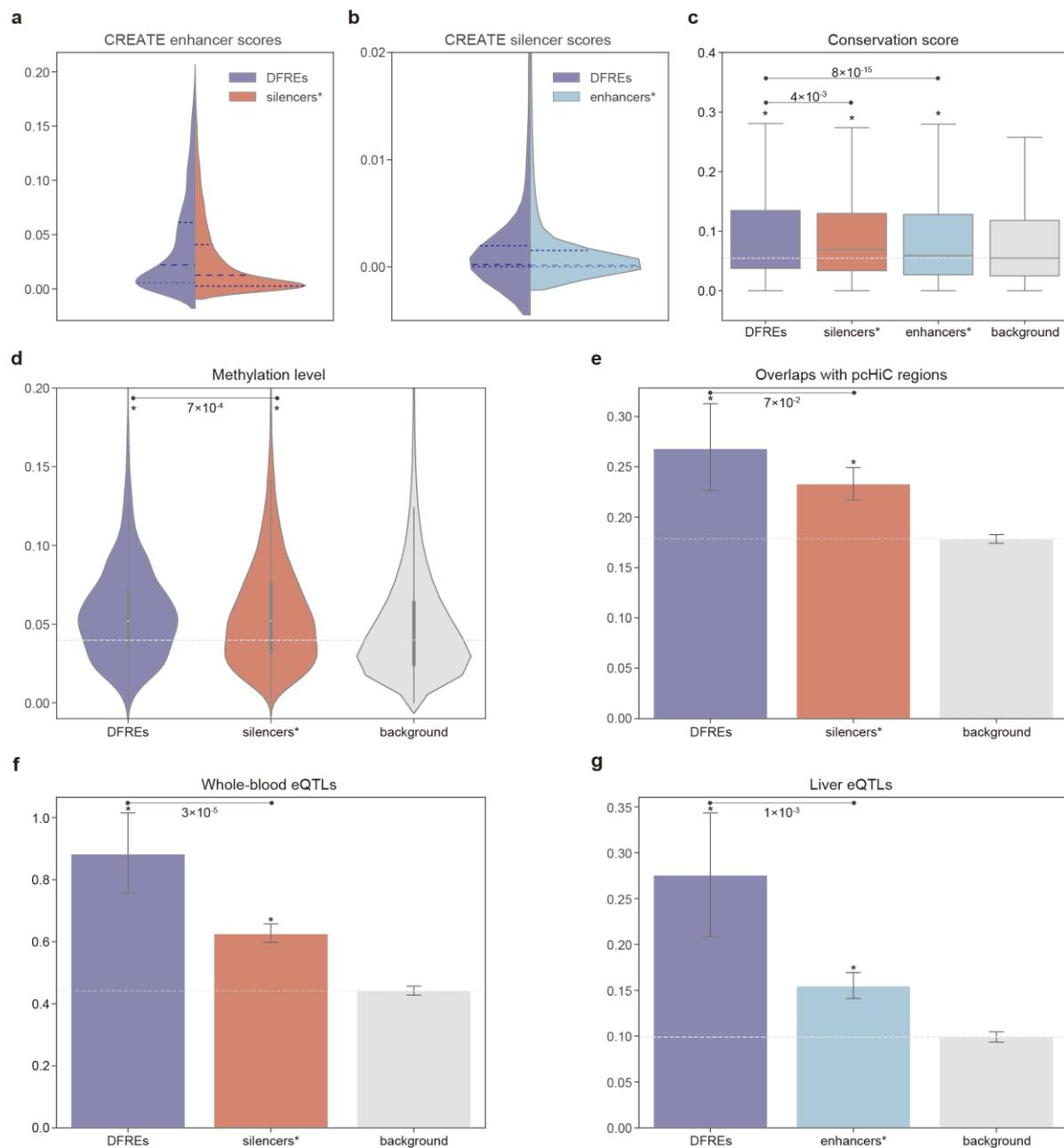


Fig. 6 | Characterization of DFREs identified by CREATE. **a**, Violin plot of the CREATE enhancer scores for DFREs and normal silencers (silencers*). **b**, Violin plot of the CREATE silencer scores for DFREs and normal enhancers (enhancers*). Each violin plot contains three horizontal dashed lines denoting the median, the upper quartile, and the lower quartile. **c**, Box plot of conservation scores at DFREs, normal silencers (silencers*), normal enhancers (enhancers*) and the background regions. The asterisks above the boxes indicate the significant enrichments compared with the background regions. (*) P -value $< 2e-31$. Each box plot ranges from the upper to lower quartiles with the median as the horizontal line, whiskers extend to 1.5 times the interquartile range. **d**, Violin plot of methylation levels at DFREs, normal silencers (silencers*) and the background regions. (*) P -value $< 2e-6$. **e**, Bar plot of overlaps between the pChIC regions and DFREs, normal silencers (silencers*) or the background regions. (*) P -value $< 2e-3$. **f**, Bar plot of overlaps between the whole-blood eQTLs and DFREs, normal silencers (silencers*) or the background regions. (*) P -value $< 5e-21$. **g**, Bar plot of overlaps between the liver eQTLs and DFREs, normal enhancers (enhancers*) and the background regions. (*) P -value $< 2e-6$. The error bars denote the 95% confidence interval, and the centers of error bars denote the average value.

399 Ultimately, CREATE's advanced capability to differentiate and interpret DFREs enriches the
400 field's ability to map regulatory landscapes and uncover the underlying mechanisms of gene
401 regulation.

402

403 **Disease-associated variations analysis and tissue-specific enrichments in CREs.** CREs play
404 crucial roles in disease susceptibility and phenotype variations, often harboring single-
405 nucleotide polymorphisms (SNPs) and eQTLs associated with various diseases and traits^{78, 79}.
406 To highlight the capacity of CREATE in uncovering disease-relevant variations within CREs,
407 we analyzed overlaps with SNPs from dbSNP⁸⁰⁻⁸² database and eQTLs from GTEx^{76, 77} for
408 both the true CREs and the predicted CREs (Methods). Our results reveal that the gene variation
409 distributions at predicted CREs align closely with those at true CREs (Fig. 7a and
410 Supplementary Fig. 9a-c). In detail, the predicted CREs are significantly enriched with more
411 rare SNPs than the background regions (one-sided Wilcoxon rank-sum tests P -values $< 2e-9$)
412 (Fig. 7b), whereas a similar significant result was not observed for common SNPs
413 (Supplementary Fig. 9d), suggesting the functional importance of identified CREs. The
414 enrichment of rare gene variants in CREs, especially silencers and enhancers, supports the
415 notions that disease-associated variants are more frequently located in gene regulatory regions⁹,
416 ^{12, 83}, and rare variants are more impactful in complex diseases compared to common gene
417 variants^{84, 85}. Similarly, compared to background regions, significant enrichment on silencers
418 and enhancers also occurred with whole-blood eQTLs (one-sided Wilcoxon rank-sum tests P -
419 values $< 4e-49$) (Supplementary Fig. 9e), but not with all eQTLs (Supplementary Fig. 9f),
420 reinforcing their tissue specificity. Besides, gene variation levels for both true CREs and
421 predicted CREs gradually decrease with an increase in CREATE background scores (Fig. 7c
422 and Supplementary Fig. 10a,d-i) but increase with an increase in CREATE silencer scores
423 (Supplementary Fig. 10b-c), explicating the ability of CREATE for quantifying the impact of
424 gene variations.

425 To further quantitatively assess the power of CREATE in unlocking the CRE-specific
426 sequence characteristics, we portrayed the correlation between the CREATE scores and the
427 motif enrichment significance at true or predicted CREs (Fig. 7d-e and Supplementary Figs.
428 11-14). We recognized the strong positive correlation between the CREATE silencer scores and

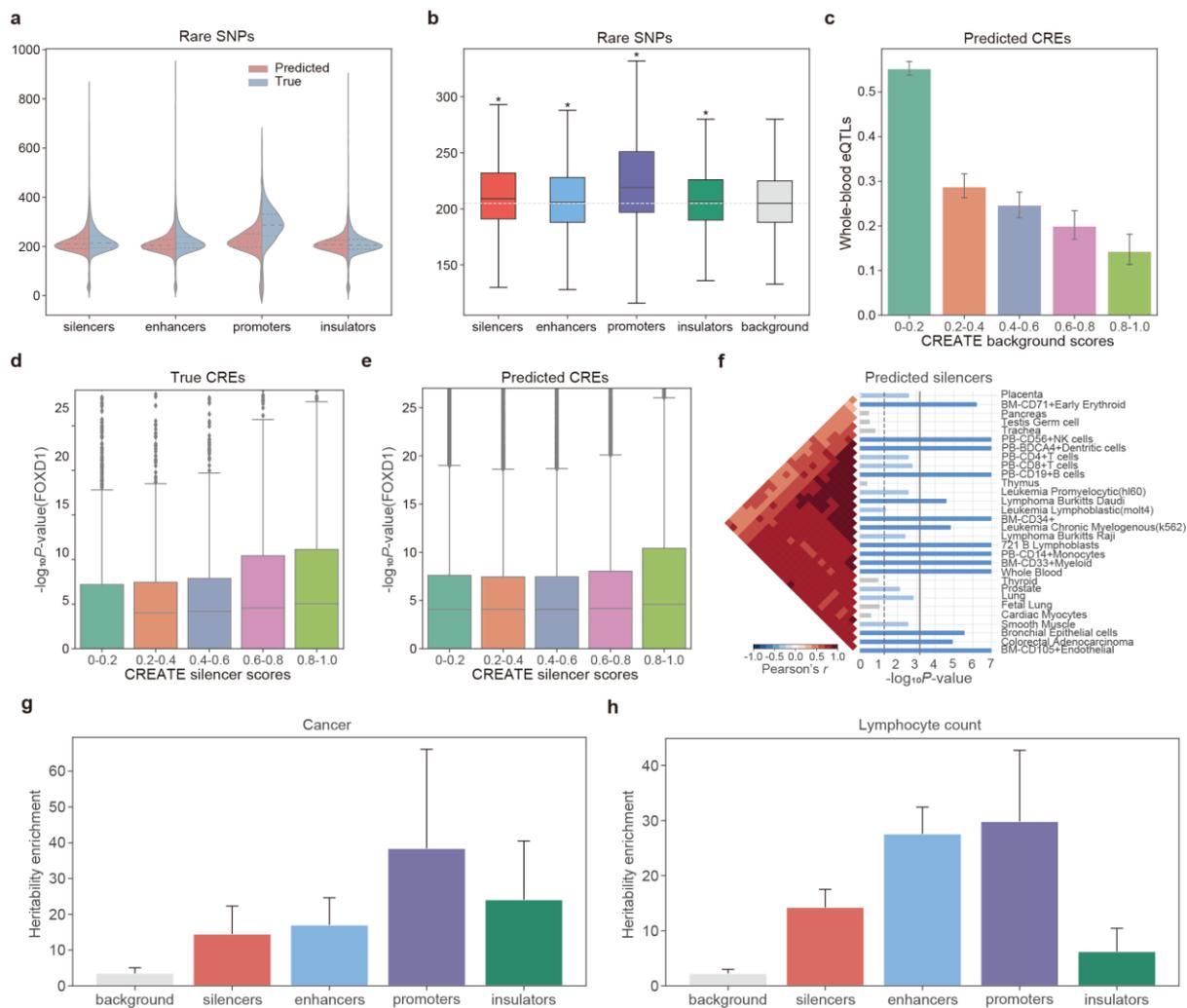


Fig. 7 | Identification of the biological variability of CREs by CREATE. **a**, Violin plot of overlaps between the rare SNPs and true CREs or predicted CREs on the K562 cell type. Each violin plot contains three horizontal dashed lines denoting the median, the upper quartile, and the lower quartile. **b**, Box plot of overlaps between the rare SNPs and the predicted CREs or background regions on the K562 cell type. The asterisks above the boxes indicate the significant enrichments compared with the background regions. (*) P -value $< 2e-9$. Each box plot ranges from the upper to lower quartiles with the median as the horizontal line, whiskers extend to 1.5 times the interquartile range. **c**, Correlation between the CREATE background scores and overlaps between the whole-blood eQTLs and predicted CREs on the K562 cell type. **d-e**, Correlation between the CREATE silencer scores and the motif enrichment significance ($-\log_{10}P$ -value) of FOXD1 at true CREs (**d**) and predicted CREs (**e**) on the K562 cell type. Each box plot ranges from the upper to lower quartiles with the median as the horizontal line, whiskers extend to 1.5 times the interquartile range, and points represent outliers. **f**, Top 30 significantly enriched tissues in SNPsea analysis on the predicted silencers of K562 cell type. The vertical dashed line represents the one-sided P -value cutoff at the 0.05 level, while the solid lines denotes the cutoff at 0.05 level for the one-sided P -value with Bonferroni correction. Each plot also contains the ordered expression profiles using hierarchical clustering with unweighted pair-group method with arithmetic means, and the Pearson correlation coefficients indicating the correlation between profiles. **g-h**, Heritability enrichments estimated by LDSC within predicted CREs and background regions identified by CREATE for blood-related traits including cancer (**g**) and lymphocyte count (**h**). The error bars denote jackknife standard errors over 200 equally sized blocks of adjacent SNPs about the estimates of enrichment, and the centers of error bars represent the average value.

430 the enrichment significance of silencer-related TFs (Fig. 5b-c), as well as between the CREATE
431 enhancer scores and the enrichment significance of enhancer-related TFs (Supplementary Fig.
432 8b-c).

433 To illustrate the ability of CREATE for revealing the tissues influenced by the identified risk
434 loci within the true CREs and the predicted CREs, we applied SNPsea⁸⁶ for tissue enrichment
435 analysis (Methods). For the silencers and enhancers predicted by CREATE, we discovered
436 more tissues related to blood than background regions (Fig. 7f and Supplementary Fig. 15),
437 which aligns with the outcomes of true CREs (Supplementary Fig. 16). Concretely,
438 Leukemia_Chronic_Myelogenous(k562) was identified as a significantly enriched tissue for
439 the predicted silencers and enhancers in the K562 cell type (P -values $< 2e-5$), confirming the
440 tissue specificity of CREs identified by CREATE.

441 To validate the competence of CREATE for studying the variations in phenotypes based on
442 the true CREs and the predicted CREs, we utilized partitioned linkage disequilibrium score
443 regression (LDSC)⁸⁷ for heritability enrichment analysis (Methods). Specifically, the
444 enrichment of heritability in the predicted CREs is higher than that in the predicted background
445 regions for the blood-related phenotypes, such as cancer, lymphocyte count and so on (Fig. 7g-
446 h and Supplementary Fig. 17). Along with the enrichment results for the true CREs and
447 background regions (Supplementary Fig. 18), CREs predicted by CREATE have inherited the
448 pattern of heritability contribution for complex traits and diseases.

449 Altogether, CREATE excels in identifying CREs with significant disease-related variations
450 and tissue-specific enrichments, providing critical insights into regulatory dynamics during
451 development and disease progression. This capability underscores the power of CREATE in
452 advancing our understanding of gene regulation and its implications for complex traits and
453 diseases.

454

455 **Discussion**

456 In this study, we introduce CREATE, a groundbreaking multimodal architecture that integrates
457 DNA sequences, cell-type-specific chromatin accessibility, and chromatin interaction features
458 for the multi-class prediction of CREs. Utilizing discrete CRE embeddings, we have verified

459 the superior performance of CREATE in accurate CRE identification compared to state-of-the-
460 art methods, as well as its effectiveness and stability across various input combinations,
461 hyperparameters and forms of latent space. One of the key strengths of CREATE lies in its
462 ability to offer improved interpretability as the CRE-specific feature spectrum, which
463 quantitatively elucidates the CRE specificity and captures CRE-specific epigenomic
464 characteristics. Moreover, CREATE has been validated the substantial potential in identifying
465 cell-type-specific CREs on a large scale and uncovering biological variabilities of CREs,
466 illustrating the ability of CREATE for unveiling the underlying regulatory dynamics that drive
467 transcriptional regulation and disease development.

468 However, despite its successes, there are several areas where CREATE could be further
469 improved. Our study identifies a few limitations and suggests several future directions for
470 enhancing the method: 1) Data imbalance and insufficient research on certain CREs. The
471 current data imbalance, particularly for silencers and certain cell types, impairs the overall
472 performance of CRE identification. The obscure understanding of general silencer
473 characteristics, the limited number of experimentally validated silencers and the restricted
474 number of cell types studied, pose challenges in both model training and the selection of
475 background regions. To address this, we plan to update our predicted silencers to the
476 SilencerDB database⁸⁸ and expand our identification to include more cell types. We also
477 anticipate that advances in biological technologies, such as HiChIP with a broader range of
478 ChIPs, will enhance multi-class CRE identification and aid in constructing a more
479 comprehensive regulatory atlas. 2) Per-base-paired input features and input-specific encoder-
480 decoder structure. While the per-base-paired input features and the input-specific encoder-
481 decoder structure are effective for extracting detailed and comprehensive CRE embeddings,
482 some epigenetic features, such as TF binding, are not well-represented in this format. To
483 improve scalability and representation, we propose integrating prior biological knowledge as
484 additional constraints directly applied to the codebook. This approach is expected to enhance
485 the model's ability of capturing complex epigenetic features and improve overall performance.
486 3) Development of a unified foundation model. Gene regulatory analysis methods typically
487 focus on specialized models for specific problems. This paradigm limits the generalizability
488 and integration of findings across different contexts. We aspire to develop a foundation model

489 for the unified characterization of key gene regulatory factors, leveraging the shareability,
490 scalability and interpretability of the discrete embedding in CREATE. We anticipate that such
491 a foundation model will facilitate a deeper understanding of gene regulation mechanisms and
492 their implications for disease development, ultimately enabling biological discoveries and
493 applications in developmental biology and precision medicine.

494 In conclusion, CREATE represents a significant advancement in the prediction and
495 interpretation of CREs, offering superior performance and insights compared to existing
496 methods. Its ability to integrate diverse data types and deliver interpretable results positions it
497 as a valuable tool for exploring gene regulation and disease mechanisms. Future improvements
498 and expansions of CREATE will continue to refine its capabilities and extend its applicability,
499 driving forward our understanding of the complex interplay between gene regulation and
500 disease.

501 **Methods**

502 **Data collection and preprocessing.** All datasets used in this study were publicly available and
503 collected from different sources. We downloaded experimentally validated silencers for K562
504 and HepG2 cell types from the SilencerDB database⁸⁸. We downloaded experimentally
505 validated enhancers for K562 and HepG2 cell types from the FANTOM5 project^{24, 26}. We
506 obtained transcription start sites (TSSs) from the EPD database⁸⁹ and defined 1kb regions
507 surrounding TSSs (500 bp upstream and 500 bp downstream) as promoters. Since CTCF
508 characterized as an insulator by blocking chromatin interactions^{90, 91}, we took as insulators the
509 CTCF Chromatin immunoprecipitation sequencing (ChIP-seq) peaks for K562 and HepG2 cell
510 types collected from the ENCODE project^{92, 93}.

511 In addition, we collected multiple histone modification ChIP-seq peaks and chromatin
512 accessibility peaks for K562 and HepG2 cell types from the Roadmap project⁹⁴ and ENCODE
513 project (Supplementary Table 1). After filtering the regions overlapped with the experimentally
514 validated CREs, known genes and consensus black list, we obtained 270,259 and 232,456
515 candidate CREs for large-scale prediction on the K562 and HepG2 cell types, respectively.

516 We randomly sampled DNA sequences from the entire human reference genome, excluding
517 the experimentally validated and candidate CREs, known genes, consensus black list. After
518 filtering overlapping regions between CREs, we obtained 6754 silencers, 10,528 enhancers,
519 15,699 promoters, 18,631 insulators and 20,000 background regions for the K562 cell type,
520 and 1456 silencers, 11,407 enhancers, 14,535 promoters, 15,650 insulators and 20,000
521 background regions for the HepG2 cell type. The input for each CRE comprises three
522 components: a one-hot encoded 1000-bp sequence from the human GRCh37/hg19 reference
523 genome, a vector containing chromatin open scores per base pair, and another vector containing
524 chromatin loop scores per base pair.

525

526 **Chromatin open score.** Chromatin accessibility is pivotal for identifying CREs, given that
527 active regulatory DNA elements are typically situated in accessible chromatin regions^{5, 6}. To
528 incorporate the information of chromatin accessibility, we adopted OpenAnnotate⁹⁵ to
529 efficiently calculate the raw read open scores of CREs and background regions per base pair.

530 We derived the chromatin open score per base pair by averaging the raw read open scores
531 across replicates for each respective cell type.

532

533 **Chromatin loop score.** Chromatin looping interactions exert a substantial influence on gene
534 regulation by establishing connections between regulatory elements and target genes^{7, 8}. We
535 incorporated cell-type-specific chromatin interaction data from HiChIP, which precisely
536 profiles both regulatory and structural interactions^{16, 96}, to enhance the identification of CREs.
537 We first calculated the number of chromatin loops per base pair for each CRE, and then
538 obtained the chromatin loop score after logarithmic transformation.

539

540 **Data augmentation.** To ensure enough training samples for our model, we applied a data
541 augmentation strategy to CREs^{21, 22, 97}. As illustrated in Supplementary Fig. 1, for each CRE
542 with length of 1000 base pairs, we shifted a window along the reference genome with a stride
543 of 10 from the midpoint towards both ends. To mitigate the impact of data imbalance, we
544 optionally incorporated data augmentation with varying augmentation ratios (5:5:3:3:1) for
545 silencers, enhancers, promoters, insulators and background regions in the training data.
546 Additionally, we augmented CREs by including the reverse complement of each original
547 sequence. To prevent information leakage, the augmentation ratios for CREs in the validation
548 and testing data are kept consistent at 5. Take the average of the predicted probabilities for all
549 augmented sequences of the input sequence as the predicted probability for that input sequence.

550

551 **The CREATE framework.** We fed CREATE with a concatenated vector $\mathbf{X}^i \in \mathbb{R}^{6 \times L}$ for the
552 i -th input sample including a one-hot encoded genomic sequence $\mathbf{S}^i \in \mathbb{R}^{4 \times L}$, a chromatin open
553 score vector $\mathbf{O}^i \in \mathbb{R}^{1 \times L}$ and a chromatin loop score vector $\mathbf{L}^i \in \mathbb{R}^{1 \times L}$, where L is the length
554 of sequence ($L=1000$). CREATE comprises encoders, a vector quantization module, and
555 decoders. The encoder module of CREATE includes encoders for multiple input-specific
556 learning and an encoder for integrating multiple inputs. Each encoder consists of a
557 convolutional layer, a max-pooling layer, a ReLU non-linear activation function and a dropout
558 layer. Correspondingly, each decoder consists of a deconvolutional layer, an upsample layer

559 and a Sigmoid or ReLU non-linear activation function. In addition, we introduced a classifier
 560 with three fully connected layers to predict CREs based on their embeddings. Specifically, the
 561 output of the encoder module is denoted as $\mathbf{e}^i \in \mathbb{R}^{L' \times D'}$ for the i -th CRE, where L' and D' are
 562 the length and dimensionality of the latent embedding respectively, and after split
 563 quantization^{32,36}, it will be split into $L' \times M$ vectors $\mathbf{e}_{l,j}^i \in \mathbb{R}^D, l \in \{1, \dots, L'\}, j \in \{1, \dots, M\}$, where
 564 M is the time of split quantization. Utilizing a shared codebook $\mathbf{v}_k \in \mathbb{R}^D, k \in \{1, \dots, K\}$ with the
 565 size of K , we obtained the quantized latent embedding $\mathbf{q}^i \in \mathbb{R}^{L' \times D'}$ for the i -th CRE by
 566 substituting the vector $\mathbf{e}_{l,j}^i$ with the nearest counterpart in the codebook as follows:

$$567 \quad \mathbf{q}_{l,j}^i = \mathbf{v}_{\underset{k \in \{1, \dots, K\}}{\operatorname{argmin}} \|\mathbf{e}_{l,j}^i - \mathbf{v}_k\|_2^2}, l \in \{1, \dots, L'\}, j \in \{1, \dots, M\}$$

568

569 **Model training.** We employed multiple update methods for different components of CREATE,
 570 mirroring the approach taken in the original studies of VQ-VAE^{34,35}. Let \mathcal{B}_0 be a mini-batch
 571 of data for training.

572 First, to optimize the decoder and encoder by reducing the distance between the original
 573 input and the reconstructed output, we integrated a hybrid reconstruction loss comprising
 574 multiple components corresponding to different inputs:

$$575 \quad L_{recon1}(\mathcal{B}_0) = -\frac{1}{P \cdot L \cdot |\mathcal{B}_0|} \sum_{i=1}^{|\mathcal{B}_0|} \sum_{l=1}^L \sum_{p=1}^P [\mathbf{s}_{lp}^i \log(\hat{\mathbf{S}}_{lp}^i) + (1 - \mathbf{s}_{lp}^i) \log(1 - \hat{\mathbf{S}}_{lp}^i)]$$

$$576 \quad L_{recon2}(\mathcal{B}_0) = \frac{1}{L \cdot |\mathcal{B}_0|} \sum_{i=1}^{|\mathcal{B}_0|} \sum_{l=1}^L \|\mathbf{o}_l^i - \hat{\mathbf{o}}_l^i\|_2^2$$

$$577 \quad L_{recon3}(\mathcal{B}_0) = \frac{1}{L \cdot |\mathcal{B}_0|} \sum_{i=1}^{|\mathcal{B}_0|} \sum_{l=1}^L \|\mathbf{l}_l^i - \hat{\mathbf{l}}_l^i\|_2^2$$

$$578 \quad L_{recon}(\mathcal{B}_0) = L_{recon1}(\mathcal{B}_0) + \beta L_{recon2}(\mathcal{B}_0) + \gamma L_{recon3}(\mathcal{B}_0)$$

579 where P represents the number of different types of bases in the DNA sequence ($P=4$), β and
 580 γ are the weights of L_{recon2} and L_{recon3} respectively ($\beta = 0.01, \gamma = 0.1$).

581 Second, to promote the encoder output to closely align with the selected codebook features
 582 and avoid excessive fluctuation, we introduced the encoder loss to aid in updating the encoder:

$$583 \quad L_{encoder}(\mathcal{B}_0) = \frac{1}{M \cdot L' \cdot |\mathcal{B}_0|} \sum_{i=1}^{|\mathcal{B}_0|} \sum_{l=1}^{L'} \sum_{j=1}^M \left\| \mathbf{e}_{l,j}^i - \operatorname{sg} \left(\mathbf{v}_{\underset{k \in \{1, \dots, K\}}{\operatorname{argmin}} \|\mathbf{e}_{l,j}^i - \mathbf{v}_k\|_2^2} \right) \right\|_2^2$$

584 where $\text{sg}(\cdot)$ denotes the stop-gradient operator with zero partial derivatives.

585 Third, we followed the recommendation from both the original studies of VQ-VAE and
 586 recent related researches⁹⁸⁻¹⁰¹ to utilize exponential moving average (EMA) for updating the
 587 codebook. Considering n_k is the number of vectors matched to \mathbf{v}_k and $\mathbf{e}_{k,m}^*$ is the m -th
 588 vector, we directly took the mean of the vectors in the set $\{\mathbf{e}_{k,m}^* | m = 1, \dots, n_k\}$ to optimize the
 589 code \mathbf{v}_k as follows:

$$590 \quad N_k^{(t)} = (1 - \mu)N_k^{(t-1)} + \mu n_k^{(t)}$$

$$591 \quad \mathbf{u}_k^{(t)} = (1 - \mu)\mathbf{u}_k^{(t-1)} + \sum_{m=1}^{n_k^{(t)}} \mu \mathbf{e}_{k,m}^{*,(t)}$$

$$592 \quad \mathbf{v}_k^{(t)} = \frac{\mathbf{u}_k^{(t)}}{N_k^{(t)}}$$

593 where μ is the update ratio of codebook. We initialized N_k as a zero vector and \mathbf{u}_k
 594 randomly from a normal distribution with a mean of 0 and a standard deviation of 1.

595 Fourth, we further incorporated a classifier based on the CRE embeddings, with the cross-
 596 entropy loss function given by:

$$597 \quad L_{class}(\mathcal{B}_0) = -\frac{1}{C \cdot |\mathcal{B}_0|} \sum_{i=1}^{|\mathcal{B}_0|} \sum_{c=1}^C y_c^i \log(\hat{y}_c^i)$$

598 where C represents the number of types of CREs ($C=5$). To sum up, we trained CREATE using
 599 EMA and the total loss function as follows:

$$600 \quad L_{CREATE}(\mathcal{B}_0) = L_{recon}(\mathcal{B}_0) + \alpha L_{encoder}(\mathcal{B}_0) + L_{class}(\mathcal{B}_0)$$

601 where α is the weights of $L_{encoder}$.

602 In this study, we implemented CREATE with “Pytorch” package¹⁰². In details, there are
 603 three one-dimensional convolutional layers (filters=256,128,128; size=8,8,8) with layer
 604 normalization in the input-specific encoder module, followed by three one-dimensional
 605 convolutional layers (filters=512,384,128; size=1,8,8). In all cases, we set the mini-batch size
 606 to 1024 and employed the Adam stochastic optimization algorithm¹⁰³ with a learning rate of
 607 5e-5. We trained CREATE with a maximum of 300 epochs and implemented early stopping if
 608 there were no reductions in validation auPRC for 20 consecutive epochs. We set the dimension
 609 of the latent embedding to 128 and trained CREATE with M of 16, K of 200, α of 0.25, and μ

610 of 0.01.

611

612 **Model evaluation.** To comprehensively evaluate the performance of CREATE for CRE
613 identification, we conducted 10-fold cross-validation experiments by dividing all CREs into
614 8:1:1 ratios for training, validation and testing data, respectively. We evenly distributed each
615 type of CRE into 10 folds. We compared the classification performance with four baseline
616 methods including DeepSEA¹⁷, DanQ¹⁸, ES-transition¹⁹ and DeepICSH²⁰, with the area under
617 the Receiver Operating Characteristic Curve (auROC), the area under the Precision-Recall
618 Curve (auPRC), F1-score, accuracy, precision and recall as evaluation metrics.

619

620 **Feature spectrum.** Supplementary Fig. 6a illustrates the process of generating the feature
621 spectrum. For the j -th codebook feature, we counted its occurrence frequency in the latent
622 embeddings of input regions, and we summed over these frequencies across all regions of the
623 i -th CRE to gain the frequency c_{ij} . We next derived a probability matrix ($\mathbf{P} \in \mathbb{R}^{C \times K}$) by the
624 following formula:

$$625 \quad t_{ij} = \frac{c_{ij}}{\sum_{k \in \{1, \dots, K\}} c_{ik}}$$

$$626 \quad p_{ij} = \frac{t_{ij}}{\sum_{c \in \{1, \dots, C\}} t_{cj}}$$

627 where C is the number of types of CREs and K is the number of codebook features. In this
628 matrix, a row corresponds to a type of CRE and a column to a codebook feature, and an element
629 p_{ij} indicates a feature probability score, representing the likelihood of the j -th codebook
630 feature appearing in the latent embeddings of the i -th CRE. For the j -th codebook feature, we
631 identified the element p_j with the highest feature probability score and the corresponding
632 CRE C_j yielding the score, as follows.

$$633 \quad p_j = \max_{i \in \{1, \dots, C\}} p_{ij}$$

$$634 \quad C_j = \operatorname{argmax}_{i \in \{1, \dots, C\}} p_{ij}$$

635 We then grouped the features corresponding to the same CRE together based on their CRE
636 indices (C_j), and further sorted these features in descending order according to their feature

637 probability scores (p_j). Finally, we attained the rearranged matrix $\mathbf{F} \in \mathbb{R}^{C \times K}$ as the
638 interpretable feature spectrum.

639

640 **Downstream analyses.**

641 *Motif enrichment analysis.* To discover enriched TF motifs for the true CREs and the predicted
642 CREs by CREATE, we applied the tool FIMO¹⁰⁴ with default settings to scan a set of input
643 sequences for searching known human TFs in the HOCOMOCO¹⁰⁵ database. For each input
644 sequence, we used Fisher's method to combine the P -values of reported binding sites for each
645 TF, and we obtained a P -value vector representing the significance that 678 human TFs
646 matched in the input sequence.

647 *Methylation levels computation.* The methylation state data at CpG in the K562 cell type was
648 obtained from ENCODE¹⁰⁶ (<https://www.encodeproject.org/files/ENCFF867JRG/>;
649 <https://www.encodeproject.org/files/ENCFF721JMB/>). Using BEDTools¹⁰⁷, we computed the
650 methylation levels for the true CREs and the predicted CREs by CREATE.

651 *Conservation scores computation.* We downloaded the phastCons^{71, 108} scores for multiple
652 alignments of 45 vertebrate genomes to the human genome from UCSC^{109, 110}
653 (<http://hgdownload.soe.ucsc.edu/goldenPath/hg19/phastCons46way/vertebrate.phastCons46way.bw>). The phastCons scores for the true CREs and the predicted CREs were calculated via
654 UCSC tool *bigWigAverageOverBed*¹¹¹.

656 *Overlaps with promoter-capture HiC regions.* The pcHiC data of K562 cell type was
657 downloaded from NCBI¹¹² under accession number "GSE236305". Using BEDTools¹⁰⁷, we
658 computed the number of overlaps between the pcHiC regions and the true CREs or the
659 predicted CREs by CREATE.

660 *Gene variation analysis.* We downloaded human all SNPs and common SNPs from dbSNP⁸⁰⁻
661 ⁸² database (https://ftp.ncbi.nih.gov/snp/organisms/human_9606_b151_GRCh37p13/VCF/),
662 all GTEx eQTLs, whole-blood eQTLs and liver eQTLs from the Genotype-Tissue Expression
663 Project^{76, 77}, GTEx database version 7 (<https://www.gtexportal.org/home/downloads/adult-gtex>). By excluding the common SNPs from all human SNPs, we obtained rare SNPs. Using
664 BEDTools¹⁰⁷, we considered the number of overlaps with SNPs or eQTLs as the corresponding
665

666 gene variation levels for the true CREs or the predicted CREs by CREATE.

667 *Tissue enrichment analysis.* To identify the tissues influenced by the identified risk loci within
668 the true CREs and the predicted CREs by CREATE, we performed SNPsea analysis⁸⁶ with
669 default settings. Based on the tissue-specific expression profiles of 17,581 genes across 79
670 human tissues (Gene Atlas¹¹³), we quantified the enrichments of these profiles on the true CREs
671 and the predicted CREs, and displayed the top 30 significantly enriched tissues in the heatmaps.

672 *Heritability enrichment analysis.* To quantify the enrichment of heritability for blood-related
673 phenotypes within the true CREs and the predicted CREs by CREATE, we conducted
674 heritability enrichment analysis using partitioned LDSC⁸⁷ with default settings. LDSC took
675 European samples from the 1000 Genomes Project as the LD reference panel. We downloaded
676 the HapMap3 SNPs and GWAS summary statistics from the Broad LD Hub
677 (<https://doi.org/10.5281/zenodo.7768714>), and then quantified the enrichment of heritability
678 for blood-related phenotypes, and displayed the results for the true CREs and the predicted
679 CREs.

680

681 **Baseline methods.** In this study, we compared CREATE to multiple baseline methods by
682 expanding them into multi-class models, including DeepSEA¹⁷, DanQ¹⁸, ES-transition¹⁹ and
683 DeepICSH²⁰. DeepSEA was implemented from their original source code
684 (<https://deepsea.princeton.edu/>). DanQ was implemented from their original source code
685 repository (<https://github.com/uci-cbcl/DanQ>). ES-transition was implemented from their
686 original source code repository (<https://github.com/ncbi/SilencerEnhancerPredict>). DeepICSH
687 was implemented from their original source code repository
688 (<https://github.com/lyli1013/DeepICSH>).

689

690 **Statistics and reproducibility**

691 No statistical method was used to predetermine sample size. No data were excluded from the
692 analyses. The experiments were not randomized. Data collection and analysis were not
693 performed blind to the conditions of the experiments.

694

695 **Data availability**

696 All datasets used in this study were obtained from public sources. We downloaded
697 experimentally validated silencers from the SilencerDB database⁸⁸
698 (<http://health.tsinghua.edu.cn/SilencerDB/>), enhancers from the FANTOM5 project^{24, 26}
699 (<https://bioinfo.vanderbilt.edu/AE/HACER/>), TSSs from the EPD database⁸⁹
700 (<https://epd.expasy.org/epd>), insulators from the ENCODE project^{92, 93}
701 (<https://www.encodeproject.org/files/ENCFF085HTY/>;
702 <https://www.encodeproject.org/files/ENCFF237OKO/>) for the K562 and HepG2 cell types. We
703 downloaded the histone modification ChIP-seq peaks and chromatin accessibility peaks from
704 the Roadmap project⁹⁴
705 (<http://egg2.wustl.edu/roadmap/data/byFileType/peaks/consolidated/narrowPeak>) and
706 ENCODE project (<https://www.encodeproject.org/files/ENCFF055NNT/>;
707 <https://www.encodeproject.org/files/ENCFF333TAT/>;
708 <https://www.encodeproject.org/files/ENCFF558BLC/>;
709 <https://www.encodeproject.org/files/ENCFF842UZU/>;
710 <https://www.encodeproject.org/files/ENCFF439EIO/>;
711 <https://www.encodeproject.org/files/ENCFF913MQB/>) for the K562 and HepG2 cell types.
712 The non-coding genes were obtained from GENCODE¹¹⁴ for human (GRCh37.p13/hg19). All
713 regions in this study are either in the genome of GRCh37/hg19 or have been converted to
714 GRCh37/hg19 by UCSC liftOver¹¹⁵ tool.

715

716 **Code availability**

717 The CREATE software, including detailed documents and tutorial, is freely available on
718 GitHub (<https://github.com/cuixj19/CREATE>).

719

720 **Acknowledgements**

721 This work was supported by the National Key Research and Development Program of China
722 grant nos. 2021YFF1200902 (R.J.), 2023YFF1204802 (R.J.), the National Natural Science
723 Foundation of China grants no. 62273194 (R.J.).

724

725 **Author contributions**

726 R.J. and W.Z. conceived the study and supervised the project. X.C. designed, implemented and
727 validated CREATE. Q.Y., Z.G., Z.L., X.C., S.C. and Q.L. helped analyze the results. X.C. and
728 W.Z. wrote the manuscript, with input from all the authors.

729

730 **Ethics declarations**

731 **Competing interests**

732 The authors declare no competing interests.

733

734

735 Reference

- 736 1. Berger, S.L. The complex language of chromatin regulation during transcription. *Nature* **447**, 407-412
737 (2007).
- 738 2. Lee, T.I. & Young, R.A. Transcriptional regulation and its misregulation in disease. *Cell* **152**, 1237-1251
739 (2013).
- 740 3. Maston, G.A., Evans, S.K. & Green, M.R. Transcriptional regulatory elements in the human genome.
741 *Annu. Rev. Genomics Hum. Genet.* **7**, 29-59 (2006).
- 742 4. Chatterjee, S. & Ahituv, N. Gene regulatory elements, major drivers of human disease. *Annu. Rev.*
743 *Genomics Hum. Genet.* **18**, 45-63 (2017).
- 744 5. Klemm, S.L., Shipony, Z. & Greenleaf, W.J. Chromatin accessibility and the regulatory epigenome. *Nat.*
745 *Rev. Genet.* **20**, 207-220 (2019).
- 746 6. Minnoye, L. et al. Chromatin accessibility profiling methods. *Nat. Rev. Methods Primers* **1**, 10 (2021).
- 747 7. Kadauke, S. & Blobel, G.A. Chromatin loops in gene regulation. *Biochimica et Biophysica Acta (BBA)-*
748 *Gene Regulatory Mechanisms* **1789**, 17-25 (2009).
- 749 8. Dekker, J., Marti-Renom, M.A. & Mirny, L.A. Exploring the three-dimensional organization of genomes:
750 interpreting chromatin interaction data. *Nat. Rev. Genet.* **14**, 390-403 (2013).
- 751 9. Pang, B., van Weerd, J.H., Hamoen, F.L. & Snyder, M.P. Identification of non-coding silencer elements
752 and their regulation of gene expression. *Nat. Rev. Mol. Cell Biol.* **24**, 383-395 (2023).
- 753 10. Sealfon, R.S., Wong, A.K. & Troyanskaya, O.G. Machine learning methods to model multicellular
754 complexity and tissue specificity. *Nature Reviews Materials* **6**, 717-729 (2021).
- 755 11. Vaishnav, E.D. et al. The evolution, evolvability and engineering of gene regulatory DNA. *Nature* **603**,
756 455-463 (2022).
- 757 12. Huang, D., Petrykowska, H.M., Miller, B.F., Elnitski, L. & Ovcharenko, I. Identification of human
758 silencers by correlating cross-tissue epigenetic profiles and gene expression. *Genome Res.* **29**, 657-667
759 (2019).
- 760 13. Lieberman-Aiden, E. et al. Comprehensive mapping of long-range interactions reveals folding principles
761 of the human genome. *Science* **326**, 289-293 (2009).
- 762 14. Song, L. & Crawford, G.E. DNase-seq: a high-resolution technique for mapping active gene regulatory
763 elements across the genome from mammalian cells. *Cold Spring Harb. Protoc.* **2010**, pdb. prot5384
764 (2010).
- 765 15. Buenrostro, J.D., Giresi, P.G., Zaba, L.C., Chang, H.Y. & Greenleaf, W.J. Transposition of native
766 chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and
767 nucleosome position. *Nat. Methods* **10**, 1213-1218 (2013).
- 768 16. Mumbach, M.R. et al. HiChIP: efficient and sensitive analysis of protein-directed genome architecture.
769 *Nat. Methods* **13**, 919-922 (2016).
- 770 17. Zhou, J. & Troyanskaya, O.G. Predicting effects of noncoding variants with deep learning-based
771 sequence model. *Nat. Methods* **12**, 931-934 (2015).
- 772 18. Quang, D. & Xie, X. DanQ: a hybrid convolutional and recurrent deep neural network for quantifying
773 the function of DNA sequences. *Nucleic Acids Res.* **44**, e107-e107 (2016).
- 774 19. Huang, D. & Ovcharenko, I. Enhancer-silencer transitions in the human genome. *Genome Res.* **32**, 437-
775 448 (2022).
- 776 20. Zhang, T., Li, L., Sun, H., Xu, D. & Wang, G. DeepICSH: a complex deep learning framework for
777 identifying cell-specific silencers and their strength from the human genome. *Brief. Bioinform.* **24**,
778 bbad316 (2023).
- 779 21. Chen, S., Gan, M., Lv, H. & Jiang, R. DeepCAPE: a deep convolutional neural network for the accurate

- 780 prediction of enhancers. *Genomics Proteomics Bioinformatics* **19**, 565-577 (2021).
- 781 22. Min, X. et al. Predicting enhancers with deep convolutional neural networks. *BMC Bioinformatics* **18**,
782 35-46 (2017).
- 783 23. Oubounyt, M., Louadi, Z., Tayara, H. & Chong, K.T. DeePromoter: robust promoter predictor using deep
784 learning. *Front. Genet.* **10**, 453150 (2019).
- 785 24. Andersson, R. et al. An atlas of active enhancers across human cell types and tissues. *Nature* **507**, 455-
786 461 (2014).
- 787 25. Khan, A. & Zhang, X. dbSUPER: a database of super-enhancers in mouse and human genome. *Nucleic
788 Acids Res.* **44**, D164-D171 (2016).
- 789 26. Wang, J. et al. HACER: an atlas of human active enhancers to interpret regulatory variants. *Nucleic Acids
790 Res.* **47**, D106-D112 (2019).
- 791 27. Zeng, W., Min, X. & Jiang, R. EnDisease: a manually curated database for enhancer-disease associations.
792 *Database* **2019**, baz020 (2019).
- 793 28. Jiang, Y. et al. SEdb: a comprehensive human super-enhancer database. *Nucleic Acids Res.* **47**, D235-
794 D243 (2019).
- 795 29. Bai, X. et al. ENdb: a manually curated database of experimentally supported enhancers for human and
796 mouse. *Nucleic Acids Res.* **48**, D51-D57 (2020).
- 797 30. Gao, T. & Qian, J. EnhancerAtlas 2.0: an updated resource with enhancer annotation in 586 tissue/cell
798 types across nine species. *Nucleic Acids Res.* **48**, D58-D64 (2020).
- 799 31. Ogbourne, S. & Antalis, T.M. Transcriptional control and the role of silencers in transcriptional regulation
800 in eukaryotes. *Biochem. J.* **331**, 1-14 (1998).
- 801 32. Cui, X. et al. Discrete latent embedding of single-cell chromatin accessibility sequencing data for
802 uncovering cell heterogeneity. *Nat. Comput. Sci.*, 1-14 (2024).
- 803 33. Zheng, A. et al. Deep neural networks identify sequence context features predictive of transcription factor
804 binding. *Nat. Mach. Intell.* **3**, 172-180 (2021).
- 805 34. Van Den Oord, A. & Vinyals, O. Neural discrete representation learning. In *Proc. Advances in Neural
806 Information Processing Systems* **30** (2017).
- 807 35. Razavi, A., Van den Oord, A. & Vinyals, O. Generating diverse high-fidelity images with vq-vae-2. In
808 *Adv. Neural Inf. Process. Syst.* **32** (2019).
- 809 36. Kobayashi, H., Cheveralls, K.C., Leonetti, M.D. & Royer, L.A. Self-supervised deep learning encodes
810 high-resolution features of protein subcellular localization. *Nat. Methods* **19**, 995-1003 (2022).
- 811 37. McInnes, L., Healy, J. & Melville, J. Umap: Uniform manifold approximation and projection for
812 dimension reduction. *arXiv preprint arXiv:1802.03426* (2018).
- 813 38. Igarashi, K. et al. Multivalent DNA binding complex generated by small Maf and Bach1 as a possible
814 biochemical basis for β -globin locus control region complex. *J. Biol. Chem.* **273**, 11783-11790 (1998).
- 815 39. Zhang, Z. et al. The LIM homeodomain transcription factor LHX6: a transcriptional repressor that
816 interacts with pituitary homeobox 2 (PITX2) to regulate odontogenesis. *J. Biol. Chem.* **288**, 2485-2500
817 (2013).
- 818 40. Gu, X. et al. PBRM1 loss in kidney cancer unbalances the proximal tubule master transcription factor
819 hub to repress proximal tubule differentiation. *Cell Rep.* **36** (2021).
- 820 41. Wang, J., Jia, Q., Jiang, S., Lu, W. & Ning, H. POU6F1 promotes ferroptosis by increasing lncRNA-
821 CASC2 transcription to regulate SOCS2/SLC7A11 signaling in gastric cancer. *Cell Biol. Toxicol.* **40**, 1-
822 17 (2024).
- 823 42. Cowling, V.H. & Cole, M.D. Mechanism of transcriptional activation by the Myc oncoproteins. In *Semin.
824 Cancer Biol.* **16**, 242-252 (2006).
- 825 43. Chittka, A., Nitarska, J., Grazini, U. & Richardson, W.D. Transcription factor positive regulatory domain

- 826 4 (PRDM4) recruits protein arginine methyltransferase 5 (PRMT5) to mediate histone arginine
827 methylation and control neural stem cell proliferation and differentiation. *J. Biol. Chem.* **287**, 42995-
828 43006 (2012).
- 829 44. Li, N., He, Y., Mi, P. & Hu, Y. ZNF582 methylation as a potential biomarker to predict cervical
830 intraepithelial neoplasia type III/worse: A meta-analysis of related studies in Chinese population.
831 *Medicine* **98**, e14297 (2019).
- 832 45. Ha, T.J. et al. Identification of novel cerebellar developmental transcriptional regulators with motif
833 activity analysis. *BMC Genomics* **20**, 1-17 (2019).
- 834 46. Della Rosa, M. & Spivakov, M. Silencers in the spotlight. *Nat. Genet.* **52**, 244-245 (2020).
- 835 47. Ngan, C.Y. et al. Chromatin interaction analyses elucidate the roles of PRC2-bound silencers in mouse
836 development. *Nat. Genet.* **52**, 264-272 (2020).
- 837 48. Cai, Y. et al. H3K27me3-rich genomic regions can function as silencers to repress gene expression via
838 chromatin interactions. *Nat. Commun.* **12**, 719 (2021).
- 839 49. Karmodiya, K., Krebs, A.R., Oulad-Abdelghani, M., Kimura, H. & Tora, L. H3K9 and H3K14
840 acetylation co-occur at many gene regulatory elements, while H3K14ac marks a subset of inactive
841 inducible promoters in mouse embryonic stem cells. *BMC Genomics* **13**, 1-18 (2012).
- 842 50. Oka, R. et al. Genome-wide mapping of transcriptional enhancer candidates using DNA and chromatin
843 features in maize. *Genome Biol.* **18**, 1-24 (2017).
- 844 51. Zhu, Y. et al. Predicting enhancer transcription and activity from chromatin modifications. *Nucleic Acids*
845 *Res.* **41**, 10032-10043 (2013).
- 846 52. Bogdanović, O. et al. Dynamics of enhancer chromatin signatures mark the transition from pluripotency
847 to cell specification during embryogenesis. *Genome Res.* **22**, 2043-2053 (2012).
- 848 53. Local, A. et al. Identification of H3K4me1-associated proteins at mammalian enhancers. *Nat. Genet.* **50**,
849 73-82 (2018).
- 850 54. Pekowska, A. et al. H3K4 tri - methylation provides an epigenetic signature of active enhancers. *EMBO*
851 *J.* **30**, 4198-4210 (2011).
- 852 55. Chen, K. et al. Broad H3K4me3 is associated with increased transcription elongation and enhancer
853 activity at tumor-suppressor genes. *Nat. Genet.* **47**, 1149-1157 (2015).
- 854 56. Whitfield, T.W. et al. Functional analysis of transcription factor binding sites in human promoters.
855 *Genome Biol.* **13**, 1-16 (2012).
- 856 57. Boeva, V. Analysis of genomic sequence motifs for deciphering transcription factor binding and
857 transcriptional regulation in eukaryotic cells. *Front. Genet.* **7**, 174397 (2016).
- 858 58. Inukai, S., Kock, K.H. & Bulyk, M.L. Transcription factor–DNA binding: beyond binding site motifs.
859 *Curr. Opin. Genet. Dev.* **43**, 110-119 (2017).
- 860 59. Polevoy, H., Malyarova, A., Fonar, Y., Elias, S. & Frank, D. FoxD1 protein interacts with Wnt and BMP
861 signaling to differentially pattern mesoderm and neural tissue. *Int. J. Dev. Biol.* **61**, 293-302 (2017).
- 862 60. Scibetta, A.G., Wong, P.-P., Chan, K.V., Canosa, M. & Hurst, H.C. Dual association by TFAP2A during
863 activation of the p21cip/CDKN1A promoter. *Cell Cycle* **9**, 4525-4532 (2010).
- 864 61. Sieweke, M.H., Tekotte, H., Frampton, J. & Graf, T. MafB is an interaction partner and repressor of Ets-
865 1 that inhibits erythroid differentiation. *Cell* **85**, 49-60 (1996).
- 866 62. Piper, M. et al. NFIA controls telencephalic progenitor cell differentiation through repression of the
867 Notch effector Hes1. *J. Neurosci.* **30**, 9127-9139 (2010).
- 868 63. Davis, C.A. et al. PRISM/PRDM6, a transcriptional repressor that promotes the proliferative gene
869 program in smooth muscle cells. *Mol. Cell. Biol.* **26**, 2626-2636 (2006).
- 870 64. Parsons, M.J. et al. The regulatory factor ZFH3 modifies circadian function in SCN via an AT motif-
871 driven axis. *Cell* **162**, 607-621 (2015).

- 872 65. Schepers, G.E., Bullejos, M., Hosking, B.M. & Koopman, P. Cloning and characterisation of the Sry-
873 related transcription factor gene Sox8. *Nucleic Acids Res.* **28**, 1473-1480 (2000).
- 874 66. Baluapuri, A., Wolf, E. & Eilers, M. Target gene-independent functions of MYC oncoproteins. *Nat. Rev.*
875 *Mol. Cell Biol.* **21**, 255-267 (2020).
- 876 67. Rubio-Alarcón, M. et al. Zfh3 transcription factor represses the expression of SCN5A gene and
877 decreases sodium current density (INa). *Int. J. Mol. Sci.* **22**, 13031 (2021).
- 878 68. Doni Jayavelu, N., Jajodia, A., Mishra, A. & Hawkins, R.D. Candidate silencer elements for the human
879 and mouse genomes. *Nat. Commun.* **11**, 1061 (2020).
- 880 69. Smith, Z.D. & Meissner, A. DNA methylation: roles in mammalian development. *Nat. Rev. Genet.* **14**,
881 204-220 (2013).
- 882 70. Moore, L.D., Le, T. & Fan, G. DNA methylation and its basic function. *Neuropsychopharmacology* **38**,
883 23-38 (2013).
- 884 71. Siepel, A. et al. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes.
885 *Genome Res.* **15**, 1034-1050 (2005).
- 886 72. Villar, D. et al. Enhancer evolution across 20 mammalian species. *Cell* **160**, 554-566 (2015).
- 887 73. Ovcharenko, I. et al. Evolution and functional classification of vertebrate gene deserts. *Genome Res.* **15**,
888 137-145 (2005).
- 889 74. Schoenfelder, S., Javierre, B.-M., Furlan-Magaril, M., Wingett, S.W. & Fraser, P. Promoter capture Hi-
890 C: high-resolution, genome-wide profiling of promoter interactions. *J. Vis. Exp.*, e57320 (2018).
- 891 75. Gisselbrecht, S.S. et al. Transcriptional silencers in Drosophila serve a dual role as transcriptional
892 enhancers in alternate cellular contexts. *Mol. Cell* **77**, 324-337. e328 (2020).
- 893 76. Lonsdale, J. et al. The genotype-tissue expression (GTEx) project. *Nat. Genet.* **45**, 580-585 (2013).
- 894 77. Consortium, G. The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science*
895 **369**, 1318-1330 (2020).
- 896 78. Maurano, M.T. et al. Systematic localization of common disease-associated variation in regulatory DNA.
897 *Science* **337**, 1190-1195 (2012).
- 898 79. Schaub, M.A., Boyle, A.P., Kundaje, A., Batzoglou, S. & Snyder, M. Linking disease associations with
899 regulatory information in the human genome. *Genome Res.* **22**, 1748-1759 (2012).
- 900 80. Sherry, S.T., Ward, M. & Sirotkin, K. dbSNP—database for single nucleotide polymorphisms and other
901 classes of minor genetic variation. *Genome Res.* **9**, 677-679 (1999).
- 902 81. Sherry, S.T. et al. dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res.* **29**, 308-311 (2001).
- 903 82. Smigielski, E.M., Sirotkin, K., Ward, M. & Sherry, S.T. dbSNP: a database of single nucleotide
904 polymorphisms. *Nucleic Acids Res.* **28**, 352-355 (2000).
- 905 83. Hawkins, R.D. et al. Global chromatin state analysis reveals lineage-specific enhancers during the
906 initiation of human T helper 1 and T helper 2 cell polarization. *Immunity* **38**, 1271-1284 (2013).
- 907 84. Raychaudhuri, S. et al. Identifying relationships among genomic disease regions: predicting genes at
908 pathogenic SNP associations and rare deletions. *PLoS Genet.* **5**, e1000534 (2009).
- 909 85. Schork, N.J., Murray, S.S., Frazer, K.A. & Topol, E.J. Common vs. rare allele hypotheses for complex
910 diseases. *Curr. Opin. Genet. Dev.* **19**, 212-219 (2009).
- 911 86. Slowikowski, K., Hu, X. & Raychaudhuri, S. SNPsea: an algorithm to identify cell types, tissues and
912 pathways affected by risk loci. *Bioinformatics* **30**, 2496-2497 (2014).
- 913 87. Finucane, H.K. et al. Partitioning heritability by functional annotation using genome-wide association
914 summary statistics. *Nat. Genet.* **47**, 1228-1235 (2015).
- 915 88. Zeng, W. et al. SilencerDB: a comprehensive database of silencers. *Nucleic Acids Res.* **49**, D221-D228
916 (2021).
- 917 89. Meylan, P., Dreos, R., Ambrosini, G., Groux, R. & Bucher, P. EPD in 2020: enhanced data visualization

- 918 and extension to ncRNA promoters. *Nucleic Acids Res.* **48**, D65-D69 (2020).
- 919 90. Bell, A.C., West, A.G. & Felsenfeld, G. The protein CTCF is required for the enhancer blocking activity
920 of vertebrate insulators. *Cell* **98**, 387-396 (1999).
- 921 91. Nora, E.P. et al. Targeted degradation of CTCF decouples local insulation of chromosome domains from
922 genomic compartmentalization. *Cell* **169**, 930-944. e922 (2017).
- 923 92. Consortium, E.P. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57
924 (2012).
- 925 93. Luo, Y. et al. New developments on the Encyclopedia of DNA Elements (ENCODE) data portal. *Nucleic
926 Acids Res.* **48**, D882-D889 (2020).
- 927 94. Consortium, R.E. Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317-330 (2015).
- 928 95. Chen, S. et al. OpenAnnotate: a web server to annotate the chromatin accessibility of genomic regions.
929 *Nucleic Acids Res.* **49**, W483-W490 (2021).
- 930 96. Zeng, W., Liu, Q., Yin, Q., Jiang, R. & Wong, W.H. HiChIPdb: a comprehensive database of HiChIP
931 regulatory interactions. *Nucleic Acids Res.* **51**, D159-D166 (2023).
- 932 97. Li, W., Wong, W.H. & Jiang, R. DeepTACT: predicting 3D chromatin contacts via bootstrapping deep
933 learning. *Nucleic Acids Res.* **47**, e60-e60 (2019).
- 934 98. Kaiser, L. et al. Fast decoding in sequence models using discrete latent variables. In *Proc. International
935 Conference on Machine Learning*, 2390-2399 (2018).
- 936 99. Peng, J., Liu, D., Xu, S. & Li, H. Generating diverse structure for image inpainting with hierarchical VQ-
937 VAE. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10775-10784 (2021).
- 938 100. Williams, W. et al. Hierarchical quantized autoencoders. In *Proc. Advances in Neural Information
939 Processing Systems* **33**, 4524-4535 (2020).
- 940 101. Takida, Y. et al. Sq-vae: Variational bayes on discrete representation with self-annealed stochastic
941 quantization. *arXiv preprint arXiv:2205.07547* (2022).
- 942 102. Paszke, A. et al. Pytorch: An imperative style, high-performance deep learning library. In *Proc. Advances
943 in Neural Information Processing Systems* **32** (2019).
- 944 103. Kingma, D.P. & Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*
945 (2014).
- 946 104. Grant, C.E., Bailey, T.L. & Noble, W.S. FIMO: scanning for occurrences of a given motif. *Bioinformatics*
947 **27**, 1017-1018 (2011).
- 948 105. Kulakovskiy, I.V. et al. HOCOMOCO: towards a complete collection of transcription factor binding
949 models for human and mouse via large-scale CHIP-Seq analysis. *Nucleic Acids Res.* **46**, D252-D259
950 (2018).
- 951 106. Zhang, J. et al. An integrative ENCODE resource for cancer genomics. *Nat. Commun.* **11**, 3696 (2020).
- 952 107. Quinlan, A.R. & Hall, I.M. BEDTools: a flexible suite of utilities for comparing genomic features.
953 *Bioinformatics* **26**, 841-842 (2010).
- 954 108. Pollard, K.S., Hubisz, M.J., Rosenbloom, K.R. & Siepel, A. Detection of nonneutral substitution rates
955 on mammalian phylogenies. *Genome Res.* **20**, 110-121 (2010).
- 956 109. Kent, W.J. et al. The human genome browser at UCSC. *Genome Res.* **12**, 996-1006 (2002).
- 957 110. Nassar, L.R. et al. The UCSC genome browser database: 2023 update. *Nucleic Acids Res.* **51**, D1188-
958 D1195 (2023).
- 959 111. Kent, W.J., Zweig, A.S., Barber, G., Hinrichs, A.S. & Karolchik, D. BigWig and BigBed: enabling
960 browsing of large distributed datasets. *Bioinformatics* **26**, 2204-2207 (2010).
- 961 112. Wheeler, D.L. et al. Database resources of the national center for biotechnology information. *Nucleic
962 Acids Res.* **36**, D13-D21 (2007).
- 963 113. Su, A.I. et al. A gene atlas of the mouse and human protein-encoding transcriptomes. *Proc. Natl. Acad.*

- 964 *Sci. U. S. A.* **101**, 6062-6067 (2004).
- 965 114. Frankish, A. et al. GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids*
966 *Res.* **47**, D766-D773 (2019).
- 967 115. Hinrichs, A.S. et al. The UCSC genome browser database: update 2006. *Nucleic Acids Res.* **34**, D590-
968 D598 (2006).
- 969
- 970