

# Single-Cell ATAC-seq analysis via Network Refinement with peaks location information

Jiating Yu<sup>1,2</sup>, Duanchen Sun<sup>3,\*</sup>, Zhichao Hou<sup>1,2</sup>, Ling-Yun Wu<sup>1,2,\*</sup>

<sup>1</sup>IAM, MADIS, NCMIS, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing 100190, China; <sup>2</sup>School of Mathematical Sciences, University of Chinese Academy of Sciences, Beijing 100049, China; <sup>3</sup>School of Mathematics, Shandong University, Jinan, Shandong 250100, China

\*To whom correspondence should be addressed. Email: [lywu@amss.ac.cn](mailto:lywu@amss.ac.cn), [dcsun@sdu.edu.cn](mailto:dcsun@sdu.edu.cn)

1   **Abstract:** Single-cell ATAC-seq (scATAC-seq) data provided new insights into the elaboration of  
2   cellular heterogeneity and transcriptional regulation. However, scATAC-seq data posed challenges  
3   for data analysis because of its near binarization, high sparsity, and ultra-high dimensionality  
4   properties. Here we proposed a novel network diffusion-based method to comprehensively analyze  
5   scATAC-seq data, named **Single-Cell ATAC-seq Analysis via Network Refinement with Peaks**  
6   Location Information (SCARP). By modeling the prior probability of co-accessibility between  
7   adjacent peaks as a decreasing function of genomic distance, SCARP is the first scATAC-seq  
8   analysis method that utilizes the genomic information of peaks, which contributed to characterizing  
9   co-accessibility of peaks. SCARP used network to model the accessible relationships between cells  
10   and peaks, aggregated information with the diffusion method, and then performed dimensionality  
11   reduction to obtain low-dimensional cell embeddings as well as peak embeddings. We have  
12   demonstrated through sufficient experiments that SCARP facilitated superior analysis of scATAC-  
13   seq data. Specifically, SCARP exhibited outstanding cell clustering performance to better elucidate  
14   cell heterogeneity, and can be used to reveal new biologically significant cell subpopulations. SCARP  
15   was also instrumental in portraying co-accessibility relationships of accessible regions and providing  
16   new insight into transcriptional regulation, and those SCARP-derived genes were involved in some  
17   key KEGG pathways related to diseases. To sum up, our studies suggested that SCARP is a  
18   promising tool to comprehensively analyze the scATAC-seq data from a new perspective.

## 19 Introduction

20 Chromatin accessibility is closely linked to the occurrence of transcriptional regulation, as accessible  
21 chromosome fragments often contain transcription factor (TF) binding sites and other important cis-  
22 regulatory elements, such as enhancers and promoters [1][2]. While single-cell transcriptome data  
23 can help us to construct gene regulatory networks from the level of transcripts and identify cellular  
24 diversity without bias, it is still challenging to reverse-engineer the mechanism of transcriptional  
25 regulation [2]. Fortunately, the development of single-cell epigenomic technology has brought new  
26 perspectives to solve this problem, such as single-cell assay for transposase-accessible chromatin  
27 using sequencing (scATAC-seq) [3] and single-cell combinatorial indexing assay for transposase  
28 accessible chromatin with sequencing (sci-ATAC-seq) [4][5], which can build chromatin accessibility  
29 profiles with single-cell resolution [6]. With the aid of scATAC-seq data, we were able to model the  
30 co-accessibility of chromosome segments from the genome level directly, which provides the basic  
31 conditions for transcriptional regulation to occur.

32 Compared with single-cell transcriptome data, scATAC-seq data exhibit more sparsity, higher  
33 dimension, and near-binariization properties, which pose many challenges to scATAC-based  
34 analysis. Because for a diploid genome, only one or two reads can be captured at each accessible  
35 chromatin site [6], resulting in the lack of sequencing data, i.e., the curse of ‘missingness’ in scATAC-  
36 seq [1][6]. In addition, accessible regions are usually called on the entire genome, so the features of  
37 scATAC-seq data can reach more than 100k dimensions or even higher. These properties make  
38 scRNA-seq-based algorithms less effective for analyzing scATAC-seq [7], although these algorithms  
39 are mathematically transferable since both genes and peaks can be understood as features to  
40 characterize cells.

41 Currently, various methods have been specially developed for scATAC-seq data to decipher their  
42 underlying information. For example, chromVAR [8] measures the accessibility gain or loss among  
43 peaks with the same motif or annotation and uses a bias-corrected deviation matrix for downstream  
44 analysis. However, the aggregation of similar peaks loses much information, making it less effective  
45 for cell clustering. SCALE [1] incorporates the variational autoencoder and Gaussian Mixture Model  
46 to extract latent features of cells. Nevertheless, it has numerous parameters and needs to be  
47 carefully tuned to get good performance [6]. Besides, cisTopic [2] uses the Latent Dirichlet Allocation  
48 model for co-optimal clustering of the cells and accessible regions, yet the collapsed Gibbs sampler  
49 makes it computationally slow. A recently proposed method, scAND [6], adopts the Katz index [9]

50 diffusion method to alleviate the high sparsity in scATAC-seq data, and using the smoothed matrix  
51 as model input can achieve satisfactory results in their following analysis. However, scAND does not  
52 consider the differences between different chromosomes in the diffusion process and neglects the  
53 location associations of peaks.

54 In genetics studies such as recombination between two loci and linkage disequilibrium [10]-[13],  
55 the genetic map distance has been considered an essential factor. Genes on different chromosomes  
56 or distantly separated on the same chromosomes are assorted independently and described as  
57 physically unlinked [12], and the plausibility of modeling chromatin contacts as a decreasing function  
58 of genomic distances has been explained and verified [14]. Therefore, as scATAC-seq data give the  
59 accessibility of genome-wide segments in cells, incorporating genomic distance information of  
60 accessible fragments may be beneficial for the downstream data analyses, such as cell  
61 subpopulation identification and cis-regulatory network construction, which has not been considered  
62 in the state-of-the-art methods. In addition, it is also critical to exploit a superior network diffusion  
63 process that can appropriately aggregate the neighborhood information to deal with the sparseness  
64 and missingness problems in scATAC-seq data. Our previously developed Network Refinement (NR)  
65 [15] diffusion method, which is a degree-normalized version of the Katz index, has been shown to  
66 achieve better performances in a series of applications [15][16]. We, therefore, expect to see that  
67 NR also handles scATAC-seq data well, with the rationale being that: 1) The diffusion of accessibility  
68 relationships can compensate for missing information and depict similarities between nodes (i.e.,  
69 cells and peaks) of the network; 2) Peaks accessible in too many cells should reduce their impact  
70 on diffusion, as their ability to portray cell similarity is overestimated (Methods).

71 In this study, we proposed a novel network-based method to comprehensively analyze scATAC-  
72 seq data, named **S**ingle-**C**ell **A**TAC-seq **A**nalysis via **N**etwork **R**efinement with **P**eaks **L**ocation  
73 **I**nformation (SCARP). Specifically, SCARP takes genetic location information of peaks into account  
74 and globally diffuses the accessibility relationships using our previously developed Network  
75 Refinement diffusion method (Methods). The output matrix derived from SCARP can be further  
76 processed by the dimension reduction method to obtain low-dimensional embeddings of cells and  
77 peaks, which can benefit the downstream analyses such as the cells clustering and cis-regulatory  
78 relationships prediction (Fig. 1).

79 We have demonstrated through sufficient experiments that SCARP facilitates superior analysis  
80 of scATAC-seq data, including improving cell clustering performance to better elucidate cell

81 heterogeneity and reveal new biologically significant cell subpopulations, and characterizing co-  
82 accessibility of genome segments to shed fresh light on transcriptional regulation. Our studies  
83 suggested that SCARP is a promising tool to comprehensively analyze the scATAC-seq data from a  
84 new perspective.  
85

## 86 **Results**

### 87 **Overview of SCARP workflow**

88 We proposed a novel network diffusion-based method, SCARP, to comprehensively analyze the  
89 scATAC-seq data. The workflow of SCARP was shown in Fig. 1. First, SCARP constructed a bipartite  
90 network to model the accessible relationships between the cells and peaks (Fig. 1a). Second, prior  
91 edge weights of adjacent peaks based on their genome distance were employed to better capture  
92 the co-accessibility of adjacent peaks (Fig. 1b). Together they constituted the input network for the  
93 next step (Fig. 1c). Third, the NR diffusion processes [15] were performed on the subnetworks  
94 corresponding to different chromosomes separately to obtain dense matrixes that can reflect  
95 accessibility similarities between the cells and peaks, as well as possible cell-peak accessible  
96 relationships (Fig. 1d). After this, diffusion matrixes obtained from different chromosomes were  
97 spliced together and the cell-cell similarities were integrated (Fig. 1e). Finally, dimensionality  
98 reduction technique was performed to map cells and peaks into low-dimensional space, and these  
99 representations were used for downstream analyses, such as cells clustering and construction of  
100 cis-regulatory networks (Fig. 1f).

101

### 102 **SCARP exhibited superior cell clustering performance on benchmarking scATAC-seq 103 datasets**

104 We selected nine benchmarking scATAC-seq datasets with reference cell type annotations. For  
105 these datasets, the number of peaks ranges from 7,000 to 140,000 with different sparsity levels (Fig.  
106 2a and Supplementary Table 1). We compared SCARP with six state-of-the-art methods (original,  
107 DCA [17], MAGIC [18], PCA, cisTopic [2], and scAND [6]) as well as SCARP without using prior  
108 weights of adjacency peaks (denoted as SCARP\*). Notably, DCA and MAGIC were especially  
109 designed for scRNA-seq analysis (Supplementary Materials). The running time of SCARP is  
110 advantageous when compared with these methods (Fig. 2b and Supplementary Table 3). Using  
111 blood2K dataset as an example, we found stronger consistency between the identified cell clusters

112 using SCARP cell embeddings and the annotated cell types when compared with other candidate  
113 methods (Fig. 2c and Supplementary Figs. 1-2). The SCARP-derived low dimensional visualization  
114 of cells also exhibits more explicit cell type boundaries and a consistent developmental trajectory  
115 with known FACS-sorted population [6][19] (Fig. 2d and Supplementary Figs. 3-6). For all  
116 benchmarking datasets, SCARP achieved the top-ranked clustering performances with different  
117 evaluation metrics (Fig. 2e and Supplementary Fig. 7). It is noteworthy that SCARP\* ranked second,  
118 meaning that both NR diffusion and prior weights play essential roles in obtaining better low-  
119 dimensional cell representations. Together with the observations that SCARP had the best averaged  
120 clustering performances regarding different metrics (Fig. 2f, g), all these pieces of evidence  
121 comprehensively demonstrated the superiority of SCARP over other methods in cell feature  
122 extraction.

123 Since cell type annotations based on RNA-seq data are more readily available, we further  
124 investigate whether the scATAC-derived cell clustering is consistent with the scRNA-based results,  
125 using the SNARE-seq data, which is paired multi-omics data that simultaneously profile gene  
126 expressions and chromatin accessibilities [20]. We used the cell clustering results on RNA-seq data  
127 processed by Seurat [21] and Signac [22] as a reference, and compared their consistency with those  
128 of SCARP, scAND, and cisTopic on scATAC-seq data. We observed that SCARP yielded higher cell  
129 clustering concordance between paired multi-omics data (Fig. 2h).

130 We next explored the robustness of SCARP in dealing with different peak numbers, missing data,  
131 and different model parameters. Since the blood2K dataset has the maximum number of peaks, we  
132 used it to test the impact of peak filtering to the performance of SCARP. We found that the number  
133 of peaks and the running time of SCARP decreased as the filtering level increased under both count-  
134 based and variance-based peaks filtering strategies (Fig. 2i). In contrast, the accuracy of cell  
135 clustering of SCARP had not been much affected, and the cell clustering performances were still  
136 better than that of two other competitive methods for most filtering levels (Fig. 2j). Besides, we  
137 randomly down-sampled the counts in Leukemia dataset, which has the minimum number of peaks,  
138 to simulate the missing values in scATAC-seq data. We observed that our method was not sensitive  
139 to missing data for different down-sampling ratios and got the best performance overall (Fig. 2k and  
140 Supplementary Fig. 8). Furthermore, SCARP has three model parameters ( $\beta$ ,  $m$ , and  $\gamma$ ; Methods).  
141 We found that SCARP is robust to the parameter selection since the cell clustering results only  
142 slightly fluctuated (Fig. 2l and Supplementary Figs. 9-11).

143 Taking together, SCARP is a fast and robust method to extract informative low-dimensional cell  
144 embeddings and has an outstanding performance of cell clustering compared with all the other  
145 methods.

146

147 **SCARP identified a new CD14 monocyte subpopulation characterized with high**  
148 **differentiation activities**

149 We applied SCARP to a 10X multiome peripheral blood mononuclear cell (PBMC) dataset [23],  
150 where scRNA-seq and scATAC-seq data can be acquired simultaneously. The cell clusters  
151 generated using SCARP-derived cell embeddings were well-separated between different cell types  
152 (Fig. 3a). Interestingly, we observed that the cells originally annotated as CD14 monocytes were  
153 divided into two groups (denoted as cluster 1 and cluster 2) in UMAP visualization (Fig. 3a and  
154 Supplementary Fig. 12), indicating the existence of possible new subpopulations of CD14  
155 monocytes. Differentially expression analysis showed that cluster 2 monocytes were characterized  
156 by the overexpression of *BCL11B*, *CD247*, and some interleukin-related genes (Fig. 3b and Fig. 3d).  
157 Several marker peaks were also identified corresponding to each monocyte subpopulations (Fig. 3c,  
158 Fig. 3e and Supplementary Figs. 13-14).

159 To ascertain whether the newly identified CD14 monocyte subpopulation is biologically  
160 meaningful, we built a CD14 monocyte signature using the top 50 marker genes of cluster 2  
161 monocytes. After this, we applied our CD14 monocyte signature to a time-series scRNA-seq dataset  
162 of human CD14 monocytes, where CD14 monocytes were stimulated by macrophage colony-  
163 stimulating factor (M-CSF) at days 0, 3, and 6 [24]. We observed that only 15% of monocytes were  
164 clustered together, and 25% of them had relatively higher CD14 monocyte signature scores at day  
165 0 (Fig. 3f, Student's *t*-test  $P = 2.84\text{e-}20$ ). This finding became more apparent at day 3, with  
166 approximately half of the monocytes had significantly higher signature scores (Fig. 3g, Student's *t*-  
167 test  $P = 4.1\text{e-}123$ ). With the continuous stimulation of M-CSF, almost all cells (86%) exhibited high  
168 signature scores, and the difference between the two clusters was less significant (Fig. 3h, Student's  
169 *t*-test  $P = 6\text{e-}08$ ). These observations can be reproduced using other replicates (Supplementary Figs.  
170 15-16). Besides, we also identified the new CD14 monocyte subpopulation in another CD14  
171 monocyte scRNA-seq dataset with higher signature scores (Supplementary Fig. 17).

172 Since monocytes are precursors to macrophages and can differentiate into dendritic cells, we  
173 hypothesized that the signature scores can serve as differentiation activities and the monocytes with

174 high signature scores are prone to differentiate into other cells. We thereby conducted GO functional  
175 enrichment analysis using the signature genes. We found that these genes were associated with  
176 mononuclear cell differentiation and its related processes, such as T cell differentiation [25] and  
177 immune response-activating signal transduction [26] (Fig. 3i and Supplementary Fig. 17). These  
178 results confirmed the ability of our gene signature to identify CD14 monocyte subtypes.

179 In conclusion, we demonstrated that the SCARP-derived low-dimensional representation of cells  
180 can refine cell types and reveal biologically meaningful cell subtypes.

181

### 182 **SCARP discovered key factors of melanoma progression by portraying co-accessibility 183 relationships of peaks**

184 Except for the low-dimensional cell embeddings, SCARP can also generate biologically meaningful  
185 peak embeddings. In this experiment, we want to investigate whether the SCARP-derived peak  
186 embeddings can uncover some potential regulatory relationships in melanoma disease progression.  
187 To achieve this, we applied SCARP on a SOX10 knockdown time series scATAC-seq dataset  
188 (SOX10KD) [2], which records changes in the accessibility of peaks at 0, 24, 48, and 72 hours after  
189 SOX10 knockdown in two patients (MM057 and MM087).

190 We first annotated two regions related to SOX10, i.e., the promoter (chr22:38380499-38380726)  
191 and 3'UTR (chr22:38364975-38365257) of it (Supplementary Fig. 18). Since the promoter initiates  
192 gene transcription while 3'UTR can decrease the expression of the corresponding gene [27] [28], it  
193 is reasonable to see that the accessibility of SOX10 promoter increases and 3'UTR decreases as  
194 knockdown time passes. By calculating the correlations between peaks' low-dimensional  
195 representations, we obtained the top 10 co-accessible regions with the highest correlation to each  
196 of the SOX10 promoter and 3'UTR. These loci have the more similar time-varying trend as the  
197 SOX10 promoter (3'UTR) than that found by using original data, indicating that SCARP can well  
198 characterize the co-accessibility of accessible segments (Fig. 4a and Supplementary Fig. 25a). To  
199 better interpret these co-accessible regions, we obtained the corresponding genes using the  
200 promoter regions (Supplementary Fig. 19; Methods) and selected the top 50 genes with the most  
201 co-accessible relationships with SOX10. By constructing a gene regulatory network (GRN) between  
202 these genes, we found abundant database-supported regulatory relationships (Fig. 4c). Notably,  
203 SCARP owned the largest number of curated edges and the second number of verified genes in the  
204 constructed GRNs compared to other candidate methods, such as cisTopic and scAND (Fig. 4b).

205 Among the genes in SCARP-derived GRN, numerous genes are strongly associated with  
206 melanoma disease. For example, the upregulations of *STMN1* and *ROR1* could contribute to tumor  
207 migration and proliferation [29][30][31]. *LTA4H* regulates the cell cycle process and can lead to skin  
208 carcinogenesis by its overexpression [32]. Notably, we uniquely identified that *HIF1A*, which is  
209 known for promoting cancer progression and causing therapy resistance [33][34], interacts with the  
210 key master regulator of melanocyte development and melanoma deterioration *MITF* [35][36] (Fig.  
211 4c). Numerous studies have reported that hypoxia is closely related to the invasiveness,  
212 angiogenesis, and therapy response in melanoma [33][37][38]. In studies of its pathogenic  
213 mechanism, *MITF* has been shown to stimulate the transcriptional activity of *HIF1A* to supply oxygen  
214 to cancer cells [39][40].

215 We further conducted the functional enrichment analysis using the genes in SCARP-derived  
216 GRN. As a result, we found many melanoma progression-related KEGG pathways, e.g.,  
217 melanogenesis [41][42], MAPK signaling pathways [43][44], and Adherens junction [45][46] (Fig. 4d  
218 and Supplementary Fig. 21). Notably, the HIF-1 signaling pathway was again exclusively identified  
219 by SCARP (Fig. 4d and Supplementary Figs. 23-25) [47][48]. In addition, GO enrichment categories  
220 contained several melanoma-related terms, such as muscle tissue development [49] and membrane  
221 raft [50] (Supplementary Fig. 22). Furthermore, we performed survival analysis using the genes in  
222 SCARP-derived GRN and found that many genes' abnormal expressions were associated with the  
223 poor prognoses, such as *TRIM69*, *MITF*, *STK17B* (Fig. 4e and Supplementary Fig. 26). In particular,  
224 the low expression of gene *STK17B*, which has been shown in the literature to play an important  
225 role in immune infiltration, was associated with skin cancer, making it a key biomarker for the  
226 diagnosis of melanoma disease [51].

227

## 228 **SCARP predicted reliable cis-regulatory interactions supported by external evidence**

229 We next explored the potential of the SCARP-derived peak embeddings and investigated their cis-  
230 regulatory relationships. We first defined the co-accessibility score of two peaks as the cosine  
231 similarity of their low-dimensional peak embeddings. This co-accessibility score can serve as the  
232 confidence of the corresponding cis-regulatory interactions. To assess the reliability of this score, we  
233 use external evidence to validate them, such as promoter-capture Hi-C (PCHi-C) [52][53] and ChIP-  
234 seq data [54]. Specifically, we selected those peaks belonging to promoter regions of certain genes  
235 and focused on the regulatory relationships between these peaks (i.e., promoters) with other peaks.

236 As PCHi-C data depicted whether there were physical interactions between the promoters (baits)  
237 and the entire genome (other ends), we can treat whether SCARP-derived cis-regulatory  
238 relationships were validated by PCHi-C data as a binary classification problem and measure its  
239 accuracy by calculating AUROC score (Methods). The result shows a significant difference between  
240 PCHi-C-validated and unvalidated SCARP-derived cis-regulatory scores (Fig. 5a, Student's *t*-test  
241  $P=1.31e-230$ ), and the receiver operating characteristic (ROC) curve demonstrates the superiority  
242 of SCARP over other methods (Fig. 5b).

243 Besides, we also used ChIP-seq data of human TF to validate SCARP-derived cis-regulatory  
244 relationships [54], which gives the experimentally verified interaction between proteins and DNA to  
245 help identify TF binding sites. To use evidence from both PCHi-C and ChIP-seq data, we selected  
246 those peaks that were both promoter regions and corresponded to a certain TF to see if the  
247 interaction with this region could be verified (Methods; Supplementary Fig. 27). As a demonstration  
248 example, we plotted the cis-regulatory interactions in the promoter region of gene ZNF384 (Fig. 5c  
249 and Fig. 5e), which is a human TF associated with B-cell acute lymphoblastic leukemia (B-ALL)  
250 disease [55][56]. Visually we see that many of the SCARP-derived interactions were validated by  
251 either Chip-seq or PCHi-C data. We further extracted the ZNF384 regulon network by mapping the  
252 SCARP-derived peaks to gene level, as shown in Fig. 5c, and those genes (yellow nodes) interacting  
253 with ZNF384 can be regarded as its target genes, since they were proven to have its binding sites  
254 (green edges). Some literature evidence can also be found to verify this result: *MLF2* is a key factor  
255 associated with leukemia disease [57], and *CHD4* is stated as a potential therapeutic target of acute  
256 myeloid leukemia [58][59]. Another example of the promoter region of gene *BCL3* shown in Fig. 5d  
257 and Fig. 5f exhibited similar results. *BCL3* is also related to chronic lymphocytic leukemia and B-cell  
258 malignancies [60][61], and the genes interacting with it were supported by external data and  
259 literature [62][63][64]. More examples can be found in Supplementary Fig. 27.

260

## 261 Discussion

262 Single-cell ATAC-seq data established genome-wide chromatin accessibility profiles with single-cell  
263 resolution that could complement single-cell transcriptome data, and together they provided new  
264 insights into transcriptional regulatory mechanisms. The accessible regions of scATAC-seq data,  
265 which can be seen as cell features, exhibited many different properties from genes of scRNA-seq  
266 data, such as higher sparsity, higher dimension, and near-binariization [1][6]. Thus, even both peaks

267 and genes can be used to characterize cellular heterogeneity, and many scRNA-seq-based  
268 algorithms are mathematically transferable to analyze scATAC-seq, we should not be limited to  
269 previous attempts and ideas for dealing with scRNA-seq data given the different nature of the two  
270 features; instead, we should develop new approaches that focus on the specific properties and  
271 information of the scATAC-seq data itself and elaborate on cellular heterogeneity from new  
272 perspectives.

273 Although the severe signal deficiencies and high sparsity of scATAC-seq data are well known,  
274 the importance of filling in missing data and denoising was not given enough attention in some  
275 commonly used scATAC methods [65]. In this study, we proposed a network diffusion-based  
276 computational approach to comprehensively analyze scATAC-seq data. Considering the high  
277 sparsity of the data, it is very appropriate and natural to use the diffusion method for aggregating the  
278 neighborhood information between two nodes. Our previously developed NR diffusion method [15]  
279 has achieved outstanding success in many other applications, such as network denoising [15] and  
280 link prediction [16], and in this study it also proved to have excellent performance in handling  
281 scATAC-seq data.

282 In addition, to mine and exploit the unique information of scATAC-seq data, we have borrowed  
283 ideas from many genomic studies and used the genetic map distance as an important prior factor  
284 in the analysis of epigenomic data. To further test the soundness of this hypothesis, we measured  
285 the co-accessibility likelihood between two peaks using the p-value of Fisher's exact test, and then  
286 investigated its relevance to genomic distance of peaks. The results shown in Supplementary Fig.  
287 28a confirm that peaks closer in the genome are more likely to be co-accessible. The peaks with  
288 network weights higher (or lower) than expected after diffusion show a certain pattern on the genome  
289 (Supplementary Fig. 28b-e), i.e., the posterior probabilities of co-accessibility in some regions are  
290 higher (or lower) than that predicted solely based on genomic distance. The enrichment of the  
291 different region types within the SCARP-derived low-dimensional peak embedding features also  
292 supported that there were clear differences in accessibility between the different region types, which  
293 were closely related to their regulatory functions (Supplementary Fig. 20). Furthermore, the results  
294 of our experiments supported that the use of prior edge weights between adjacent peaks facilitate  
295 SCARP to mine more valuable information from scATAC-seq data, and also helped to obtain better  
296 downstream analysis performance.

297 We have demonstrated through sufficient experiments that SCARP comprehensively promotes

298 superior analysis of scATAC-seq data compared to other state-of-the-art methods, such as scAND  
299 [6] and cisTopic [2]. Of all methods evaluated, SCARP had the best averaged clustering  
300 performances on benchmarking scATAC-seq data, and was robust to different peaks filter strategies  
301 and parameter selections. The running time of SCARP is also advantageous since we used some  
302 appropriate tricks to reduce the time cost. For example, to make those peaks on the same  
303 chromosome better learn from each other, we performed NR diffusion on subgraphs in parallel and  
304 spliced them back together.

305 Apart from the benchmarking dataset, we have designed three innovative case studies to  
306 explore what kind of analyses could be done with the scATAC-seq data to contribute to our further  
307 understanding of biological regulatory mechanisms. Specifically, we used SCARP-derived cell  
308 embeddings to identify a new CD14 monocyte subpopulation. Differential analysis between the two  
309 cell clusters obtained by SCARP identified a CD14 monocyte signature, which was later proved to  
310 be associated with mononuclear cell differentiation. We also applied SCARP on a melanoma-related  
311 scATAC-seq data to uncover some potential regulatory relationships in melanoma disease  
312 progression. Our analyses discovered several key factors of melanoma using SCARP-derived peak  
313 embeddings, such as genes *MITF*, *HIF1A*, and *STMN1*. Altogether, this showed how SCARP can  
314 be used to reveal disease mechanism by portraying co-accessibility relationships of peaks.  
315 Furthermore, the comparison of co-accessibility scores between SCARP-derived cis-regulatory  
316 interactions of peaks with external evidence, such as PCHi-C and CHIP-seq data, demonstrated the  
317 reliability of SCARP's peak low-dimensional embeddings.

318 In the future, we expect to explore more unique properties of the scATAC-seq data to better  
319 handle it, and consider doing joint analyses with other multi-omics data, such as data alignment and  
320 integration.

321

## 322 **Methods**

### 323 **SCARP workflow**

324 Given scATAC-seq data matrix  $A_{c \times p}$  with  $c$  cells and  $p$  peaks, we first binarized  $A$  to better cater  
325 to its biological interpretation of depicting accessibility. Specifically, we set all values greater than 1  
326 to 1, making  $A$  a Boolean matrix with entries from the Boolean domain  $\{0, 1\}$ , reflecting whether a  
327 cell is accessible on a certain peak. Then a symmetric square root graph normalization  $\tilde{A} =$

328  $D_L^{-\frac{1}{2}}AD_R^{-\frac{1}{2}}$  was applied to eliminate the library size differences between cells and peaks, where  $D_L$   
329 and  $D_R$  are diagonal degree matrixes with  $(D_L)_{ii} = \sum_j A_{ij}$  and  $(D_R)_{jj} = \sum_i A_{ij}$  [6].

330 After the pre-processing step, SCARP constructed a bipartite network with an adjacency matrix  
331  $B_{N \times N}$  to model the accessible relationship between the cells and peaks, where  $N = c + p$  is the  
332 number of nodes of the network, i.e., the peaks and cells are all treated as nodes of the network. To  
333 incorporate genetic location information for a better depiction of the co-accessibility of adjacent  
334 peaks, we computed the prior weights between adjacent peaks using the negative Haldane's genetic  
335 map function. That is to say, we can now represent the network as:

336 
$$B_{N \times N} = \begin{bmatrix} 0 & \tilde{A} \\ \tilde{A}^T & H \end{bmatrix}$$

337 where  $H$  is a sub-diagonal matrix of size  $p \times p$ , whose element represents the prior weight between  
338 the adjacent peaks in the same chromosome.

339 NR diffusion method was then applied to compute accessibility similarities between the cells and  
340 peaks, as well as predict possible cell-peak accessible relationships. Considering that those peaks  
341 on the same chromosome can better learn from each other, the NR diffusion process was calculated  
342 on the subnetwork of  $B$ , which is induced from nodes of cells and those peaks on the same  
343 chromosome, and then spliced back together. Specifically, the graph  $B$  was divided into a number  
344 of subgraphs  $B_{N_1 \times N_1}^{(1)}, \dots, B_{N_M \times N_M}^{(M)}$ , where  $N_k = c + p_k$  and  $p_k$  is the number of peaks in the  
345 subgraph  $B^{(k)}$ , and NR was performed on  $B^{(k)}$  to get a diffused matrix. This diffusion process can  
346 be performed in parallel, which reduces the time cost. Notably, when splitting peaks, we divided  
347 them according to whether they were on the same chromosome or not by default. However, given  
348 the sparsity of the scATAC-seq data, when the number of peaks is too small, we merged the adjacent  
349 chromosomes, and set the threshold parameter  $\gamma$  (default as 3000) to determine whether we need  
350 to merge the adjacent chromosomes.

351 Notice that after parallel diffusion, we got  $M$  dense matrixes,  $S^{(k)} = \begin{bmatrix} S_{c \times c}^{(k)(11)} & S_{c \times p_k}^{(k)(12)} \\ S_{p_k \times c}^{(k)(21)} & S_{p_k \times p_k}^{(k)(22)} \end{bmatrix}$ , each of  
352 which has a submatrix of size  $c \times c$ ,  $S_{c \times c}^{(k)(11)}$ , representing the cell similarity calculated based on the  
353 peaks of the current subgraph, and when splicing them back into matrix  $D_{N \times N}$ :

354

$$D_{N \times N} = \begin{bmatrix} D_{c \times c}^{(11)} & D_{c \times p}^{(12)} \\ D_{p \times c}^{(21)} & D_{p \times p}^{(22)} \end{bmatrix} = \begin{bmatrix} D_{c \times c}^{(11)} & S_{c \times p_1}^{(1)(12)} & \dots & S_{c \times p_M}^{(M)(12)} \\ S_{p_1 \times c}^{(1)(21)} & S_{p_1 \times p_1}^{(1)(22)} & & \\ \vdots & & \ddots & \\ S_{p_M \times c}^{(M)(21)} & & & S_{p_M \times p_M}^{(M)(22)} \end{bmatrix}$$

355  $D_{c \times c}^{(11)}$  was the average over  $M$  matrixes  $S_{c \times c}^{(k)(11)}$  (Fig. 1).

356 Finally, to get cell and peak embeddings, the dimensional reduction method, such as UCPCA,  
357 was performed on the matrix  $[D_{c \times c}^{(11)}, D_{c \times p}^{(12)}]$  and  $D_{p \times c}^{(21)}$  separately. The number of kept components  
358 was determined based on the proportion of variance that can be explained, as detailed below.

359

### 360 **Prior weights between adjacent peaks**

361 The prior edge weights between the adjacent peaks in this study were calculated using the negative  
362 Haldane's genetic map function, motivated by the traditional use of gene map function [10] to  
363 estimate the relation connecting recombination fractions  $r$  and genomic distances  $d$ :

$$364 \quad r = e^{-d} + e^{-d} \frac{d^3}{3!} + \dots = \frac{1}{2}(1 - e^{-2d})$$

365 We assumed that the probability of chromosome segments being co-accessible should be a  
366 decreasing function of its genomic distance, as opposed to the probability of recombination [52][54].  
367 Specifically, the distance between two peaks  $\text{peak}_1 = [s_1, e_1]$  and  $\text{peak}_2 = [s_2, e_2]$  (any two peaks  
368 do not intersect with each other) was  $d(\text{peak}_1, \text{peak}_2) = \max(s_2 - e_1, s_1 - e_2)$ , where  $s_i$  stands for  
369 the start position and  $e_i$  stands for the end position on the genome, and the prior weight of the two  
370 peaks was computed as:

$$371 \quad w(\text{peak}_1, \text{peak}_2) = e^{-\frac{2d(\text{peak}_1, \text{peak}_2)}{\beta}}$$

372 here we set a parameter  $\beta$  (default as 5000) to control the extent to which prior edge weight decays  
373 as peak distance increases.

374

### 375 **NR diffusion method**

376 We used our previously proposed method, Network Refinement (NR) [15], to cope with the sparsity  
377 of scATAC-seq data. Specifically, NR takes the adjacency matrix of a network as input, uses a  
378 diffusion process defined by a random walk on the graph to enhance the self-organization properties  
379 of complex networks, and outputs a network with adjusted edge weights. The matrix obtained by NR  
380 reflects both the similarity of cells and peaks, and gives a confidence score for the possible cell-peak

381 accessible relationship.

382 The graph operator  $F_m$  of NR is defined based on three operators:  $f_m$ ,  $g$ , and  $h$ . The operator  
383  $f_m$  transforms a transition matrix  $P$  to another by adding the probability of all paths of different  
384 lengths joining two nodes, with a smaller weight coefficient  $1/m^k$  for a longer path of length  $k$ :

385 
$$f_m: \mathcal{P} \rightarrow \mathcal{P}$$

386 
$$f_m(P) = \frac{1}{\sum_{k=1}^{\infty} (1/m^k)} \sum_{k=1}^{\infty} \frac{P^k}{m^k} = (m-1)P(mI - P)^{-1} \quad (1)$$

387 where  $\sum_k (1/m^k)$  is a normalization factor, and  $m > 1$  to ensure that the series converges when  $k$   
388 approaches infinity [15]. The parameter  $m$  controls the diffusion intensity. The smaller the  $m$ , the  
389 higher the degree of diffusion (default as  $m = 1.5$ ). Notably, the operator  $f_m$  has described a signal  
390 enhancement process by accumulating the transition probabilities of all lengths of paths between  
391 the two vertexes, which can also be regarded as a process of signal diffusion.

392 Then we defined two auxiliary operators  $g$  and  $h$  to help us map the diffusion process of the  
393 random walk defined by the operator  $f_m$  to the diffusion process on the graph [15]. Denote  $D_W$  as  
394 the diagonal degree matrix of the input network  $W$ , then  $g(W)$  defines a random walk on the graph  
395 whose (weighted) adjacency matrix is  $W$ :

396 
$$g: \mathcal{W} \rightarrow \mathcal{P}$$

397 
$$g(W) = D_W^{-1}W \quad (2)$$

398 The operator  $h$  has the opposite effect of  $g$ , which recovers the underlying graph of the random  
399 walk defined by the transition matrix  $P$ :

400 
$$h: \mathcal{P} \rightarrow \mathcal{W}$$

401 
$$h(P) = \alpha \cdot \text{diag}(\pi(P))P \quad (3)$$

402 where  $\pi(P) = (\pi_1, \dots, \pi_n)$  is the stationary distribution of transition matrix  $P$  such that  $\pi P = \pi$ .  
403 When  $P$  defines an irreducible and aperiodic random walk,  $\pi(P)$  exists and can be guaranteed to  
404 be unique (that is generally true in practical applications).  $\text{diag}(x)$  means the diagonal matrix whose  
405 diagonal element is the vector  $x$ , and  $\alpha$  is a constant which controls the sum of weight matrix  $h(P)$ .  
406 The operator  $h$  multiplies the transition probability  $P_{ij}$  by the stationary distribution of node  $i$   
407 which reflects the degree information of the graph.

408 Formally, to map the diffusion process of random walk on the graph defined by  $f_m$  onto the  
409 diffusion process on the graph, NR wraps the operator  $f_m$  in operators  $g$  and  $h$  to get the  
410 composite operator  $F_m$ :

411

$$F_m: \mathcal{W} \rightarrow \mathcal{W}$$

412

$$F_m(W) = h(f_m(g(W))) \quad (4)$$

413

where operators  $g$  and  $h$  realize the conversion between the graph and random walk on the graph, and  $f_m$  realize the diffusion of random walk on the graph.

415

#### 416 The difference between NR and Katz Index

417

We have proven in previous work that the graph operator of NR can be further written as [15]:

418

$$F_m(W) = \frac{1}{\sum_{k=1}^{\infty} (1/m^k)} \sum_{k=1}^{\infty} \frac{D_W(D_W^{-1}W)^k}{m^k} = \frac{m-1}{m} D_W(D_W^{-1}W) + \frac{m-1}{m^2} D_W(D_W^{-1}W)^2 + \dots$$

419

And the Katz index was calculated by [9]:

420

$$Katz(W) = \sum_{k=1}^{\infty} \left(\frac{\beta}{r}\right)^k W^k = \frac{\beta}{r} W + \frac{\beta^2}{r^2} W^2 + \dots$$

421

The difference between the two diffusion methods can be understood from two perspectives.

422

First, as we described in our previous work [15], NR is a degree-normalized version of the Katz index, which makes the path with higher intermediate node degrees less important because their information is dispersed through more adjacent edges. Second, we can compute NR by computing the Katz index with a few more processing steps. When  $W$  is an arbitrary matrix, we cannot construct a direct connection between the two methods. But when we take a pre-processing step

427

$W = D_W^{-\frac{1}{2}} W D_W^{-\frac{1}{2}}$ , the spectral radius of the input matrix can be guaranteed to be one, which means the  $r$  in the formula of Katz index is one and thus have:

429

$$Katz(D_W^{-\frac{1}{2}} W D_W^{-\frac{1}{2}}) = \sum_{k=1}^{\infty} \beta^k \left(D_W^{-\frac{1}{2}} W D_W^{-\frac{1}{2}}\right)^k = \beta D_W^{-\frac{1}{2}} W D_W^{-\frac{1}{2}} + \beta^2 \left(D_W^{-\frac{1}{2}} W D_W^{-\frac{1}{2}}\right)^2 + \beta^3 \left(D_W^{-\frac{1}{2}} W D_W^{-\frac{1}{2}}\right)^3 \dots$$

430

Note that the formula of NR can be further expressed as:

431

$$F_m(W) = \frac{m-1}{m} W + \frac{m-1}{m^2} W D_W^{-1} W + \frac{m-1}{m^3} W D_W^{-1} W D_W^{-1} W + \dots$$

432

It can be seen that when  $\beta = \frac{1}{m}$ , we have:

434

$$D_W^{-\frac{1}{2}} F_m(W) D_W^{-\frac{1}{2}} = (m-1) * \left( Katz\left(D_W^{-\frac{1}{2}} W D_W^{-\frac{1}{2}}\right) \right)$$

435

#### Dimensionality reduction

436

We used the Uncentred Principal Component Analysis (UCPCA) [66] to obtain the low-dimensional

437 embeddings of the corresponding cells and peaks in this study. The number of kept components  
438 was determined based on the proportion of variance that can be explained [21]. Specifically, we kept  
439 the first 50 (if peaks number less than 50,000) or 100 (if peaks number more than 50,000) principal  
440 components (PCs) and plotted the standard deviation (SD) of each PC. The kept number of PCs  
441 was determined based on the SD variance of PCs that were not selected, i.e., if the SD of the  
442 remaining PCs does not change much, we did not keep them. We give in Supplementary Fig. 29 the  
443 SD plots and the number of PCs retained on all scATAC-seq datasets in this study. Finally, we  
444 applied L2-normalization to each of the low-dimensional representations, which was a classical  
445 processing step [6][21][22] after dimensionality reduction.

446

#### 447 **Clustering and visualization**

448 We used the ‘scanpy.tl.louvain’ function from the Scanpy [67] python package to implement the  
449 Louvain algorithm for clustering cells, and the ‘scanpy.pl.tsne’ and ‘scanpy.pl.umap’ functions of  
450 Scanpy python package to visualize cells in two-dimensional coordinates, using the t-SNE basis and  
451 UMAP basis.

452

#### 453 **Evaluation metrics**

454 We used the Normalized Mutual Information (NMI) and the Adjusted Rand Index (ARI) to evaluate  
455 the clustering performance. Assuming that there are two partitions  $X = \{X_1, \dots, X_r\}$  and  $Y = \{Y_1, \dots, Y_s\}$   
456 for  $N$  vertexes, ARI is computed as follows:

$$457 \quad ARI = \frac{\sum_{ij} \binom{n_{ij}}{2} - \frac{\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}}{\binom{n}{2}}}{\frac{1}{2} \left( \sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right) - \frac{\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}}{\binom{n}{2}}}$$

458 where  $n_{ij}$  is the number of vertexes in partitions  $X_i$  and  $Y_j$ , and  $a_i$  is the number of vertexes in  
459 partition  $X_i$ ,  $b_j$  is the number of vertexes in partition  $Y_j$ . ARI takes values from -1 to 1, and a higher  
460 value of ARI indicates better performance. Besides, NMI is computed as:

$$461 \quad NMI(X, Y) = \frac{MI(X, Y)}{\sqrt{H(X)H(Y)}}$$

462 where the mutual information of two partitions and the entropy of a partition are computed as:

$$463 \quad MI(X, Y) = \sum_{ij} \frac{n_{ij}}{N} \log \frac{n_{ij}/N}{a_i/N \times b_j/N}$$

464

$$H(X) = \sum_i \frac{a_i}{N} \log \frac{a_i}{N}, H(Y) = \sum_j \frac{b_j}{N} \log \frac{b_j}{N}$$

465 NMI takes values from 0 to 1, and a higher value of NMI indicates better performance.

466 We also used the silhouette coefficient to evaluate the distance between clusters, given a  
467 partition for cells. Specifically, for cell  $x$ :

468

$$\text{silhouette}(x) = \frac{b(x) - a(x)}{\max\{a(x), b(x)\}}$$

469 where  $a(x)$  is the average distance between  $x$  and other cells in the same cluster, and  $b(x)$  is  
470 the average distance between  $x$  and cells in the closest different cluster [65]. The distance is  
471 measured by their Euclidean distance on a two-dimensional visualization (such as UMAP and t-  
472 SNE). The silhouette coefficient of the entire partition was the average over all cells. Silhouette  
473 coefficient takes values from -1 to 1, and a higher value indicates better performance.

474

#### 475 Peaks annotation

476 Peaks in the scATAC-seq datasets were annotated using the R package ‘ChIPseeker’ [68]. In this  
477 study, the upstream and downstream regions of the gene Transcription Start Site (TSS) were defined  
478 as  $\text{TSS} \pm 3,000$  base pairs.

479

#### 480 Differential genes in the SOX10KD datasets

481 After annotating the region type of peaks, we screened 14,266 peaks (promoters of genes) in total,  
482 which were referred as background genes. So, when we talk about background genes, we are  
483 actually talking about the peaks annotated as their promoter regions. The top 50 genes with the  
484 highest correlation to SOX10 were selected to validate the regulatory relationships with it, and we  
485 refer to them as differential genes since their accessibility changes over time. Specifically, we  
486 calculated the correlations between the SOX10 promoter (3'UTR) and all other background genes  
487 using their low-dimensional features, and then selected the highest 50 from them. That’s to say,  
488 those differential genes may have the same (indicating positive regulatory relationship) or opposite  
489 (indicating negative regulatory relationship) co-accessibility with SOX10.

490

#### 491 Validation of cis-regulatory between genes using Cytoscape

492 The application Reactome [69] in the software Cytoscape [70] was used to validate the gene  
493 interactions. Specifically, we input the differential genes set, and then use the database to validate

494 direct or indirect (by linker genes) regulatory relationships between them.

495

## 496 **Functional enrichment analysis**

497 In the SOX10KD experiment, given a gene set, we used the R package ‘clusterProfiler’ [71][72] for  
498 gene functional annotation, including GO enrichment analysis (molecular function, cellular  
499 components, and biological process) and KEGG pathway enrichment analysis. Specifically, we  
500 treated the genes in the GRN obtained by Cytoscape (including those differential genes and linker  
501 genes) as foreground genes set, the background genes (14,266 genes) and linker genes as the  
502 background genes set, using the ‘enrichGO’ and ‘enrichKEGG’ functions to perform gene enrichment  
503 analysis.

504

## 505 **Survival analysis**

506 Survival analysis was performed with R package ‘survival’, and the external gene expression data  
507 as well as clinical data were downloaded using the online tool Xena [73].

508

## 509 **Differential expression (accessibility) analysis in PBMC dataset**

510 We used scanpy.tl.rank\_genes\_groups function with default settings from python package ‘Scanpy’  
511 [67] for the standard differential expression analysis as well as differential accessibility analysis of  
512 cell subgroups in PBMC dataset.

513

## 514 **Gene set variation analysis (GSVA)**

515 We built our SCARP gene signatures using the top 50 marker genes for cell subgroup 2, as shown  
516 in Fig. 3b. We used the R package ‘GSVA’ to implement the Gene set variation analysis (GSVA) with  
517 default settings for the calculation of the SCARP signature score.

518

## 519 **Validation of cis-regulatory using PCHi-C**

520 PCHi-C data depicted physical interactions between the promoters (baits) and the entire genome  
521 (the other end). Since PCHi-C data were cell type-specific, we only selected those cells annotated  
522 as CD4 naive cells in the 10X Multiome data, so that there were overlapping cell types for both data  
523 types, and used SCARP to get peak embeddings. After selecting those peaks corresponding to  
524 promoters of specific genes, we got promoter-peak pairs with confidence scores calculated by cos

525 similarity of their low-dimensional embeddings. A promoter-peak pair was considered validated by  
526 PCHi-C evidence if (1) the overlap area between promoter (from 10X Multiome) and bait (from PCHi-  
527 C) reaches half the length of one of them, or the distance between the centers of the two is less than  
528 5kb, and (2) the overlap area between peak (from 10X Multiome) and other end (from PCHi-C)  
529 reaches half the length of one of them, or the distance between the centers of the two is less than  
530 5kb [52]. We can then treat whether promoter-peak pairs were validated by PCHi-C as a binary  
531 classification problem and measure its accuracy by calculating AUROC score.

532 We used t-test to verify whether there was a significant difference in SCARP-derived cis-  
533 regulatory scores between PCHi-C validated and unvalidated promoter-peak pairs.

534

### 535 **Validation of cis-regulatory using ChIP-Seq**

536 TF ChIP-seq data downloaded from ENCODE [74] was also used as another external evidence to  
537 validate cis-regulatory relationships. We took the intersections of all known human TF, promoters of  
538 PCHi-C, and the peaks annotated as promoter regions in scATAC-seq data, so that we can verify  
539 the third one with the other two. The interactions between chromosome segments were plotted using  
540 the R package ‘Cicero’ [54]. The human GRCh38 reference genome and annotation files were  
541 downloaded from ENSEMBL [75].

542

### 543 **Data availability**

544 All datasets analyzed in this study were publicly available. The detailed information for the scATAC-  
545 seq and the multi-omics datasets used in this study can be found in Supplementary Table 1. The  
546 gene expression and phenotype data of the Skin Cutaneous Melanoma (SKCM) project were  
547 downloaded from TCGA. For promoter-capture Hi-C data, it can be obtained from the original  
548 publication [52][53]. For TF ChIP-seq data, it can be downloaded from the ENCODE project  
549 (<https://www.encodeproject.org/>) and further processed using GLUE’s script (<https://github.com/gao-lab/GLUE>).  
550

551

### 552 **Code availability**

553 SCARP is available at <https://github.com/Wu-Lab/SCARP>.

554

### 555 **Reference**

- 556 [1] Xiong, L. et al. SCALE method for single-cell ATAC-seq analysis via latent feature extraction. *Nat. Commun.* **10**, 1–10 (2019).
- 557
- 558 [2] Bravo González-Blas, C. et al. cisTopic: cis-regulatory topic modeling on single-cell ATAC-seq
- 559 data. *Nat. Methods* **16**, 397–400 (2019).
- 560 [3] Buenrostro, J. D. et al. Single-cell chromatin accessibility reveals principles of regulatory variation.
- 561 *Nature* **523**, 486–490 (2015).
- 562 [4] Cusanovich, D. A. et al. Multiplex single-cell profiling of chromatin accessibility by combinatorial
- 563 cellular indexing. *Science* **348**, 910–914 (2015).
- 564 [5] Cusanovich, D. A. et al. The cis-regulatory dynamics of embryonic development at single-cell
- 565 resolution. *Nature* **555**, 538–542 (2018).
- 566 [6] Dong, K. & Zhang, S. Network diffusion for scalable embedding of massive single-cell ATAC-seq
- 567 data. *Sci. Bull.* **66**, 2271–2276 (2021).
- 568 [7] Liu, Y., Zhang, J., Wang, S., Zeng, X. & Zhang, W. Are dropout imputation methods for scRNA-
- 569 seq effective for scATAC-seq data? *Brief. Bioinform.* **23**, 1–12 (2022).
- 570 [8] Schep, A. N., Wu, B., Buenrostro, J. D. & Greenleaf, W. J. ChromVAR: Inferring transcription-
- 571 factor-associated accessibility from single-cell epigenomic data. *Nat. Methods* **14**, 975–978
- 572 (2017).
- 573 [9] Katz, L. A new status index derived from sociometric analysis. *Psychometrika* **18**, 39–43 (1953).
- 574 [10] Speed, T. P. Genetic Map Functions. *Encycl. Biostat.* (2005).
- 575 [11] Altshuler, D., Daly, M. J. & Lander, E. S. Genetic mapping in human disease. *Science* **322**, 881–
- 576 888 (2008).
- 577 [12] Robinson, M. A. Linkage Disequilibrium. *Encycl. Immunol.* 1586–1588 (1998).
- 578 [13] Teare, M. D. & Barrett, J. H. Genetic linkage studies. *Lancet* **366**, 1036–1044 (2005).
- 579 [14] Dekker, J., Marti-Renom, M. & Mirny, L. Exploring the three-dimensional organization of
- 580 genomes: interpreting chromatin interaction data. *Nat. Rev. Genet.* **14**, 390–403 (2013).
- 581 [15] Yu, J., Leng, J. & Wu, L.-Y. Network Refinement: A unified framework for enhancing signal or
- 582 removing noise of networks. Preprint at arXiv <https://arxiv.org/abs/2109.09119> (2021).
- 583 [16] Yu, J. & Wu, L. Y. Multiple Order Local Information model for link prediction in complex networks.
- 584 *Phys. A Stat. Mech. its Appl.* **600**, 127522 (2022).
- 585 [17] Eraslan, G., Simon, L. M., Mircea, M., Mueller, N. S. & Theis, F. J. Single-cell RNA-seq denoising
- 586 using a deep count autoencoder. *Nat. Commun.* **10**, 1–14 (2019).

- 587 [18] van Dijk, D. et al. Recovering Gene Interactions from Single-Cell Data Using Data Diffusion.  
588 *Cell* **174**, 716-729.e27 (2018).
- 589 [19] Buenrostro, J. D. et al. Integrated Single-Cell Analysis Maps the Continuous Regulatory  
590 Landscape of Human Hematopoietic Differentiation. *Cell* **173**, 1535-1548.e16 (2018).
- 591 [20] Chen, S., Lake, B. B. & Zhang, K. High-throughput sequencing of the transcriptome and  
592 chromatin accessibility in the same cell. *Nat. Biotechnol.* **37**, 1452–1457 (2019).
- 593 [21] Stuart, T. et al. Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902.e21 (2019).
- 594 [22] Stuart, T., Srivastava, A., Madad, S., Lareau, C. A. & Satija, R. Single-cell chromatin state  
595 analysis with Signac. *Nat. Methods* **18**, 1333–1341 (2021).
- 596 [23] PBMC from a healthy donor, single cell multiome ATAC gene expression demonstration data by  
597 Cell Ranger ARC 1.0.0. 10X Genomics [https://support.10xgenomics.com/single-cell-multiome-  
598 atac-gex/datasets/1.0.0/pbmc\\_granulocyte\\_sorted\\_10k](https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/1.0.0/pbmc_granulocyte_sorted_10k) (2020).
- 599 [24] Hie, B., Bryson, B. & Berger, B. Efficient integration of heterogeneous single-cell transcriptomes  
600 using Scanorama. *Nat. Biotechnol.* **37**, 685–691 (2019).
- 601 [25] Zimmermann, H. W. et al. Bidirectional transendothelial migration of monocytes across hepatic  
602 sinusoidal endothelium shapes monocyte differentiation and regulates the balance between  
603 immunity and tolerance in liver. *Hepatology* **63**, 233–246 (2016).
- 604 [26] Wierda, R. J. et al. A role for KMT1c in monocyte to dendritic cell differentiation: Epigenetic  
605 regulation of monocyte differentiation. *Hum. Immunol.* **76**, 431–437 (2015).
- 606 [27] <https://www.genome.gov/genetics-glossary/Promoter>.
- 607 [28] [https://en.wikipedia.org/wiki/Three\\_prime\\_untranslated\\_region](https://en.wikipedia.org/wiki/Three_prime_untranslated_region).
- 608 [29] Chen, J. et al. Stathmin 1 is a potential novel oncogene in melanoma. *Oncogene* **32**, 1330–  
609 1337 (2012).
- 610 [30] Fernández, N. B. et al. ROR1 contributes to melanoma cell growth and migration by regulating  
611 N-cadherin expression via the PI3K/Akt pathway. *Mol. Carcinog.* **55**, 1772–1785 (2016).
- 612 [31] Hojjat-Farsangi, M. et al. Inhibition of the Receptor Tyrosine Kinase ROR1 by Anti-ROR1  
613 Monoclonal Antibodies and siRNA Induced Apoptosis of Melanoma Cells. *PLoS One* **8**, e61167  
614 (2013).
- 615 [32] Oi, N. et al. LTA4H regulates cell cycle and skin carcinogenesis. *Carcinogenesis* **38**, 728–737  
616 (2017).
- 617 [33] D'Aguanno, S., Mallone, F., Marenco, M., Del Bufalo, D. & Moramarco, A. Hypoxia-dependent

- 618 drivers of melanoma progression. *J. Exp. Clin. Cancer Res.* **40**, 159 (2021).
- 619 [34] Jing, X. et al. Role of hypoxia in cancer therapy by regulating the tumor microenvironment. *Mol.*  
620 *Cancer* **18**, 1–15 (2019).
- 621 [35] Levy, C., Khaled, M. & Fisher, D. E. MITF: master regulator of melanocyte development and  
622 melanoma oncogene. *Trends Mol. Med.* **12**, 406–414 (2006).
- 623 [36] Hartman, M. L. & Czyz, M. MITF in melanoma: Mechanisms behind its expression and activity.  
624 *Cell. Mol. Life Sci.* **72**, 1249–1260 (2015).
- 625 [37] Hwang, H. W. et al. Distinct microRNA expression signatures are associated with melanoma  
626 subtypes and are regulated by HIF1A. *Pigment Cell Melanoma Res.* **27**, 777–787 (2014).
- 627 [38] Lakhter, A. J., Lahm, T., Broxmeyer, H. E. & Naidu, S. R. Golgi Associated HIF1a Serves as a  
628 Reserve in Melanoma Cells. *J. Cell. Biochem.* **117**, 853–859 (2016).
- 629 [39] Feige, E. et al. Hypoxia-induced transcriptional repression of the melanoma-associated  
630 oncogene MITF. *Proc. Natl. Acad. Sci. U. S. A.* **108**, E924–E933 (2011).
- 631 [40] Buscà, R. et al. Hypoxia-inducible factor 1 $\alpha$  is a new target of microphthalmia-associated  
632 transcription factor (MITF) in melanoma cells. *J. Cell Biol.* **170**, 49–59 (2005).
- 633 [41] Slominski, A. et al. The role of melanogenesis in regulation of melanoma behavior:  
634 Melanogenesis leads to stimulation of HIF-1 $\alpha$  expression and HIF-dependent attendant  
635 pathways. *Arch. Biochem. Biophys.* **563**, 79–93 (2014).
- 636 [42] Slominski, A., Zbytek, B. & Slominski, R. Inhibitors of melanogenesis increase toxicity of  
637 cyclophosphamide and lymphocytes against melanoma cells. *Int. J. Cancer* **124**, 1470–1477  
638 (2009).
- 639 [43] Sullivan, R. J. & Flaherty, K. MAP kinase signaling and inhibition in melanoma. *Oncogene* **32**,  
640 2373–2379 (2012).
- 641 [44] Fecher, L. A., Amaravadi, R. K. & Flaherty, K. T. The MAPK pathway in melanoma. *Curr. Opin.*  
642 *Oncol.* **20**, 183–189 (2008).
- 643 [45] Lee, D. J. et al. Peroxiredoxin-2 represses melanoma metastasis by increasing E-cadherin/ $\beta$ -  
644 catenin complexes in adherens junctions. *Cancer Res.* **73**, 4744–4757 (2013).
- 645 [46] Korla, P. K. et al. Somatic mutational landscapes of adherens junctions and their functional  
646 consequences in cutaneous melanoma development. *Theranostics* **10**, 12026–12043 (2020).
- 647 [47] Slominski, A. et al. The role of melanogenesis in regulation of melanoma behavior:  
648 Melanogenesis leads to stimulation of HIF-1 $\alpha$  expression and HIF-dependent attendant

- 649 pathways. *Arch. Biochem. Biophys.* **563**, 79–93 (2014).
- 650 [48] Malekan, M., Ebrahimzadeh, M. A. & Sheida, F. The role of Hypoxia-Inducible Factor-1alpha  
651 and its signaling in melanoma. *Biomed. Pharmacother.* **141**, 111873 (2021).
- 652 [49] Moss, A. L. H. & Rees, M. J. W. Metastatic malignant melanoma in muscle. *Br. J. Plast. Surg.*  
653 **37**, 250–252 (1984).
- 654 [50] Baruthio, F., Quadroni, M., Rüegg, C. & Mariotti, A. Proteomic analysis of membrane rafts of  
655 melanoma cells identifies protein patterns characteristic of the tumor progression stage.  
656 *Proteomics* **8**, 4733–4747 (2008).
- 657 [51] Shi, X. et al. Prognostic and immune-related value of STK17B in skin cutaneous melanoma.  
658 *PLoS One* **17**, 1–21 (2022).
- 659 [52] Cao, Z. J. & Gao, G. Multi-omics single-cell data integration and regulatory inference with graph-  
660 linked embedding. *Nat. Biotechnol.* **40**, 1458–1466 (2022).
- 661 [53] Javierre, B. M. et al. Lineage-Specific Genome Architecture Links Enhancers and Non-coding  
662 Disease Variants to Target Gene Promoters. *Cell* **167**, 1369–1384.e19 (2016).
- 663 [54] Pliner, H. A. et al. Cicero Predicts cis-Regulatory DNA Interactions from Single-Cell Chromatin  
664 Accessibility Data. *Mol. Cell* **71**, 858–871.e8 (2018).
- 665 [55] McClure, B. J. et al. Pre-B acute lymphoblastic leukaemia recurrent fusion, EP300-ZNF384, is  
666 associated with a distinct gene expression. *Br. J. Cancer* **118**, 1000–1004 (2018).
- 667 [56] Zhang, X. Y. et al. MRD-Negative Remission Induced in EP300-ZNF384 Positive B-ALL Patients  
668 by Tandem CD19/CD22 CAR T-Cell Therapy Bridging to Allogeneic Stem Cell Transplantation.  
669 *Onco. Targets. Ther.* **14**, 5197–5204 (2021).
- 670 [57] Kuefer, M. U. et al. cDNA Cloning, Tissue Distribution, and Chromosomal Localization of  
671 Myelodysplasia/Myeloid Leukemia Factor 2 (MLF2). *Genomics* **35**, 392–396 (1996)
- 672 [58] Heshmati, Y. et al. The chromatin-remodeling factor CHD4 is required for maintenance of  
673 childhood acute myeloid leukemia. *Haematologica* **103**, 1169–1181 (2018).
- 674 [59] Heshmati, Y. et al. Identification of CHD4 As a Potential Therapeutic Target of Acute Myeloid  
675 Leukemia. *Blood* **128**, 1648 (2016).
- 676 [60] Yabumoto, K. et al. Involvement of the BCL3 gene in two patients with chronic lymphocytic  
677 leukemia. *Int. J. Hematol.* **59**, 211–218 (1994).
- 678 [61] McKeithan, T. W. et al. BCL3 Rearrangements and t(14;19) in Chronic Lymphocytic Leukemia  
679 and Other B-Cell Malignancies: A Molecular and Cytogenetic Study. *Genes Chromosom. Cancer*

- 680           **20**, 64–72 (1997).
- 681       [62] Weinberg, J. B. et al. Apolipoprotein E (APOE) Genotype as a Determinant of Survival in Women  
682           with Chronic Lymphocytic Leukemia. *Blood* **110**, 3081–3081 (2007).
- 683       [63] Weinberg, J. B. et al. Apolipoprotein E genotype as a determinant of survival in chronic  
684           lymphocytic leukemia. *Leukemia* **22**, 2184–2192 (2008).
- 685       [64] Chun, E. M. et al. Expression of the apolipoprotein C-II gene during myelomonocytic  
686           differentiation of human leukemic cells. *J. Leukoc. Biol.* **69**, 645–650 (2001).
- 687       [65] Li, Z. et al. Chromatin-accessibility estimation from single-cell ATAC-seq data with scOpen. *Nat.*  
688           *Commun.* **12**, 6386 (2021).
- 689       [66] Cadima, J. & Jolliffe, I. On relationships between uncentred and column-centred principal  
690           component analysis. *Pakistan J. Stat.* **25**, 473–503 (2009).
- 691       [67] Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: Large-scale single-cell gene expression data  
692           analysis. *Genome Biol.* **19**, 1–5 (2018).
- 693       [68] Yu, G., Wang, L. G. & He, Q. Y. ChIPseeker: an R/Bioconductor package for ChIP peak  
694           annotation, comparison and visualization. *Bioinformatics* **31**, 2382–2383 (2015).
- 695       [69] Croft, D. et al. Reactome: a database of reactions, pathways and biological processes. *Nucleic*  
696           *Acids Res.* **39**, D691–D697 (2011).
- 697       [70] Shannon, P. et al. Cytoscape: A Software Environment for Integrated Models of Biomolecular  
698           Interaction Networks. *Genome Res.* **13**, 2498–2504 (2003).
- 699       [71] Yu, G., Wang, L. G., Han, Y. & He, Q. Y. ClusterProfiler: An R package for comparing biological  
700           themes among gene clusters. *Omi. A J. Integr. Biol.* **16**, 284–287 (2012).
- 701       [72] Wu, T. et al. clusterProfiler 4.0: A universal enrichment tool for interpreting omics data. *Innov.* **2**,  
702           100141 (2021).
- 703       [73] <https://xena.ucsc.edu/>.
- 704       [74] Davis, C. A. et al. The Encyclopedia of DNA elements (ENCODE): data portal update. *Nucleic*  
705           *Acids Res.* **46**, D794–D801 (2018).
- 706       [75] <https://asia.ensembl.org/index.html>.
- 707
- 708       **Acknowledgements**
- 709       This work has been supported by the National Key Research and Development Program of China  
710           (No. 2020YFA0712402) and the National Natural Science Foundation of China (No. 12231018).

711

712 **Author contributions**

713 **Jiating Yu:** Conceptualization, Methodology, Software, Investigation, Formal analysis, Writing –  
714 original draft. **Duanchen Sun:** Analysis, Writing – review & editing. **Zhichao Hou:** Writing – review  
715 & editing. **Ling-Yun Wu:** Conceptualization, Methodology, Supervision, Writing – review & editing,  
716 Funding acquisition.

717

718 **Competing interests**

719 The authors declare no competing interests.

720

721 **Additional information**

722 Supplementary materials are available online.

723

724 **Abbreviations**

725 PBMC: Peripheral blood mononuclear cell

726 HSC: Hematopoietic Stem Cells

727 MPP: Multipotent Progenitor

728 CMP: Common Myeloid Progenitor

729 LMPP: Lymphoid-primed Multipotent Progenitors

730 MEP: Megakaryocyte-erythroid Progenitor

731 GMP: Granulocyte-monocyte Progenitor

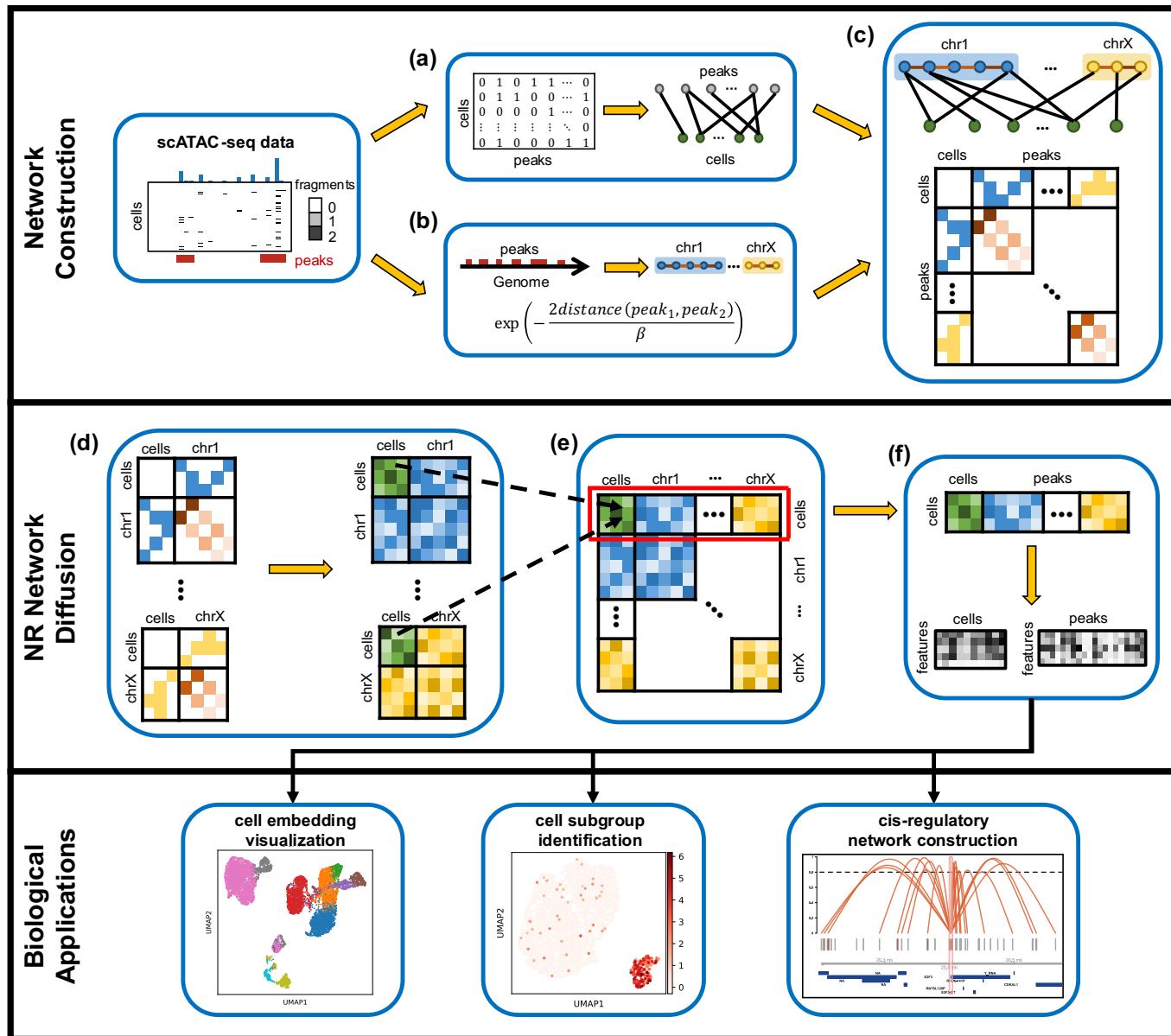
732 pDC: Plasmacytoid Dendritic Cells

733 CLP: Common Lymphoid Progenitor

734 Mono: Monocytes

735

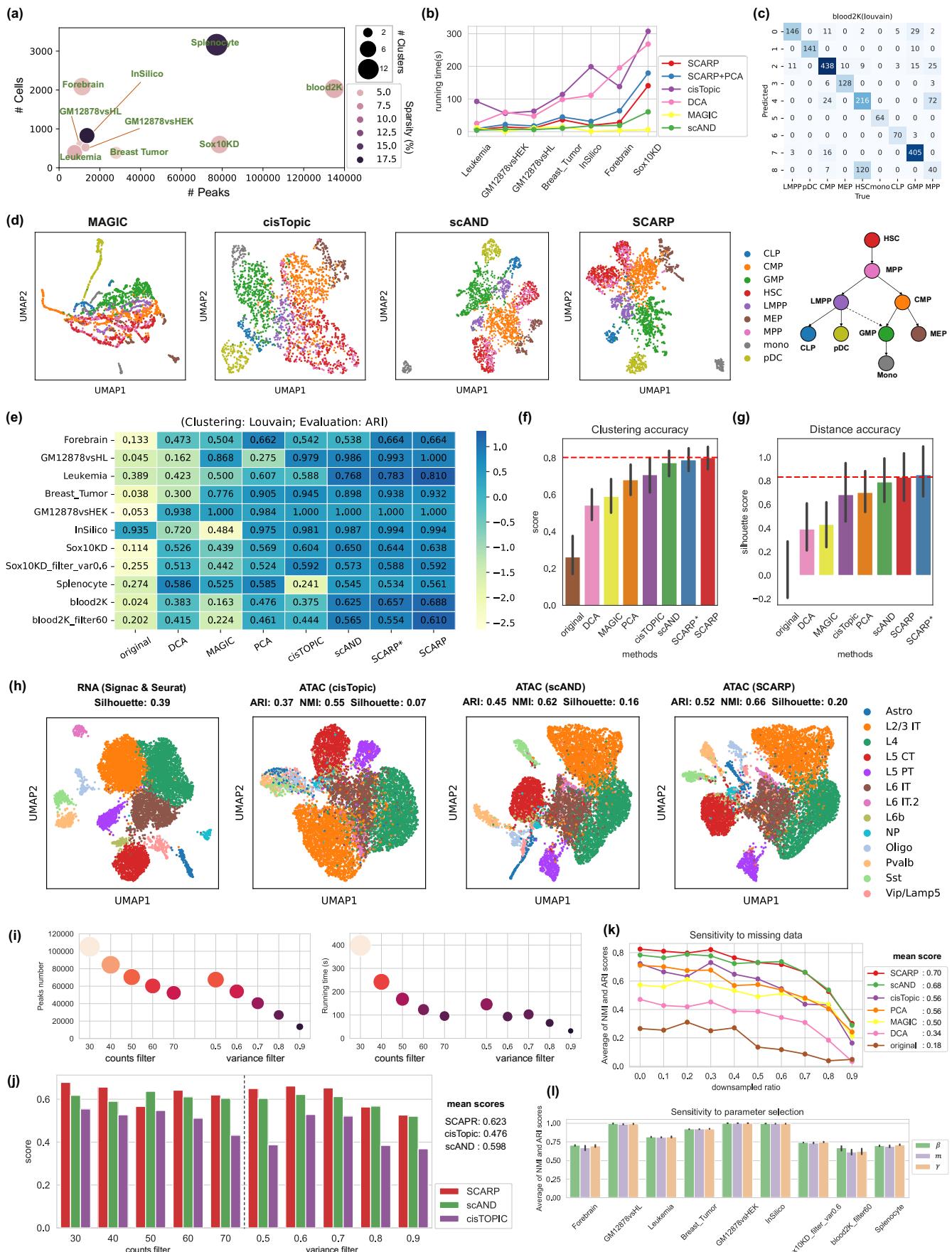
736



737

738 **Fig. 1. The workflow of SCARP.** SCARP consists of two important modules, the first is constructing the  
 739 network from scATAC-seq data, with the nodes being cells or peaks, the edges between cells and peaks  
 740 reflecting accessibility, and the edges between peaks reflecting prior information about their genome position;  
 741 the second is using NR diffusion to compensate for the effects of high sparsity of scATAC-seq data as well  
 742 as characterize the similarity between cells. The output matrix of SCARP with the low-dimensional embedding  
 743 of cells and peaks can further benefit the downstream analyses.

744

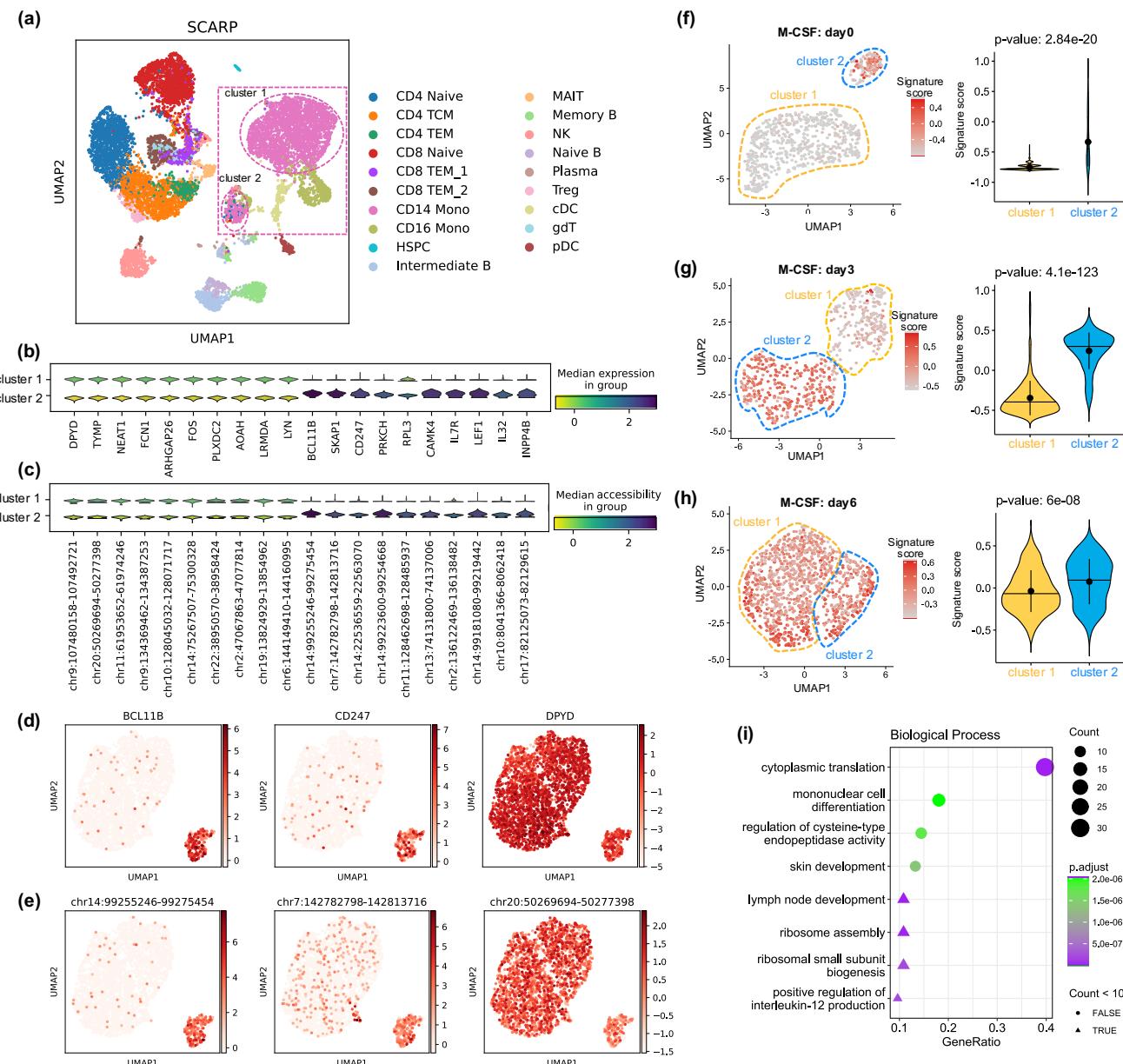


745

746 **Fig. 2. Cell clustering performance on benchmarking scATAC-seq datasets. (a)** Data presentation.  
747 Each dot represents a scATAC-seq dataset, with the x-axis (y-axis) showing its number of peaks (cells), the

748 size of the dots showing the number of annotated cell types, and the color of the dots showing the data  
749 sparsity. **(b)** Running times (y-axis) of various methods on seven benchmarking datasets (x-axis). One  
750 method is shown by one line (color). **(c)** The confusion matrix of SCARP-derived cell embeddings clustered  
751 by Louvain (rows) and the annotated cell types (columns) on the blood2K dataset. The matrix elements  
752 represent the number of overlapping cells in a predicted cluster and a real cluster. **(d)** The UMAP plots of four  
753 methods MAGIC, cisTopic, scAND, and SCARP on the blood2K dataset with 2034 cells, colored by the  
754 annotated cell types. The rightest legend was the known FACS-sorted population of origin, which shows the  
755 process of cell development and differentiation. **(e)** The cell clustering results of eight methods (x-axis) on  
756 eleven datasets (y-axis) under ARI evaluation. Colors represent the Z-scores calculated from ARI score of  
757 the corresponding dataset. **(f)** Cell clustering accuracy (y-axis) of various methods (x-axis). Scores were  
758 calculated by averaging the ARI scores and NMI scores for a method over all datasets. **(g)** Cell distance  
759 accuracy (y-axis) of various methods (x-axis). Scores were calculated by averaging the silhouette coefficient  
760 scores for a method over all datasets. **(h)** Four UMAP visualizations showing the cell clustering results of  
761 Seurat & Signac based on scRNA-seq data, as well as cisTopic, scAND, and SCARP based on scATAC-seq  
762 data, and their clustering accuracy scores (title). The plots were colored by the annotated cell types obtained  
763 by Seurat. **(i)** The x-axis represents different peak filtering strategies, and the y-axis on the left panel  
764 represents number of peaks in the blood2K dataset, the y-axis on the right panel represents running time of  
765 SCARP. **(j)** Barplots of cell clustering performance (y-axis) for SCARP, scAND and cisTopic under different  
766 strategies of peaks filtering (x-axis). Scores (y-axis) were calculated by averaging the ARI score and NMI  
767 score for one method. Mean score (right legend) for each method was calculated by averaging the scores for  
768 all peak filtering strategies. **(k)** The average of NMI and ARI scores (y-axis) of various methods on the  
769 Leukemia dataset with different downsampling ratios (x-axis). **(l)** The average of NMI and ARI scores (y-axis)  
770 when taking different values of parameters  $\beta$ ,  $m$ , and  $\gamma$  on various scATAC-seq dataset (x-axis). Error bars  
771 represent the standard deviation (SD) of the scores.

772

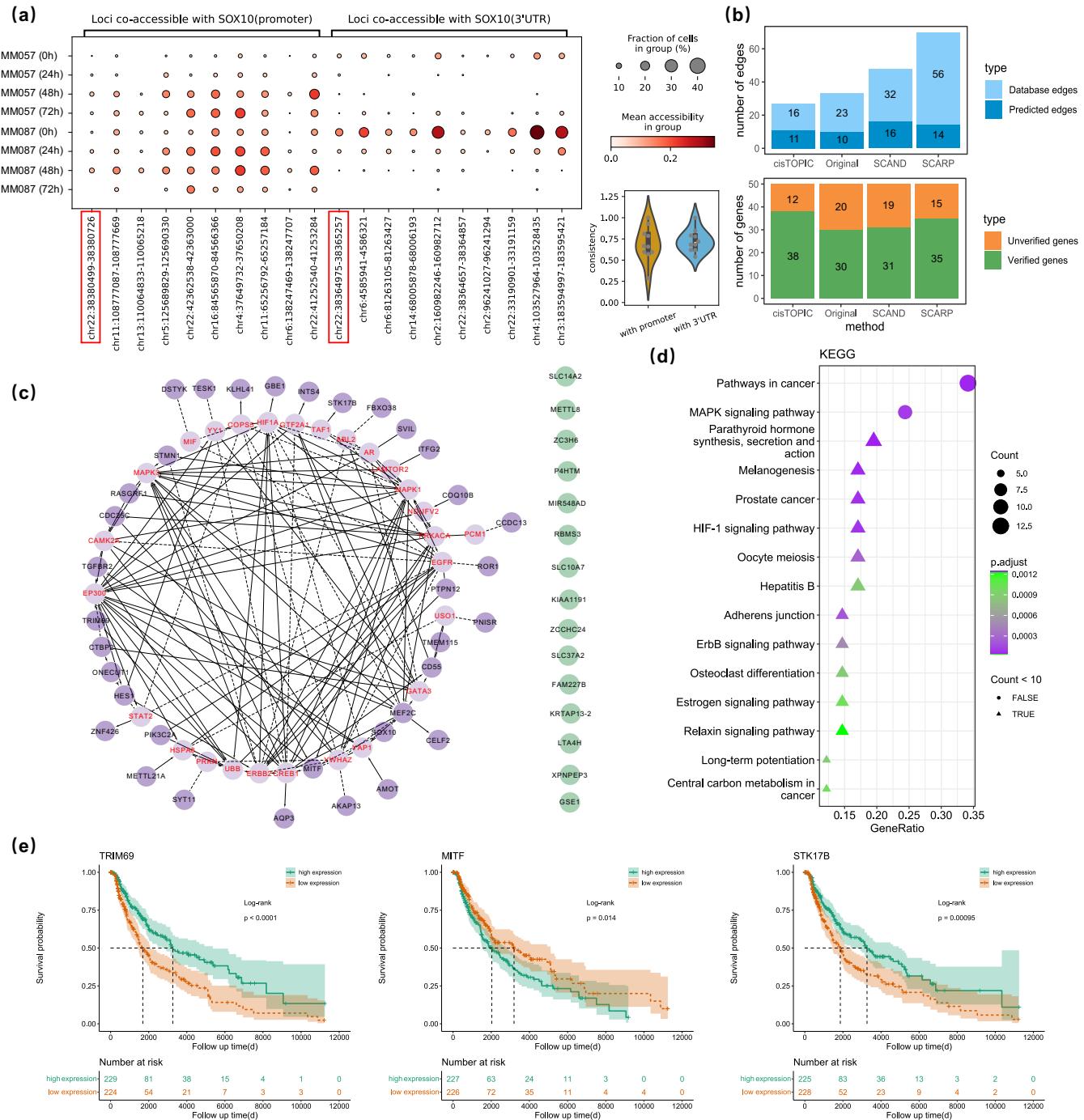


773

774 **Fig. 3. SCARP identified a new CD14 monocyte subpopulation.** (a) UMAP visualization of SCARP  
 775 -derived cell embeddings colored by the cell types annotated by Seurat. Those CD14 monocytes were framed  
 776 in pink, and two cell subgroups were detected, denoted as cluster 1 and cluster 2. (b) Violin plots showing  
 777 the expression levels of the top 10 marker genes (x-axis) for two cell subpopulations (y-axis). (c) Violin plots  
 778 showing the accessibility levels of the top 10 marker peaks (x-axis) for two cell subpopulations (y-axis). (d)  
 779 The UMAP visualizations of cells colored by the expressions of cluster 2 marker genes *BCL11B*, *CD247*, and  
 780 the expressions of cluster 1 marker gene *DPYD*. (e) The UMAP visualizations of cells colored by the  
 781 accessibility of cluster 2 marker peaks chr14:99255246-99275454 and chr7:142782798-142813716, and the  
 782 accessibility of cluster 1 marker peak chr20:50269694-50277398. (f)(g)(h): The UMAP visualizations of cells  
 783 in external M-CSF scRNA-seq dataset, colored by gene signature scores for CD14 monocytes (donor 1)  
 784 under M-CSF stimulation at day 0 (f), day 3 (g) and day 6 (h), along with the violin plots showing the difference

785 between the two cell clusters. **(i)** Top eight biological processes that gene signatures are involved in, using  
786 all genes in this scRNA-seq dataset as the background gene set. The y-axis represents the name of KEGG  
787 pathway, and the x-axis represents the ratio of genes enriched in a pathway. The number of genes enriched  
788 in a pathway is related to the size of the circle (gene count less than 10) or triangle (gene count greater than  
789 10), and the color represents the corresponding adjusted p-value.

790

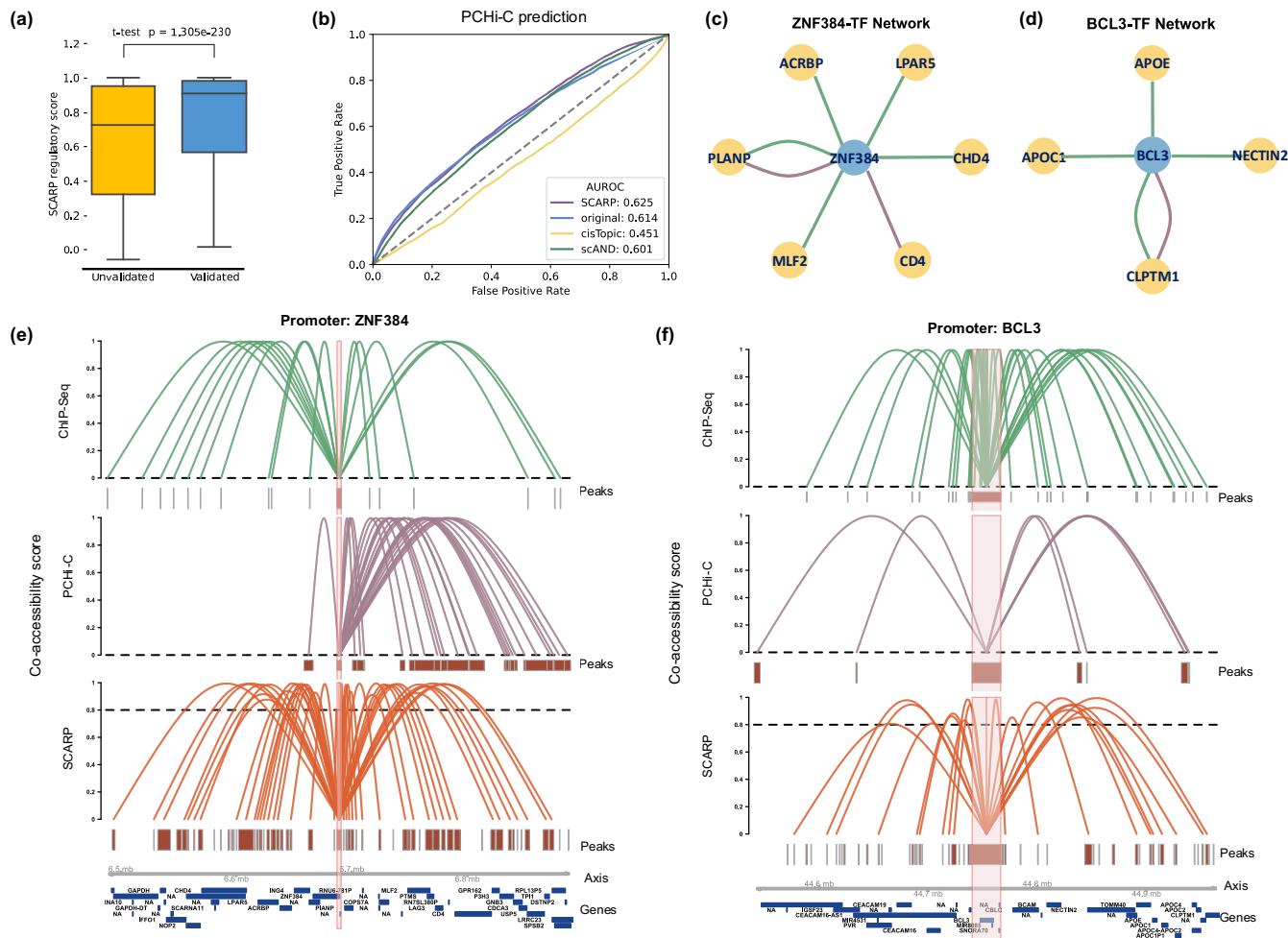


791

792 **Fig. 4. SCARP portrayed the co-accessibility of peaks on SOX10KD scATAC-seq dataset. (a)** Changes  
 793 in accessibility for those loci (x-axis) that were highly correlated with SOX10 promoter and 3'UTR as  
 794 knockdown time walked by (y-axis). The SOX10 promoter and 3'UTR were framed in red. The size of the dots  
 795 showing the fraction of cells in cell groups, and the color of the dots showing the mean accessibility in cell  
 796 groups. The violin plot at the right bottom showed the consistency of the trend for the top ten peaks  
 797 coaccessible with the promoter (3'UTR) of SOX10, with consistency calculated from the correlation of mean  
 798 accessibility within eight cell groups. **(b)** For each of the differential gene sets (50 genes) obtained by four  
 799 methods (x-axis), we compared the number of genes (y-axis of lower panel) that were verified (green bars)

800 as directly or indirectly related by the database, and the number of genes that were not verified (orange bars),  
801 as well as the number of edges (y-axis of upper panel) in the corresponding GRNs, including predicted edges  
802 (dark blue), i.e., those with a functional interaction (FI) type of “predicted”, and database edges (light blue),  
803 i.e., those with a FI type other than “predicted”, such as “activate” and “inhibited by”. **(c)** The transcriptional  
804 regulatory relationships validated by the Reactome database. Solid lines indicate experimentally validated  
805 regulatory relationships, while dashed lines indicate database-predicted regulatory relationships. The dark  
806 purple nodes, light purple nodes, and green nodes represent the SCARP-derived genes, linker genes, and  
807 unvalidated genes, respectively. **(d)** Results of KEGG pathway analysis. The legends are the same as in Fig.  
808 3*i*. **(e)** The results of survival analysis on genes *TRIM69*, *MITF*, and *STK17B*.

809



810

811 **Fig. 5. SCARP-derived cis-regulatory interactions are significant and validated by external evidence.**

812 **(a)** The boxplot showing the difference between PCHi-C validated (blue) and unvalidated (yellow) SCARP-  
813 derived transcriptional regulatory scores (y-axis), followed by t-test to test the difference. **(b)** ROC curve of  
814 various methods, treating whether SCARP-derived transcriptional regulatory relationships were validated by  
815 PCHi-C data as a binary classification problem. **(c), (d)**, SCARP-derived TF-regulon networks of genes  
816 ZNF384 **(c)** and BCL3 **(d)**. The green and purple edges mean those interactions are also supported by ChIP-  
817 seq or PCHi-C data, respectively. The blue and yellow nodes mean the TF genes and target genes  
818 respectively. **(e), (f)**, Visualizations of co-accessibility scores (y-axis) of SCARP-derived cis-regulatory  
819 interactions (orange), PCHi-C evidence (purple), and ChIP-seq evidence (green) around the genes **(e)**  
820 ZNF384 locus (x-axis) and **(f)** BCL3 locus (x-axis).

821