

1 **APEC: an accesson-based method for single-cell chromatin**  
2 **accessibility analysis**

3

4 Bin Li<sup>1,3</sup>, Young Li<sup>1,3</sup>, Kun Li<sup>1</sup>, Lianbang Zhu<sup>1</sup>, Qiaoni Yu<sup>1</sup>, Pengfei Cai<sup>1</sup>, Jingwen Fang<sup>1</sup>,  
5 Wen Zhang<sup>1</sup>, Pengcheng Du<sup>1</sup>, Chen Jiang<sup>1</sup>, Kun Qu<sup>1,2,\*</sup>

6

7 **Affiliation**

8 <sup>1</sup>Division of Molecular Medicine, Hefei National Laboratory for Physical Sciences at  
9 Microscale, The CAS Key Laboratory of Innate Immunity and Chronic Disease,  
10 Department of Oncology of the First Affiliated Hospital, Division of Life Sciences and  
11 Medicine, University of Science and Technology of China

12 <sup>2</sup>CAS Center for Excellence in Molecular Cell Sciences, University of Science and  
13 Technology of China

14 <sup>3</sup>Cofirst authors

15

16 \*Correspondence: Kun Qu ([gukun@ustc.edu.cn](mailto:gukun@ustc.edu.cn))

17

18 **Contact Information**

19 Kun Qu, Ph.D.

20 Division of Molecular Medicine, Hefei National Laboratory for Physical Sciences at  
21 Microscale, The CAS Key Laboratory of Innate Immunity and Chronic Disease, School of  
22 Life Sciences, University of Science and Technology of China

23 Hefei, Anhui, China, 230027

24 Email: [gukun@ustc.edu.cn](mailto:gukun@ustc.edu.cn)

25 Phone: +86-551-63606257

26

27 **ABSTRACT**

28

29 The development of sequencing technologies has promoted the survey of genome-wide  
30 chromatin accessibility at single-cell resolution; however, comprehensive analysis of  
31 single-cell epigenomic profiles remains a challenge. Here, we introduce an accessibility  
32 pattern-based epigenomic clustering (APEC) method, which classifies each individual cell  
33 by groups of accessible regions with synergistic signal patterns termed “accessions”. By  
34 integrating with other analytical tools, this python-based APEC package greatly improves  
35 the accuracy of unsupervised single-cell clustering for many different public data sets.  
36 APEC also predicts gene expressions, identifies significant differential enriched motifs,  
37 discovers super enhancers, and projects pseudotime trajectories. Furthermore, we  
38 adopted a fluorescent tagmentation-based single-cell ATAC-seq technique (ftATAC-seq)  
39 to investigate the per cell regulome dynamics of mouse thymocytes. Associated with  
40 ftATAC-seq, APEC revealed a detailed epigenomic heterogeneity of thymocytes,  
41 characterized the developmental trajectory and predicted the regulators that control the  
42 stages of maturation process. Overall, this work illustrates a powerful approach to study  
43 single-cell epigenomic heterogeneity and regulome dynamics.

44

## 45 INTRODUCTION

46

47 As a technique for probing genome-wide chromatin accessibility in a small number of cells  
48 *in vivo*, the assay for transposase-accessible chromatin with high-throughput sequencing  
49 (ATAC-seq) has been widely applied to investigate the cellular regulomes of many  
50 important biological processes <sup>1</sup>, such as hematopoietic stem cell (HSC) differentiation <sup>2</sup>,  
51 embryonic development <sup>3</sup>, neuronal activity and regeneration <sup>4,5</sup>, tumor cell metastasis<sup>6</sup>,  
52 and patient responses to anticancer drug treatment <sup>7</sup>. Recently, several experimental  
53 schemes have been developed to capture chromatin accessibility at single-cell/nucleus  
54 resolution, i.e., single-cell ATAC-seq (scATAC-seq) <sup>8</sup>, single-nucleus ATAC-seq (snATAC-  
55 seq) <sup>9,10</sup>, and single-cell combinatorial indexing ATAC-seq (sci-ATAC-seq) <sup>11,12</sup>, which  
56 significantly extended researchers' ability to uncover cell-to-cell epigenetic variation and  
57 other fundamental mechanisms that generate heterogeneity from identical DNA  
58 sequences. By contrast, the in-depth analysis of single-cell chromatin accessibility profiles  
59 for this purpose remains a challenge. Numerous efficient algorithms have been developed  
60 to accurately normalize, cluster and visualize cells from single-cell transcriptome  
61 sequencing profiles, including but not limited to Seurat <sup>13</sup>, SC3 <sup>14</sup>, SIMLR <sup>15</sup>, and SCANPY  
62 <sup>16</sup>. However, most of these algorithms are not directly compatible with a single-cell ATAC-  
63 seq dataset, for which the signal matrix is much sparser. To characterize scATAC-seq  
64 data, the Greenleaf lab developed an algorithm named chromVAR <sup>17</sup>, which aggregates  
65 mapped reads at accessible sites based on annotated motifs of known transcription  
66 factors (TFs) and thus projects the sparse per accessible peak per cell matrix to a bias-  
67 corrected deviation per motif per cell matrix and significantly stabilizes the data matrix for  
68 downstream clustering analysis. Other mathematical tools, such as the latent semantic  
69 indexing (LSI) <sup>11,12</sup>, Cicero <sup>18</sup>, cisTopic <sup>19</sup>, and SnapATAC <sup>20</sup> have also been applied to  
70 process single-cell/nucleus ATAC-seq data <sup>10,12</sup>. However, great challenges still remain  
71 for current algorithms to precisely cluster large number of cells and predict gene  
72 expressions from single cell chromatin accessibility profiles. Therefore, a refined algorithm  
73 is needed to better categorize cell subgroups with minor differences, thereby providing a  
74 deeper mechanistic understanding of single-cell epigenetic heterogeneity and regulation.

75 Here, we introduce a new single-cell chromatin accessibility analysis toolkit named  
76 APEC (accessibility pattern-based epigenomic clustering), which combines peaks with the  
77 same signal fluctuation among all single cells into peak groups, termed "accessions", and  
78 converts the original sparse cell-peak matrix to a much denser cell-accession matrix for

79 cell type categorization (Fig. 1a). In contrast to previous methods, this accesson-based  
80 reduction scheme naturally groups synergistic accessible regions genome-wide together  
81 without a priori knowledge of genetic information (such as TF motifs or genomic distance)  
82 and provides an efficient, accurate, robust and rapid cell clustering from single-cell ATAC-  
83 seq profiles. APEC was also integrated into a head-to-toe program package that has been  
84 made available on GitHub (<https://github.com/QuKunLab/APEC>).

85

86

## 87 RESULTS

88

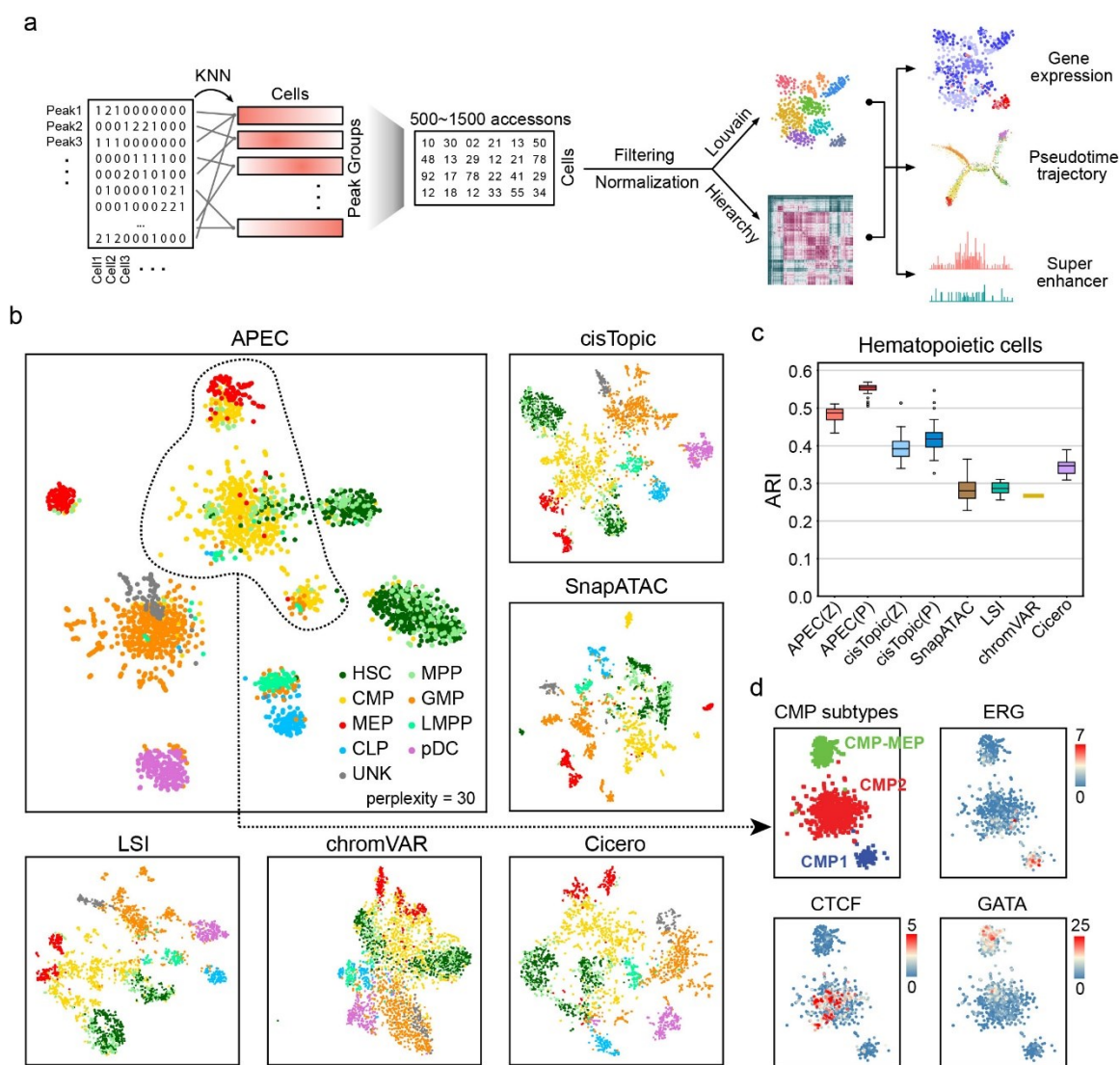
### 89 Accesson-based algorithm improves single-cell clustering

90

91 To test the performance of APEC, we first obtained data from previous publications  
92 that performed scATAC-seq on a variety of cell types with known identity during  
93 hematopoietic stem cell (HSC) differentiation<sup>21</sup> as a gold standard. Compared to other  
94 state-of-the-art single cell chromatin accessibility analysis methods, this new accesson-  
95 based algorithm can clearly cluster cells into their corresponding identities according to  
96 the Adjusted Rand Index (ARI) (Fig. 1b & c). Here, we adopted the Louvain clustering  
97 algorithm for all the methods, and assumed that the number of cell types was unknown to  
98 ensure that the clustering analysis was unsupervised. To calculate ARI values, we ignored  
99 all “unknown” cells, and sampled the tunable parameter (e.g., accesson number in APEC,  
100 random seed in cisTopic, etc.) multiple times for each tool (see Methods). On average,  
101 67% of cells were correctly classified by APEC with ARI=0.483/0.552 (using matrices  
102 normalized by Z-score or probability, respectively), while cisTopic was the second most  
103 accurate method to predict cell identities (average ARI=0.392/0.418) (Supplementary  
104 Table 1).

105 Moreover, APEC identified 3 sub-clusters of CMP cells that were not discovered  
106 by any other algorithms, namely CMP1, CMP2 and CMP-MEP (Fig. 1d). CMP1 cells are  
107 early stage of CMPs that enriched TFs associated with stem cell self-renewal, such as Erg  
108<sup>22</sup>; CMP2 cells are enriched with CTCF motif, suggesting that these cells are at the fate  
109 decision stage with CTCF associated chromatin remodeling<sup>23</sup>; CMP-MEP cells are  
110 considered as MEP committed CMPs, and are strongly enriched with crucial regulators  
111 for MEP differentiation, such as GATA1<sup>24</sup>. However, these 3 CMP sub-clusters were not  
112 as clearly distinguishable and cells were more scattered on the tSNE maps generated by

113 other methods (Supplementary Fig. 1). More details about the distribution of these 3 sub-  
 114 clusters of CMP cells on the development trajectory will be discussed later in the section  
 115 of pseudotime prediction.



116

117 **Figure 1. The accession matrix constructed from the sparse fragment count matrix improved**  
 118 **the clustering of scATAC-seq data.** (a) Step-by-step workflow of APEC. Peaks were grouped  
 119 into accessions by their accessibility pattern among cells with the K nearest neighbors (KNN)  
 120 method. (b) t-Distributed Stochastic Neighbor Embedding (tSNE) diagrams of the hematopoietic  
 121 single cells dataset based on the dimension-transformed matrices from different algorithms, i.e.,  
 122 APEC: accession matrix; cisTopic: topic matrix; LSI: LSI matrix; chromVAR: bias-corrected  
 123 deviation matrix; Cicero: aggregated model matrix. The cells are FACS-indexed human  
 124 hematopoietic cells, including HSCs (hematopoietic stem cells), MPPs (multipotent progenitors),  
 125 LMPPs (lymphoid-primed multipotential progenitors), CMPs (common myeloid progenitors), CLPs  
 126 (common lymphoid progenitors), pDCs (plasmacytoid dendritic cells), GMPs (granulocyte-  
 127 macrophage progenitors), MEPs (megakaryocyte-erythroid progenitors), and UNK (unknown type)  
 128 cells. (c) The ARI (Adjusted Rand Index) values for the clustering of the human hematopoietic cells  
 129 by different algorithms. Same as the two normalization methods applied in cisTopic, we normalized

130 the accession matrix in APEC based on probability (P) and z-score (Z). Center line, median; box  
131 limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers. (d) Three CMP  
132 subtypes identified in APEC and the motifs enriched in each cell subtype.

133

134 To further confirm the superiority of APEC, we performed the same comparison  
135 analysis with another scATAC-seq dataset on three distinct cell types, namely, lymphoid-  
136 primed multipotent progenitors (LMPPs), monocytes, and HL-60 lymphoblastoid cells  
137 (HL60), and four similar cell types, namely, blast cells and leukemic stem cells (LSCs)  
138 from two acute myeloid leukemia (AML) patients<sup>17</sup>. We found that both APEC and cisTopic  
139 were tied for best to classify these cells ([Supplementary Fig. 2a](#)). Interestingly, APEC,  
140 cisTopic and LSI were all capable of almost perfectly separating the three distinct cell  
141 types (LMPPs, monocytes, and HL60), with ARI = 1.000/0.988, 0.987/0.987 and 0.969,  
142 respectively. However, in terms of clustering the four similar cell types from AML patients,  
143 APEC (average ARI=0.575/0.564) outperformed other tools ([Supplementary Fig. 2b](#)),  
144 suggesting that APEC was the most sensitive among all the tools. Since each method can  
145 generate varying numbers of clusters depending on the parameters used, we  
146 benchmarked the performance of all the methods using ARI across a wide range of  
147 tunable parameters to ensure the reliability of their predictions ([Supplementary Fig. 2c](#),  
148 [Supplementary Table 2](#)). To further test the robustness of APEC at low sequencing depth,  
149 we randomly selected reads from the raw data and calculated the ARI values for each  
150 down-sampled data. APEC exhibits better performance at sequencing depths as low as  
151 20% of the original data ([Supplementary Fig. 2d](#)), confirming the sensitivity of the algorithm.

152 Compare with chromVAR, the contribution of the minor differences between similar  
153 cells is aggregated in accessions but diluted in motifs. For example, APEC identified  
154 prominent super-enhancers around the genes *N4BP1*<sup>25</sup> and *GPHN*<sup>26</sup> in the LSC cells  
155 from AML patient 1 (P1-LSC) but not the other cell types ([Supplementary Fig. 3a & b](#)).  
156 These two loci were also confirmed as super enhancers by ROSE<sup>27</sup> (the top 2 candidates  
157 in [Supplementary Table 3](#)). We noticed that all peaks in these super-enhancers were  
158 classified into one accession that was critical for distinguishing P1-LSCs from P2-LSCs,  
159 P1-blast cells and P2-blast cells. However, these peaks were distributed in multiple TF  
160 motifs, which significantly diluted the contributions of the minor differences  
161 ([Supplementary Fig. 3c & d](#)).

162 In contrast to Cicero, which aggregates peaks based on their cis-co-accessibilities  
163 networks (CCAN) within a certain range of genomic distance<sup>18</sup>, APEC combine synergistic  
164 peaks genome-wide. Take the human hematopoietic cells dataset as an example, 600

165 accessions were built from the 54,212 peaks, and each accession contained ~40 peaks  
166 (median number) compare with ~4 peaks in each CCAN ([Supplementary Fig. 4a & b](#)). The  
167 average distance between peaks in a same accession is ~50 million base pairs (compared  
168 with ~0.2 million bps from CCAN), and over 57% of accessions contain peaks from more  
169 than 15 different chromosomes. From the same dataset, Cicero identified 732,306 pairs  
170 of site links from 25,102 peaks, and information from the remaining peaks were simply  
171 discarded. APEC identified more than 9.2 million pairs of site links from all the 54,212  
172 peaks, within which only 3080 site links were identified by both methods ([Supplementary  
173 Fig. 4c](#)), therefore, APEC and Cicero are two completely different approaches.  
174 Furthermore, Buenrostro et al. showed that the covariation of the accessible sites across  
175 all the cells may reflect the spatial distance between the corresponding peaks <sup>8</sup>. By  
176 integrating the chromatin conformation profiles from Hi-C experiments with the scATAC-  
177 seq profile for the same cells, we found that peaks in the same accession are spatially  
178 much closer to each other than randomly selected peaks ([Supplementary Fig. 5a](#), P-  
179 value<10<sup>-7</sup>), suggesting that they may belong to the same topologically associated  
180 domains (TADs) ([Supplementary Fig. 5b](#)).

181 Speed and scalability are now extremely important for single-cell analytical tools  
182 due to the rapid growth in the number of cells sequenced in each experiment. We  
183 benchmarked APEC and all the other tools based on a random sampling of the mouse *in  
184 vivo* single-cell chromatin accessibility atlas dataset <sup>28</sup>, which contains 81,173 high quality  
185 cells. Taking into account of all the 436,206 peaks, it took APEC 310 min to cluster 80,000  
186 cells with 1 CPU thread ([Supplementary Fig. 5c](#)). We also randomly select 100,000 peaks  
187 from the entire dataset to test the computer time spent by these tools ([Supplementary Fig.  
188 5d](#)). In addition, APEC is very stable for a wide range of parameter values used in the  
189 algorithm, such as the number of accessions, nearest neighbors and principle components  
190 ([Supplementary Fig. 5e-g](#)).

191

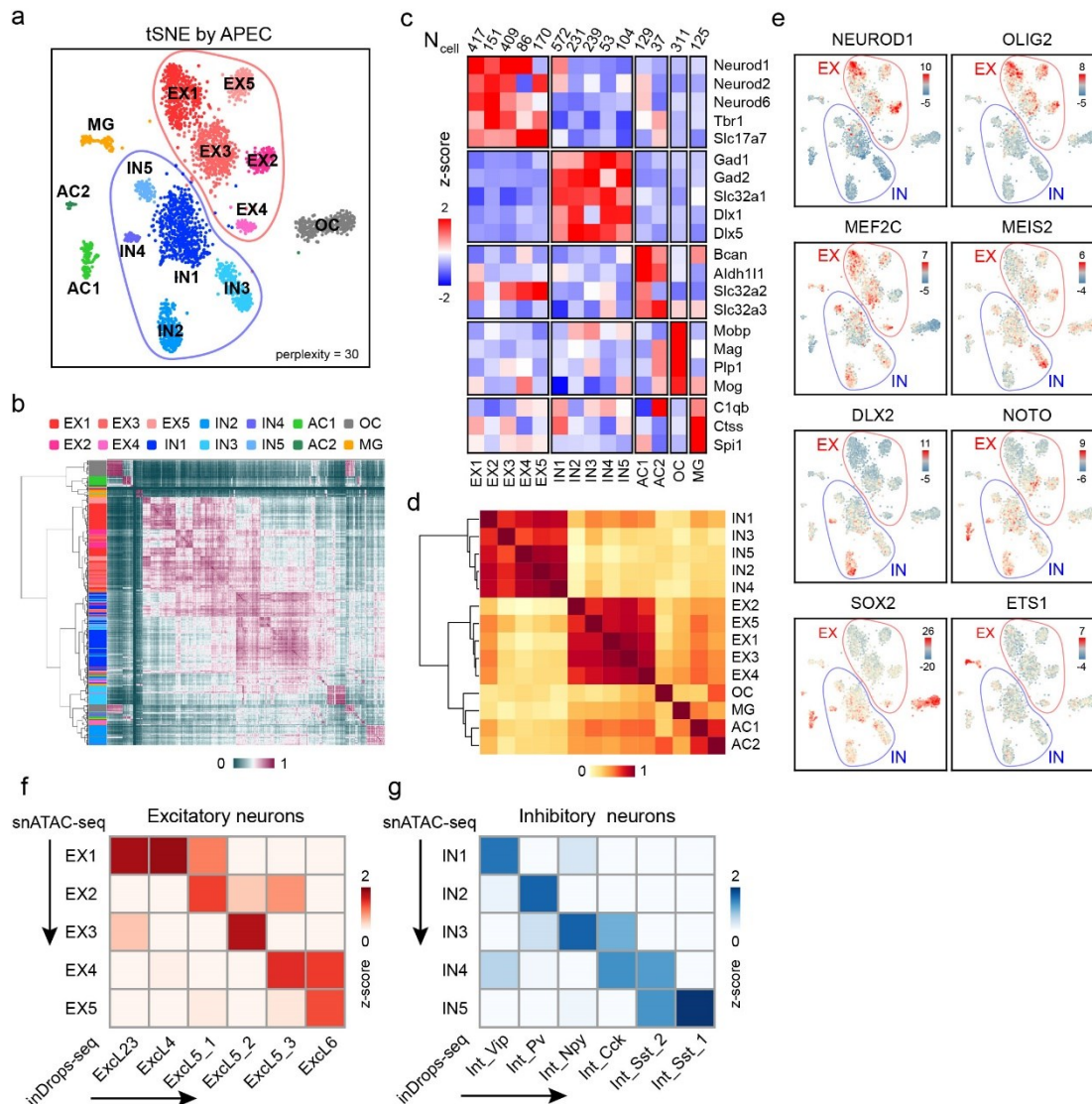
## 192 **APEC is applicable to other single-cell chromatin detection techniques**

193

194 To evaluate the compatibility and performance of APEC with other single-cell  
195 chromatin accessibility detection techniques, such as snATAC-seq <sup>10</sup>, transcript-indexed  
196 scATAC-seq <sup>29</sup> and sciATAC-seq <sup>11</sup>, APEC was also tested with the data sets generated  
197 by those experiments. For example, APEC discovered 14 cell subpopulations in adult  
198 mouse forebrain snATAC-seq data <sup>10</sup>, including four clusters of excitatory neurons (EX1-

199 5), five groups of inhibitory neurons (IN1-5), astroglia cells (AC1&2), oligodendrocyte cells  
200 (OC), and microglial cells (MG; Fig. 2a & b). To quantify gene expression level, we defined  
201 a gene's score as the average signal of the peaks close to its TSS region (Fig. 2c; see  
202 Methods). With that, we identified 5 excitatory subpopulations and 5 distinct inhibitory  
203 subpopulations, and all the cell groups were clearly distinguished from each other by  
204 hierarchical clustering (Fig. 2d). In contrast, more than 29.7% (946 out of 3034) of cells  
205 were unable to be correctly assigned into any subpopulation of interest in previous study<sup>10</sup>.  
206 Besides, scores for several marker genes such as *Neurod6* and *Aldh111* were not available  
207 in Cicero and cisTopic, respectively, making it difficult to identify the corresponding cell  
208 clusters (i.e. c0~c2 in Supplementary Fig. 6a & b). The SnapATAC algorithms however,  
209 mis-clustered the AC1/2 with excitatory neurons in the correlation matrix (Supplementary  
210 Fig. 6c). On the other hand, the motif enrichment analysis module in APEC identified cell  
211 type-specific regulators that are also consistent with previous publications<sup>10</sup>. For example,  
212 the NEUROD1 and OLIG2 motifs were generally enriched on excitatory clusters  
213 (EX1\2\3\5); the MEF2C motif was more enriched on EX1/3/4/5; the motifs of MEIS2 and  
214 DLX2 were differentially enriched on different subtypes of inhibitory neurons (IN3 and  
215 IN2/4, respectively); and the NOTO, SOX2, and ETS1 motifs were enriched on the AC1,  
216 OC, and MG clusters, respectively (Fig. 2e). These results confirm that APEC can  
217 distinguish and categorize single cells with great sensitivity and reliability.





218

219 **Figure 2. APEC improved the cell type classification of adult mouse forebrain snATAC-seq**  
 220 **data.** (a) A tSNE diagram demonstrates the APEC clustering of forebrain cells. (b) Hierarchical  
 221 clustering of the cell-cell correlation matrix. The side bar denotes cell clusters from APEC. (c)  
 222 Average scores of the marker genes for each cell cluster generated by the method mentioned in  
 223 the data source paper<sup>10</sup>, and normalized by the standard score (z-score). The top row lists the  
 224 number of cells in each cluster. (d) Hierarchical clustering of the cluster-cluster correlation matrix.  
 225 (e) Differential enrichments of cell type-specific motifs in each cluster. (f, g) Fisher's exact test  
 226 between the differential genes of the excitatory (EX) and inhibitory (IN) neurons from snATAC-seq  
 227 and the signature genes of the excitatory (Excl) and inhibitory (Int) neurons defined in inDrops-  
 228 seq<sup>30</sup>. The heatmaps show the  $-\log(P\text{-value})$  normalized by calculating the z-scores through each  
 229 row and column.

230

231 Since single-cell transcriptome analysis is also capable to identify novel cell  
 232 subpopulations, it is critical to anchor the cell types identified from scATAC-seq to those  
 233 from scRNA-seq. Hrvatin et al. identified multiple excitatory and inhibitory neuronal

234 subtypes in the mouse visual cortex using single cell inDrops sequencing<sup>30</sup> and provided  
235 top 20 signature genes that distinguished these cell subtypes. However, due to the  
236 sparseness of snATAC-seq matrix, scores of many signature genes were not strong  
237 enough to distinguish the cell sub-clusters. To overcome this, we developed a gene sets  
238 overlap algorithm to associate cell clusters from scATAC-seq and scRNA-seq profiles (see  
239 [Methods](#)). We found that sub-cluster EX1~5 and IN1~5 in snATAC-seq can nicely  
240 correspond to the neuron subtypes classified by Hrvatin et al ([Fig. 2f-g](#)). These results  
241 highlight the potential advantages of the accesson-based approach for the integrative  
242 analysis of scATAC-seq and scRNA-seq data.

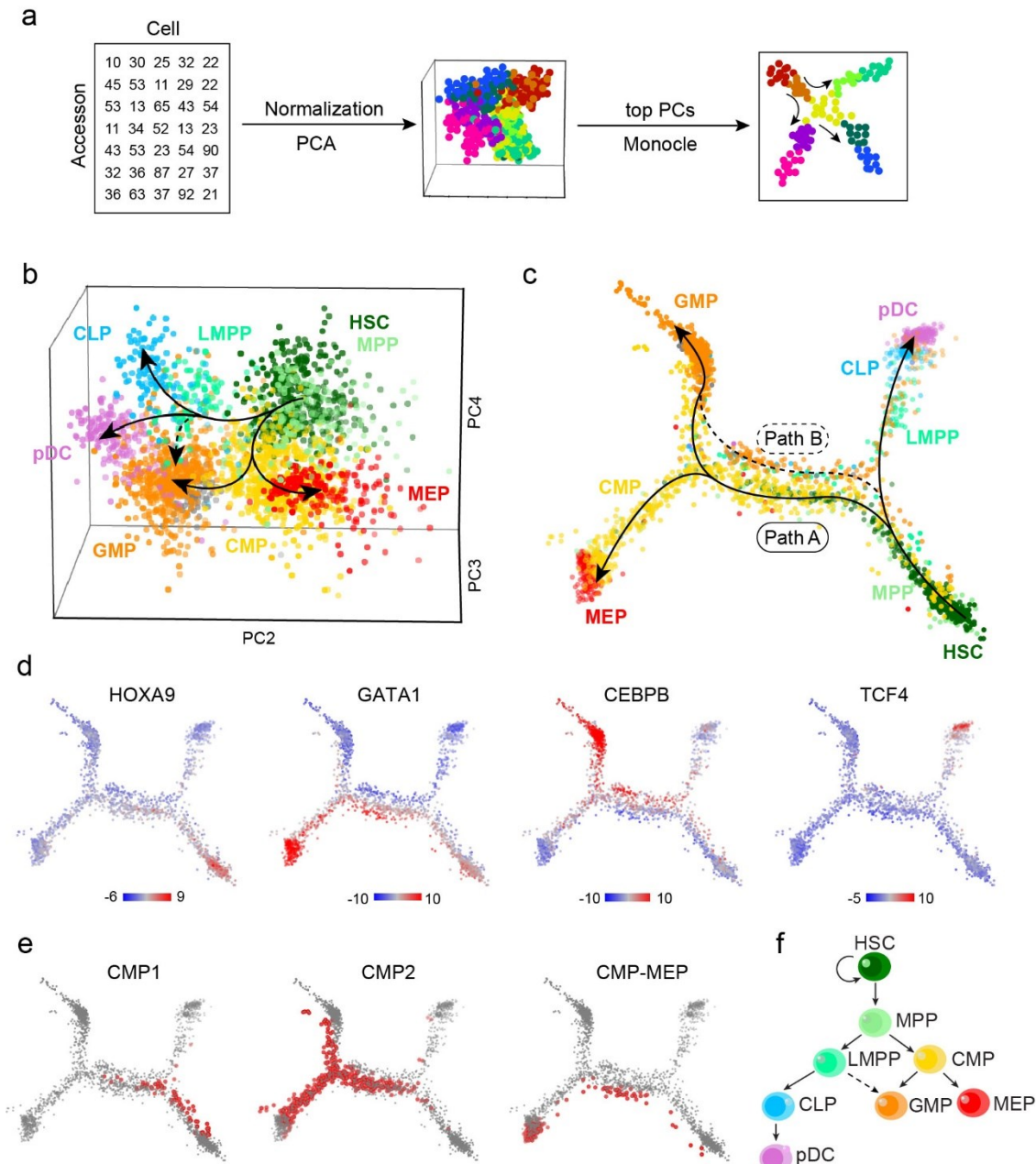
243

### 244 **APEC constructs a pseudotime trajectory that predicts cell differentiation lineage**

245

246 Cells are not static but dynamic entities, and they have a history, particularly a  
247 developmental history. Although single-cell experiments often profile a momentary  
248 snapshot, a number of remarkable computational algorithms have been developed to  
249 pseudo-order cells based on the different points they were assumed to occupy in a  
250 trajectory, thereby leveraging biological asynchrony<sup>31,32</sup>. For instance, Monocle<sup>32,33</sup>  
251 constructs the minimum spanning tree, and Wishbone<sup>34</sup> and SPRING<sup>35</sup> construct the  
252 nearest neighbor graph from single-cell transcriptome profiles. These tools have been  
253 widely used to depict neurogenesis<sup>36</sup>, hematopoiesis<sup>37,38</sup> and reprogramming<sup>39</sup>. APEC  
254 integrates the Monocle algorithm into the accesson-based method and enables  
255 pseudotime prediction from scATAC-seq data<sup>21</sup> and was applied to investigate HSC  
256 differentiation lineages ([Fig. 3a](#)). Principal component analysis (PCA) of the accesson  
257 matrix revealed multiple stages of the lineage during HSC differentiation ([Fig. 3b](#)) and was  
258 consistent with previous publications<sup>2,21</sup>. After utilizing the Monocle package, APEC  
259 provided more precise pathways from HSCs to the differentiated cell types ([Fig. 3c](#)). In  
260 addition to the differentiation pathways to MEP cells through the CMP state and to CLP  
261 cells through the LMPP state, MPP cells may differentiate into GMP cells through two  
262 distinct trajectories: Path A through the CMP state and Path B through the LMPP state,  
263 which is consistent with the composite model of HSC and blood lineage commitment<sup>40</sup>.  
264 Notably, pDCs from bone marrow are CD34<sup>+</sup> ([Supplementary Fig. 7](#)), indicative of  
265 precursors of pDCs. APEC suggested that pDC precursors were derived from CLP cells  
266 on the pseudotime trajectory ([Fig. 3c](#)), which also agrees with previous reports<sup>41</sup>.  
267 Furthermore, APEC incorporated the chromVAR algorithm to determine the regulatory

268 mechanisms during HSC differentiation by evaluating the deviation of each TF along the  
269 single-cell trajectory. As expected, the HOX motif is highly enriched in the accessible sites  
270 of HSCs/MPP cells, as are the GATA1, CEBPB and TCF4 motifs, which exhibit gradients  
271 that increase along the erythroid, myeloid and lymphoid differentiation pathways,  
272 respectively<sup>21</sup> (Fig. 3d). We also noticed that the TF regulatory strategies of the two paths  
273 from MPP towards GMP cells were very different. In addition, the 3 CMP sub-clusters  
274 identified in Figure 1 were differentially distributed along the developmental trajectory (Fig.  
275 3e). CMP1 cells that close to HSCs and MPPs are early stage CMPs; CMP2 cells are  
276 distributed in both the GMP and MEP branches; CMP-MEP cells are MEP committed  
277 CMPs and are dominantly distributed in the MEP differentiation branch. The distributions  
278 of these CMP sub-clusters are also consistent with the functions of their enriched motifs  
279 mentioned in the first section (Fig. 1d)<sup>22-24</sup>. Finally, we generated a hematopoiesis tree  
280 based on the APEC analysis (Fig. 3f).



281

282 **Figure 3. APEC constructed a differentiation pathway from scATAC-seq data from human**  
 283 **hematopoietic cells. (a)** The pseudotime trajectory construction scheme based on the accession  
 284 matrix and Monocle. **(b)** Principal component analysis (PCA) of the accession matrix for human  
 285 hematopoietic cells. The first principal component is not shown here because it was highly  
 286 correlated with sequencing depth <sup>21</sup>. HSC, hematopoietic stem cell; MPP, multipotent progenitor;  
 287 LMPP, lymphoid-primed multipotential progenitor; CMP, common myeloid progenitor; CLP,  
 288 common lymphoid progenitor; pDC, plasmacytoid dendritic cell; GMP, granulocyte-macrophage  
 289 progenitor; MEP, megakaryocyte-erythroid progenitor. **(c)** Pseudotime trajectory for the same data  
 290 constructed by applying Monocle on the accession matrix. Paths A and B represent different  
 291 pathways for GMP cell differentiation. **(d)** The deviations of significant differential motifs (HOXA9,  
 292 GATA1, CEBPB, and TCF4) plotted on the pseudotime trajectory. **(e)** Distributions of the CMP sub-  
 293 clusters on the trajectory. **(f)** Modified schematic of human hematopoietic differentiation.

294

295 Furthermore, we benchmarked the performance of APEC and of all the other tools  
296 in constructing a pseudotime trajectory from the scATAC-seq profile on the same dataset.  
297 We found that (1) when the raw peak count matrix was invoked into Monocle, almost none  
298 developmental pathways were constructed ([Supplementary Fig. 8a](#)), suggesting that the  
299 peak aggregation step in APEC greatly improves the pseudotime estimation; (2) APEC +  
300 Monocle provides the most precise pathways from HSCs to differentiated cells, compared  
301 to other methods, including cisTopic, SnapATAC, LSI, Cicero, and chromVAR  
302 ([Supplementary Fig. 8b-f](#)); and (3) when we applied other pseudotime trajectory  
303 construction methods, such as SPRING<sup>35</sup>, after APEC, a similar though less clear cell  
304 differentiation diagram was also obtained, suggesting the reliability of our prediction  
305 ([Supplementary Fig. 8g](#)).

306

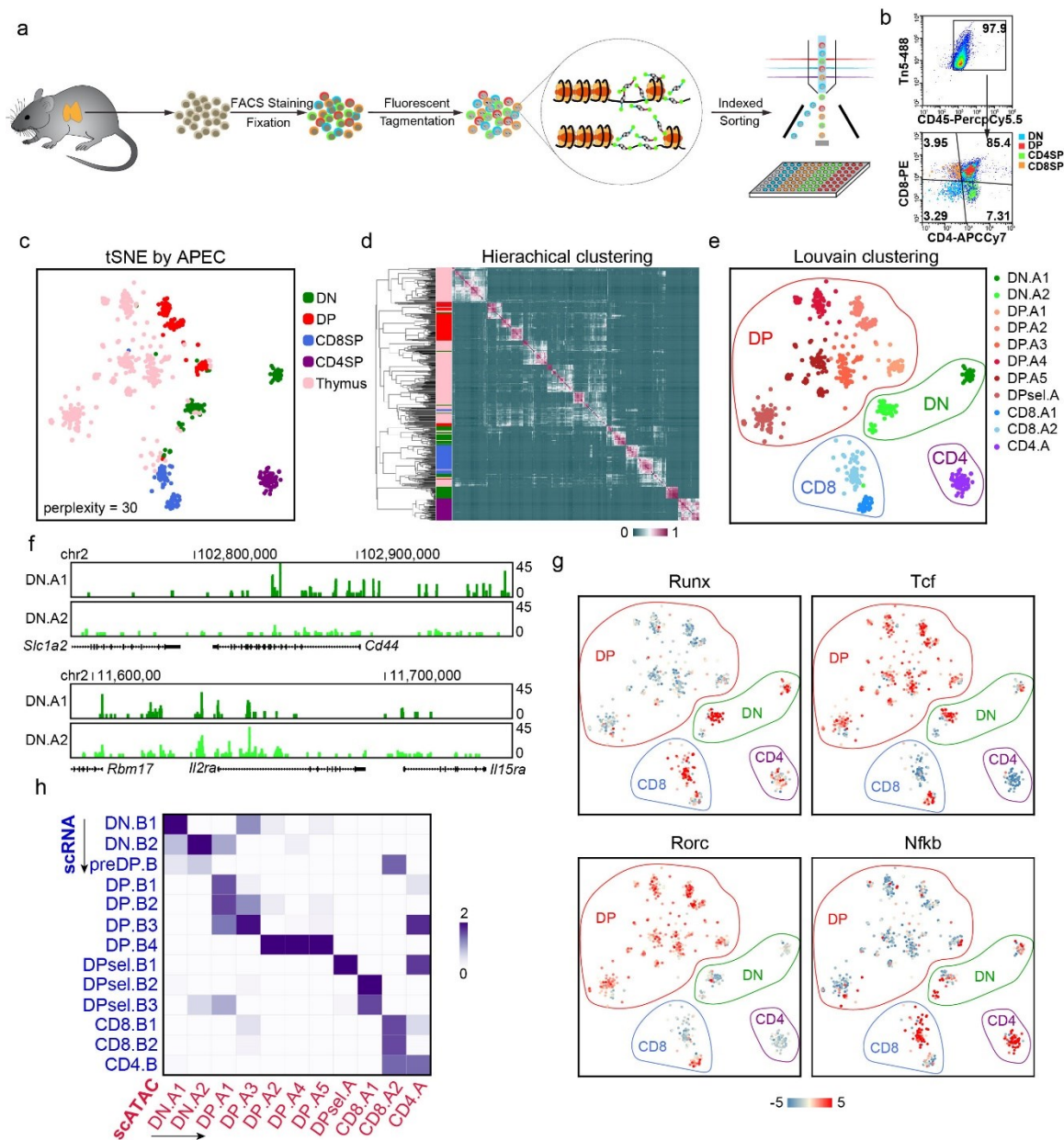
### 307 **APEC reveals the single-cell regulatory heterogeneity of thymocytes**

308

309 T cells generated in the thymus play a critical role in the adaptive immune system,  
310 and the development of thymocytes can be divided into 3 main stages based on the  
311 expression of the surface markers CD4 and CD8, namely, CD4 CD8 double-negative (DN),  
312 CD4 CD8 double-positive (DP) and CD4 or CD8 single-positive (CD4SP or CD8SP,  
313 respectively) stages<sup>42</sup>. However, due to technical limitations, our genome-wide  
314 understanding of thymocyte development at single-cell resolution remains unclear.  
315 Typically, more than 80% of thymocytes stay in the DP stage in the thymus, whereas DN  
316 cells account for only approximately 3% of the thymocyte population. To eliminate the  
317 impacts of great differences in proportion, we developed a fluorescent tagmentation- and  
318 FACS-sorting-based scATAC-seq strategy (ftATAC-seq), which combined the advantages  
319 of ATAC-seq<sup>43</sup> and Pi-ATAC-seq<sup>44</sup> to manipulate the desired number of target cells by  
320 indexed sorting ([Fig. 4a](#)). Tn5 transposomes were fluorescently labeled in each cell to  
321 evaluate the tagmentation efficiency so that cells with low ATAC signals could be gated  
322 out easily ([Fig. 4b](#), [Supplementary Fig. 9a](#)). With ftATAC-seq, we acquired high-quality  
323 chromatin accessibility data for 352 index-sorted DN, DP, CD4SP, and CD8SP single cells  
324 and 352 mixed thymocytes ([Supplementary Fig. 9b-d](#)). Correlation analysis with the  
325 published bulk ATAC-seq data of thymocytes<sup>45</sup> indicates that the cells we sorted in  
326 ftATAC-seq were correctly labeled ([Supplementary Fig. 9e](#)). We then applied APEC on  
327 this dataset to investigate the chromatin accessibility divergence during developmental

328 process and to reveal refined regulome heterogeneity of mouse thymocytes at single-cell  
329 resolution. Taking into account of all the 130685 peaks called from the raw sequencing  
330 data, APEC aggregated 600 accessions and successfully assigned over 82% of index-  
331 sorted DN, DP, CD4SP and CD8SP cells into the correct subpopulations (Fig. 4c & 4d).  
332 As expected, the majority of randomly sorted and mixed thymocytes were classified into  
333 DP subtypes based on hierarchical clustering of cell-cell correlation matrix, which was  
334 consistent with the cellular subtype proportions in the thymus. APEC further classified all  
335 thymocytes into 11 subpopulations, including 2 DN, 6 DP, 1 CD4SP, 2 CD8SP, suggesting  
336 that extensive epigenetic heterogeneity exists among cells with the same CD4 and CD8  
337 surface markers (Fig. 4e). For instance, there are four main subtypes of DN cells,  
338 according to the expression of the surface markers CD44 and CD25<sup>46</sup>, while two clusters  
339 were identified in ftATAC-seq. The accessibility signals around the *Ii2ra* (Cd25) and *Cd44*  
340 gene loci demonstrated that DN.A1 comprised CD44<sup>+</sup>CD25<sup>-</sup> and CD44<sup>+</sup>CD25<sup>+</sup> DN  
341 subtypes (DN1 and DN2), and DN.A2 cells comprised CD44<sup>-</sup>CD25<sup>+</sup> and CD44<sup>-</sup>CD25<sup>-</sup>  
342 subtypes (DN3 and DN4), suggesting significant chromatin changes between DN2 and  
343 DN3 cell development (Fig. 4f).

344



345

346 **Figure 4. APEC accurately identified cell subtypes based on scATAC-seq data from *Mus***  
 347 ***musculus* thymocytes. (a)** Experimental workflow of the fluorescent tagmentation- and FACS-  
 348 **sorting-based scATAC-seq strategy (ftATAC-seq). (b)** Indexed sorting of double-negative (DN),  
 349 **double-positive (DP), CD4<sup>+</sup> single-positive (CD4SP), and CD8<sup>+</sup> single-positive (CD8SP) cells with**  
 350 **strong tagmentation signals. (c)** The tSNE of thymocyte single-cell ftATAC-seq data based on the  
 351 **accession matrix, in which the cells are labeled by the sorting index. (d)** Hierarchical clustering of  
 352 **the cell-cell correlation matrix. On the sidebar, each cell was colored by the sorting index. (e)** The  
 353 **accession-based Louvain method clustered thymocytes into 11 subtypes. DN.A1 (dark green) & A2**  
 354 **(light green), double-negative clusters; DP.A1~A5 and DPsel.A, double-positive clusters; CD8.A1**  
 355 **(dark blue) & A2 (light blue), CD8<sup>+</sup> single-positive clusters; CD4.A (purple), CD4<sup>+</sup> single-positive**  
 356 **cluster. (f)** Average fragment counts of two DN clusters around the marker genes *Cd44* and *Il2ra*.  
 357 **(g)** Differential enrichment of the motifs Runx, Tcf, Rorc, and Nfkb in the cell clusters. (h) Z-score  
 358 **of the Fisher exact test -log(p-value) of the common differential genes between the cell clusters**

359 from different experiments. The row and column clusters were identified by data from single-cell  
360 transcriptome (SMART-seq) and chromatin accessibility (ftATAC-seq) analysis respectively.

361

362 Many TFs have been reported to be essential in regulating thymocyte development,  
363 and we found that their motifs were remarkably enriched at different stages during the  
364 process (Fig. 4g). For instance, Runx3 is well known for regulating CD8SP cells <sup>47</sup>, and  
365 we observed significant enrichment of the RUNX motif on DN cells and a group of CD8SP  
366 cells. Similarly, the TCF <sup>48,49</sup>, RORC <sup>50</sup> and NFkB <sup>51</sup> family in regulating the corresponding  
367 stages during this process. More enriched TF motifs in each cell subpopulation were also  
368 observed, suggesting significant regulatory divergence in thymocytes (Supplementary Fig.  
369 10). Interestingly, two clusters of CD8SP cells appear to be differentially regulated based  
370 on motif analysis, in which CD8.A1 cells are closer to DP cells, while CD8.A2 cells are  
371 more distant at the chromatin level, suggesting that CD8.A2 cells are more mature CD8SP  
372 cells, and CD8.A1 cells are in a transitional state between DP and SP cells.

373 APEC is capable of integrating single-cell transcriptional and epigenetic  
374 information by scoring gene sets of interest based on their nearby peaks from scATAC-  
375 seq, thereby converting the chromatin accessibility signals to values that are comparable  
376 to gene expression profiles (see Methods). To test the performance of this integrative  
377 analysis approach and to evaluate the accuracy of thymocyte classification by APEC, we  
378 assayed the transcriptomes of single thymocytes and obtained 357 high-quality scRNA-  
379 seq profiles using the SMART-seq2 protocol <sup>52</sup>. Unsupervised analysis of gene expression  
380 profiles clustered these thymocytes into 13 groups in Seurat <sup>13</sup> (Supplementary Fig. 11a  
381 & b), and each subpopulation was identified based on known feature genes  
382 (Supplementary Fig. 11c & d). We then adopted fisher's exact test on the shared  
383 differential genes in cell clusters identified from scATAC-seq and scRNA-seq profiles (see  
384 Methods), and observed a strong correlation between the subtypes identified from the  
385 transcriptome and those from chromatin accessibility (Fig. 4h), confirming the reliability  
386 and stability of cellular classification using APEC.

387

### 388 **APEC reconstructs the thymocyte developmental trajectory**

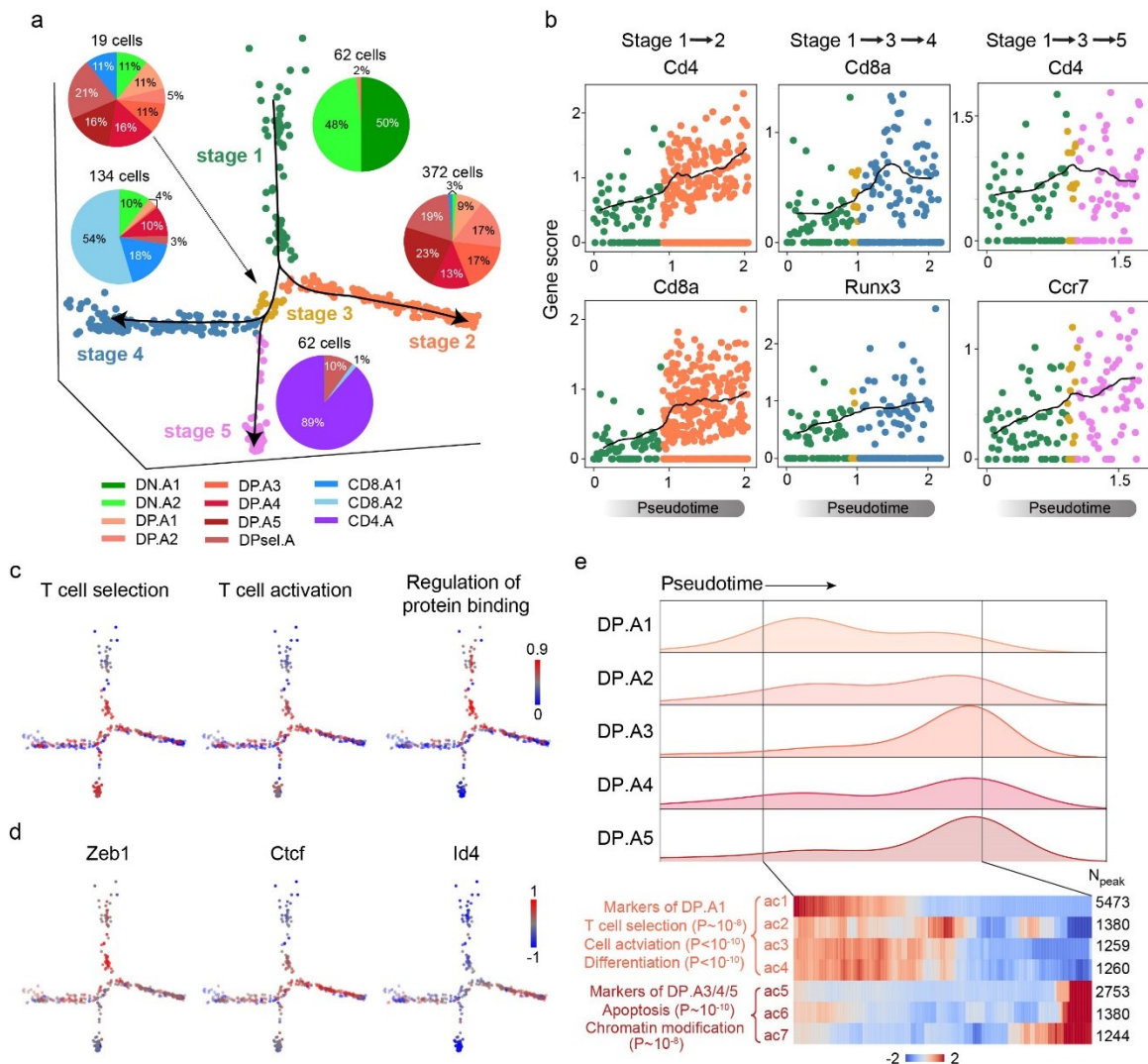
389

390 APEC is capable of constructing a pseudotime trajectory and then predicting the  
391 cell differentiation lineage from a "snapshot" of single-cell epigenomes (Fig. 3). We applied  
392 APEC to recapitulate the developmental trajectory and thereby reveal the single-cell  
393 regulatory dynamics during the maturation of thymocytes. Pseudotime analysis based on



394 single-cell ATAC-seq data shaped thymocytes into 5 developing stages (Fig. 5a,  
395 [Supplementary Fig. 12a & b](#)), where most of the cells in stages 1, 2, 4, and 5 were DN,  
396 DP, CD8SP and CD4SP cells, respectively. APEC also identified a transitional stage 3,  
397 which was mainly consisted of last stages of DP cells. Besides Monocle, a similar  
398 developmental pathways can also be constructed by SPRING<sup>35</sup> based on the accesson  
399 matrix ([Supplementary Fig. 12c](#)). Interestingly, the pseudotime trajectory suggests three  
400 developmental pathways for this process, one of which started with stage 1 (DN) and  
401 ended in stage 2 (DP), and the other two of which started with stage 1 (DN), went through  
402 a transitional stage 3 and then bifurcated into stage 4 (CD8SP) and 5 (CD4SP). The  
403 predicted developmental trajectory could also be confirmed by the gene expression of  
404 surface markers, such as Cd4, Cd8, Runx3 and Ccr7 (Fig. 5b). To evaluate the gene  
405 ontology (GO) enrichments over the entire process, we implemented an accesson-based  
406 GO module in APEC, which highlights the significance of the association between cells  
407 and biological function (Fig. 5c). For instance, T cells selections, including  $\beta$ -selection,  
408 positive selection and negative selection, are initiated in the late DN stage. Consistent with  
409 this process, we observed a strong “T cell selection” GO term on the trajectory path after  
410 DN.A1 ([Supplementary Fig. 12d](#)). Since TCR signals are essential for T cell selection, we  
411 also observed the “T cell activation” GO term accompanied by “T cell selection”.  
412 Meanwhile, the signal for regulation of protein binding was found decreased at SP stages,  
413 indicating the necessity of weak TCR signal for the survival of SP T cells during negative  
414 selection.

415



416

417 **Figure 5. APEC depicted the developmental pathways of *Mus musculus* thymocytes by**  
 418 **pseudotime analysis. (a)** Pseudotime trajectory based on the accession matrix of thymocyte  
 419 ftATAC-seq data. Cell colors were defined by the developmental stages along pseudotime. Pie  
 420 charts show the proportion of cell clusters at each stage. **(b)** APEC scores of important marker  
 421 genes (*Cd8a*, *Cd4*, *Runx3*, and *Ccr7*) along each branch of the pseudotime trajectory. **(c)**  
 422 Weighted scores of important functional GO terms along each branch of the pseudotime trajectory. **(d)**  
 423 Enrichment of specific motifs searched from the differential accessions of each cell subtype. **(e)**  
 424 On the stage 2 branch, the cell number distribution of clusters DP.A1~A5 along pseudotime (upper  
 425 panel) and the intensity of marker accessions of DP.A1 and DP.A3/4/5 (lower right panel) with top  
 426 enriched GO terms with significance (lower left panel).  
 427

428 To further uncover the regulatory mechanism underlying this developmental  
 429 process, APEC was implemented to identify stage-specific enriched TFs along the  
 430 trajectory and pinpoint the “pseudotime” at which the regulation occurs. In addition to the  
 431 well-studied TFs mentioned above (Fig. 4g), APEC also identified Zeb1<sup>53</sup>, Ctcf<sup>54</sup> and Id4

432 as potential stage-specific regulators (Fig. 5d). Interestingly, the Id4 motif enriched on DP  
433 cells was also reported to regulate apoptosis in other cell types<sup>55,56</sup>. Associated with the  
434 fact that the vast majority of DP thymocytes die because of a failure of positive selection  
435<sup>57</sup>, we hypothesize that stage 2 may be the path towards DP cell apoptosis. We then  
436 checked the distribution of DP cells along the stage 2 trajectory and found that most DP.A1  
437 cells were scattered in “early” stage 2, and they were enriched with GO terms such as “T  
438 cell selection”, “cell activation” and “differentiation” (Fig. 5e, Supplementary Fig. 12e).  
439 However, most DP.A3/4/5 cells were distributed at the end of stage 2, and their principle  
440 accessions were enriched with GO terms such as “apoptosis” and “chromatin modification”.  
441 Although it is believed that more than 95% of DP thymocytes die during positive selection,  
442 only a small proportion of apoptotic cells could be detected in a snapshot of the thymus,  
443 which in our data are the cells at the end of stage 2. By comparing the number of cells  
444 near stage 3 with all the cells in stage 2, we estimated that ~3-5% of cells would survive  
445 positive selection, which is consistent with the findings reported in previous publications  
446<sup>58,59</sup>. Our data suggest that before entering the final apoptotic stage, DP thymocytes under  
447 selection could have already been under apoptotic stress at the chromatin level, which  
448 explains why DP cells are more susceptible to apoptosis than other thymocyte subtypes  
449<sup>60</sup>.

450

## 451 **DISCUSSION**

452

453 Here, we introduced an accession-based algorithm for single-cell chromatin  
454 accessibility analysis. Without relying on any prior information (such as bulk sequencing  
455 data or known cell types), this approach generated more refined cell groups with reliable  
456 biological functions and properties. Integrating the new algorithm with all necessary  
457 chromatin sequencing data processing tools, APEC provides a comprehensive solution  
458 for transforming raw experimental single-cell data into final visualized results. In addition  
459 to improving the clustering of subtle cell subtypes, APEC is also capable of locating  
460 potential specific super-enhancers, searching enriched motifs, estimating gene activities,  
461 and constructing time-dependent cell developmental trajectories, and it is compatible with  
462 many existing single-cell accessibility datasets. Compared with all the other state-of-the-  
463 art single cell chromatin accessibility analysis methods, APEC clearly shows superiority in  
464 correctly predicting cell identities and precisely constructing developmental trajectories,  
465 and provides new biological insights. APEC is also very robust and stable and is scalable

466 to clustering a large number of cells using limited computational resources. Despite these  
467 advantages, the biological implications of accessions are still obscure, especially for those  
468 that involve only a small number of peaks. Although we noticed peaks in the same  
469 accession may belong to the same TADs, further investigations are still required to fully  
470 uncover the biology that underlies accessions.

471 To evaluate the performance of this approach in the context of the immune system,  
472 we adopted APEC with scATAC-seq technology to investigate the regulome dynamics of  
473 the thymic development process. Coordinated with essential cell surface markers, APEC  
474 provided a much more in-depth classification of thymocytes than the conventional DN, DP,  
475 CD4SP and CD8SP stages based on single-cell chromatin status. By reconstructing the  
476 developmental pseudotime trajectory, APEC discovered a transitional stage before  
477 thymocytes bifurcate into CD4SP and CD8SP cells and inferred that one of the stages  
478 leads to cell apoptosis. Considering that more than 95% of DP cells undergo apoptosis as  
479 a programmed cell death process, our data suggested that before DP cells enter the final  
480 apoptotic state, there would already be some intracellular changes towards apoptosis at  
481 the chromatin level. However, further studies are still needed to fully understand the  
482 regulatory mechanism of this process.

483

484

## 485 **METHODS**

486

### 487 **ftATAC-seq on mouse thymocytes**

488 Alexa fluor 488-labeled adaptor oligonucleotides were synthesized at Sangon Biotech as  
489 follows: Tn5ME, 5'-[phos]CTGTCTCTTATACACATCT-3'; AF488-R1, 5'-AF488-  
490 TCGTCGGCAGCGTCAGATGTGTATAAGAGACAG-3'; and AF488-R2, 5'-AF488-  
491 GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAG-3'. Then, 50  $\mu$ M of AF488-  
492 R1/Tn5ME and AF488-R2/Tn5ME were denatured separately in TE buffer (Qiagen) at  
493 95 °C for 5 min and cooled down to 22 °C at 0.1 °C/s. AF488-labeled adaptors were  
494 assembled onto Robust Tn5 transposase (Robustnique) according to the user manual to  
495 form fluorescent transposomes.

496 Thymus tissues isolated from 6- to 8-week-old male mice were gently ground in 1 mL of  
497 RPMI-1640. Thymocytes in a single-cell suspension were counted after passing through  
498 a 40  $\mu$ m nylon mesh. A total of  $1 \times 10^6$  thymocytes were stained with PerCP-Cy5.5-anti-  
499 CD45, PE-anti-CD8a and APC-Cy7-anti-CD4 antibodies (Biolegend) and then fixed in  $1 \times$   
500 PBS containing 1% methanal at room temperature for 5 min. After washing twice with  $1 \times$   
501 PBS, the cells were counted again. A total of  $1 \times 10^5$  fixed cells were resuspended in 40  
502  $\mu$ L of  $1 \times$  TD buffer (5 mM Tris-HCl, pH 8.0, 5 mM  $MgCl_2$ , and 10% DMF) containing 0.1%  
503 NP-40. Then, 10  $\mu$ L of fluorescent transposomes were added and mixed gently.  
504 Fluorescent tagmentation was conducted at 55 °C for 30 min and stopped by adding 200  
505  $\mu$ L of 100 mM EDTA directly to the reaction mixture. The cells were loaded on a Sony  
506 SH800S sorter, and single cells of the CD45<sup>+</sup>/AF488-Tn5<sup>hi</sup> population were index-sorted  
507 into each well of 384-well plates. The 384-well plates used to acquire sorted cells were  
508 loaded with 2  $\mu$ L of release buffer (50 mM EDTA, 0.02% SDS) before use. After sorting,  
509 the cells in the wells were incubated for 1 min. Plates that were not processed immediately  
510 were preserved at -80 °C.

511 To prepare a single-cell ATAC-seq library, plates containing fluorescently tagmented cells  
512 were incubated at 55 °C for 30 min. Then, 4.2  $\mu$ L of PCR round 1 buffer (1  $\mu$ L of 100  $\mu$ M  
513  $MgCl_2$ , 3  $\mu$ L of  $2 \times$  I-5 PCR mix [MCLAB], and 0.1  $\mu$ L each of 10  $\mu$ M R1 and R2 primers)  
514 were added to each well, followed by PCR: 72 °C for 10 min; 98 °C for 3 min; 10 cycles of  
515 98 °C for 10 s, 63 °C for 30 s and 72 °C for 1 min; 72 °C for 3 min; and holding at 4 °C.  
516 Thereafter, each well received 4  $\mu$ L of PCR round 2 buffer (2  $\mu$ L of I-5 PCR Mix, 0.5  $\mu$ L  
517 each of Ad1 and barcoded Ad2 primers, and 1  $\mu$ L of ddH<sub>2</sub>O), and final PCR amplification  
518 was carried out: 98 °C for 3 min; 12 cycles of 98 °C for 10 s, 63 °C for 30 s and 72 °C for

519 1 min; 72 °C for 3 min; and holding at 4 °C. Wells containing different Ad2 barcodes were  
520 collected together and purified with a QIAquick PCR purification kit (Qiagen). Libraries  
521 were sequenced on an Illumina HiSeq X Ten system.

522

### 523 **SMART-seq on thymocytes**

524 Thymocytes were stained and sorted directly into 384-well plates without fixation. SMART-  
525 seq was performed as described with some modifications <sup>61</sup>. Reverse transcription and  
526 the template-switch reaction were performed at 50 °C for 1 hr with Maxima H Minus  
527 Reverse Transcriptase (Thermo Fisher); for library construction, 0.5-1 ng of cDNA was  
528 fragmented with 0.05 µL of Robust Tn5 transposome in 20 µL of TD buffer at 55 °C for 10  
529 min, then purified with 0.8× VAHTS DNA Clean Beads (Vazyme Biotech), followed by PCR  
530 amplification with Ad1 and barcoded Ad2 primers and purification with 0.6× VAHTS DNA  
531 Clean Beads. Libraries were sequenced on an Illumina HiSeq X Ten system.

532

### 533 **Data source**

534 All experimental raw data used in this paper are available online. The single-cell data for  
535 mouse thymocytes captured by the ftATAC-seq experiment can be obtained from the  
536 Genome Sequence Archive at BIG Data Center with the accession number CRA001267  
537 and is available via <http://bigd.big.ac.cn/gsa/s/yp1164Et>. Other published data sets used  
538 in this study are available from NIH GEO: (1) scATAC-seq data for LSCs and leukemic  
539 blast cells from patients SU070 and SU353, LMPP cells, and monocytes from GSE74310  
540 <sup>2</sup>; (2) scATAC-seq data for HL-60 cells from GSE65360 <sup>8</sup>; and (3) scATAC-seq data for  
541 hematopoietic development (HSCs, MPPs, CMPs, LMPPs, GMPs, EMPs, CLPs and  
542 pDCs) from GSE96772 <sup>21</sup>. (4) APEC is also compatible with a preprocessed fragment  
543 count matrix from the snATAC-seq data for the forebrain of adult mice (p56) from  
544 GSE100033 <sup>10</sup>. (5) The computational efficiency of APEC and other methods was tested  
545 using data from the single-cell atlas of mouse chromatin accessibility (sciATAC-seq) from  
546 GSE111586 <sup>28</sup>. (6) The scATAC-seq (GSE65360 <sup>8</sup>) and Hi-C (GSE63525 <sup>62</sup>) data of  
547 GM12878 cells were used to generate the spatial correlation of peaks in the same or in  
548 different accessions.

549

### 550 **Preparing the fragment count matrix from the raw data**

551 APEC adopted the general mapping, alignment, peak calling and motif searching  
552 procedures to process the scATAC-seq data from ATAC-pipe <sup>63</sup>. We also implemented

553 the python script in ATAC-pipe<sup>63</sup> to trim the adapters in the raw data (in paired-end fastq  
554 format files for each single-cell sample). APEC used BOWTIE2 to map the trimmed  
555 sequencing data to the corresponding genome index and used PICARD for the sorting,  
556 duplicate removal, and fragment length counting of the aligned data.

557 Unlike several previous methods that call peaks from bulk ATAC-seq data or  
558 aggregated cell populations according to cell type<sup>20,21,64</sup>, the APEC pipeline calls peaks  
559 from the merged single-cell profiles of all cells using MACS2 to ensure that the entire  
560 analysis is unsupervised. We then ranked and filtered out the low quality peaks based on  
561 the false discovery rate (Q-value). Genomic locations of the peaks were annotated by  
562 HOMER, and motifs searched by FIMO. APEC calculates the number of fragments and  
563 the percent of reads mapped to the TSS region ( $\pm 2000$  BP) for each cell, and filters out  
564 high quality cells for downstream analysis. All required files for the hg19 and mm10  
565 assembly have been integrated into the pipeline. If users want to process data from other  
566 species, they can also download corresponding reference files from the UCSC website.  
567 By combining existing tools, APEC made it possible to finish all of the above data  
568 processing steps by one command line, and generate a fragment count matrix for  
569 subsequent cell clustering and differential analysis. APEC has been made available on  
570 GitHub (<https://github.com/QuKunLab/APEC>).

571

## 572 **Accession-based clustering algorithm**

573 We define accession as a set of peaks with similar accessibility patterns across all single  
574 cells, similar to the definition of gene module for RNA-seq data. After preprocessing, a  
575 filtered fragment count matrix  $\mathbf{B}$  is obtained, and APEC groups peaks to construct  
576 accessions and then performs cell clustering analysis as follows:

577 (1) *Normalization of the fragment count matrix.* Each matrix element  $B_{ij}$  represents the  
578 number of raw reads in cell  $i$  and peak  $j$ , and element  $B_{ij}$  was then normalized by the  
579 total number of reads in each cell  $i$ , as if there are 10,000 reads in each cell.

$$580 \quad B'_{ij} = \log_2 \left( \frac{B_{ij} \times 10000}{\sum_{j'} B_{ij'}} + 1 \right)$$

581 (2) *Constructing accessions.* The top 40 principal components of the normalized matrix  
582 were used to construct the connectivity matrix ( $\mathbf{C}_{\text{peak}}$ ) of peaks by the K-nearest-  
583 neighbor (KNN) method with  $K=10$ . The grouping of peaks is insensitive to the  
584 number of principal components and the number of nearest neighbors, so it is usually  
585 not necessary to change these two parameters for different data sets. Based on the

586 matrix  $\mathbf{C}_{\text{peak}}$ , all peaks were grouped by agglomerative clustering with Euclidean  
587 distance and Ward linkage method, and the sum of one peak group was an accesson.  
588 For most data sets, we recommend setting the number of accessons to a value  
589 between 500 and 1500, and the default was set to 600, however the cell clustering  
590 result is not sensitive to the choice of accesson number within this range. We then  
591 built the accesson count matrix  $\mathbf{M}$  by summing the fragment count of all peaks in one  
592 accesson. Thus, each column of matrix  $\mathbf{M}$  is an accesson, each row is a cell, and  
593 each element represents the cumulative fragment count of each accesson in each  
594 cell. The accesson matrix was then normalized by calculating the z-score or  
595 probability of the fragment count for each row (i.e., each cell) to generate normalized  
596 matrix  $\mathbf{M}_a$  for the next step of cell clustering.

597 (3) *Cell clustering.* From the normalized accesson matrix  $\mathbf{M}_a$ , APEC established the  
598 connectivity matrix by computing the k-neighbor graph of all cells. Since the Louvain  
599 algorithm was proven to be a reliable single-cell clustering method in Seurat<sup>13</sup> and  
600 Scanpy<sup>16</sup>, we adopted it in APEC to automatically predict the number of clusters from  
601 the connectivity matrix and defined each Louvain community as a cell cluster. APEC  
602 uses the Louvain algorithm to predict cluster number to ensure that the clustering  
603 analysis is unsupervised and then performs cell clustering as default. Meanwhile, if  
604 users want to artificially define the number of cell clusters, APEC can also perform  
605 KNN clustering on the connectivity matrix.

606 (4) *Compare the performance of APEC with that of other methods on cells with known*  
607 *identity.* To investigate the accuracy of the cell clusters predicted by different  
608 algorithms, we used the ARI value, which evaluates the similarity of clustering results  
609 with all known types of cells<sup>14</sup>. The ARI value can be calculated as follows:

$$610 \quad ARI = \frac{\sum_{ij} \binom{n_{ij}}{2} - \sum_i \binom{a_i}{2} \cdot \sum_j \binom{b_j}{2} / \binom{n}{2}}{\left[ \sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right] / 2 - \sum_i \binom{a_i}{2} \cdot \sum_j \binom{b_j}{2} / \binom{n}{2}}$$

611 where  $n_{ij}$  is the element from the contingency matrix (i.e. the number of type  $i$  cells  
612 that were classified into cluster  $j$ ),  $a_i$  and  $b_j$  are the sums of the  $i$ th row and  $j$ th  
613 column, respectively, and  $\binom{x}{y}$  denotes a binomial coefficient. A higher ARI value  
614 indicates more accurate classification of cell types.

615 (5) *Characteristics of accesson.* The peaks of a same accesson can be distant from each  
616 other on the genome, and sometimes even on multiple chromosomes. The average



617 number of peaks per accession depends on the total number of peaks in the dataset  
618 and the number of accessions set in the program (default 600). Usually the total  
619 number of peaks can vary between ~40,000-150,000 depending on the total number  
620 of cells and the sequencing depth for each cell, thereby the average number of peaks  
621 per accession is around ~60-250. Beside, we chose the top 40 principle components  
622 (PCs) of the normalized matrix to construct the connectivity matrix since the first 3~5  
623 PCs are usually not sufficient to capture the detailed features of a single cell dataset,  
624 as described in Seurat and many other single cell analysis tools. Although default  
625 values were chosen to provide better clustering results based on analysis of multiple  
626 datasets, users can adjust these parameters as needed.

627

### 628 **Sampling of accession number**

629 To test if the APEC clustering result is sensitive to the choice of accession numbers, we  
630 sampled 100 different accession numbers from 500 to 1500 in steps of 10 and clustered  
631 the cells of each dataset 100 times (Fig. 1c and Supplementary Fig. 2c). APEC generated  
632 stable clustering results in terms of the average ARI on these datasets, with a wide range  
633 of different accession numbers (Supplementary Fig. 5e). We used 600 as the default  
634 number of accessions in APEC.

635

### 636 **Parameter settings for other algorithms**

637 To quantify the cell clustering performance of APEC, we compared APEC with other state-  
638 of-the-art single-cell epigenomic algorithms on the same datasets with gold standards,  
639 including cisTopic<sup>19</sup>, LSI<sup>11,12</sup>, chromVAR<sup>17</sup> and Cicero<sup>18</sup>. Since most of them have no  
640 cell clustering algorithm within their original codes, we applied the Louvain clustering  
641 algorithm on their transformed matrices to fairly compare their performance. We adopted  
642 the default settings of these tools for most of the comparisons in this paper, except for  
643 some parameters that were manually defined as necessary, such as the random seed in  
644 cisTopic, the number of top components in LSI, and the peak aggregation distance in  
645 Cicero. Therefore, we sampled these parameters multiple times to obtain the average ARI  
646 and ratio of correctly classified cells of the clustering results for each tool (Fig. 1c and  
647 Supplementary Fig. 2c), just as we sampled the accession number for APEC. We set the  
648 same parameters for all the datasets as follows:

649 (1) *cisTopic*. The scanning range of the topic number was set to [10, 40], the number  
650 of parallel CPUs was set to 5, and the random seed was sampled 100 times from

651 100 to 600 in steps of 5. We kept the topic matrices normalized by z-score and  
652 probability and provided the performances based on both normalization methods.  
653 We then applied the Louvain algorithm as we did in APEC to cluster cells from the  
654 normalized topic matrix generated by cisTopic.

655 (2) *SnapATAC*. As SnapATAC uses a mapping procedure totally different with other  
656 tools, we adopt it with default parameters (binsize=5k) to build its own fragment  
657 count matrix from the merged bam file. We called the function “runJDA” with  
658 “bin.cov.zscore.lower=-2, bin.cov.zscore.upper=2, pc.num=50, norm.method  
659 =‘normOVE’, max.var=5000, do.par=TRUE, ncell.chunk=1000, num.cores=1,  
660 seed.use=10, tmp.folder=tempdir()”. We then called the function “runKNN” and  
661 sampled the number of principal components from 20 to 40 (step-size=2), and the  
662 number of nearest neighbors from 10 to 20 (step-size=1). Therefore, we sampled  
663 a total of 100 times (10×10) to calculate the ARI values. Finally, we called the  
664 function “runCluster” with “louvain.lib=‘leiden’, seed.use=10, resolution=1”.

665 (3) *LSI*. We performed truncated SVD (singular value decomposition) analysis on the  
666 TF-IDF (term frequency-inverse document frequency) matrix and chose  $N_{SVD}$  top  
667 components to generate the LSI matrix.  $N_{SVD}$  was sampled 6 times from 6 to 11.  
668 The first component was ignored since it is always related to read depth, and the  
669 LSI scores were capped at  $\pm 1.5$ . Then, we used the Louvain algorithm to cluster  
670 the cells of the LSI-processed matrix.

671 (4) *chromVAR*. The number of background iterations was set to 50, and the number  
672 of parallel CPUs was set to 1. We then used the Louvain algorithm to cluster cells  
673 based on the bias corrected deviation matrix generated by chromVAR.

674 (5) *Cicero*. The genome window was set to 500k BPs, the normalization method was  
675 set to “log”, the number of sample regions was set to 100, the number of  
676 dimensions was set to 40, and the peak aggregation distance was sampled at 20  
677 values from 1k to 20k BPs in steps of 1k BPs. Then, we used the Louvain algorithm  
678 to cluster cells based on the aggregated model matrix generated by Cicero.

679 To test the robustness of each algorithm, we randomly sampled 20%~90% of the raw  
680 sequence reads from the dataset of AML cells and 3 cell lines (LMPP, HL60 and monocyte)  
681 and calculated the ARI accordingly. This random sampling experiment was performed 50  
682 times for each method and average ARIs were reported ([Supplementary Fig. 2d](#)). The  
683 manually defined parameters for each method were set to: APEC, 600 accessions;

684 cisTopic, random seed 100; SnapATAC, 20 PCs and 15 nearest neighbors; LSI, top 2-6  
685 principle components; Cicero, 10k BPs aggregation distance.

686

### 687 **Gene scores and differential analysis**

688 APEC scores a gene by the peaks around its TSS region, which is similar to the algorithm  
689 used by Preissl et al. <sup>10</sup>. We calculate the average read counts of all peaks around a gene's  
690 TSS ( $\pm 20000$  BP by default) as its raw score ( $S_{ij}$  for cell  $i$  and gene  $j$ ), then define the  
691 gene expression by normalizing the raw score by ( $S'_{ij} = S_{ij} * 10000 / \sum_i S_{ij}$ ), making it in a  
692 range comparable to the gene expression from scRNA-seq data. After scoring genes,  
693 APEC uses the Student's t-test to estimate the significance of each genes between cell  
694 clusters and therefore obtained a list of differential genes filtered by P-values and fold  
695 changes. We also tested the gene scoring algorithms of other tools such as Cicero,  
696 cisTopic, and SnapATAC, and the parameters set in those algorithms were the same as  
697 described in the previous section. For scATAC-seq with very sparse signals, APEC first  
698 search for differential accessions between cell clusters, and extracted all the peaks in the  
699 differential accessions. APEC then defines genes close to these peaks ( $\pm 20000$  BP around  
700 TSS) as the differential genes.

701

### 702 **Association of cell clusters from scATAC-seq and scRNA-seq data**

703 To determine the association between cell clusters from the epigenomics and  
704 transcriptomic sequencing, we calculated the P-values of Fisher's exact test of the  
705 differential/non-differential genes between each pair of cell clusters from scATAC-seq and  
706 scRNA-seq data. For example, for cell cluster  $a$  from ftATAC-seq and cell cluster  $b$  from  
707 SMART-seq (Fig. 4h), if the number of consensus differential genes in both cluster  $a$  and  
708  $b$  is  $G_{11}$ , and the number of differential genes in either cluster  $a$  (or cluster  $b$ ) is  $G_{12}$  (or  
709  $G_{21}$ ), and the number of all the other genes is  $G_{22}$ , then the  $2 \times 2$  matrix  $\mathbf{G}$  can be directly  
710 used for Fisher's exact test to evaluate the P-value  $A_{ab}$  between cluster  $a$  and  $b$ . After  
711 constructing a matrix  $\mathbf{A}$  filled with  $-\log(A_{ab})$  for ftATAC-seq cluster  $a$  and SMART-seq  
712 cluster  $b$ , we calculated the z-score for each row and column of  $\mathbf{A}$  to determine the  
713 correlation between cell clusters from different sequencing experiments. Same algorithm  
714 was also used to show the correspondence between the cell sub-clusters (EX1~5 and  
715 IN1~5) generated from snATAC-seq data and the cell subtypes identified from inDrops-  
716 seq data (Fig. 2f & g).

717

## 718 **Potential super-enhancers**

719 Here, we defined a super-enhancer as a long continuous genomic area containing many  
720 accessible regions and have the same accessibility pattern in different cells. The  
721 accession-based algorithm can group most peaks in one super-enhancer to one accession  
722 since they always present the same accessibility pattern between cells. APEC identified  
723 super-enhancers by counting the number of peaks in a 1 million BP genomic area that  
724 belong to a same accession. It also requires that the percentage of the putative peaks in  
725 one super-enhancer is one of the highest among all genomic areas (P-value<0.01). The  
726 pipeline can also aggregate bam files by cell types/clusters and convert them to BigWig  
727 format for users to upload to the UCSC genome browser for visualization. ROSE <sup>27</sup> was  
728 used to confirm the super enhancers called from APEC, with command line “python  
729 ROSE\_main.py -g hg19 -i top\_filtered\_peaks.gff -r P1-LSC.bam -s 12500 -t 2500 -o P1-  
730 LSC-SuperEnhancer” applied to the merged bam files of all the P1-LSC cells.

731

## 732 **Spatial correlation of peaks in the same accession**

733 To test whether peaks in the same accession are closer in space, we integrated the Hi-C  
734 <sup>62</sup> data (GSE63525) and scATAC-seq <sup>8</sup> data (GSE65360) on GM12878 cells. The spatial  
735 correlation of different windows, both intra- and inter-chromosomal, can be directly  
736 extracted by Juicer <sup>65</sup>. The Pearson's correlation matrix of the intra-chromosomal or inter-  
737 chromosomal windows can be calculated from the corresponding observed/expected  
738 matrix. We constructed 600 accessions by grouping peaks in the GM12878 scATAC-seq  
739 data in APEC, and removed the accessions that contained more than 1000 peaks or less  
740 than 5 peaks. The width of the window was set to 500k BPs, and peaks were then  
741 assigned to each window. Next, we collected Hi-C correlations between windows that  
742 contained peaks in the same accession, termed “Accession” correlations. For comparison,  
743 we also shuffled all peaks in different accessions to make fake accessions and re-collected  
744 the Hi-C correlations between windows that contained peaks in each fake accession,  
745 termed “Shuffled” correlations. Meanwhile, we also collected Hi-C correlations between  
746 windows that contained no peaks, termed “Non-accession” correlations. We made  
747 boxplots for these three types of correlations for intra-chromosomal and inter-  
748 chromosomal peaks and found that co-accessible peaks are spatially closer to each other  
749 than random peaks.

750

## 751 **Pseudotime trajectory constructed by APEC**

752 As a tool to simulate the time-dependent variation of gene expression and the cell  
753 development pathway, Monocle has been widely used for the analysis of single-cell RNA-  
754 seq experiments<sup>32,66</sup>. APEC reduced the dimension of the accession count matrix **M** by  
755 PCA, and then performed pseudotime analysis using the Monocle program. For complex  
756 datasets, it is necessary to limit the number of principal components, since too many  
757 features will cause too many branches on the pseudotime trajectory, and makes it difficult  
758 for a user to identify the biological significance of each branch. For the hematopoietic  
759 single cell data and thymocyte data, we used the top 5 principal components of the  
760 accession matrix to construct the developmental and differentiation trajectories.

761

### 762 **Pseudotime trajectory constructed by other algorithms**

763 To check whether other algorithms can provide solutions to construct cell developmental  
764 pathways, we combined their transformed count matrix with Monocle to build the  
765 pseudotime trajectory from scATAC-seq data. A similar preprocessing method was  
766 applied to ensure the fairness of the comparison:

- 767 (1) Raw fragment count matrix. We normalized the raw count matrix **B** exactly as in  
768 the first step of APEC, i.e.,  $B'_{ij} = \log_2 \left( \frac{B_{ij} \times 10000}{\sum_{j'} B_{ij'}} + 1 \right)$ , and performed PCA analysis  
769 on the normalized matrix **B'**. Only the top 5 PCs were subjected to Monocle to  
770 construct the trajectory.
- 771 (2) cisTopic. The topic matrix generated by cisTopic was normalized by making the  
772 sum of each row the same (i.e., the probability). Then, we performed PCA analysis  
773 on the normalized topic matrix and subjected the top 5 PCs to Monocle to build the  
774 trajectory.
- 775 (3) SnapATAC. We run PCA analysis on the normalized fragment count matrix  
776 generated by SnapATAC, and subjected the top 5 PCs to Monocle to build the  
777 trajectory.
- 778 (4) LSI. We chose the 2nd~6th principle components of the SVD transformation of  
779 the LSI matrix and subjected them to Monocle to construct the trajectory. As the  
780 dimensions had been reduced by the LSI method, we skipped the PCA analysis.
- 781 (5) ChromVAR. After the PCA analysis of the bias corrected deviation matrix  
782 generated by chromVAR, the top 5 PCs were combined with Monocle to construct  
783 the trajectory.

784 (6) Cicero. We performed PCA analysis on the aggregated matrix generated by  
785 Cicero and used the top 5 PCs in Monocle to build the trajectory.

786 In addition, to confirm the reliability of the APEC + Monocle prediction of the  
787 developmental pathway, we applied another pseudotime trajectory constructing method,  
788 SPRING<sup>35</sup>, to the accession count matrix **M** from APEC to reconstruct the pathways for  
789 the hematopoietic differentiation dataset and thymocyte developmental dataset. We  
790 performed PCA analysis of the accession matrix **M** and subjected the top 5 PCs to SPRING  
791 to generate the trajectories. The number of edges per node in SPRING was set to 5.

792

### 793 **Parameter settings for each dataset**

794 In the quality control (QC) step, cells are filtered by two constraints: the percentage of the  
795 fragments in peaks ( $P_f$ ) and the total number of valid fragments ( $N_f$ ). However, there is no  
796 fixed cutoff for these two parameters since the quality of different cell types and/or  
797 experiment batches are completely different. The total number of peaks is usually limited  
798 to approximately 50000 to reduce computer time, but we recommend using all peaks if the  
799 users want to obtain better cell clusters. (1) For the data set from hematopoietic cells, the  
800  $-\log(Q\text{-value})$  threshold of high-quality peaks was set to 35 to retain 54212 peaks, and the  
801 cutoff values of  $P_f$  and  $N_f$  were 0.1 and 1000, respectively. (2) For the scATAC-seq data  
802 on the two types of cells from 2 AML patients (P1-LSC, P1-Blast, P2-LSC, P2-Blast), the  
803 threshold of  $-\log(Q\text{-value})$  was set to 5 to retain 38683 high-quality peaks for subsequent  
804 processing. When LMPPs, HL60 and monocytes were added to this dataset with the AML  
805 cells, the threshold of  $-\log(Q\text{-value})$  was set to 8 to retain 42139 peaks. In the QC step,  
806 we set the  $P_f$  cutoff to 0.05 and the  $N_f$  cutoff to 800. (3) For the snATAC-seq data from  
807 the adult mouse forebrain, all peaks and the raw count matrix obtained from the original  
808 data source were adopted in the analysis. (4) For the ftATAC-seq data from thymocytes,  
809 all 130685 peaks called by MACS2 were reserved for the fragment count matrix ( $Q\text{-}$   
810  $\text{value} < 0.05$ ), and we retained cells with  $P_f > 0.2$  and  $N_f > 2000$ .

811

### 812 **SMART-seq data analysis with Seurat**

813 For the analysis of SMART-seq data from mouse thymocytes, we employed STAR  
814 (version 2.5.2a) with the ratio of mismatches to mapped length  
815 (`outFilterMismatchNoverLmax`) less than or equal to 0.05, translated output alignments  
816 into transcript coordinates (i.e., `quantMode TranscriptomeSAM`) for mapping<sup>67</sup> and used

817 RSEM<sup>68</sup> to calculate the TPM of genes. For QC, we excluded cells in which fewer than  
818 2000 genes were detected and genes that were expressed in only 3 or fewer cells. Seurat  
819 filtered cells with several specific parameters to limit the number of genes detected in each  
820 cell to 2000~6000 and the proportion of mitochondrial genes in each cell was set to less  
821 than 0.4 (i.e., low.thresholds=c(2000, Inf), high.thresholds=c(6000, 0.4)). Additionally, the  
822 top 12 principal components were used for dimension reduction with a resolution of 3.2  
823 (dims.use =1:12, resolution=3.2), followed by cell clustering and differential expressed  
824 gene analysis<sup>69</sup>.

825

### 826 **GO term analysis of cells along pseudotime trajectory**

827 We defined the functional characteristics of each accession by the GO terms and motifs  
828 enriched on its peaks. The GO terms of an accession were obtained by submitting all of  
829 its peaks to the GREAT website<sup>70</sup>. The negative logarithm of the P-value of each GO term  
830 in each accession was filled into a (GO terms) × (accessions) matrix **L**. The significance of  
831 each GO term on each cell was evaluated by the product of the matrix **L** and the accession  
832 count matrix **M**, i.e.

$$833 \quad GO_{ij} = \sum_k L_{ik} \cdot M_{kj}$$

834 where *i* is the *i*th GO term, *j* is the *j*th cell, and *k* is the *k*th accession. Then we calculated  
835 the z-score for each row of this product matrix, and plotted the z-score as the GO-term  
836 score on the trajectory diagram.

837

### 838 **Motif enrichment of cells along pseudotime trajectory**

839 To assess the motif enrichment of the accessions, we used the Centrimo tool of the MEME  
840 suite<sup>71</sup> to search for the enriched motifs for the peaks of each accession and applied the  
841 same algorithm as to the GO term score to obtain the motif score. The negative logarithm  
842 of the E-value (product of adjusted P-value and motif number)<sup>71</sup> of each motif in each  
843 accession was used to construct a (motifs) × (accessions) matrix **F**. The enrichment of each  
844 motif on each cell was evaluated by the product of the matrix **F** and the accession count  
845 matrix **M**, i.e.,

$$846 \quad Motif_{ij} = \sum_k F_{ik} \cdot M_{kj}$$

847 where  $i$  is the  $i$ th motif,  $j$  is the  $j$ th cell, and  $k$  is the  $k$ th accession. Then, we calculated  
848 the z-score for each row of this product matrix and plotted the z-score as the motif score  
849 on the trajectory diagram.

850

851

## 852 **DATA AND CODE AVAILABILITY**

853 Mouse thymocytes ftATAC-seq data can be obtained from the Genome Sequence Archive  
854 at BIG Data Center with the accession number CRA001267 and is available via  
855 <http://bigd.big.ac.cn/gsa/s/yp1164Et>. Other published data sets used in this study are  
856 available from NIH GEO with accession numbers GSE74310<sup>2</sup>, GSE65360<sup>8</sup>, GSE96772  
857<sup>21</sup>, GSE100033<sup>10</sup>, GSE111586<sup>28</sup>, and GSE63525<sup>62</sup>. APEC pipeline can be downloaded  
858 from the GitHub website (<https://github.com/QuKunLab/APEC>).

859

## 860 **ACKNOWLEDGMENTS**

861

862 This work was supported by the National Key R&D Program of China (2017YFA0102900  
863 to K.Q.) and by National Natural Science Foundation of China grants (81788101,  
864 91640113, 91940306, 31771428 to K.Q.). It was also supported by the Fundamental  
865 Research Funds for the Central Universities (to K.Q.) and the Anhui Provincial Natural  
866 Science Foundation grant BJ2070000097 (to B.L.) and 1908085QH326 (to Y.L.). We  
867 thank the Howard Chang lab at Stanford University for helpful discussion. We thank the  
868 USTC supercomputing center and the School of Life Science Bioinformatics Center for  
869 providing supercomputing resources for this project.

870

## 871 **AUTHOR CONTRIBUTIONS**

872

873 KQ, BL and YL conceived the project, BL developed the APEC software and performed  
874 all data analysis with helps from KL, QY, PC, JF, WZ, PD and CJ. YL developed ftATAC-  
875 seq technique and performed all scATAC-seq and scRNA-seq experiments with helps  
876 from LZ. KL analyzed scRNA-seq data. BL, YL and KQ wrote the manuscript with inputs  
877 from all other authors.

878

## 879 **DECLARATION OF INTERESTS**



880

881 The authors declare no competing interests.

882

## 883 REFERENCES

- 884 1 Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. & Greenleaf, W. J. Transposition of  
885 native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-  
886 binding proteins and nucleosome position. *Nat Methods* **10**, 1213-1218,  
887 doi:10.1038/nmeth.2688 (2013).
- 888 2 Corces, M. R. *et al.* Lineage-specific and single-cell chromatin accessibility charts human  
889 hematopoiesis and leukemia evolution. *Nat Genet* **48**, 1193-1203, doi:10.1038/ng.3646  
890 (2016).
- 891 3 Wu, J. *et al.* The landscape of accessible chromatin in mammalian preimplantation  
892 embryos. *Nature* **534**, 652-657, doi:10.1038/nature18606 (2016).
- 893 4 Jorstad, N. L. *et al.* Stimulation of functional neuronal regeneration from Muller glia in  
894 adult mice. *Nature* **548**, 103-107, doi:10.1038/nature23283 (2017).
- 895 5 Su, Y. *et al.* Neuronal activity modifies the chromatin accessibility landscape in the adult  
896 brain. *Nat Neurosci* **20**, 476-483, doi:10.1038/nn.4494 (2017).
- 897 6 Denny, S. K. *et al.* Nfib Promotes Metastasis through a Widespread Increase in Chromatin  
898 Accessibility. *Cell* **166**, 328-342, doi:10.1016/j.cell.2016.05.052 (2016).
- 899 7 Qu, K. *et al.* Chromatin Accessibility Landscape of Cutaneous T Cell Lymphoma and  
900 Dynamic Response to HDAC Inhibitors. *Cancer Cell* **32**, 27-41 e24,  
901 doi:10.1016/j.ccell.2017.05.008 (2017).
- 902 8 Buenrostro, J. D. *et al.* Single-cell chromatin accessibility reveals principles of regulatory  
903 variation. *Nature* **523**, 486-490, doi:10.1038/nature14590 (2015).
- 904 9 Lake, B. B. *et al.* Integrative single-cell analysis of transcriptional and epigenetic states in  
905 the human adult brain. *Nat Biotechnol* **36**, 70-80, doi:10.1038/nbt.4038 (2018).
- 906 10 Preissl, S. *et al.* Single-nucleus analysis of accessible chromatin in developing mouse  
907 forebrain reveals cell-type-specific transcriptional regulation. *Nat Neurosci* **21**, 432-439,  
908 doi:10.1038/s41593-018-0079-3 (2018).
- 909 11 Cusanovich, D. A. *et al.* Multiplex single cell profiling of chromatin accessibility by  
910 combinatorial cellular indexing. *Science* **348**, 910-914, doi:10.1126/science.aab1601  
911 (2015).
- 912 12 Cusanovich, D. A. *et al.* The cis-regulatory dynamics of embryonic development at single-  
913 cell resolution. *Nature* **555**, 538-542, doi:10.1038/nature25981 (2018).
- 914 13 Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell  
915 transcriptomic data across different conditions, technologies, and species. *Nat Biotechnol*  
916 **36**, 411-420, doi:10.1038/nbt.4096 (2018).
- 917 14 Kiselev, V. Y. *et al.* SC3: consensus clustering of single-cell RNA-seq data. *Nat Methods* **14**,  
918 483-486, doi:10.1038/nmeth.4236 (2017).
- 919 15 Wang, B., Zhu, J., Pierson, E., Ramazzotti, D. & Batzoglou, S. Visualization and analysis of  
920 single-cell RNA-seq data by kernel-based similarity learning. *Nat Methods* **14**, 414-416,  
921 doi:10.1038/nmeth.4207 (2017).
- 922 16 Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data  
923 analysis. *Genome Biol* **19**, 15, doi:10.1186/s13059-017-1382-0 (2018).

- 924 17 Schep, A. N., Wu, B., Buenrostro, J. D. & Greenleaf, W. J. chromVAR: inferring  
925 transcription-factor-associated accessibility from single-cell epigenomic data. *Nat*  
926 *Methods* **14**, 975-978, doi:10.1038/nmeth.4401 (2017).
- 927 18 Pliner, H. A. *et al.* Cicero Predicts cis-Regulatory DNA Interactions from Single-Cell  
928 Chromatin Accessibility Data. *Mol Cell* **71**, 858-871 e858,  
929 doi:10.1016/j.molcel.2018.06.044 (2018).
- 930 19 Bravo González-Blas, C. *et al.* *Cis-topic* modelling of single cell epigenomes.  
931 *bioRxiv* (2018).
- 932 20 Fang, R. *et al.* Fast and Accurate Clustering of Single Cell Epigenomes Reveals  
933 *Cis*-Regulatory Elements in Rare Cell Types. *bioRxiv* (2019).
- 934 21 Buenrostro, J. D. *et al.* Integrated Single-Cell Analysis Maps the Continuous Regulatory  
935 Landscape of Human Hematopoietic Differentiation. *Cell* **173**, 1535-1548 e1516,  
936 doi:10.1016/j.cell.2018.03.074 (2018).
- 937 22 Ng, A. P. *et al.* Erg is required for self-renewal of hematopoietic stem cells during stress  
938 hematopoiesis in mice. *Blood* **118**, 2454-2461, doi:10.1182/blood-2011-03-344739 (2011).
- 939 23 Ong, C. T. & Corces, V. G. CTCF: an architectural protein bridging genome topology and  
940 function. *Nature Reviews Genetics* **15**, 234-246, doi:10.1038/nrg3663 (2014).
- 941 24 Dore, L. C. & Crispino, J. D. Transcription factor networks in erythroid cell and  
942 megakaryocyte development. *Blood* **118**, 231-239, doi:10.1182/blood-2011-04-285981  
943 (2011).
- 944 25 Oberst, A. *et al.* The Nedd4-binding partner 1 (N4BP1) protein is an inhibitor of the E3  
945 ligase Itch. *Proc Natl Acad Sci U S A* **104**, 11280-11285, doi:10.1073/pnas.0701773104  
946 (2007).
- 947 26 Eguchi, M. *et al.* GPHN, a novel partner gene fused to MLL in a leukemia with  
948 t(11;14)(q23;q24). *Genes Chromosomes Cancer* **32**, 212-221 (2001).
- 949 27 Whyte, W. A. *et al.* Master transcription factors and mediator establish super-enhancers  
950 at key cell identity genes. *Cell* **153**, 307-319, doi:10.1016/j.cell.2013.03.035 (2013).
- 951 28 Cusanovich, D. A. *et al.* A Single-Cell Atlas of In Vivo Mammalian Chromatin Accessibility.  
952 *Cell* **174**, 1309+, doi:10.1016/j.cell.2018.06.052 (2018).
- 953 29 Satpathy, A. T. *et al.* Transcript-indexed ATAC-seq for precision immune profiling. *Nat*  
954 *Med* **24**, 580-590, doi:10.1038/s41591-018-0008-8 (2018).
- 955 30 Hrvatin, S. *et al.* Single-cell analysis of experience-dependent transcriptomic states in the  
956 mouse visual cortex. *Nat Neurosci* **21**, 120-129, doi:10.1038/s41593-017-0029-5 (2018).
- 957 31 Bendall, S. C. *et al.* Single-cell trajectory detection uncovers progression and regulatory  
958 coordination in human B cell development. *Cell* **157**, 714-725,  
959 doi:10.1016/j.cell.2014.04.005 (2014).
- 960 32 Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by  
961 pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, 381-386,  
962 doi:10.1038/nbt.2859 (2014).
- 963 33 Qiu, X. *et al.* Reversed graph embedding resolves complex single-cell trajectories. *Nat*  
964 *Methods* **14**, 979-982, doi:10.1038/nmeth.4402 (2017).
- 965 34 Setty, M. *et al.* Wishbone identifies bifurcating developmental trajectories from single-  
966 cell data. *Nat Biotechnol* **34**, 637-645, doi:10.1038/nbt.3569 (2016).
- 967 35 Weinreb, C., Wolock, S. & Klein, A. M. SPRING: a kinetic interface for visualizing high  
968 dimensional single-cell expression data. *Bioinformatics* **34**, 1246-1248,  
969 doi:10.1093/bioinformatics/btx792 (2018).
- 970 36 Habib, N. *et al.* Div-Seq: Single-nucleus RNA-Seq reveals dynamics of rare adult newborn  
971 neurons. *Science* **353**, 925-928, doi:10.1126/science.aad7038 (2016).

- 972 37 Olsson, A. *et al.* Single-cell analysis of mixed-lineage states leading to a binary cell fate  
973 choice. *Nature* **537**, 698-+, doi:10.1038/nature19348 (2016).
- 974 38 Zhou, F. *et al.* Tracing haematopoietic stem cell formation at single-cell resolution. *Nature*  
975 **533**, 487-+, doi:10.1038/nature17997 (2016).
- 976 39 Treutlein, B. *et al.* Dissecting direct reprogramming from fibroblast to neuron using single-  
977 cell RNA-seq. *Nature* **534**, 391-+, doi:10.1038/nature18323 (2016).
- 978 40 Adolfsson, J. *et al.* Identification of Flt3<sup>+</sup> lympho-myeloid stem cells lacking erythro-  
979 megakaryocytic potential a revised road map for adult blood lineage commitment. *Cell*  
980 **121**, 295-306, doi:10.1016/j.cell.2005.02.013 (2005).
- 981 41 Chistiakov, D. A., Orekhov, A. N., Sobenin, I. A. & Bobryshev, Y. V. Plasmacytoid dendritic  
982 cells: development, functions, and role in atherosclerotic inflammation. *Front Physiol* **5**,  
983 279, doi:10.3389/fphys.2014.00279 (2014).
- 984 42 Germain, R. N. T-cell development and the CD4-CD8 lineage decision. *Nat Rev Immunol* **2**,  
985 309-322, doi:10.1038/nri798 (2002).
- 986 43 Chen, X. *et al.* ATAC-seq reveals the accessible genome by transposase-mediated imaging  
987 and sequencing. *Nat Methods* **13**, 1013-1020, doi:10.1038/nmeth.4031 (2016).
- 988 44 Chen, X. *et al.* Joint single-cell DNA accessibility and protein epitope profiling reveals  
989 environmental regulation of epigenomic heterogeneity. *Nat Commun* **9**, 4590,  
990 doi:10.1038/s41467-018-07115-y (2018).
- 991 45 Yoshida, H. *et al.* The cis-Regulatory Atlas of the Mouse Immune System. *Cell* **176**, 897-  
992 912 e820, doi:10.1016/j.cell.2018.12.036 (2019).
- 993 46 Godfrey, D. I., Kennedy, J., Suda, T. & Zlotnik, A. A developmental pathway involving four  
994 phenotypically and functionally distinct subsets of CD3-CD4-CD8- triple-negative adult  
995 mouse thymocytes defined by CD44 and CD25 expression. *J Immunol* **150**, 4244-4252  
996 (1993).
- 997 47 Taniuchi, I. *et al.* Differential requirements for Runx proteins in CD4 repression and  
998 epigenetic silencing during T lymphocyte development. *Cell* **111**, 621-633 (2002).
- 999 48 Ioannidis, V., Beermann, F., Clevers, H. & Held, W. The beta-catenin--TCF-1 pathway  
1000 ensures CD4(+)CD8(+) thymocyte survival. *Nat Immunol* **2**, 691-697, doi:10.1038/90623  
1001 (2001).
- 1002 49 Yu, S. *et al.* The TCF-1 and LEF-1 transcription factors have cooperative and opposing roles  
1003 in T cell development and malignancy. *Immunity* **37**, 813-826,  
1004 doi:10.1016/j.immuni.2012.08.009 (2012).
- 1005 50 Sun, Z. *et al.* Requirement for RORgamma in thymocyte survival and lymphoid organ  
1006 development. *Science* **288**, 2369-2373, doi:10.1126/science.288.5475.2369 (2000).
- 1007 51 Gerondakis, S., Fulford, T. S., Messina, N. L. & Grumont, R. J. NF-kappaB control of T cell  
1008 development. *Nat Immunol* **15**, 15-25, doi:10.1038/ni.2785 (2014).
- 1009 52 Picelli, S. *et al.* Smart-seq2 for sensitive full-length transcriptome profiling in single cells.  
1010 *Nat Methods* **10**, 1096-1098, doi:10.1038/nmeth.2639 (2013).
- 1011 53 Higashi, Y. *et al.* Impairment of T cell development in deltaEF1 mutant mice. *J Exp Med*  
1012 **185**, 1467-1479 (1997).
- 1013 54 Heath, H. *et al.* CTCF regulates cell cycle progression of alphabeta T cells in the thymus.  
1014 *EMBO J* **27**, 2839-2850, doi:10.1038/emboj.2008.214 (2008).
- 1015 55 Andres-Barquin, P. J., Hernandez, M. C. & Israel, M. A. Id4 expression induces apoptosis  
1016 in astrocytic cultures and is down-regulated by activation of the cAMP-dependent signal  
1017 transduction pathway. *Exp Cell Res* **247**, 347-355, doi:10.1006/excr.1998.4360 (1999).

1018 56 Carey, J. P., Knowell, A. E., Chinaranagari, S. & Chaudhary, J. Id4 promotes senescence and  
1019 sensitivity to doxorubicin-induced apoptosis in DU145 prostate cancer cells. *Anticancer*  
1020 *Res* **33**, 4271-4278 (2013).

1021 57 Surh, C. D. & Sprent, J. T-cell apoptosis detected in situ during positive and negative  
1022 selection in the thymus. *Nature* **372**, 100-103, doi:10.1038/372100a0 (1994).

1023 58 Huesmann, M., Scott, B., Kisielow, P. & von Boehmer, H. Kinetics and efficacy of positive  
1024 selection in the thymus of normal and T cell receptor transgenic mice. *Cell* **66**, 533-540  
1025 (1991).

1026 59 Shortman, K., Vremec, D. & Egerton, M. The kinetics of T cell antigen receptor expression  
1027 by subgroups of CD4+8+ thymocytes: delineation of CD4+8+3(2+) thymocytes as post-  
1028 selection intermediates leading to mature T cells. *J Exp Med* **173**, 323-332 (1991).

1029 60 Chow, S. C., Snowden, R., Orrenius, S. & Cohen, G. M. Susceptibility of different subsets  
1030 of immature thymocytes to apoptosis. *Febs Lett* **408**, 141-146 (1997).

1031 61 Picelli, S. *et al.* Full-length RNA-seq from single cells using Smart-seq2. *Nat Protoc* **9**, 171-  
1032 181, doi:10.1038/nprot.2014.006 (2014).

1033 62 Rao, S. S. *et al.* A 3D map of the human genome at kilobase resolution reveals principles  
1034 of chromatin looping. *Cell* **159**, 1665-1680, doi:10.1016/j.cell.2014.11.021 (2014).

1035 63 Zuo, Z. *et al.* ATAC-pipe: general analysis of genome-wide chromatin accessibility.  
1036 *Briefings in Bioinformatics*, bby056-bby056, doi:10.1093/bib/bby056 (2018).

1037 64 Chen, H. *et al.* Assessment of computational methods for the analysis of single-cell ATAC-  
1038 seq data. *Genome Biol* **20**, 241, doi:10.1186/s13059-019-1854-5 (2019).

1039 65 Durand, N. C. *et al.* Juicer Provides a One-Click System for Analyzing Loop-Resolution Hi-  
1040 C Experiments. *Cell Syst* **3**, 95-98, doi:10.1016/j.cels.2016.07.002 (2016).

1041 66 Qiu, X. *et al.* Single-cell mRNA quantification and differential analysis with Census. *Nat*  
1042 *Methods* **14**, 309-315, doi:10.1038/nmeth.4150 (2017).

1043 67 Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21,  
1044 doi:10.1093/bioinformatics/bts635 (2013).

1045 68 Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with or  
1046 without a reference genome. *Bmc Bioinformatics* **12**, 323, doi:10.1186/1471-2105-12-323  
1047 (2011).

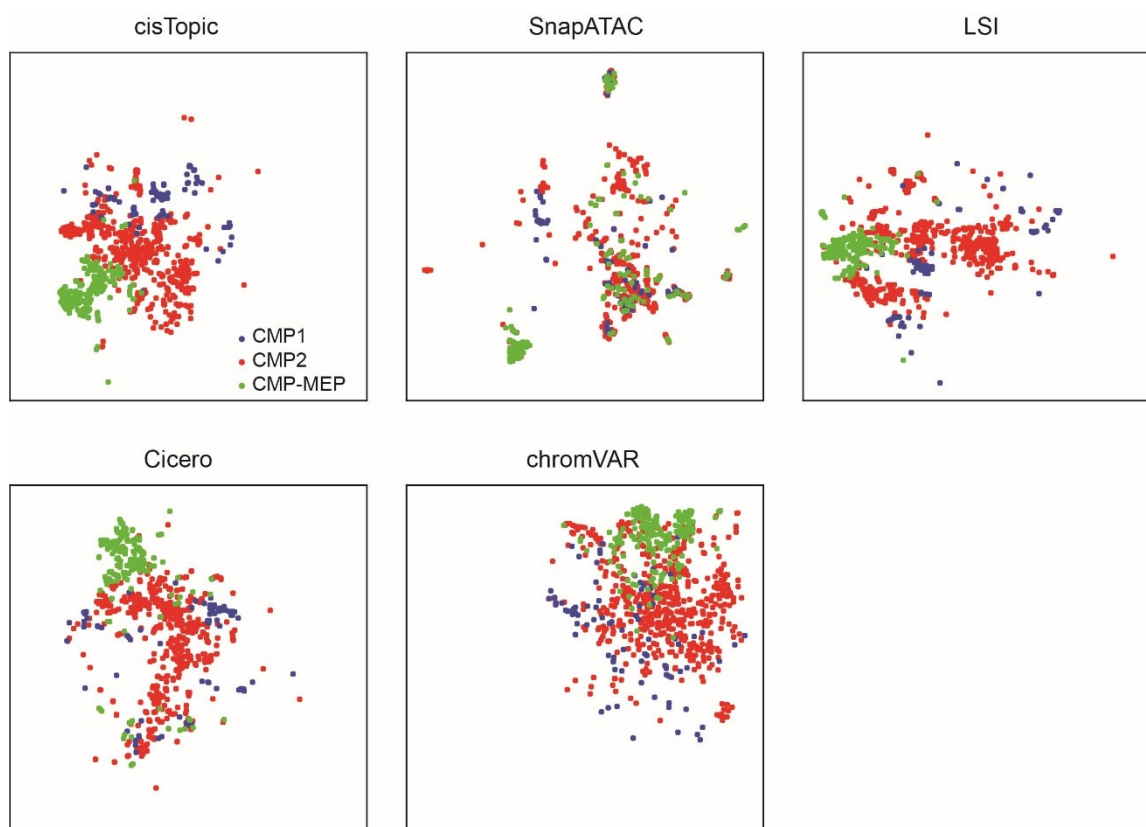
1048 69 Macosko, E. Z. *et al.* Highly Parallel Genome-wide Expression Profiling of Individual Cells  
1049 Using Nanoliter Droplets. *Cell* **161**, 1202-1214, doi:10.1016/j.cell.2015.05.002 (2015).

1050 70 McLean, C. Y. *et al.* GREAT improves functional interpretation of cis-regulatory regions.  
1051 *Nat Biotechnol* **28**, 495-501, doi:10.1038/nbt.1630 (2010).

1052 71 Bailey, T. L. *et al.* MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res*  
1053 **37**, W202-208, doi:10.1093/nar/gkp335 (2009).

1054

1055 **Supplementary Figures**

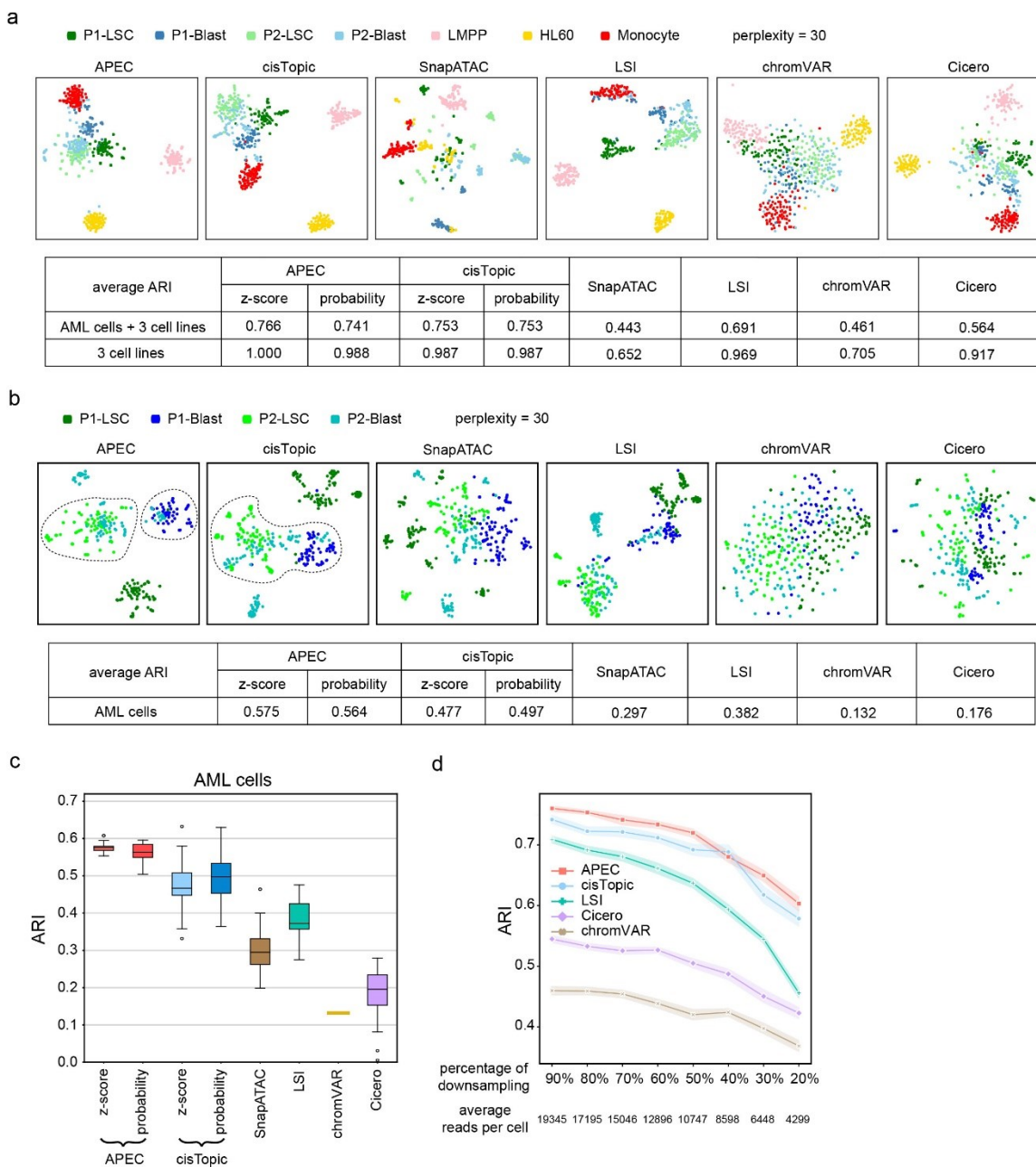


1056

1057 **Supplementary Fig. 1.** The 3 subtypes of CMP cells in the tSNE maps generated by other  
1058 methods.

1059

1060

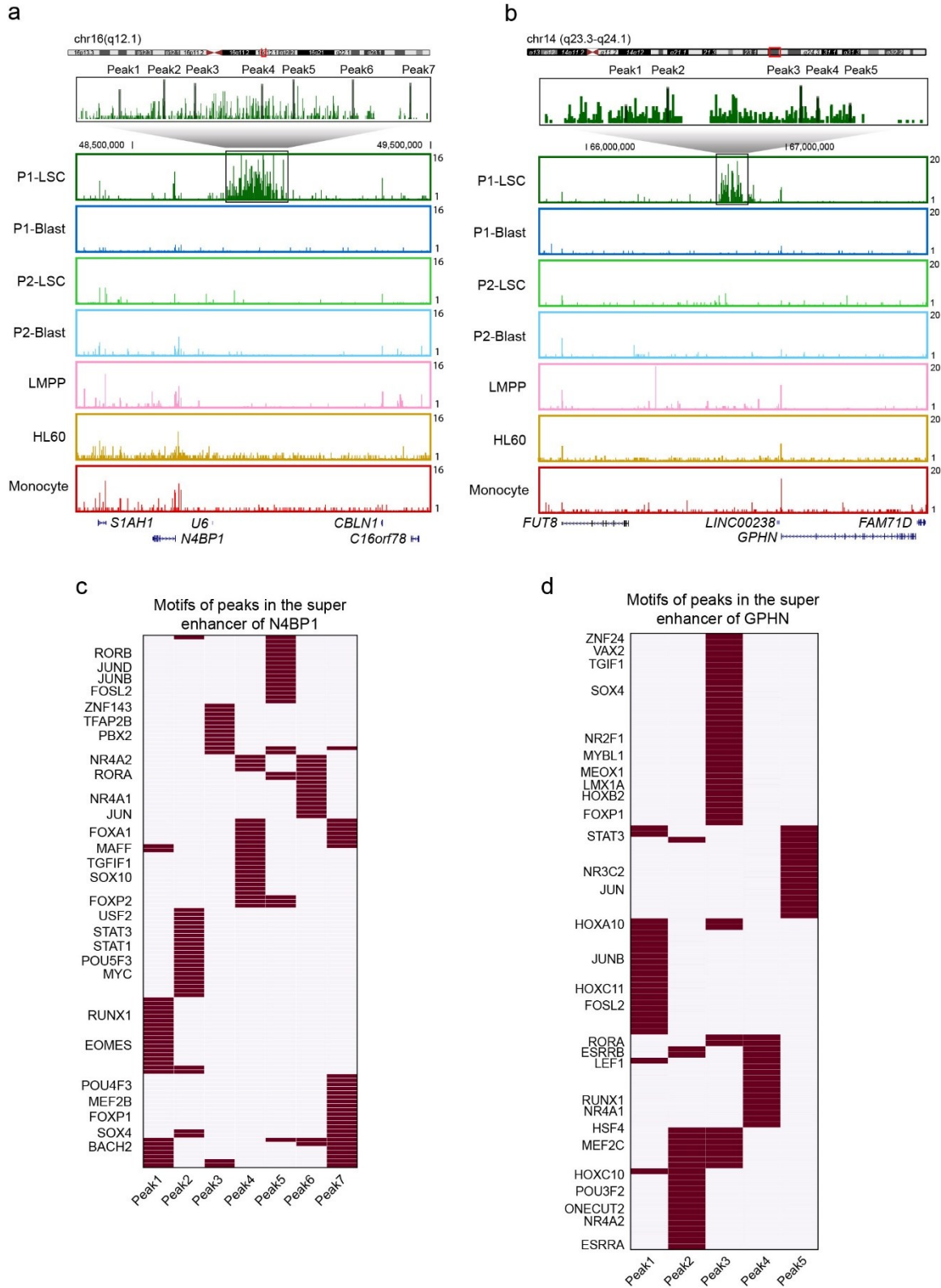


1061

1062 **Supplementary Fig. 2.** Clustering performance of the dimension-transformed matrices  
 1063 generated by different algorithms. (a) The tSNE diagrams of the cells from AML patients  
 1064 and three distinct cell lines (LMPP, monocyte and HL60). Different algorithms provided  
 1065 different dimension-transformed matrices for tSNE analysis, i.e., APEC: accession matrix;  
 1066 cisTopic: topic matrix; LSI: LSI matrix; chromVAR: bias corrected deviation matrix; Cicero:  
 1067 aggregated model matrix. The table below the diagrams contains the average ARI of the  
 1068 cell clustering results for each algorithm. (b) The tSNE diagrams and ARI table for the  
 1069 leukemic stem cells (LSCs) and blast cells from 2 different AML patients only, as in (a). (c)

1070 Box-plots showing the ARI values for the clustering of the blast and LSC cells from two  
1071 AML patients. We sampled different tunable parameters for different algorithms. APEC:  
1072 the accession number; cisTopic: the random seed; SnapATAC: the number of principal  
1073 components and the number of nearest neighbors; LSI: the number of top SVD  
1074 components; Cicero: the peak aggregation distance; chromVAR: no sampling. Z-score  
1075 and probability denote different methods of normalizing the dimension-transformed  
1076 matrices. Center line, median; box limits, upper and lower quartiles; whiskers, 1.5x  
1077 interquartile range; points, outliers. (d) The average ARI values calculated by down-  
1078 sampling 50 times from the raw data of the AML cells and three cell lines for each method.  
1079 The X-axis represents the percentage of down-sampled sequencing reads. Shaded error  
1080 band: 95% confidence interval.

1081



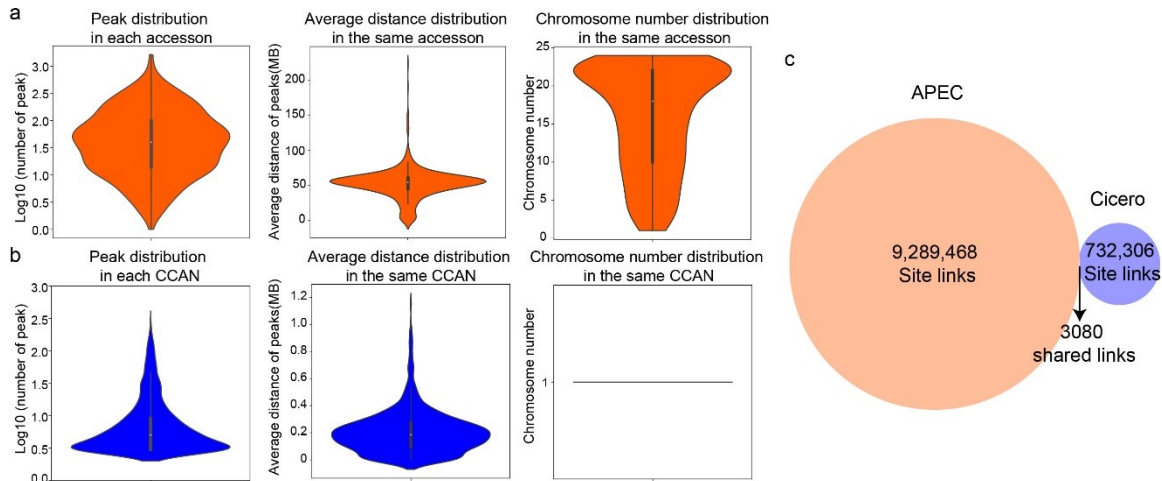
1082

1083 **Supplementary Fig. 3.** Super-enhancers predicted by APEC for the scATAC-seq data of  
 1084 cells from AML patients. (a, b) The genome browser track shows the aggregated scATAC-



1085 seq signal of the super-enhancer of P1-LSC cells upstream of *N4BP1* (a) and *GPHN* (b).  
1086 (c, d) The motifs associated with peaks in the super-enhancer upstream of *N4BP1* (c) and  
1087 *GPHN* (d).

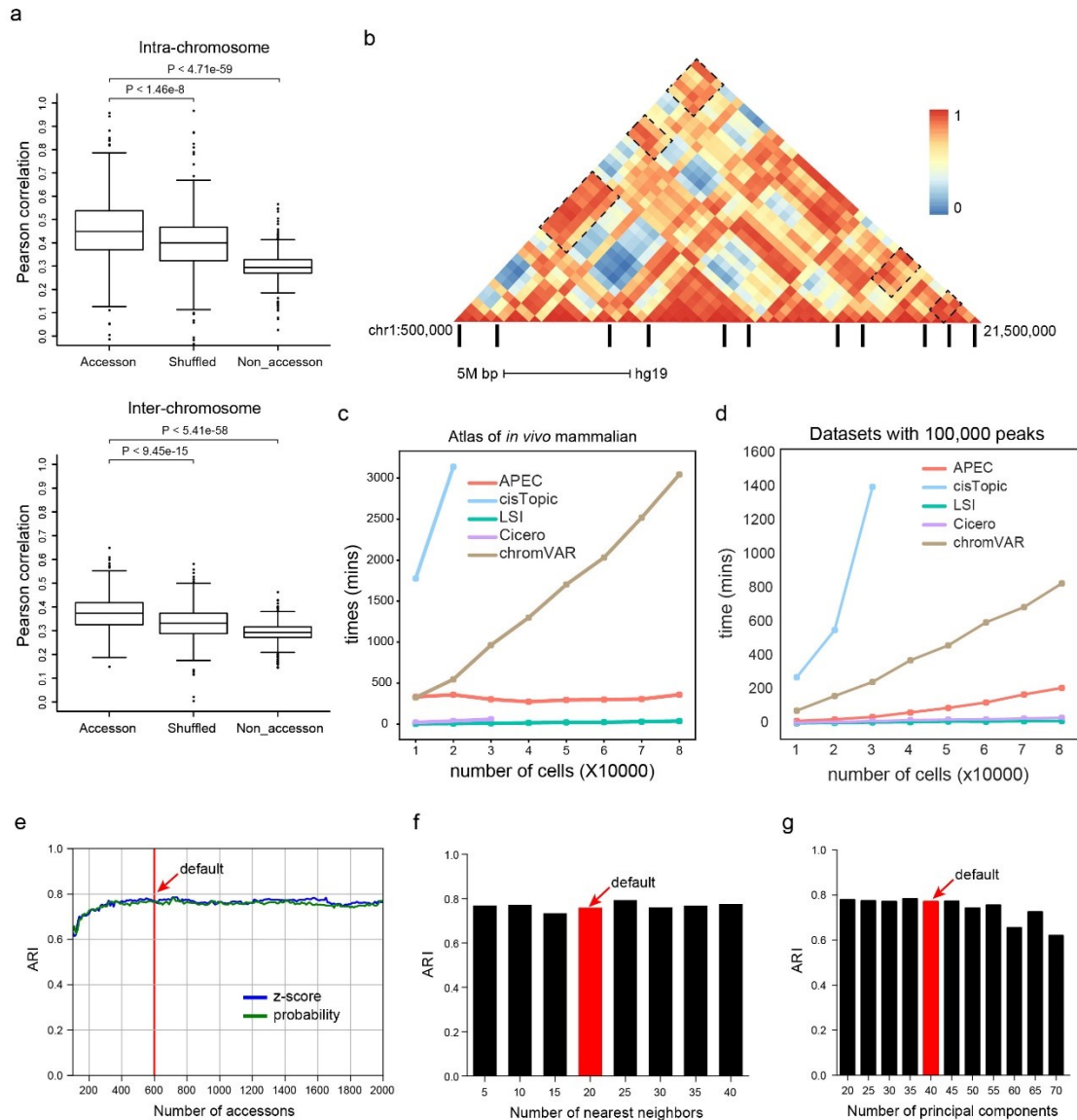
1088



1089

1090 **Supplementary Fig. 4.** Comparison of the peak grouping algorithms used by APEC and  
1091 Cicero on the hematopoietic dataset. (a) The characteristics of accessions in APEC. Left  
1092 panel: distribution of peaks in each accession; middle panel: genomic distances of peaks  
1093 belong to the same accession; right panel: number of chromosomes with peaks belong to  
1094 the same accession. (b) The characteristics of CCAN (defined by Cicero), as in (a). (c)  
1095 Site links discovered by APEC and Cicero.

1096



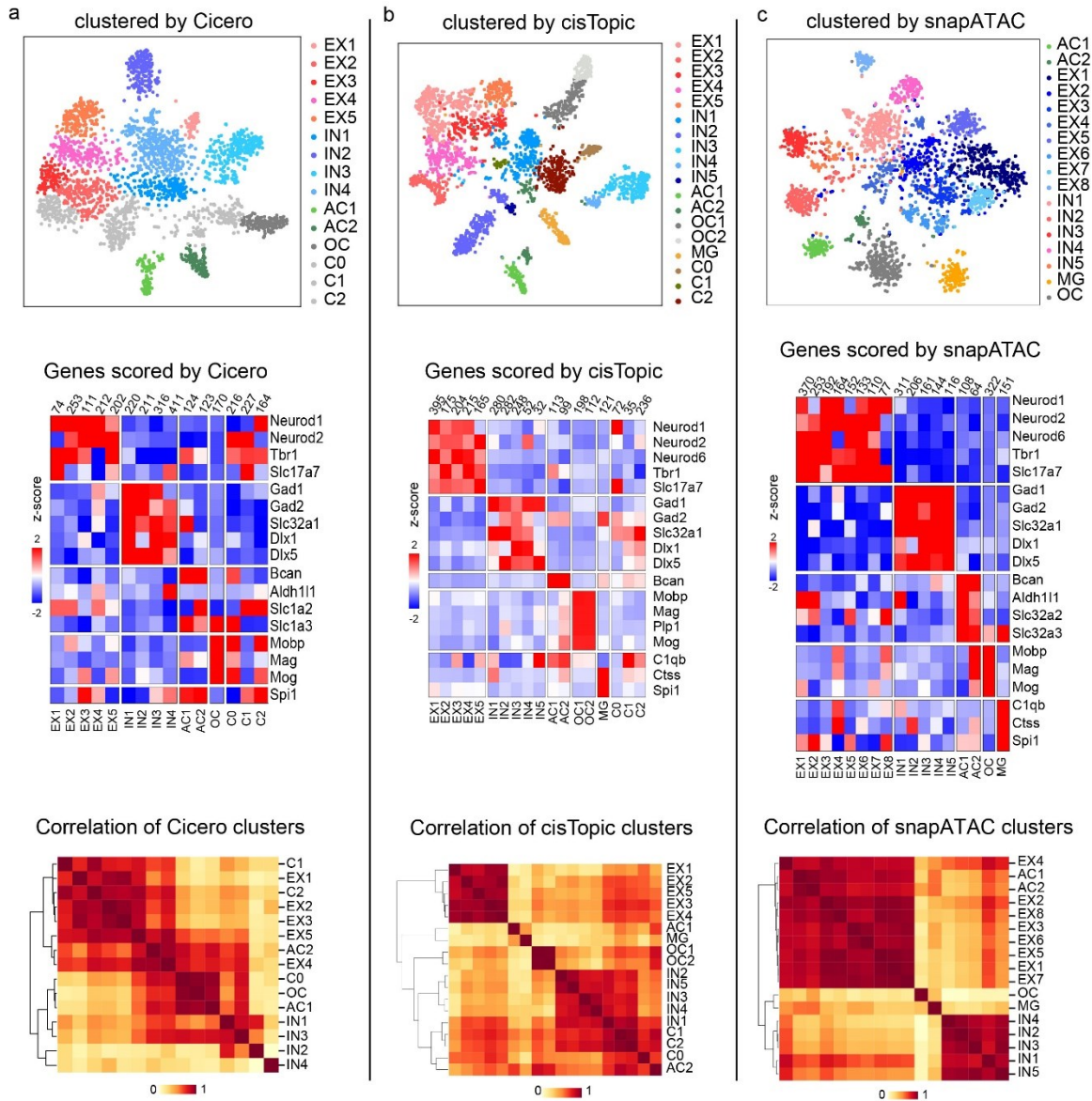
1097

1098 **Supplementary Fig. 5.** Biological insights of the accession and stability and scalability  
 1099 analysis of APEC. **(a)** Box plot showing the average spatial distance between peaks in the  
 1100 same accession (from scATAC-seq) versus randomly shuffled peaks versus non-  
 1101 accessible genomic regions. Spatial distance was estimated from chromosome  
 1102 conformation capture (Hi-C) technology. Both Hi-C and scATAC-seq data were generated  
 1103 from the same cell line GM12878. Upper panel: intra-chromosomal correlation of windows  
 1104 in the Hi-C data; Lower panel: inter-chromosomal correlation of windows in the Hi-C data.  
 1105 Accession: The correlation between two windows that contain peaks in the same accession;  
 1106 Shuffled: The correlation calculated by randomly shuffling peaks in each accession; Non-  
 1107 accession: The correlation between two windows that contain no peaks. Center line,  
 1108 median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points,  
 1109 outliers. **(b)** The Hi-C profile of windows between chr1:500,000-21,500,000. The black  
 1110 bars below the Hi-C track denote peaks in the same accession from APEC. Dotted boxes

1111 indicate examples of peaks in the same accession that are distant in genomic position but  
1112 close in space. **(c, d)** The computing time required for different algorithms to cluster cell  
1113 numbers from 10,000 to 80,000 with all peaks (c) and 100,000 peaks (d). The data were  
1114 sampled from the single-cell atlas of *in vivo* mammalian chromatin accessibility. CisTopic  
1115 was performed using 8 CPU threads and all the other tools with 1 CPU thread. **(e-g)** The  
1116 ARI values of the clustering results that used different numbers of accessions (e), nearest  
1117 neighbors (f), and principle components (g). The dataset includes the cells from two AML  
1118 patients and three cell lines. Default values are noted in red.

1119

1120



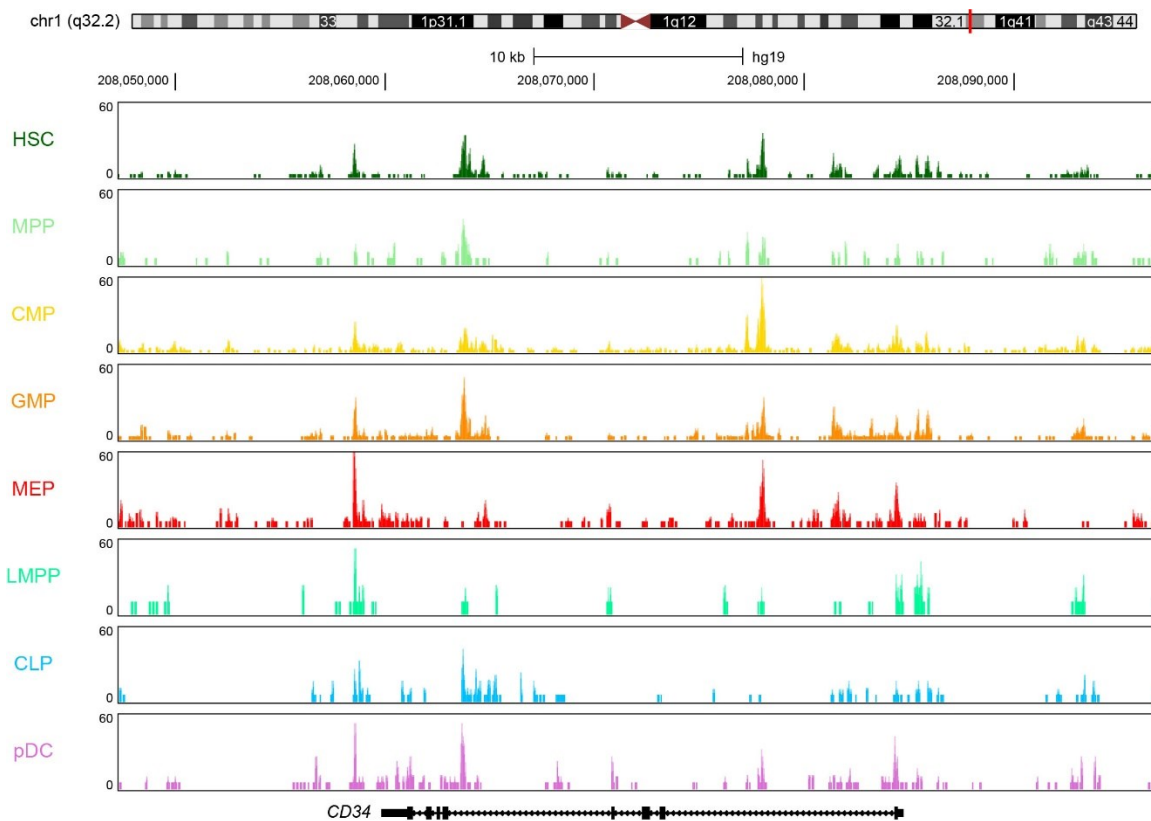
1121

1122 **Supplementary Fig. 6. (a)** The clustering and cell-type classification of the mouse  
 1123 forebrain dataset by Cicero. Upper panel: cell clusters obtained by Cicero, illustrated in  
 1124 the tSNE diagram. Middle panel: the z-scores of the average gene scores of cell clusters,  
 1125 obtained by Cicero. Lower panel: the hierarchical clustering of the Pearson correlations  
 1126 between cell clusters identified by Cicero. **(b, c)** The clustering and cell-type classification  
 1127 of the same dataset by cisTopic and SnapATAC respectively, as in (a).

1128

1129

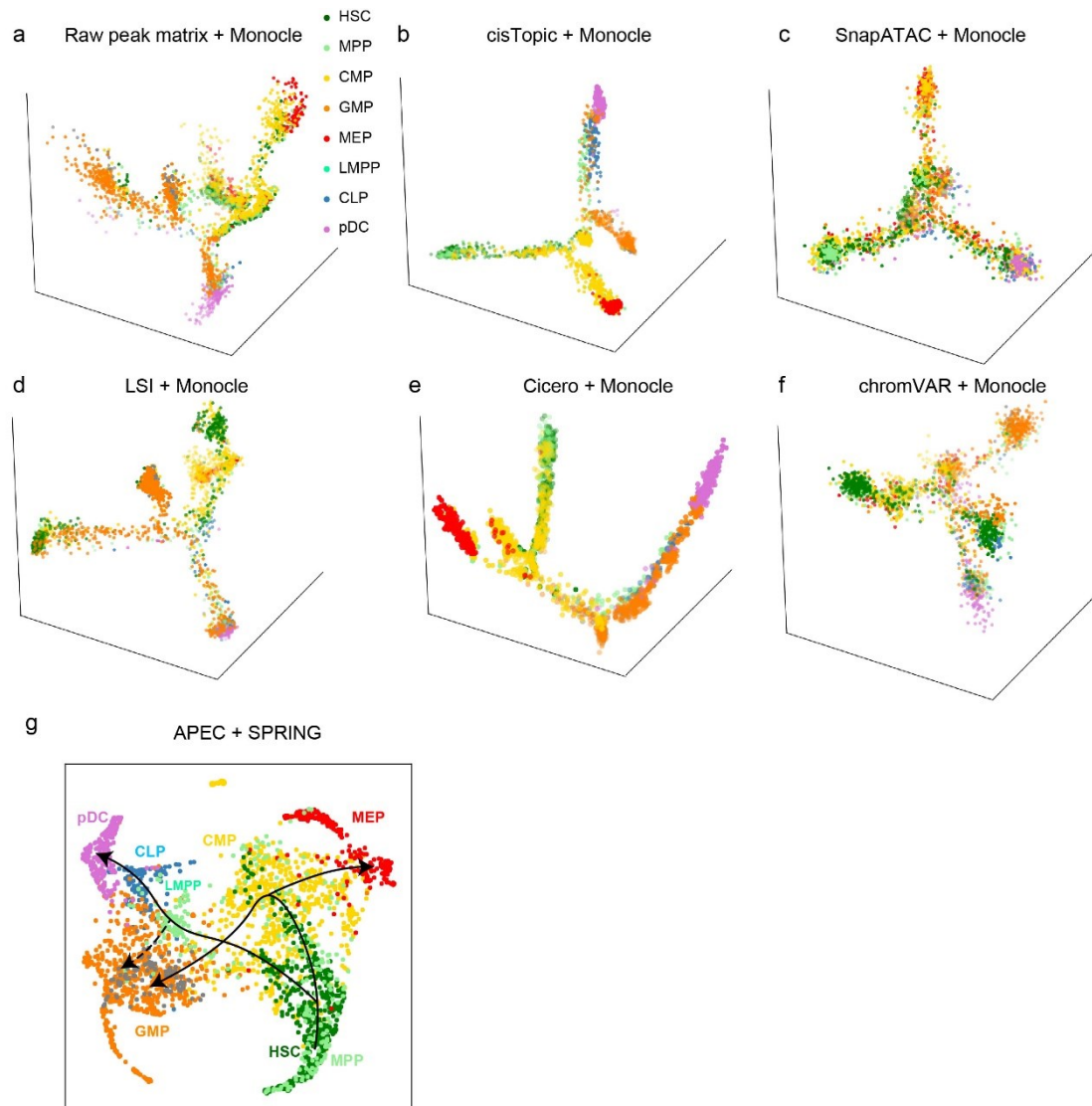
1130



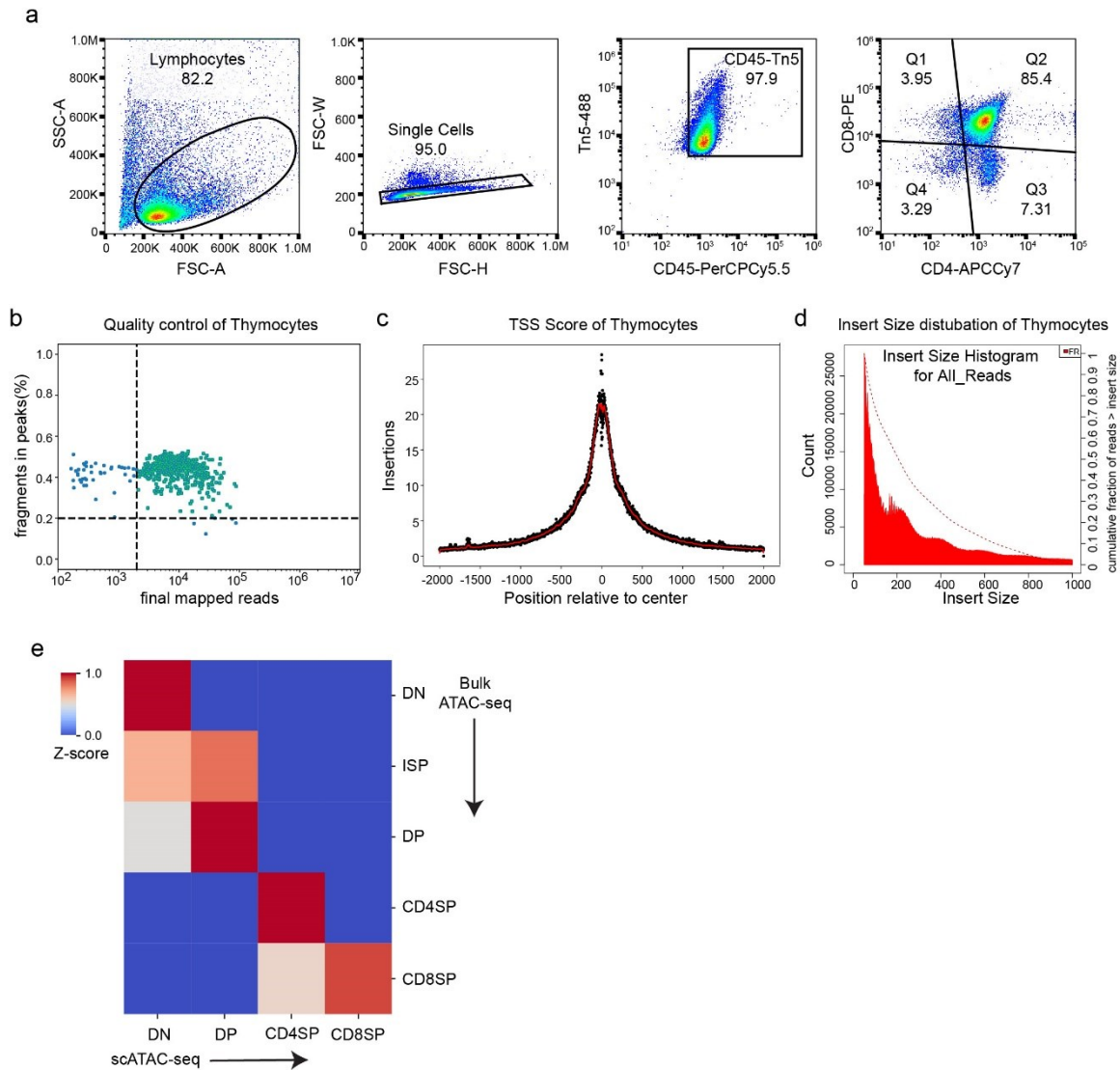
1131

1132 **Supplementary Fig. 7.** UCSC genome browser track diagram of the normalized fragment  
1133 count around gene *CD34* for each hematopoietic cell type.

1134



**Supplementary Fig. 8.** Cell differentiation trajectories of the human hematopoietic dataset constructed by different algorithms. (a-f) The pseudotime trajectories constructed by the combination of Monocle and the raw peak count matrix, the topic matrix from cisTopic, the normalized count matrix from SnapATAC, the LSI matrix, the aggregated model matrix from Cicero, and the bias corrected deviation matrix from chromVAR, respectively. (g) The pseudotime trajectory constructed by the combination of SPRING and the accession matrix from APEC.

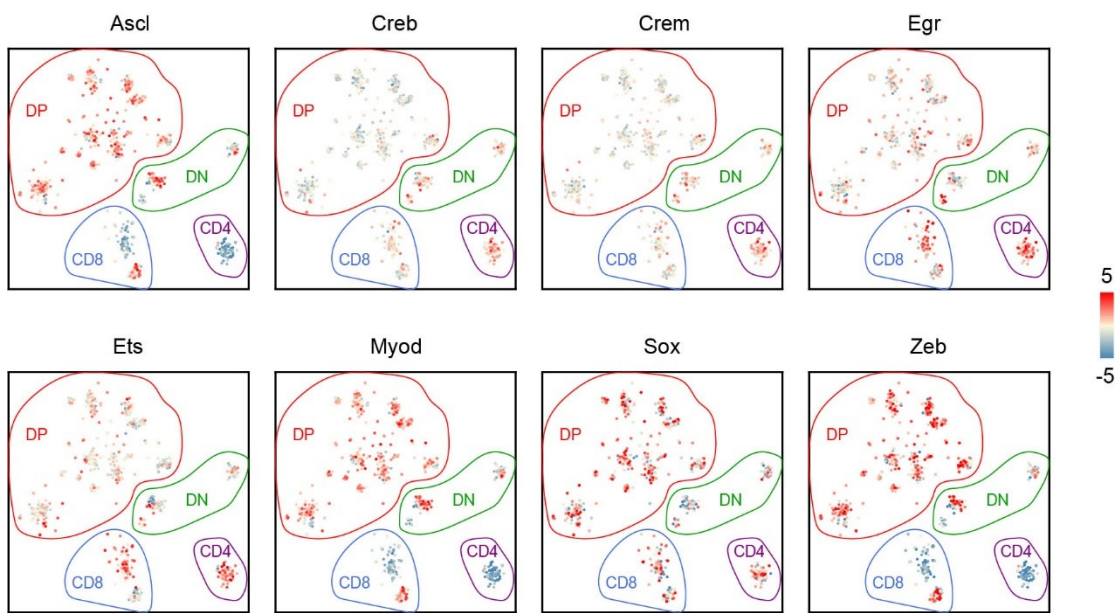


1144

1145 **Supplementary Fig. 9.** (a) Gating strategy of the mouse thymocytes in ftATAC-seq. (b-d)  
 1146 Quality control diagrams for the mouse thymocyte data, including reads numbers and  
 1147 percentage of fragments in peaks for each cell (b), average count of scATAC-seq  
 1148 insertions around TSS regions (c), and statistical distribution of fragment lengths (d). (e)  
 1149 The z-score of correlation between the cell types from ftATAC-seq and bulk ATAC-seq  
 1150 data.

1151

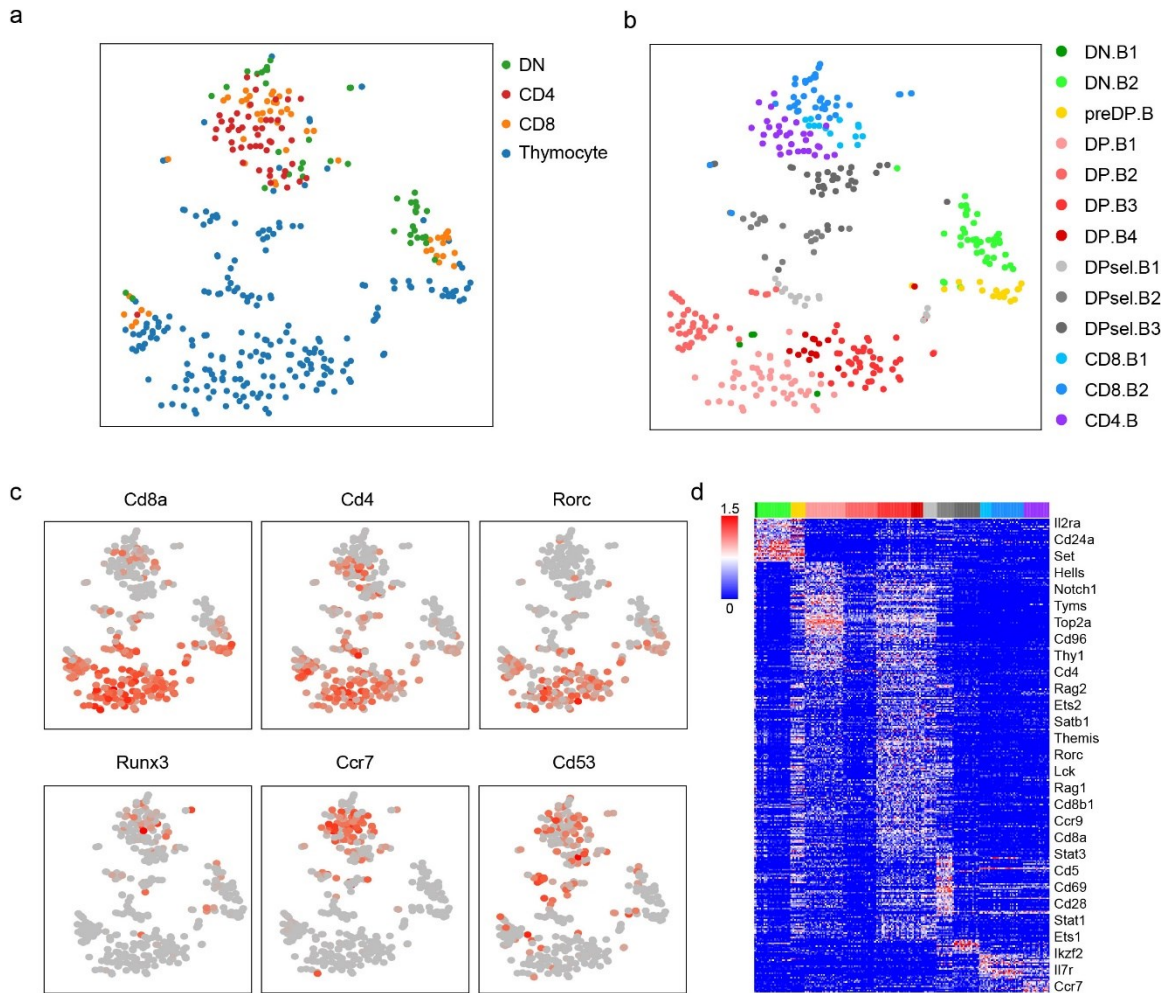




1152

1153 **Supplementary Fig. 10.** Selected significant motifs enriched in different thymocyte  
1154 subtypes obtained by the APEC algorithm.

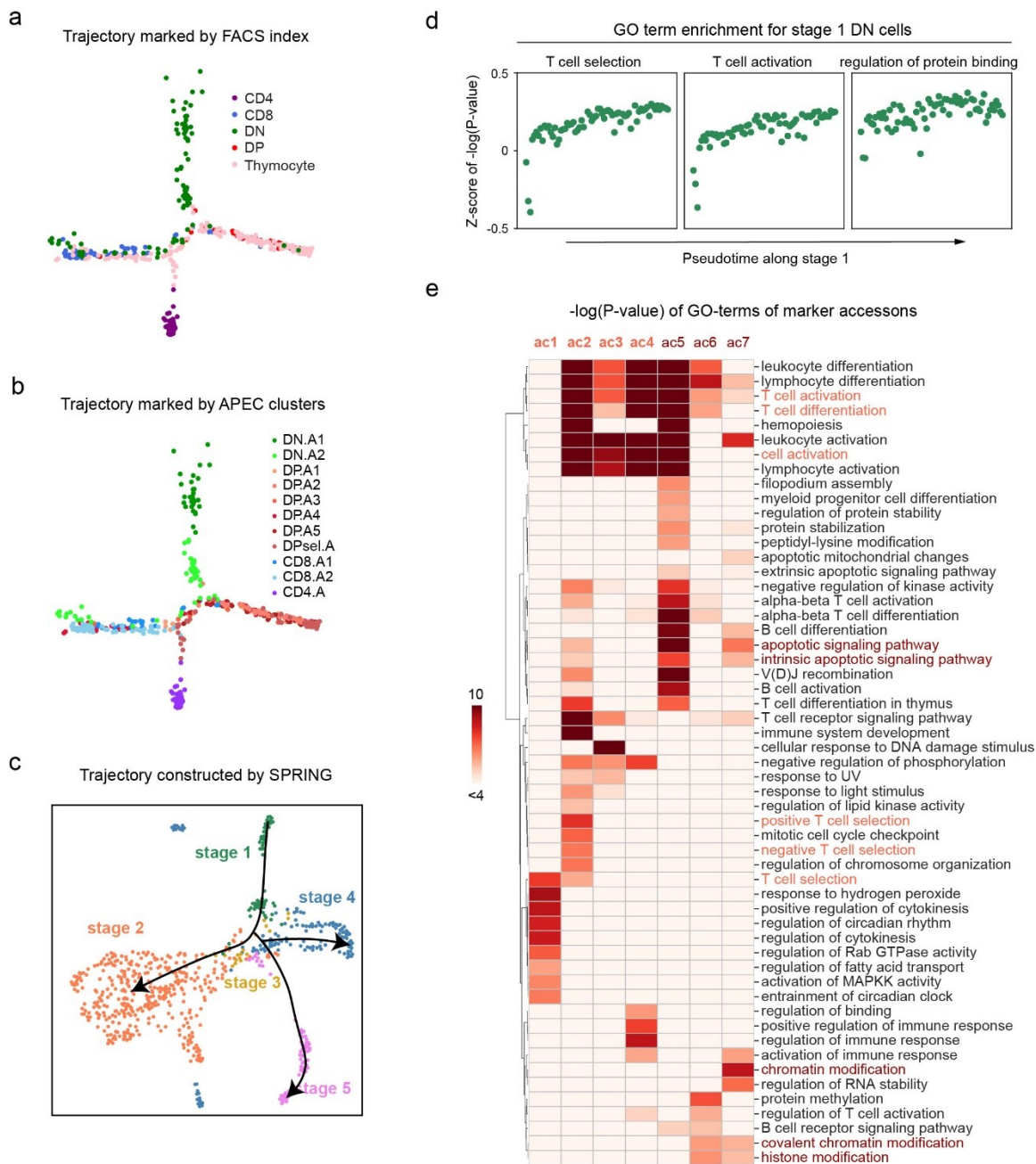
1155



1156

1157 **Supplementary Fig. 11.** Single-cell transcriptome analysis of *Mus musculus* thymocytes  
 1158 from SMART-seq. **(a)** tSNE diagram of the single-cell expression matrix of *Mus musculus*  
 1159 thymocytes, labeled by the FACS index of each cell. **(b)** Louvain clustering of the same  
 1160 single-cell dataset obtained by Seurat. The cell types of these clusters were classified by  
 1161 the expression of corresponding marker genes. **(c)** Important marker genes were  
 1162 differentially expressed in different cell clusters. **(d)** Heatmap of the expressions of all  
 1163 genes significantly differentially expressed between cell clusters. The top color bar used  
 1164 the same scheme described in (b) to render cells of different clusters.

1165



1166

1167 **Supplementary Fig. 12.** Developmental characteristics of single-cell samples captured  
 1168 by APEC. (a, b) Pseudotime trajectory of scATAC-seq data from *Mus musculus*  
 1169 thymocytes labeled with the FACS index and APEC cluster index. (c) Pseudotime  
 1170 trajectory constructed by applying SPRING to the accession matrix. The colors of cells  
 1171 denote their stages in the APEC trajectory results. (d) Z-scores of the  $-\log(P\text{-value})$  of the  
 1172 GO terms along the pseudotime trajectory of stage 1 cells. (e) Logarithm of the P-value of  
 1173 GO terms searched from peaks in accessions ac1~ac7, which are the marker accessions  
 1174 of cluster DP. A1 and DP. A3/4/5 of thymocytes.