

APEC: An accesson-based method for single-cell chromatin accessibility analysis

Bin Li^{1,3}, Young Li^{1,3}, Kun Li¹, Lianbang Zhu¹, Qiaoni Yu¹, Pengfei Cai¹, Jingwen Fang¹, Wen Zhang¹, Pengcheng Du¹, Chen Jiang¹, Kun Qu^{1,2,*}

Affiliation

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

²CAS Center for Excellence in Molecular Cell Sciences, University of Science and Technology of China

³Cofirst authors

*Correspondence: Kun Qu (gukun@ustc.edu.cn)

Contact Information

Kun Qu, Ph.D.

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

Hefei, Anhui, China, 230027

Email: gukun@ustc.edu.cn

Phone: +86-551-63606257

Abstract:

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

INTRODUCTION

As a technique for probing genome-wide chromatin accessibility in a small number of cells *in vivo*, the assay for transposase-accessible chromatin with high-throughput sequencing (ATAC-seq) has been widely applied to investigate the cellular regulomes of many important biological processes¹, such as hematopoietic stem cell (HSC) differentiation², embryonic development³, neuronal activity and regeneration^{4,5}, tumor cell metastasis⁶, and patient responses to anticancer drug treatment⁷. Recently, several experimental schemes have been developed to capture chromatin accessibility at single-cell/nucleus resolution, i.e., single-cell ATAC-seq (scATAC-seq)⁸, single-nucleus ATAC-seq (snATAC-seq)^{9, 10}, and single-cell combinatorial indexing ATAC-seq (sci-ATAC-seq)^{11, 12}, which significantly extended researchers' ability to uncover cell-to-cell epigenetic variation and other fundamental mechanisms that generate heterogeneity from identical DNA sequences. By contrast, the in-depth analysis of single-cell chromatin accessibility profiles for this purpose remains a challenge. Numerous efficient algorithms have been developed to accurately normalize, cluster and visualize cells from single-cell transcriptome sequencing profiles, including but not limited to Seurat¹³, SC3¹⁴, SIMLR¹⁵, and SCANPY¹⁶. However, most of these algorithms are not directly compatible with a single-cell ATAC-seq dataset, for which the signal matrix is much sparser. To characterize scATAC-seq data, the Greenleaf lab developed an algorithm named chromVAR¹⁷, which aggregates mapped reads at accessible sites based on annotated motifs of known transcription factors (TFs) and thus projects the sparse per accessible peak per cell matrix to a bias-corrected deviation per motif per cell matrix and significantly stabilizes the data matrix for downstream clustering analysis. Other mathematical tools, such as the latent semantic indexing (LSI)^{11, 12}, Cicero¹⁸, and cisTopic¹⁹ have also been applied to process single-cell/nucleus ATAC-seq data^{10, 12}. However, great challenges still remain for current algorithms to accurately cluster large number of cells and precisely predict gene expressions from single cell chromatin accessibility profiles. Therefore, a refined algorithm is urgently needed to better categorize cell subgroups with minor differences, thereby providing a deeper mechanistic understanding of single-cell epigenetic heterogeneity and regulation.

RESULTS

Accession-based algorithm improves single-cell clustering

Here, we introduce a new single-cell chromatin accessibility analysis toolkit named APEC (accessibility pattern-based epigenomic clustering), which combines peaks with the same signal fluctuation among all single cells into peak groups, termed "accessions", and converts the original sparse cell-peak matrix to a much denser cell-accession matrix for cell type categorization (Fig. 1a). In contrast to previous methods (e.g., chromVAR¹⁷, LSI^{11, 12}, Cicero¹⁸, and cisTopic¹⁹), this accession-based reduction scheme naturally groups synergistic accessible regions genome-wide together without a priori knowledge of genetic information (such as TF motifs or genomic distance) and provides more efficient, accurate, robust and rapid cell clustering from single-cell ATAC-seq profiles. More conveniently, APEC integrates all necessary procedures, from raw sequence trimming, alignment, and quality control (Supplementary Fig. 1) to cell clustering, motif enrichment, and pseudotime trajectory prediction into a head-to-toe program package that has been made available on GitHub (<https://github.com/QuKunLab/APEC>).

To test the performance of APEC, we first obtained data from previous publications that performed scATAC-seq on a variety of cell types with known identity during hematopoietic stem cell (HSC) differentiation²⁰ as a gold standard. Compared to other state-of-the-art single cell chromatin accessibility analysis methods such as chromVAR^{17, 21}, LSI^{11, 12}, Cicero¹⁸ and cisTopic¹⁹, this new accession-based algorithm more precisely and clearly clustered cells into their corresponding identities according to the Adjusted Rand Index (ARI) (Fig. 1b-c). On average, 67% of cells were correctly classified by APEC with ARI=0.480/0.522 (normalized by Z-score or probability, respectively). CisTopic was the second most accurate method to predict cell identities (ARI=0.392/0.418), and correctly classified ~56% of cells. However, given 5 times more CPU threads, the cisTopic algorithm was still 15-50 times slower than other methods. Moreover, APEC identified 3 sub-clusters of CMP cells that were not discovered by any other algorithms, namely CMP1, CMP2 and CMP-MEP (Fig. 1d). CMP1 cells are early stage of CMPs that enriched TFs associated with stem cell self-renewal, such as Erg²²; CMP2 cells are enriched with CTCF motif, suggesting that these cells are at the fate decision stage with CTCF associated chromatin remodeling²³; CMP-MEP cells are considered as MEP committed CMPs, and are strongly enriched with crucial regulators for MEP differentiation, such as GATA1²⁴. More details about the distribution of these 3 sub-clusters of CMP cells on the development trajectory will be discussed later in the section of pseudotime prediction. Another advantage of APEC over all the other tools is the capability to precisely evaluate gene expressions from single cell chromatin accessibility information. For instance, genes *FOXO1*²⁵, *CEBPA*²⁶, *CD86*²⁷, *IKZF1*²⁸, *GFI1B*²⁹, and *AQP1*³⁰ are well studied marker genes for HSC, GMP, CLP, and MEP cells respectively, and were also predicted exclusivity highly expressed in the correct cell types by APEC (Fig. 1e). On the contrary,

although cisTopic and Cicero can also assess gene expression, they both failed to present the enrichment of these genes in the corresponding cell types ([Supplementary Fig. 2](#)).

To further confirm the superiority of APEC, we performed the same comparison analysis with another scATAC-seq dataset on three distinct cell types, namely, lymphoid-primed multipotent progenitors (LMPPs), monocytes, and HL-60 lymphoblastoid cells (HL60), and four similar cell types, namely, blast cells and leukemic stem cells (LSCs) from two acute myeloid leukemia (AML) patients¹⁷. Again, we see that APEC outperformed all the other tools in cell clustering with an overall ARI=0.770/0.763 ([Supplementary Fig. 3a](#)). Interestingly, APEC, cisTopic and LSI were all capable of almost perfectly separating the three distinct cell types (LMPPs, monocytes, and HL60), with ARI = 1.000/0.994, 0.987/0.987 and 0.969, respectively, while the other two Cicero (ARI=0.917) and chromVAR (ARI=0.705) were not as good. However, in terms of clustering the four similar cell types from AML patients, APEC (ARI=0.572/0.551, 80% of cells correctly classified) clearly outperformed cisTopic (ARI=0.477/0.497, 61% of cells correctly classified) and LSI (ARI=0.382, 53% of cells correctly classified) ([Supplementary Fig. 3b](#)), suggesting that APEC was the most sensitive among all the tools. Since each method can generate varying numbers of clusters depending on the parameters used, we benchmarked the performance of all the methods using ARI across a wide range of parameters to ensure the reliability of their predictions ([Supplementary Fig. 3c-d](#)). To further test the robustness of APEC at low sequencing depth, we randomly selected reads from the raw data and calculated the ARI values for each down-sampled data. APEC exhibits better performance at sequencing depths as low as 20% of the original data ([Supplementary Fig. 3e](#)), confirming the sensitivity of the algorithm.

Compare with chromVAR, the contribution of the minor differences between similar cells is aggregated in accessions but diluted in motifs. For example, APEC identified prominent super-enhancers around the E3 ligase inhibitor gene *N4BP1*³¹ and the MLL fusion gene *GPHN*³² in the LSC cells from AML patient 1 (P1-LSC) but not in the other cell types ([Supplementary Fig. 4a-b](#)). We noticed that all peaks in these super-enhancers were classified into one accession that was critical for distinguishing P1-LSCs from P2-LSCs, P1-blast cells and P2-blast cells. However, these peaks were distributed in multiple TF motifs, which significantly diluted the contributions of the minor differences ([Supplementary Fig. 4c-d](#)). In contrast to Cicero, which aggregates peaks based on their cis-co-accessibilities networks (CCAN) within a certain range of genomic distance¹⁸, APEC combine synergistic peaks genome-wide. Take the human hematopoietic cells dataset as an example, 600 accessions were built from the 54,212 peaks, and each accession contained ~40 peaks (median number) compare with ~4 peaks in each CCAN ([Supplementary Fig. 5a~b](#)). The average distance between peaks in a same accession is ~50 million base pairs

(compared with ~0.2 million bps from CCAN), and over 57% of accessions contain peaks from more than 15 different chromosomes. From the same dataset, Cicero identified 732,306 pairs of site links from 25,102 peaks, and information from the remaining peaks were simply discarded. APEC identified more than 9.2 million pairs of site links from all the 54,212 peaks, within which only 3080 site links were identified by both methods ([Supplementary Fig. 5c](#)), therefore, APEC and Cicero are two completely different approaches. Furthermore, Buenrostro et al. showed that the covariation of the accessible sites across all the cells may reflect the spatial distance between the corresponding peaks⁸. By integrating the chromatin conformation profiles from Hi-C experiments with the scATAC-seq profile for the same cells, we found that peaks in the same accession are spatially much closer to each other than randomly selected peaks ([Supplementary Fig. 6a](#), $P\text{-value} < 10^{-7}$), suggesting that they may belong to the same topologically associated domains (TADs) ([Supplementary Fig. 6b](#)).

Speed and scalability are now extremely important for single-cell analytical tools due to the rapid growth in the number of cells sequenced in each experiment. We benchmarked APEC and all the other tools based on a random sampling of the mouse *in vivo* single-cell chromatin accessibility atlas dataset³³, which contains 81,173 high quality cells. Taking into account of all the 436,206 peaks, it took APEC 310 min to cluster 80,000 cells with 1 CPU thread, ~9 times faster than chromVAR and ~8 times slower than LSI ([Fig. 1f](#)). CisTopic took 5,000 min to categorize 20,000 cells using 5 CPU threads, and we estimated it might take 30 days to finish clustering 80,000 cells if no parallel computing is applied. On the other hand, Cicero led to memory overflow to cluster 40,000 cells (on 512 GM RAM computer). Therefore, neither cisTopic nor Cicero may be applicable to analyze large scale datasets using limited computing resources. We also randomly select 100,000 peaks from the entire dataset to test the performance of these tools where similar conclusions can be obtained ([Supplementary Fig. 6c](#)). In addition, APEC is very stable for a wide range of parameter values used in the algorithm, such as the number of accessions, nearest neighbors and principle components ([Supplementary Fig. 6d-f](#)).

APEC is applicable to multiple single-cell chromatin detection techniques

To evaluate the compatibility and performance of APEC with other single-cell chromatin accessibility detection techniques, such as snATAC-seq¹⁰, transcript-indexed scATAC-seq³⁴ and sciATAC-seq¹¹, APEC was also tested with the data sets generated by those experiments. For example, APEC discovered 13 cell subpopulations in adult mouse forebrain snATAC-seq data¹⁰, including four clusters of excitatory neurons (EX1-4), five groups of inhibitory neurons (IN1-5),

astroglia cells (AC1&2), oligodendrocyte cells (OC), and microglial cells (MG; Fig. 2a-b), defined by the chromatin accessibilities at the loci of cell type-specific genes (Fig. 2c). Compared to published results¹⁰, APEC identified 4 rather than 3 excitatory subpopulations and 5 rather than 2 distinct inhibitory subpopulations, and all the cell groups were clearly distinguished from each other by hierarchical clustering (Fig. 2d). The motif enrichment analysis module in APEC identified cell type-specific regulators that are also consistent with previous publications¹⁰. For example, the NEUROD1 and OLIG2 motifs were generally enriched on excitatory clusters (EX1\3\4); the MEF2C motif was more enriched on EX1/2 and the left part of EX3 than other excitatory neurons; the motifs of MEIS2 and DLX2 were differentially enriched on two subtypes of inhibitory neurons (IN2 and IN4, respectively); and the NOTO, SOX2, and ETS1 motifs were enriched on the AC1, OC, and MG clusters, respectively (Fig. 2e). These results confirm that APEC is capable of identifying cell subtype-specific regulators.

Since single-cell transcriptome analysis is also capable to identify novel cell subpopulations, it is critical to anchor the cell types identified from scATAC-seq to those from scRNA-seq. Hrvatin et al. identified dozens of excitatory and inhibitory neuronal subtypes in the mouse visual cortex using single cell inDrops sequencing³⁵ and provided tens of signature genes that distinguished these cell types. Interestingly, the accessions that represent these signature genes were also distinctly enriched at corresponding clusters of neurons. For example, *Enpp2* is a marker gene for cluster Excl23 defined in the inDrops-seq data, and accessions represent this gene were also enriched in the EX1 cluster in snATAC-seq, suggesting an anchor between EX1 and Excl23 (Fig. 2f, Supplementary Fig. 7a). Cluster EX1 in scATAC-seq also matched with Excl5_3 in inDrops-seq (marked by *Deptor* and *Fam3c*). Similarly, cluster EX2 matched with Excl4 (marked by *Rorb* and *Tshz1*), EX3 matched with Excl5_1 (marked by *Bmp3*), and EX4 matched with Excl6 (marked by *Col5a1*, *Foxp2* and *Pcp4*). The same method also works to anchor inhibitory neurons, as the IN1 cells in the snATAC-seq data corresponded to the Int_Vip and Int_Npy clusters in the inDrops-seq data (marked by *Vip* and *Npy*). Cluster IN2 matched with Int_Cck (marked by *Cck*), IN3 matched with Int_Pv (marked by *Pvalb*), IN4 matched with Int_Sst_2 (marked by *Crhbp*), and IN5 matched with Int_Sst_1 (marked by *Lypd6b*) (Fig. 2g, Supplementary Fig. 7b). These results highlight the potential advantages of the accession-based approach for the integrative analysis of scATAC-seq and scRNA-seq data.

In addition, due to the sparser fragment count matrix (~1200 reads per cell), more than 29.7% (946 out of 3034) of cells were previously unable to be correctly assigned into any subpopulation of interest¹⁰, but APEC successfully categorized all cells into their corresponding

subtypes, confirming its high sensitivity. In contrast, cisTopic identified 5 EX, 5 IN, 2 AC, 2 OC and 1 MG clusters, however 11% of cells in clusters C0~C2 cannot be clearly classified, and it mis-clustered AC1 and AC2 in the heatmap (Supplementary Fig. 8a). LSI identified 5 EX, 3 IN, 2 AC, 1 OC and 1 MG cell clusters, but failed to classify 13% of cells in the C0 cluster (Supplementary Fig. 8b). Cicero was also applied to cluster these cells, however since it failed to predict the activities of several critical genes (e.g., *Neurod6*, *C1qb* and *Ctss*), the algorithm was not able to distinguish 20% of cells in clusters C0~C2 (Supplementary Fig. 8c). ChromVAR generated 11 cell clusters, but it mis-clustered AC and EX4 with inhibitory neurons (Supplementary Fig. 8d). These results confirm that APEC can better distinguish and categorize single cells with great sensitivity and reliability.

APEC constructs a pseudotime trajectory that predicts cell differentiation lineage

Cells are not static but dynamic entities, and they have a history, particularly a developmental history. Although single-cell experiments often profile a momentary snapshot, a number of remarkable computational algorithms have been developed to pseudo-order cells based on the different points they were assumed to occupy in a trajectory, thereby leveraging biological asynchrony^{36, 37}. For instance, Monocle^{37, 38} constructs the minimum spanning tree, and Wishbone³⁹ and SPRING⁴⁰ construct the nearest neighbor graph from single-cell transcriptome profiles. These tools have been widely used to depict neurogenesis⁴¹, hematopoiesis^{42, 43} and reprogramming⁴⁴. APEC integrates the Monocle algorithm into the accesson-based method and enables pseudotime prediction from scATAC-seq data²⁰ and was applied to investigate HSC differentiation lineages (Fig. 3a). Principal component analysis (PCA) of the accesson matrix revealed multiple stages of the lineage during HSC differentiation (Fig. 3b) and was consistent with previous publications^{2, 20}. After utilizing the Monocle package, APEC provided more precise pathways from HSCs to the differentiated cell types (Fig. 3c). In addition to the differentiation pathways to MEP cells through the CMP state and to CLP cells through the LMPP state, MPP cells may differentiate into GMP cells through two distinct trajectories: Path A through the CMP state and Path B through the LMPP state, which is consistent with the composite model of HSC and blood lineage commitment⁴⁵. Notably, pDCs from bone marrow are CD34⁺ (Supplementary Fig. 9), indicative of precursors of pDCs. APEC suggested that pDC precursors were derived from CLP cells on the pseudotime trajectory (Fig. 3c), which also agrees with previous reports⁴⁶. Furthermore, APEC incorporated the chromVAR algorithm to determine the regulatory mechanisms during HSC differentiation by evaluating the deviation of each TF along the single-

cell trajectory. As expected, the HOX motif is highly enriched in the accessible sites of HSCs/MPP cells, as are the GATA1, CEBPB and TCF4 motifs, which exhibit gradients that increase along the erythroid, myeloid and lymphoid differentiation pathways, respectively²⁰ (Fig. 3d). We also noticed that the TF regulatory strategies of the two paths from MPP towards GMP cells were very different. In addition, the 3 CMP sub-clusters identified in Fig. 1 were differentially distributed along the developmental trajectory (Fig. 3e). CMP1 cells that close to HSCs and MPPs are early stage CMPs; CMP2 cells are distributed in both the GMP and MEP branches; CMP-MEP cells are MEP committed CMPs and are dominantly distributed in the MEP differentiation branch. The distributions of these CMP sub-clusters are also consistent with the functions of their enriched motifs mentioned in the first section (Fig. 1d)²²⁻²⁴. Finally, we generated a hematopoiesis tree based on the APEC analysis (Fig. 3f).

Furthermore, we benchmarked the performance of APEC and of all the other tools in constructing a pseudotime trajectory from the scATAC-seq profile on the same dataset. We found that (1) when the raw peak count matrix was invoked into Monocle, almost none developmental pathways were constructed (Supplementary Fig. 10a), suggesting that the peak aggregation step in APEC greatly improves the pseudotime estimation; (2) APEC + Monocle provides the most precise pathways from HSCs to differentiated cells, compared to those of other peak aggregation methods, such as chromVAR, Cicero, LSI and cisTopic (Supplementary Fig. 10b-e); and (3) when we applied other pseudotime trajectory construction methods, such as SPRING⁴⁰, after APEC, a similar though less clear cell differentiation diagram was also obtained, suggesting the reliability of our prediction (Supplementary Fig. 10f).

APEC reveals the single-cell regulatory heterogeneity of thymocytes

T cells generated in the thymus play a critical role in the adaptive immune system, and the development of thymocytes can be divided into 3 main stages based on the expression of the surface markers CD4 and CD8, namely, CD4 CD8 double-negative (DN), CD4 CD8 double-positive (DP) and CD4 or CD8 single-positive (CD4SP or CD8SP, respectively) stages⁴⁷. However, due to technical limitations, our genome-wide understanding of thymocyte development at single-cell resolution remains unclear. Typically, more than 80% of thymocytes stay in the DP stage in the thymus, whereas DN cells account for only approximately 3% of the thymocyte population. To eliminate the impacts of great differences in proportion, we developed a fluorescent tagmentation- and FACS-sorting-based scATAC-seq strategy (ftATAC-seq), which combined the advantages of ATAC-seq⁴⁸ and Pi-ATAC-seq⁴⁹ to manipulate the desired number of target cells by indexed

sorting (Fig. 4a). Tn5 transposomes were fluorescently labeled in each cell to evaluate the tagmentation efficiency so that cells with low ATAC signals could be gated out easily (Fig. 4b, Supplementary Fig. 11a). With ftATAC-seq, we acquired high-quality chromatin accessibility data for 352 index-sorted DN, DP, CD4SP, and CD8SP single cells and 352 mixed thymocytes (Supplementary Fig. 11b-d). Correlation analysis with the published bulk ATACs-eq data of thymocytes⁵⁰ indicates that the cells we sorted in ftATAC-seq were correctly labeled (Supplementary Fig. 11e). We then applied APEC on this dataset to investigate the chromatin accessibility divergence during developmental process and to reveal refined regulome heterogeneity of mouse thymocytes at single-cell resolution. Taking into account of all the 130685 peaks called from the raw sequencing data, APEC aggregated 600 accessions and successfully assigned over 82% of index-sorted DN, DP, CD4SP and CD8SP cells into the correct subpopulations (Fig. 4c-d). As expected, the majority of randomly sorted and mixed thymocytes were classified into DP subtypes based on hierarchical clustering of cell-cell correlation matrix, which was consistent with the cellular subtype proportions in the thymus. APEC further classified all thymocytes into 11 subpopulations, including 2 DN, 6 DP, 1 CD4SP, 2 CD8SP, suggesting that extensive epigenetic heterogeneity exists among cells with the same CD4 and CD8 surface markers (Fig. 4e). For instance, there are four main subtypes of DN cells, according to the expression of the surface markers CD44 and CD25⁵¹, while two clusters were identified in ftATAC-seq. The accessibility signals around the *Il2ra* (*Cd25*) and *Cd44* gene loci demonstrated that DN.A1 comprised CD44⁺CD25⁻ and CD44⁺CD25⁺ DN subtypes (DN1 and DN2), and DN.A2 cells comprised CD44⁻CD25⁺ and CD44⁻CD25⁻ subtypes (DN3 and DN4), suggesting significant chromatin changes between DN2 and DN3 cell development (Fig. 4f).

Many TFs have been reported to be essential in regulating thymocyte development, and we found that their motifs were remarkably enriched at different stages during the process (Fig. 4g). For instance, Runx3 is well known for regulating CD8SP cells⁵², and we observed significant enrichment of the RUNX motif on DN cells and a group of CD8SP cells. Similarly, the TCF^{53, 54}, RORC⁵⁵ and NFkB⁵⁶ family in regulating the corresponding stages during this process. More enriched TF motifs in each cell subpopulation were also observed, suggesting significant regulatory divergence in thymocytes (Supplementary Fig. 12a). Interestingly, two clusters of CD8SP cells appear to be differentially regulated based on motif analysis, in which CD8.A1 cells are closer to DP cells, while CD8.A2 cells are more distant at the chromatin level, suggesting that CD8.A2 cells are more mature CD8SP cells, and CD8.A1 cells are in a transitional state between DP and SP cells.

APEC is capable of integrating single-cell transcriptional and epigenetic information by scoring gene sets of interest based on their nearby peaks from scATAC-seq, thereby converting the chromatin accessibility signals to values that are comparable to gene expression profiles (see **Methods**). To test the performance of this integrative analysis approach and to evaluate the accuracy of thymocyte classification by APEC, we assayed the transcriptomes of single thymocytes and obtained 357 high-quality scRNA-seq profiles using the SMART-seq2 protocol⁵⁷. Unsupervised analysis of gene expression profiles clustered these thymocytes into 13 groups in Seurat¹³ ([Supplementary Fig. 13a-b](#)), and each subpopulation was identified based on known feature genes ([Supplementary Fig. 13c-d](#)). We then adopted fisher's exact test on the shared differential genes in cell clusters identified from scATAC-seq and scRNA-seq profiles (see **Methods**), and observed a strong correlation between the subtypes identified from the transcriptome and those from chromatin accessibility ([Fig. 4h](#)), confirming the reliability and stability of cellular classification using APEC.

APEC reconstructs the thymocyte developmental trajectory

APEC is capable of constructing a pseudotime trajectory and then predicting the cell differentiation lineage from a “snapshot” of single-cell epigenomes ([Fig. 3](#)). We applied APEC to recapitulate the developmental trajectory and thereby reveal the single-cell regulatory dynamics during the maturation of thymocytes. Pseudotime analysis based on single-cell ATAC-seq data shaped thymocytes into 5 developing stages ([Fig. 5a](#), [Supplementary Fig. 14a-b](#)), where most of the cells in stages 1, 2, 4, and 5 were DN, DP, CD8SP and CD4SP cells, respectively. APEC also identified a transitional stage 3, which was mainly consisted of last stages of DP cells. Besides Monocle, a similar developmental pathways can also be constructed by SPRING⁴⁰ based on the accession matrix ([Supplementary Fig 14c](#)). Interestingly, the pseudotime trajectory suggests three developmental pathways for this process, one of which started with stage 1 (DN) and ended in stage 2 (DP), and the other two of which started with stage 1 (DN), went through a transitional stage 3 and then bifurcated into stage 4 (CD8SP) and 5 (CD4SP). The predicted developmental trajectory could also be confirmed by the gene expression of surface markers, such as Cd4, Cd8, Runx3 and Ccr7 ([Fig. 5b](#)). To evaluate the gene ontology (GO) enrichments over the entire process, we implemented an accession-based GO module in APEC, which highlights the significance of the association between cells and biological function ([Fig. 5c](#)). For instance, T cells selections, including β -selection, positive selection and negative selection, are initiated in the late DN stage. Consistent with this process, we observed a strong “T cell selection” GO term on the

trajectory path after DN.A1 ([Supplemental Fig. 14d](#)). Since TCR signals are essential for T cell selection, we also observed the “T cell activation” GO term accompanied by “T cell selection”. Meanwhile, the signal for regulation of protein binding was found decreased at SP stages, indicating the necessity of weak TCR signal for the survival of SP T cells during negative selection.

To further uncover the regulatory mechanism underlying this developmental process, APEC was implemented to identify stage-specific enriched TFs along the trajectory and pinpoint the “pseudotime” at which the regulation occurs. In addition to the well-studied TFs mentioned above ([Fig. 4g](#), [Supplemental Fig. 12a](#)), APEC also identified Zeb1⁵⁸, Ctf⁵⁹ and Id4 as potential stage-specific regulators ([Fig. 5d](#)). Interestingly, the Id4 motif enriched on DP cells was also reported to regulate apoptosis in other cell types^{60, 61}. Associated with the fact that the vast majority of DP thymocytes die because of a failure of positive selection⁶², we hypothesize that stage 2 may be the path towards DP cell apoptosis. We then checked the distribution of DP cells along the stage 2 trajectory and found that most DP.A1 cells were scattered in “early” stage 2, and they were enriched with GO terms such as “T cell selection”, “cell activation” and “differentiation” ([Fig. 5e](#), [Supplementary Fig. 14e](#)). However, most DP.A3/4/5 cells were distributed at the end of stage 2, and their principle accessions were enriched with GO terms such as “apoptosis” and “chromatin modification”. Although it is believed that more than 95% of DP thymocytes die during positive selection, only a small proportion of apoptotic cells could be detected in a snapshot of the thymus, which in our data are the cells at the end of stage 2. By comparing the number of cells near stage 3 with all the cells in stage 2, we estimated that ~3-5% of cells would survive positive selection, which is consistent with the findings reported in previous publications^{63, 64}. Our data suggest that before entering the final apoptotic stage, DP thymocytes under selection could have already been under apoptotic stress at the chromatin level, which explains why DP cells are more susceptible to apoptosis than other thymocyte subtypes⁶⁵.

DISCUSSION

Here, we introduced an accession-based algorithm for single-cell chromatin accessibility analysis. Without any prior information (such as motifs), this approach generated more refined cell groups with reliable biological functions and properties. Integrating the new algorithm with all necessary chromatin sequencing data processing tools, APEC provides a comprehensive solution for transforming raw experimental single-cell data into final visualized results. In addition to improving the clustering of subtle cell subtypes, APEC is also capable of locating potential specific super-enhancers, searching enriched motifs, estimating gene activities, and constructing

time-dependent cell developmental trajectories, and it is compatible with many existing single-cell accessibility datasets. Compared with all the other state-of-the-art single cell chromatin accessibility analysis methods, APEC clearly shows superiority in correctly predicting cell identities and precisely constructing developmental trajectories, and provides new biological insights that no other tools can. APEC is also very robust and stable and is scalable to clustering a large number of cells using limited computational resources. Despite these advantages, the biological implications of accessions are still obscure, especially for those that involve only a small number of peaks. Although we noticed peaks in the same accession may belong to the same TADs, further investigations are still required to fully uncover the biology that underlies accessions. Another caveat of APEC is that Monocle can be very sensitive to the input data, and thereby pseudotime trajectory predictions from Monocle are better to be confirmed by multiple algorithms with similar function.

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

Acknowledgments

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

Authors' contributions

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

Data and code availability

Mouse thymocytes ftATAC-seq data can be obtained from the Genome Sequence Archive at BIG Data Center with the accession number CRA001267 and is available via <http://bigd.big.ac.cn/gsa/s/yp1164Et>. Other published data sets used in this study are available from NIH GEO with accession numbers GSE74310², GSE65360⁸, GSE96772²⁰, GSE100033¹⁰, GSE111586³³, and GSE63525⁶⁶. APEC pipeline can be downloaded from the GitHub website (<https://github.com/QuKunLab/APEC>).

REFERENCES

1. Buenrostro, J.D., Giresi, P.G., Zaba, L.C., Chang, H.Y. & Greenleaf, W.J. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nat Methods* **10**, 1213-1218 (2013).
2. Corces, M.R. et al. Lineage-specific and single-cell chromatin accessibility charts human hematopoiesis and leukemia evolution. *Nat Genet* **48**, 1193-1203 (2016).
3. Wu, J. et al. The landscape of accessible chromatin in mammalian preimplantation embryos. *Nature* **534**, 652-657 (2016).
4. Jorstad, N.L. et al. Stimulation of functional neuronal regeneration from Muller glia in adult mice. *Nature* **548**, 103-107 (2017).
5. Su, Y. et al. Neuronal activity modifies the chromatin accessibility landscape in the adult brain. *Nat Neurosci* **20**, 476-483 (2017).
6. Denny, S.K. et al. Nfib Promotes Metastasis through a Widespread Increase in Chromatin Accessibility. *Cell* **166**, 328-342 (2016).
7. Qu, K. et al. Chromatin Accessibility Landscape of Cutaneous T Cell Lymphoma and Dynamic Response to HDAC Inhibitors. *Cancer Cell* **32**, 27-41 e24 (2017).
8. Buenrostro, J.D. et al. Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature* **523**, 486-490 (2015).
9. Lake, B.B. et al. Integrative single-cell analysis of transcriptional and epigenetic states in the human adult brain. *Nat Biotechnol* **36**, 70-80 (2018).
10. Preissl, S. et al. Single-nucleus analysis of accessible chromatin in developing mouse forebrain reveals cell-type-specific transcriptional regulation. *Nat Neurosci* **21**, 432-439 (2018).
11. Cusanovich, D.A. et al. Multiplex single cell profiling of chromatin accessibility by combinatorial cellular indexing. *Science* **348**, 910-914 (2015).

12. Cusanovich, D.A. et al. The cis-regulatory dynamics of embryonic development at single-cell resolution. *Nature* **555**, 538-542 (2018).
13. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat Biotechnol* **36**, 411-420 (2018).
14. Kiselev, V.Y. et al. SC3: consensus clustering of single-cell RNA-seq data. *Nat Methods* **14**, 483-486 (2017).
15. Wang, B., Zhu, J., Pierson, E., Ramazzotti, D. & Batzoglou, S. Visualization and analysis of single-cell RNA-seq data by kernel-based similarity learning. *Nat Methods* **14**, 414-416 (2017).
16. Wolf, F.A., Angerer, P. & Theis, F.J. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol* **19**, 15 (2018).
17. Schep, A.N., Wu, B., Buenrostro, J.D. & Greenleaf, W.J. chromVAR: inferring transcription-factor-associated accessibility from single-cell epigenomic data. *Nat Methods* **14**, 975-978 (2017).
18. Pliner, H.A. et al. Cicero Predicts cis-Regulatory DNA Interactions from Single-Cell Chromatin Accessibility Data. *Mol Cell* **71**, 858-871 e858 (2018).
19. Bravo González-Blas, C. et al. Cis-topic modelling of single cell epigenomes. *bioRxiv* (2018).
20. Buenrostro, J.D. et al. Integrated Single-Cell Analysis Maps the Continuous Regulatory Landscape of Human Hematopoietic Differentiation. *Cell* **173**, 1535-1548 e1516 (2018).
21. Zamanighomi, M. et al. Unsupervised clustering and epigenetic classification of single cells. *Nat Commun* **9**, 2410 (2018).
22. Ng, A.P. et al. Erg is required for self-renewal of hematopoietic stem cells during stress hematopoiesis in mice. *Blood* **118**, 2454-2461 (2011).
23. Ong, C.T. & Corces, V.G. CTCF: an architectural protein bridging genome topology and function. *Nature Reviews Genetics* **15**, 234-246 (2014).
24. Dore, L.C. & Crispino, J.D. Transcription factor networks in erythroid cell and megakaryocyte development. *Blood* **118**, 231-239 (2011).
25. Tothova, Z. et al. FoxOs are critical mediators of hematopoietic stem cell resistance to physiologic oxidative stress. *Cell* **128**, 325-339 (2007).
26. Akashi, K., Traver, D., Miyamoto, T. & Weissman, I.L. A clonogenic common myeloid progenitor that gives rise to all myeloid lineages. *Nature* **404**, 193-197 (2000).
27. Ryncarz, R.E. & Anasetti, C. Expression of CD86 on human marrow CD34(+) cells identifies immunocompetent committed precursors of macrophages and dendritic cells. *Blood* **91**, 3892-3900 (1998).
28. Georgopoulos, K. et al. The Ikaros Gene Is Required for the Development of All Lymphoid Lineages. *Cell* **79**, 143-156 (1994).
29. van der Meer, L.T., Jansen, J.H. & van der Reijden, B.A. Gfi1 and Gfi1b: key regulators of hematopoiesis. *Leukemia* **24**, 1834-1843 (2010).
30. Blanc, L. et al. The water channel aquaporin-1 partitions into exosomes during reticulocyte maturation: implication for the regulation of cell volume. *Blood* **114**, 3928-3934 (2009).
31. Oberst, A. et al. The Nedd4-binding partner 1 (N4BP1) protein is an inhibitor of the E3 ligase Itch. *Proc Natl Acad Sci U S A* **104**, 11280-11285 (2007).
32. Eguchi, M. et al. GPHN, a novel partner gene fused to MLL in a leukemia with t(11;14)(q23;q24). *Genes Chromosomes Cancer* **32**, 212-221 (2001).
33. Cusanovich, D.A. et al. A Single-Cell Atlas of In Vivo Mammalian Chromatin Accessibility. *Cell* **174**, 1309-+ (2018).
34. Satpathy, A.T. et al. Transcript-indexed ATAC-seq for precision immune profiling. *Nat Med* **24**, 580-590 (2018).

35. Hrvatin, S. et al. Single-cell analysis of experience-dependent transcriptomic states in the mouse visual cortex. *Nat Neurosci* **21**, 120-129 (2018).
36. Bendall, S.C. et al. Single-cell trajectory detection uncovers progression and regulatory coordination in human B cell development. *Cell* **157**, 714-725 (2014).
37. Trapnell, C. et al. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, 381-386 (2014).
38. Qiu, X. et al. Reversed graph embedding resolves complex single-cell trajectories. *Nat Methods* **14**, 979-982 (2017).
39. Setty, M. et al. Wishbone identifies bifurcating developmental trajectories from single-cell data. *Nat Biotechnol* **34**, 637-645 (2016).
40. Weinreb, C., Wolock, S. & Klein, A.M. SPRING: a kinetic interface for visualizing high dimensional single-cell expression data. *Bioinformatics* **34**, 1246-1248 (2018).
41. Habib, N. et al. Div-Seq: Single-nucleus RNA-Seq reveals dynamics of rare adult newborn neurons. *Science* **353**, 925-928 (2016).
42. Olsson, A. et al. Single-cell analysis of mixed-lineage states leading to a binary cell fate choice. *Nature* **537**, 698-+ (2016).
43. Zhou, F. et al. Tracing haematopoietic stem cell formation at single-cell resolution. *Nature* **533**, 487-+ (2016).
44. Treutlein, B. et al. Dissecting direct reprogramming from fibroblast to neuron using single-cell RNA-seq. *Nature* **534**, 391-+ (2016).
45. Adolfsson, J. et al. Identification of Flt3+ lympho-myeloid stem cells lacking erythro-megakaryocytic potential a revised road map for adult blood lineage commitment. *Cell* **121**, 295-306 (2005).
46. Chistiakov, D.A., Orekhov, A.N., Sobenin, I.A. & Bobryshev, Y.V. Plasmacytoid dendritic cells: development, functions, and role in atherosclerotic inflammation. *Front Physiol* **5**, 279 (2014).
47. Germain, R.N. T-cell development and the CD4-CD8 lineage decision. *Nat Rev Immunol* **2**, 309-322 (2002).
48. Chen, X. et al. ATAC-seq reveals the accessible genome by transposase-mediated imaging and sequencing. *Nat Methods* **13**, 1013-1020 (2016).
49. Chen, X. et al. Joint single-cell DNA accessibility and protein epitope profiling reveals environmental regulation of epigenomic heterogeneity. *Nat Commun* **9**, 4590 (2018).
50. Yoshida, H. et al. The cis-Regulatory Atlas of the Mouse Immune System. *Cell* **176**, 897-912 e820 (2019).
51. Godfrey, D.I., Kennedy, J., Suda, T. & Zlotnik, A. A developmental pathway involving four phenotypically and functionally distinct subsets of CD3-CD4-CD8- triple-negative adult mouse thymocytes defined by CD44 and CD25 expression. *J Immunol* **150**, 4244-4252 (1993).
52. Taniuchi, I. et al. Differential requirements for Runx proteins in CD4 repression and epigenetic silencing during T lymphocyte development. *Cell* **111**, 621-633 (2002).
53. Ioannidis, V., Beermann, F., Clevers, H. & Held, W. The beta-catenin--TCF-1 pathway ensures CD4(+)-CD8(+) thymocyte survival. *Nat Immunol* **2**, 691-697 (2001).
54. Yu, S. et al. The TCF-1 and LEF-1 transcription factors have cooperative and opposing roles in T cell development and malignancy. *Immunity* **37**, 813-826 (2012).
55. Sun, Z. et al. Requirement for RORgamma in thymocyte survival and lymphoid organ development. *Science* **288**, 2369-2373 (2000).
56. Gerondakis, S., Fulford, T.S., Messina, N.L. & Grumont, R.J. NF-kappaB control of T cell development. *Nat Immunol* **15**, 15-25 (2014).
57. Picelli, S. et al. Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat Methods* **10**, 1096-1098 (2013).

58. Higashi, Y. et al. Impairment of T cell development in deltaEF1 mutant mice. *J Exp Med* **185**, 1467-1479 (1997).
59. Heath, H. et al. CTCF regulates cell cycle progression of alphabeta T cells in the thymus. *EMBO J* **27**, 2839-2850 (2008).
60. Andres-Barquin, P.J., Hernandez, M.C. & Israel, M.A. Id4 expression induces apoptosis in astrocytic cultures and is down-regulated by activation of the cAMP-dependent signal transduction pathway. *Exp Cell Res* **247**, 347-355 (1999).
61. Carey, J.P., Knowell, A.E., Chinaranagari, S. & Chaudhary, J. Id4 promotes senescence and sensitivity to doxorubicin-induced apoptosis in DU145 prostate cancer cells. *Anticancer Res* **33**, 4271-4278 (2013).
62. Surh, C.D. & Sprent, J. T-cell apoptosis detected in situ during positive and negative selection in the thymus. *Nature* **372**, 100-103 (1994).
63. Huesmann, M., Scott, B., Kisielow, P. & von Boehmer, H. Kinetics and efficacy of positive selection in the thymus of normal and T cell receptor transgenic mice. *Cell* **66**, 533-540 (1991).
64. Shortman, K., Vremec, D. & Egerton, M. The kinetics of T cell antigen receptor expression by subgroups of CD4+8+ thymocytes: delineation of CD4+8+3(2+) thymocytes as post-selection intermediates leading to mature T cells. *J Exp Med* **173**, 323-332 (1991).
65. Chow, S.C., Snowden, R., Orrenius, S. & Cohen, G.M. Susceptibility of different subsets of immature thymocytes to apoptosis. *Febs Lett* **408**, 141-146 (1997).
66. Rao, S.S. et al. A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell* **159**, 1665-1680 (2014).

FIGURES

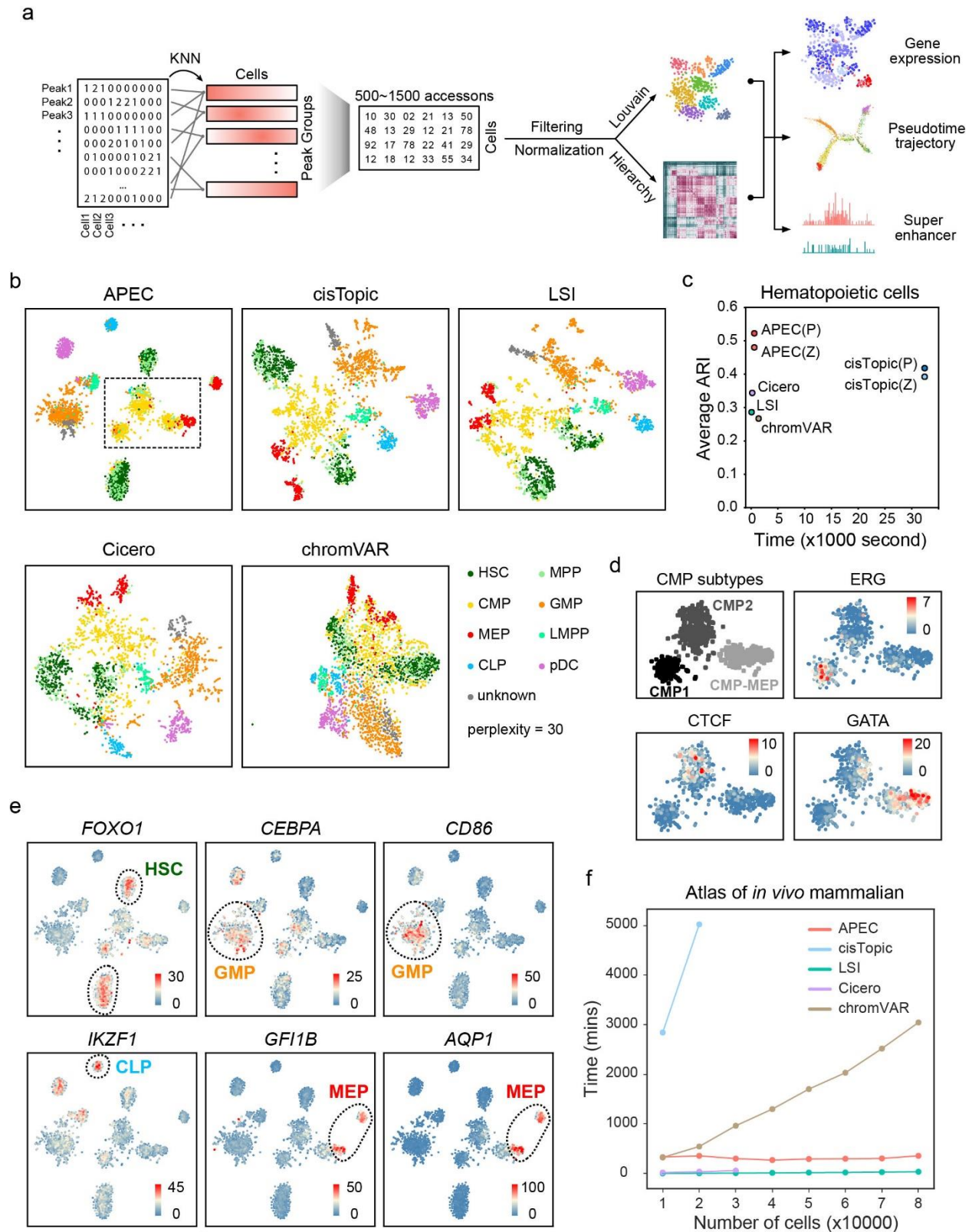


Figure 1. The accession matrix constructed from the sparse fragment count matrix improved the clustering of scATAC-seq data. **(a)** Step-by-step workflow of APEC. Peaks were grouped into accessions by their accessibility pattern among cells with the K nearest neighbors (KNN) method. **(b)** t-Distributed Stochastic Neighbor Embedding (tSNE) diagrams of the hematopoietic single cells dataset based on the dimension-transformed matrices from different algorithms, i.e., APEC: accession matrix; cisTopic: topic matrix; LSI: LSI matrix; chromVAR: bias-corrected deviation matrix; Cicero: aggregated model matrix. The cells are FACS-indexed human hematopoietic cells, including HSCs (hematopoietic stem cells), MPPs (multipotent progenitors), LMPPs (lymphoid-primed multipotential progenitors), CMPs (common myeloid progenitors), CLPs (common lymphoid progenitors), pDCs (plasmacytoid dendritic cells), GMPs (granulocyte-macrophage progenitors), MEPs (megakaryocyte-erythroid progenitors), and unknown cells. **(c)** The average ARI (Adjusted Rand Index) scores and computing time for the clustering of the human hematopoietic cells by different algorithms. Same as the two normalization methods applied in cisTopic, we normalized the accession matrix in APEC based on probability (P) and z-score (Z). CisTopic was performed using 5 CPU threads and all other tools with 1 CPU thread. **(d)** Three CMP subtypes identified in APEC and the motifs enriched in each cell subtype. **(e)** Gene expressions predicted by APEC for all cells. Predicted expressions of marker genes for each cell type including *FOXO1* (marker for HSC), *CEBPA/CD86* (markers for GMP), *IKZF1* (marker for CLP), and *GFI1B/AQP1* (markers for MEP) were shown. **(f)** The computing time required for different algorithms to cluster the data with cell numbers from 10,000 to 80,000, sampled from the mouse *in vivo* single-cell chromatin accessibility atlas dataset.

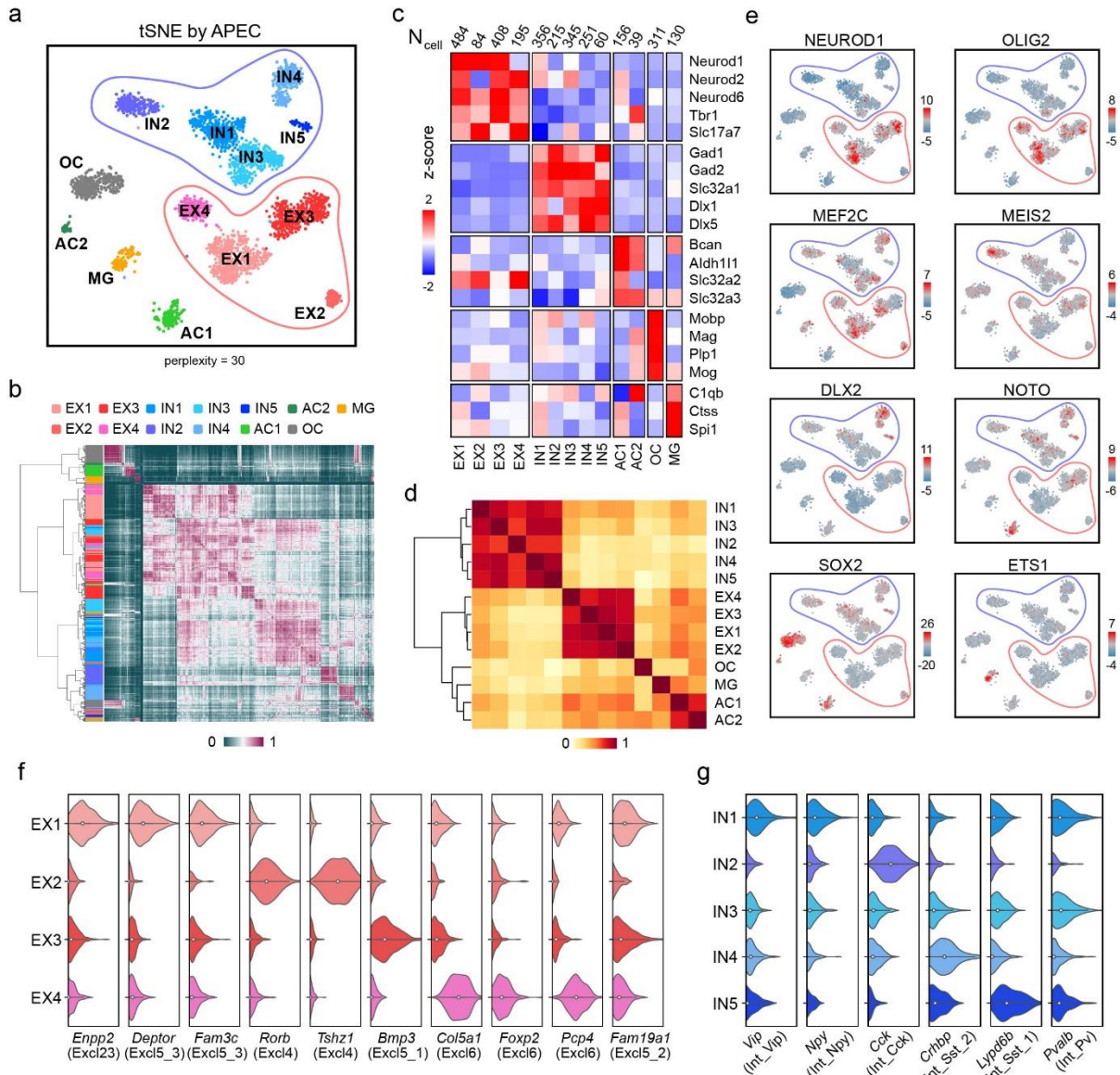


Figure 2. APEC improved the cell type classification of adult mouse forebrain snATAC-seq data. **(a)** A tSNE diagram demonstrates the APEC clustering of forebrain cells. **(b)** Hierarchical clustering of the cell-cell correlation matrix. The side bar denotes cell clusters from APEC. **(c)** Average scores of the marker genes for each cell cluster generated by the method mentioned in the data source paper¹⁰, and normalized by the standard score (z-score). The top row lists the number of cells in each cluster. **(d)** Hierarchical clustering of the cluster-cluster correlation matrix. **(e)** Differential enrichments of cell type-specific motifs in each cluster. **(f, g)** Intensities of representative accessions of the excitatory (EX) and inhibitory (IN) neuron subtypes from snATAC-seq associated with the activities of signature genes of the excitatory (Excl) and inhibitory (Int) neuron subtypes defined in inDrops-seq³⁵.

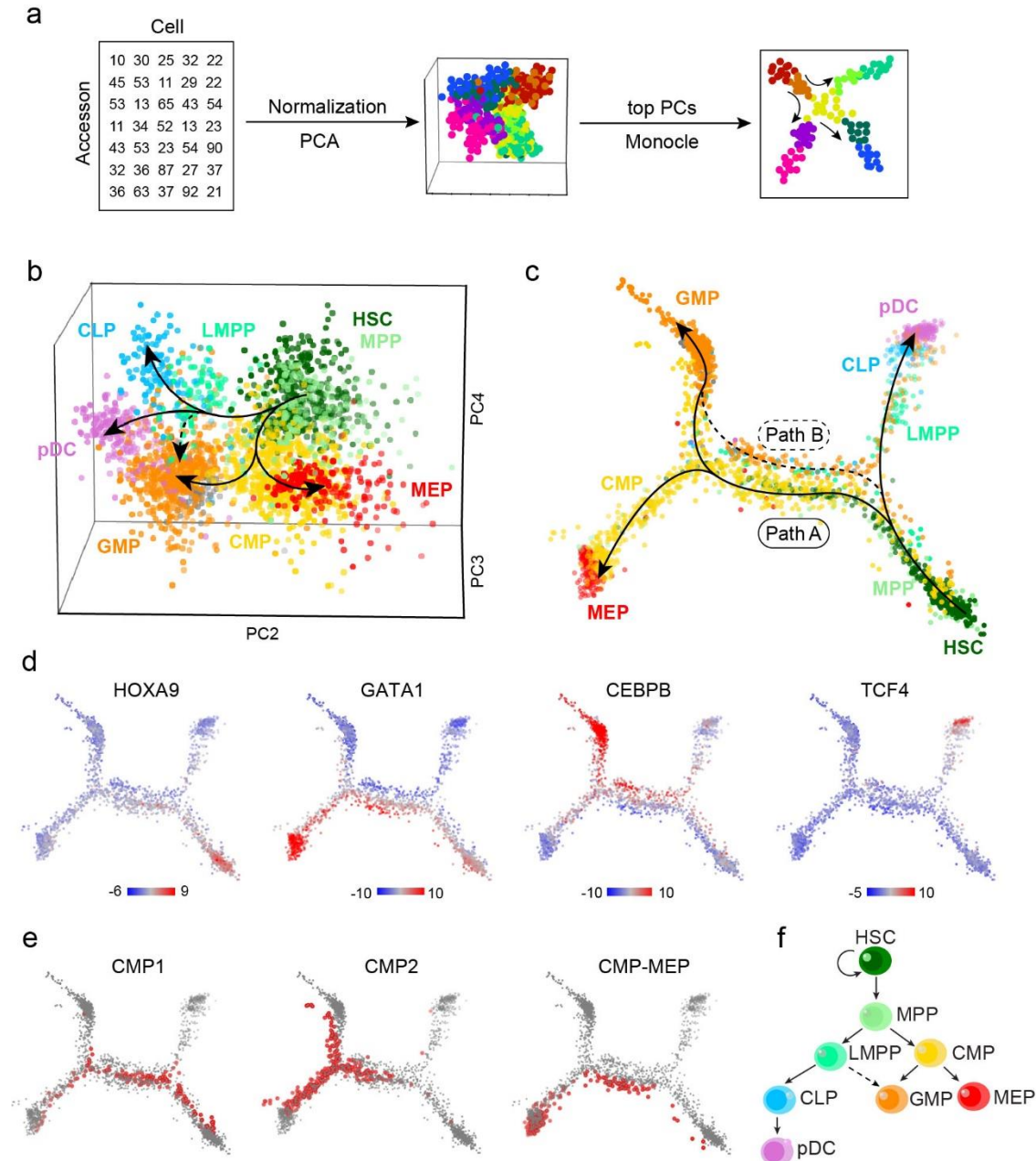


Figure 3. APEC constructed a differentiation pathway from scATAC-seq data from human hematopoietic cells. **(a)** The pseudotime trajectory construction scheme based on the accession matrix and Monocle. **(b)** Principal component analysis (PCA) of the accession matrix for human hematopoietic cells. The first principal component is not shown here because it was highly correlated with sequencing depth²⁰. HSC, hematopoietic stem cell; MPP, multipotent progenitor; LMPP, lymphoid-primed multipotential progenitor; CMP, common myeloid progenitor; CLP, common lymphoid progenitor; pDC, plasmacytoid dendritic cell; GMP, granulocyte-macrophage progenitor; MEP, megakaryocyte-erythroid progenitor; unknown, unlabeled cells. **(c)** Pseudotime

trajectory for the same data constructed by applying Monocle on the accession matrix. Paths A and B represent different pathways for GMP cell differentiation. **(d)** The deviations of significant differential motifs (HOXA9, GATA1, CEBPB, and TCF4) plotted on the pseudotime trajectory. **(e)** Distributions of the CMP sub-clusters on the trajectory. **(f)** Modified schematic of human hematopoietic differentiation.

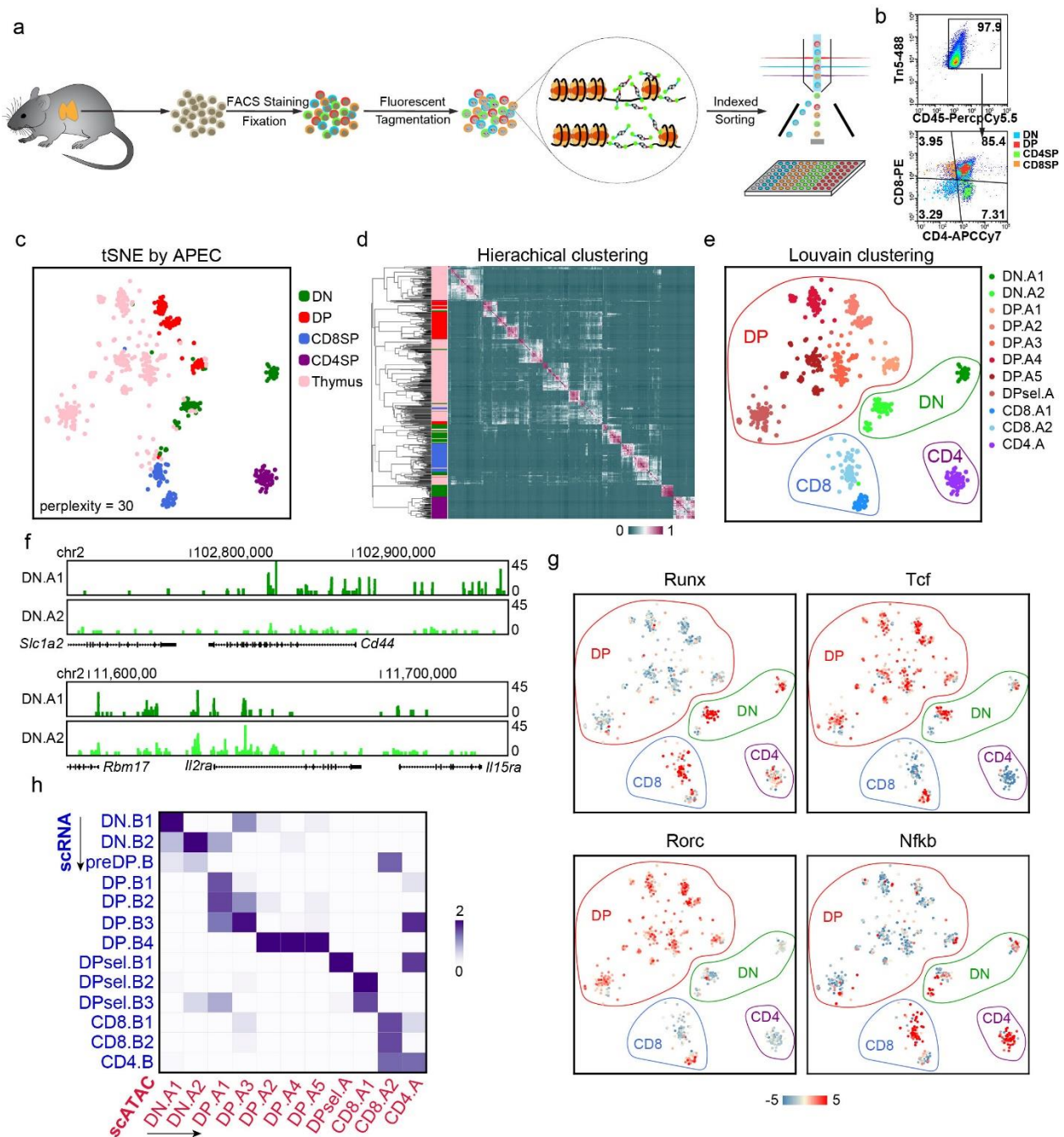


Figure 4. APEC accurately identified cell subtypes based on scATAC-seq data from *Mus musculus* thymocytes. **(a)** Experimental workflow of the fluorescent tagmentation- and FACS-sorting-based scATAC-seq strategy (ftATAC-seq). **(b)** Indexed sorting of double-negative (DN), double-positive (DP), CD4⁺ single-positive (CD4SP), and CD8⁺ single-positive (CD8SP) cells with strong tagmentation signals. **(c)** The tSNE of thymocyte single-cell ftATAC-seq data based on the accession matrix, in which the cells are labeled by the sorting index. **(d)** Hierarchical clustering of the cell-cell correlation matrix. On the sidebar, each cell was colored by the sorting index. **(e)** The

accession-based Louvain method clustered thymocytes into 11 subtypes. DN.A1 (dark green) & A2 (light green), double-negative clusters; DP.A1~A5 and DPsel.A, double-positive clusters; CD8.A1 (dark blue) & A2 (light blue), CD8⁺ single-positive clusters; CD4.A (purple), CD4⁺ single-positive cluster. **(f)** Average fragment counts of two DN clusters around the marker genes *Cd44* and *Il2ra*. **(g)** Differential enrichment of the motifs Runx, Tcf, Rorc, and Nfkb in the cell clusters. **(h)** Z-score of the Fisher exact test $-\log(p\text{-value})$ of the common differential genes between the cell clusters from different experiments. The row and column clusters were identified by data from single-cell transcriptome (SMART-seq) and chromatin accessibility (ftATAC-seq) analysis respectively.

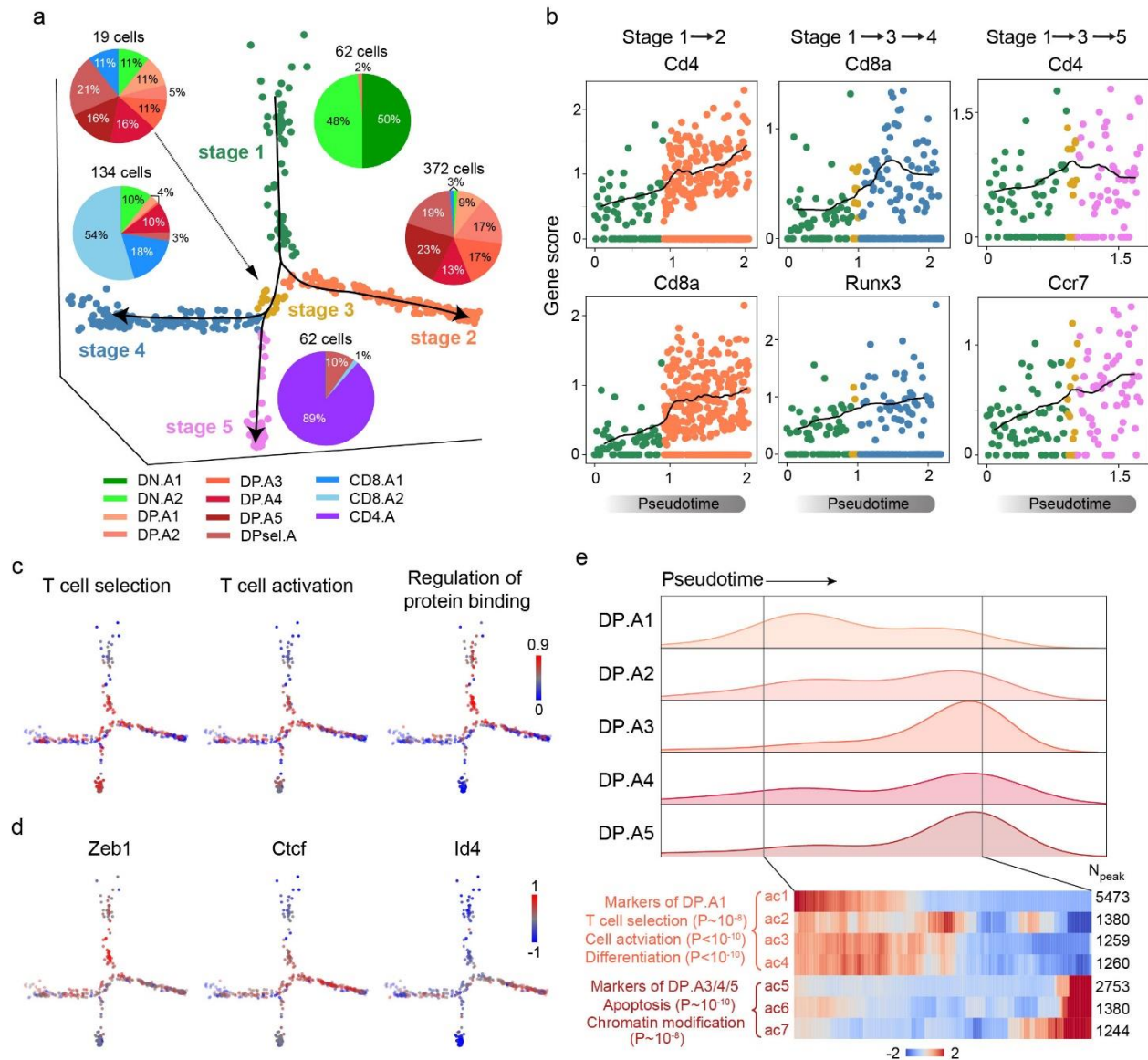


Figure 5. APEC depicted the developmental pathways of *Mus musculus* thymocytes by pseudotime analysis. **(a)** Pseudotime trajectory based on the accession matrix of thymocyte ftATAC-seq data. Cell colors were defined by the developmental stages along pseudotime. Pie charts show the proportion of cell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

METHODS

Mice. C57BL/6 mice were purchased from Beijing Vital River Laboratory Animal Technology and maintained under specific pathogen-free conditions until the time of experiments. All mouse experiments in this study were reviewed and approved by the Institutional Animal Care and Use Committee of the University of Science and Technology of China.

ftATAC-seq on mouse thymocytes. Alexa fluor 488-labeled adaptor oligonucleotides were synthesized at Sangon Biotech as follows: Tn5ME, 5'-[phos]CTGTCTCTTATACACATCT-3'; AF488-R1, 5'-AF488-TCGTCCGCGAGCGTCAGATGTGTATAAGAGACAG-3'; and AF488-R2, 5'-AF488-GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAG-3'. Then, 50 μ M of AF488-R1/Tn5ME and AF488-R2/Tn5ME were denatured separately in TE buffer (Qiagen) at 95 °C for 5 min and cooled down to 22 °C at 0.1 °C/s. AF488-labeled adaptors were assembled onto Robust Tn5 transposase (Robustnique) according to the user manual to form fluorescent transposomes.

Thymus tissues isolated from 6- to 8-week-old male mice were gently ground in 1 mL of RPMI-1640. Thymocytes in a single-cell suspension were counted after passing through a 40 μ m nylon mesh. A total of 1×10^6 thymocytes were stained with PerCP-Cy5.5-anti-CD45, PE-anti-CD8a and APC-Cy7-anti-CD4 antibodies (Biolegend) and then fixed in 1 \times PBS containing 1% methanal at room temperature for 5 min. After washing twice with 1 \times PBS, the cells were counted again. A total of 1×10^5 fixed cells were resuspended in 40 μ L of 1 \times TD buffer (5 mM Tris-HCl, pH 8.0, 5 mM MgCl₂, and 10% DMF) containing 0.1% NP-40. Then, 10 μ L of fluorescent transposomes were added and mixed gently. Fluorescent tagmentation was conducted at 55 °C for 30 min and stopped by adding 200 μ L of 100 mM EDTA directly to the reaction mixture. The cells were loaded on a Sony SH800S sorter, and single cells of the CD45⁺/AF488-Tn5^{hi} population were index-sorted into each well of 384-well plates. The 384-well plates used to acquire sorted cells were loaded with 2 μ L of release buffer (50 mM EDTA, 0.02% SDS) before use. After sorting, the cells in the wells were incubated for 1 min. Plates that were not processed immediately were preserved at -80 °C.

To prepare a single-cell ATAC-seq library, plates containing fluorescently tagmented cells were incubated at 55 °C for 30 min. Then, 4.2 μ L of PCR round 1 buffer (1 μ L of 100 μ M MgCl₂, 3 μ L of 2 \times I-5 PCR mix [MCLAB], and 0.1 μ L each of 10 μ M R1 and R2 primers) were added to each well, followed by PCR: 72 °C for 10 min; 98 °C for 3 min; 10 cycles of 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. Thereafter, each well received 4 μ L of PCR round 2 buffer (2 μ L of I-5 PCR Mix, 0.5 μ L each of Ad1 and barcoded Ad2 primers, and

1 μL of ddH₂O), and final PCR amplification 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 1 min; 72 °C for 3 min; and holding at 4 °C. Wells containing different Ad2 barcodes were collected together and purified with a QIAquick PCR purification kit (Qiagen). Libraries were sequenced on an Illumina HiSeq X Ten system.

SMART-seq on thymocytes. Thymocytes were stained and sorted directly into 384-well plates without fixation. SMART-seq was performed as described with some modifications⁶⁷. Reverse transcription and the template-switch reaction were performed at 50 °C for 1 hr with Maxima H Minus Reverse Transcriptase (Thermo Fisher); for library construction, 0.5-1 ng of cDNA was fragmented with 0.05 μL of Robust Tn5 transposome in 20 μL of TD buffer at 55 °C for 10 min, then purified with 0.8x VAHTS DNA Clean Beads (Vazyme Biotech), followed by PCR amplification with Ad1 and barcoded Ad2 primers and purification with 0.6x VAHTS DNA Clean Beads. Libraries were sequenced on an Illumina HiSeq X Ten system.

Data source. All experimental raw data used in this paper are available online. The single-cell data for mouse thymocytes captured by the ftATAC-seq experiment can be obtained from the Genome Sequence Archive at BIG Data Center with the accession number CRA001267 and is available via <http://bigd.big.ac.cn/gsa/s/yp1164Et>. Other published data sets used in this study are available from NIH GEO: (1) scATAC-seq data for LSCs and leukemic blast cells from patients SU070 and SU353, LMPP cells, and monocytes from GSE74310²; (2) scATAC-seq data for HL-60 cells from GSE65360⁸; and (3) scATAC-seq data for hematopoietic development (HSCs, MPPs, CMPs, LMPPs, GMPs, EMPs, CLPs and pDCs) from GSE96772²⁰. (4) APEC is also compatible with a preprocessed fragment count matrix from the snATAC-seq data for the forebrain of adult mice (p56) from GSE100033¹⁰. (5) The computational efficiency of APEC and other methods was tested using data from the single-cell atlas of mouse chromatin accessibility (sciATAC-seq) from GSE111586³³. (6) The scATAC-seq (GSE65360⁸) and Hi-C (GSE63525⁶⁶) data of GM12878 cells were used to generate the spatial correlation of peaks in the same or in different accessions.

Preparing the fragment count matrix from the raw data. APEC adopted the general mapping, alignment, peak calling and motif searching procedures to process the scATAC-seq data from ATAC-pipe⁶⁸. We also implemented the python script in ATAC-pipe⁶⁸ to trim the adapters in the raw data (in paired-end fastq format files for each single-cell sample). APEC used BOWTIE2 to map the trimmed sequencing data to the corresponding genome index and used PICARD for the sorting, duplicate removal, and fragment length counting of the aligned data. The pipeline called peaks from the merged file of all cells by MACS2, ranked and filtered out the low quality peaks

based on the false discovery rate (Q-value). Genomic locations of the peaks were annotated by HOMER, and motifs searched by FIMO. APEC calculates the number of fragments and the percent of reads mapped to the TSS region (± 2000 BP) for each cell, and filters out high quality cells for downstream analysis. All required files for the hg19 and mm10 assembly have been integrated into the pipeline. If users want to process data from other species, they can also download corresponding reference files from the UCSC website. By combining existing tools, APEC made it possible to finish all of the above data processing steps by one command line, and generate a fragment count matrix for subsequent cell clustering and differential analysis. APEC has been made available on GitHub (<https://github.com/QuKunLab/APEC>).

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

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

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

- (2) *Constructing accessions.* The top 40 principal components of the normalized matrix were used to construct the connectivity matrix (\mathbf{C}_{peak}) of peaks by the K-nearest-neighbor (KNN) method with $K=10$. The grouping of peaks is insensitive to the number of principal components and the number of nearest neighbors, so it is usually not necessary to change these two parameters for different data sets. Based on the matrix \mathbf{C}_{peak} , all peaks were grouped by agglomerative clustering with Euclidean distance and Ward linkage method, and the sum of one peak group was an accession. For most data sets, we recommend setting the number of accessions to a value between 500 and 1500, and the default was set to 600, however the cell clustering result is not sensitive to the choice of accession number within this range. We then built the accession count matrix \mathbf{M} by summing the fragment count of all peaks in one accession. Thus, each column of matrix \mathbf{M} is an accession, each row is a cell, and each element represents the cumulative fragment count of each accession in each cell.
- (3) *Accession filtering and normalization.* Not all accessions were used for cell clustering, and those with low dispersion were filtered out to improve cell clustering. The Gini coefficient^{69, 70}

was used to measure the dispersion/inequality of the fragment count numbers of each accession among cells, i.e.,

$$Gini_j = \frac{Mean(|\bar{M}_j \otimes \bar{M}_j|)}{Mean(\bar{M}_j)}, \quad j = 1 \sim N_{accession}$$

where \bar{M}_j is the j th column of the accession count matrix \mathbf{M} . Since the Gini coefficient increases with the mean count of the low dispersion accessions, we fitted the Gini coefficients and the mean counts of all accessions into a linear equation, i.e.,

$$Gini_j = a \cdot Mean(\bar{M}_j) + b, \quad j = 1 \sim N_{accession}$$

and selected the accessions above the line, i.e., the accessions with high dispersion, from the accession count matrix \mathbf{M} . The filtered accession matrix was then normalized by calculating the z-score or probability of the fragment count for each row (i.e., each cell) to generate normalized matrix \mathbf{M}_a for the next step of cell clustering.

(4) *Cell clustering.* From the filtered and normalized accession matrix \mathbf{M}_a , APEC established the connectivity matrix by computing the k-neighbor graph of all cells. Since the Louvain algorithm was proven to be a reliable single-cell clustering method in Seurat¹³ and Scanpy¹⁶, we adopted it in APEC to automatically predict the number of clusters from the connectivity matrix and defined each Louvain community as a cell cluster. APEC uses the Louvain algorithm to predict cluster number and perform cell clustering as default. Meanwhile, if users want to artificially define the number of cell clusters, APEC can also perform KNN clustering on the connectivity matrix.

(5) *Compare the performance of APEC with that of other methods on cells with known identity.* To investigate the accuracy of the cell clusters predicted by different algorithms, we used the ARI value, which evaluates the similarity of clustering results with all known types of cells¹⁴. The ARI value can be calculated as follows:

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

where n_{ij} is the element from the contingency matrix (i.e. the number of type i cells that were classified into cluster j), a_i and b_j are the sums of the i th row and j th column, respectively, and $\binom{x}{y}$ denotes a binomial coefficient. The higher the ARI value, the more accurate the classification. In addition, we define the largest element of each row in the contingency matrix as the “number of correctly classified cells” of the corresponding cell type. The ratio of the

correctly classified cells was also used to assess the clustering accuracy. When all cells are correctly clustered, this ratio is 1, and if most cells are evenly distributed into all clusters, this ratio is close to 0.

(6) *Characteristics of accesson*. The peaks of a same accesson can be distant from each other on the genome, and sometimes even on multiple chromosomes. The average number of peaks per accesson depends on the total number of peaks in the dataset and the number of accessions set in the program (default 600). Usually the total number of peaks can vary between ~40,000-150,000 depending on the total number of cells and the sequencing depth for each cell, thereby the average number of peaks per accesson is around ~60-250. Beside, we chose the top 40 principle components (PCs) of the normalized matrix to construct the connectivity matrix since the first 3~5 PCs are usually not sufficient to capture the detailed features of a single cell dataset, as described in Seurat and many other single cell analysis tools. As described on the scikit-learn website (<https://scikit-learn.org/stable/modules/clustering.html>), the use of the KNN and then 'Euclidean' distance and 'Ward' linkage methods to build the connectivity matrix between cells and perform the agglomerative clustering for the connectivity matrix usually provides better clustering results for different types of datasets. Although default values were chosen to provide better clustering results based on analysis of multiple datasets, users can adjust these parameters as needed.

Sampling of accesson number. To test if the APEC clustering result is sensitive to the choice of accesson numbers, we sampled 100 different accesson numbers from 500 to 1500 in steps of 10 and clustered the cells of each dataset 100 times. APEC generated stable clustering results in terms of the average ARI on these datasets, with a wide range of different accesson numbers ([Supplementary Fig. 4d](#)). We used accesson number = 600 as the default in APEC.

Parameter settings for other algorithms. To quantify the cell clustering performance of APEC, we compared APEC with other state-of-the-art single-cell epigenomic algorithms on the same datasets with gold standards, including cisTopic¹⁹, LSI^{11,12}, chromVAR¹⁷ and Cicero¹⁸. Since most of them have no cell clustering algorithm within their original codes, we applied the Louvain clustering algorithm on their transformed matrices to fairly compare their performance. We adopted the default settings of these tools for most of the comparisons in this paper, except for some parameters that were manually defined as necessary, such as the random seed in cisTopic, the number of top components in LSI, and the peak aggregation distance in Cicero. Therefore,

we sampled these parameters multiple times to obtain the average ARI and ratio of correctly classified cells of the clustering results for each tool, just as we sampled the accession number for APEC. We set the same parameters for all the datasets as follows:

- (1) *cisTopic*. The scanning range of the topic number was set to [10, 40], the number of parallel CPUs was set to 5, and the random seed was sampled 100 times from 100 to 600 in steps of 5. We kept the topic matrices normalized by z-score and probability and provided the performances based on both normalization methods. We then applied the Louvain algorithm as we did in APEC to cluster cells from the normalized topic matrix generated by *cisTopic*.
- (2) *LSI*. We performed truncated SVD (singular value decomposition) analysis on the TF-IDF (term frequency-inverse document frequency) matrix and chose N_{SVD} top components to generate the LSI matrix. N_{SVD} was sampled 6 times from 6 to 11. The first component was ignored since it is always related to read depth, and the LSI scores were capped at ± 1.5 . Then, we used the Louvain algorithm to cluster the cells of the LSI-processed matrix.
- (3) *chromVAR*. The number of background iterations was set to 50, and the number of parallel CPUs was set to 1. We then used the Louvain algorithm to cluster cells based on the bias corrected deviation matrix generated by *chromVAR*.
- (4) *Cicero*. The genome window was set to 500k BPs, the normalization method was set to “log”, the number of sample regions was set to 100, the number of dimensions was set to 40, and the peak aggregation distance was sampled at 20 values from 1k to 20k BPs in steps of 1k BPs. Then, we used the Louvain algorithm to cluster cells based on the aggregated model matrix generated by *Cicero*.

To test the robustness of each algorithm, we randomly sampled 20%~90% of the raw sequence reads from the dataset of AML cells and 3 cell lines (LMPP, HL60 and monocyte) and calculated the ARI accordingly. This random sampling experiment was performed 50 times for each method and average ARIs were reported. The manually defined parameters for each method were set to: APEC, 600 accessions; *cisTopic*, random seed 100; *LSI*, top 2-6 principle components; *Cicero*, 10k BPs aggregation distance.

Gene expression predicted by APEC. To evaluate the gene expressions from scATAC-seq data, we integrated two gene score assessment algorithms in APEC. The first algorithm is based on related accessions. Within a certain genomic distance around a peak i (1 Mbp by default), we calculate the odds of the all the nearby peaks (including peak i itself) that belong to the same accession by Fisher’s exact test. If the P value of the test is less than 0.01, we defined the

expression of the first downstream gene of peak i as $E = -\log_{10}(P_i) * M_i$, where M_i is the read counts of the accession that contains peak i . If multiple peaks are located upstream of a gene, then the expression of this gene is defined as the average of the E values of these peaks. In some cases, when the quality of scATAC-seq data is not as good, it is difficult to estimate the expressions of many genes by the above algorithm. Therefore, APEC offers another gene assessment method, by scoring a gene by the peaks around its TSS region, which is similar to the algorithm used by Preissl et al.¹⁰. We calculate the average read counts of all peaks around a gene's TSS (± 20000 BP) as its raw score (S_{ij} for cell i and gene j), then define the gene expression by normalizing the raw score by ($S'_{ij} = S_{ij} * 10000 / \sum_i S_{ij}$), making it in a range comparable to the gene expression from scRNA-seq data. Both the mouse forebrain and the human hematopoietic cells datasets were used to compare the performance of gene expressions evaluated by APEC and other algorithms such as cisTopic and Cicero, and the parameters set in those algorithms were the same as described in the previous section.

Significant differential peaks, genes and motifs. APEC used the Student's t-test to estimate the significance of the fragment count differences between cell clusters, with P-value and fold changes, and one can determine the thresholds to identify significant differential peaks for each cluster. The significant differential genes of each cell cluster can also be acquired from the gene score (\bar{S}_{kj}) by the same method. To accurately quantify the enrichment of motifs on each cell, APEC applied the bias-corrected deviation algorithm from chromVAR¹⁷; thus, the chromVAR algorithm has been embedded into the pipeline to facilitate the calculation of the corrected deviation of the motifs. In this python version of chromVAR, permuted sampling and background deviation calculation can be run in parallel on multiple processors to reduce the computer time. The differentially enriched motifs were defined by an absolute fold change >1 in the average motif deviations between one cluster and another.

Potential super-enhancers. Here, we defined a super-enhancer as a long continuous genomic area containing many accessible regions and have the same accessibility pattern in different cells. Many different motifs appear in one super-enhancer, therefore, the motif-based clustering method cannot reflect the critical contributions from super-enhancers for cell clustering. However, the accession-based algorithm can group most peaks in one super-enhancer to one accession since they always present the same accessibility pattern between cells. APEC identified super-enhancers by counting the number of peaks in a 1 million BP genomic area that belong to a same accession. It also requires that more than 3/4 of the putative peaks in one super-enhancer be adjacent on the initial peak list. The pipeline can also aggregate bam files by cell types/clusters

and convert them to BigWig format for users to upload to the UCSC genome browser for visualization.

Spatial correlation of peaks in the same accesson. To test whether peaks in the same accesson are closer in space, we integrated the Hi-C⁶⁶ data (GSE63525) and scATAC-seq⁸ data (GSE65360) on GM12878 cells. The spatial correlation of different windows, both intra- and inter-chromosomal, can be directly extracted by Juicer⁷¹. The Pearson's correlation matrix of the intra-chromosomal or inter-chromosomal windows can be calculated from the corresponding observed/expected matrix. We constructed 600 accessons by grouping peaks in the GM12878 scATAC-seq data in APEC, and removed the accessons that contained more than 1000 peaks or less than 5 peaks. The width of the window was set to 500k BPs, and peaks were then assigned to each window. Next, we collected Hi-C correlations between windows that contained peaks in the same accesson, termed “Accesson” correlations. For comparison, we also shuffled all peaks in different accessons to make fake accessons and re-collected the Hi-C correlations between windows that contained peaks in each fake accesson, termed “Shuffled” correlations. Meanwhile, we also collected Hi-C correlations between windows that contained no peaks, termed “Non-accesson” correlations. We made boxplots for these three types of correlations for intra-chromosomal and inter-chromosomal peaks and found that co-accessible peaks are spatially closer to each other than random peaks.

Pseudotime trajectory constructed by APEC. As a tool to simulate the time-dependent variation of gene expression and the cell development pathway, Monocle has been widely used for the analysis of single-cell RNA-seq experiments^{37, 72}. APEC reduced the dimension of the accesson count matrix **M** by PCA, and then performed pseudotime analysis using the Monocle program. For complex datasets, it is necessary to limit the number of principal components, since too many features will cause too many branches on the pseudotime trajectory, and makes it difficult for a user to identify the biological significance of each branch. For the hematopoietic single cell data and thymocyte data, we used the top 5 principal components of the accesson matrix to construct the developmental and differentiation trajectories.

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

- (1) Raw fragment count matrix. We normalized the raw count matrix **B** exactly as in 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 on the normalized matrix **B'**. Only the top 5 PCs were subjected to Monocle to construct the trajectory.
- (2) cisTopic. The topic matrix generated by cisTopic was normalized by making the sum of each row the same (i.e., the probability). Then, we performed PCA analysis on the normalized topic matrix and subjected the top 5 PCs to Monocle to build the trajectory.
- (3) LSI. We chose the 2nd–6th principle components of the SVD transformation of the LSI matrix and subjected them to Monocle to construct the trajectory.
- (4) ChromVAR. After the PCA analysis of the bias corrected deviation matrix generated by chromVAR, the top 5 PCs were combined with Monocle to construct the trajectory.
- (5) Cicero. We performed PCA analysis on the aggregated matrix generated by Cicero and used the top 5 PCs in Monocle to build the trajectory.

In addition, to confirm the reliability of the APEC + Monocle prediction of the developmental pathway, we applied another pseudotime trajectory constructing method, SPRING⁴⁰, to the accesson count matrix **M** from APEC to reconstruct the pathways for the hematopoietic differentiation dataset and thymocyte developmental dataset. We performed PCA analysis of the accesson matrix **M** and subjected the top 5 PCs to SPRING to generate the trajectories. The number of edges per node in SPRING was set to 5.

Parameter settings for each dataset. In the quality control (QC) step, cells are filtered by two constraints: the percentage of the fragments in peaks (P_f) and the total number of valid fragments (N_f). However, there is no fixed cutoff for these two parameters since the quality of different cell types and/or experiment batches are completely different. The total number of peaks is usually limited to approximately 50000 to reduce computer time, but we recommend using all peaks if the users want to obtain better cell clusters. (1) For the data set from hematopoietic cells, the $-\log(Q\text{-value})$ threshold of high-quality peaks was set to 35 to retain 54212 peaks, and the cutoff values of P_f and N_f were 0.1 and 1000, respectively. (2) For the scATAC-seq data on the two types of cells from 2 AML patients (P1-LSC, P1-Blast, P2-LSC, P2-Blast), the threshold of $-\log(Q\text{-value})$ was set to 5 to retain 38683 high-quality peaks for subsequent processing. When LMPPs, HL60 and monocytes were added to this dataset with the AML cells, the threshold of $-\log(Q\text{-value})$ was set to 8 to retain 42139 peaks. In the QC step, we set the P_f cutoff to 0.05 and the N_f cutoff to 800. (3) For the snATAC-seq data from the adult mouse forebrain, all peaks and the raw count

matrix obtained from the original data source were adopted in the analysis. (4) For the ftATAC-seq data from thymocytes, all 130685 peaks called by MACS2 were reserved for the fragment count matrix ($Q\text{-value} < 0.05$), and we retained cells with $P_f > 0.2$ and $N_f > 2000$.

SMART-seq data analysis with Seurat. For the analysis of SMART-seq data from mouse thymocytes, we employed STAR (version 2.5.2a) with the ratio of mismatches to mapped length (`outFilterMismatchNoverLmax`) less than or equal to 0.05, translated output alignments into transcript coordinates (i.e., `quantMode TranscriptomeSAM`) for mapping⁷³ (Dobin et al., 2013) and used RSEM⁷⁴ (Bo et al., 2011) to calculate the TPM of genes. For QC, we excluded cells in which fewer than 2000 genes were detected and genes that were expressed in only 3 or fewer cells. Seurat filtered cells with several specific parameters to limit the number of genes detected in each cell to 2000~6000 and the proportion of mitochondrial genes in each cell was set to less than 0.4 (i.e., `low.thresholds=c(2000, Inf)`, `high.thresholds=c(6000, 0.4)`). Additionally, the top 12 principal components were used for dimension reduction with a resolution of 3.2 (`dims.use = 1:12`, `resolution=3.2`), followed by cell clustering and differential expressed gene analysis⁷⁵.

Association of cell clusters from scATAC-seq and scRNA-seq data. To determine the association between cell clusters from the epigenomics and transcriptomic sequencing, we calculated the P-values of Fisher's exact test of marker/non-marker genes between each pair of cell clusters from scATAC-seq and scRNA-seq data. For example, for cell cluster a from ftATAC-seq and cell cluster b from SMART-seq, if the number of consensus marker genes in both cluster a and b is G_{11} , the number of genes that are not markers in either cluster a or b is G_{22} , and the number of markers in only cluster a (or cluster b) is G_{12} (or G_{21}), then the 2×2 matrix \mathbf{G} can be directly used for Fisher's exact test to evaluate the P-value A_{ab} between cluster a and b . After constructing a matrix \mathbf{A} filled with the negative logarithm of A_{ab} for ftATAC-seq cluster a and SMART-seq cluster b , we calculated the z-score for each row and column of \mathbf{A} to determine the correlation between cell clusters from different sequencing experiments.

GO term analysis of cells along pseudotime trajectory. We defined the functional characteristics of each accession by the GO terms and motifs enriched on its peaks. The GO terms of an accession were obtained by submitting all of its peaks to the GREAT website⁷⁶. The negative logarithm of the P-value of each GO term in each accession was filled into a (GO terms) \times (accessions) matrix \mathbf{L} . The significance of each GO term on each cell was evaluated by the product of the matrix \mathbf{L} and the accession count matrix \mathbf{M} , i.e.

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

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

Motif enrichment of cells along pseudotime trajectory. To assess the motif enrichment of the accessions, we used the Centrimo tool of the MEME suite⁷⁷ to search for the enriched motifs for the peaks of each accession and applied the same algorithm as to the GO term score to obtain the motif score. The negative logarithm of the E-value (product of adjusted P-value and motif number)⁷⁷ of each motif in each accession was used to construct a (motifs) \times (accessions) matrix **F**. The enrichment of each motif on each cell was evaluated by the product of the matrix **F** and the accession count matrix **M**, i.e.,

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

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

REFERENCES

67. Picelli, S. et al. Full-length RNA-seq from single cells using Smart-seq2. *Nat Protoc* **9**, 171-181 (2014).
68. Zuo, Z. et al. ATAC-pipe: general analysis of genome-wide chromatin accessibility. *Briefings in Bioinformatics*, bby056-bby056 (2018).
69. Shaffer, S.M. et al. Rare cell variability and drug-induced reprogramming as a mode of cancer drug resistance. *Nature* **546**, 431-435 (2017).
70. Huang, M. et al. SAVER: gene expression recovery for single-cell RNA sequencing. *Nat Methods* **15**, 539-542 (2018).
71. Durand, N.C. et al. Juicer Provides a One-Click System for Analyzing Loop-Resolution Hi-C Experiments. *Cell Syst* **3**, 95-98 (2016).
72. Qiu, X. et al. Single-cell mRNA quantification and differential analysis with Census. *Nat Methods* **14**, 309-315 (2017).
73. Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21 (2013).
74. Li, B. & Dewey, C.N. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *Bmc Bioinformatics* **12**, 323 (2011).
75. Macosko, E.Z. et al. Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* **161**, 1202-1214 (2015).

76. McLean, C.Y. et al. GREAT improves functional interpretation of cis-regulatory regions. *Nat Biotechnol* **28**, 495-501 (2010).
77. Bailey, T.L. et al. MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res* **37**, W202-208 (2009).

Supplementary Figures

Supplementary Figure 1. Quality control diagrams generated by APEC for the scATAC-seq datasets. **(a)** Number of final mapped reads and percentage of fragments in peaks for each cell in the human hematopoietic dataset. **(b)** Average count of scATAC-seq insertions around TSS regions in the same dataset. **(c)** Statistical distribution of fragment lengths in the same dataset. **(d-f)** Quality control diagrams, as in (a-c), for the blast and LSC cells from two AML patients. **(g-i)** Quality control diagrams, as in (a-c), for the cells from AML patients and three cell lines (LMPP, HL60 and Monocyte).

Supplementary Figure 2. Predicted gene expressions from scATAC-seq data of hematopoietic cells by cisTopic and Cicero. **(a)** Expressions of marker genes (*FOXO1*, *CEBPA*, *CD86*, *IKZF1*, *GFI1B*, and *AQP1*) evaluated by cisTopic. **(b)** Expressions of these marker genes evaluated by Cicero.

Supplementary Figure 3. Clustering performance of the dimension-transformed matrices generated by different algorithms. **(a)** The tSNE diagrams of the cells from AML patients and three distinct cell lines (LMPP, monocyte and HL60). Different algorithms provided different dimension-transformed matrices for tSNE analysis, i.e., APEC: accesson matrix; cisTopic: topic matrix; LSI: LSI matrix; chromVAR: bias corrected deviation matrix; Cicero: aggregated model matrix. The table below the diagrams contains the average ARI of the cell clustering results for each algorithm. **(b)** The tSNE diagrams and ARI table for the leukemic stem cells (LSCs) and blast cells from 2 different AML patients only, as in (a). **(c-d)** Box-plots showing the ARI values for the clustering of human hematopoietic cells (c) and the blast and LSC cells from two AML patients (d). We sampled different parameters for different algorithms. APEC: 100 different accesson numbers from 500 to 1500 in steps of 10; cisTopic: 100 different random seeds from 100 to 600 in steps of 5; LSI: 6 different numbers of top SVD components from 6 to 11; Cicero: 20 different genomic distances from 1k BPs to 20k BPs in steps of 1k BPs; chromVAR: no sampling. Z-score and probability denote different methods of normalizing the dimension-transformed matrices. Center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers. **(e)** The average ARI values calculated by down-sampling 50 times from the raw data of the AML cells and three cell lines for each method. The X-axis represents the percentage of down-sampled sequencing reads. Shaded error band: 95% confidence interval.

Supplementary Figure 4. Super-enhancers predicted by APEC for the scATAC-seq data of cells from AML patients. **(a-b)** The genome browser track shows the aggregated scATAC-seq signal

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

Supplementary Figure 5. Comparison of the peak grouping algorithms used by APEC and Cicero on the hematopoietic dataset. (a) The characteristics of accessions in APEC. Left panel: distribution of peaks in each accession; middle panel: genomic distance of peaks belong to the same accession; right panel: number of chromosomes with peaks belong to the same accession. (b) The characteristics of CCAN (defined by Cicero), as in (a). (c) Site links discovered by APEC and Cicero.

Supplementary Figure 6. Biological insights of the accession and stability and scalability analysis of APEC. (a) Box plot showing the average spatial distance between peaks in the same accession (from scATAC-seq) versus randomly shuffled peaks versus non-accessible genomic regions. Spatial distance was estimated from chromosome conformation capture (Hi-C) technology. Both Hi-C and scATAC-seq data were generated from the same cell line GM12878. Left panel: intra-chromosomal correlation of windows in the Hi-C data; right panel: inter-chromosomal correlation of windows in the Hi-C data. Accession: The correlation between two windows that contain peaks in the same accession; Shuffled: The correlation calculated by randomly shuffling peaks in each accession; Non-accession: The correlation between two windows that contain no peaks. Center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers. (b) The Hi-C profile of windows between chr1:500,000-21,500,000. The black bars below the Hi-C track denote peaks in the same accession from APEC. Dotted boxes indicate examples of peaks in the same accession that are distant in genomic position but close in space. (c) The computing time required for different algorithms to cluster cell numbers from 10,000 to 80,000 with peak number=100,000 randomly sampled from the original dataset. The data were sampled from the single-cell atlas of *in vivo* mammalian chromatin accessibility. CisTopic was performed using 5 CPU threads and all the other tools with 1 CPU thread. (d-f) The ARI values of the clustering results that used different numbers of accessions (d), nearest neighbors (e), and principle components (f). The dataset includes the cells from two AML patients and three cell lines. Default values are noted in red.

Supplementary Figure 7. (a-b) Intensity of the accessions associated with signature genes of excitatory (Excl) and inhibitory (Int) subtypes. The subtypes listed in parentheses were defined by the signature genes in the results from inDrops-seq data³⁵.

Supplementary Figure 8. (a) The clustering and cell-type classification of the mouse forebrain dataset by cisTopic. Left panel: cell clusters obtained by cisTopic, illustrated in the tSNE diagram. Middle panel: the z-scores of the average gene scores obtained from cisTopic clusters. Right panel: the hierarchical clustering of the correlations between cell clusters defined by cisTopic. (b-d) The clustering and cell-type classification of the same dataset by LSI, Cicero and chromVAR respectively, as in (a).

Supplementary Figure 9. UCSC genome browser track diagram of the normalized fragment count around gene *CD34* for each hematopoietic cell type.

Supplementary Figure 10. Cell differentiation trajectories of the human hematopoietic dataset constructed by different algorithms. (a-e) The pseudotime trajectories constructed by the combination of Monocle and the raw peak count matrix, the topic matrix from cisTopic, the LSI matrix, the aggregated model matrix from Cicero, and the bias corrected deviation matrix from chromVAR, respectively. (f) The pseudotime trajectory constructed by the combination of SPRING and the accession matrix from APEC.

Supplementary Figure 11. (a) Gating strategy of the mouse thymocytes in ftATAC-seq. (b-d) Quality control diagrams for the mouse thymocyte data, similar to Supplementary Fig. (1a-1c). (e) The z-score of correlation between the cell types from ftATAC-seq and bulk ATAC-seq data.

Supplementary Figure 12. (a) Selected significant motifs enriched in different thymocyte subtypes obtained by the APEC algorithm.

Supplementary Figure 13. Single-cell transcriptome analysis of *Mus musculus* thymocytes from SMART-seq. (a) tSNE diagram of the single-cell expression matrix of *Mus musculus* thymocytes, labeled by the FACS index of each cell. (b) Louvain clustering of the same single-cell dataset obtained by Seurat. The cell types of these clusters were classified by the expression of corresponding marker genes. (c) Important marker genes were differentially expressed in different cell clusters. (d) Heatmap of the expressions of all genes significantly differentially expressed between cell clusters. The top color bar used the same scheme described in (b) to render cells of different clusters.

Supplementary Figure 14. Developmental characteristics of single-cell samples captured by APEC. (a, b) Pseudotime trajectory of scATAC-seq data from *Mus musculus* thymocytes labeled with the FACS index and APEC cluster index. (c) Pseudotime trajectory constructed by applying SPRING to the accession matrix. The colors of cells denote their stages in the APEC trajectory

results. **(d)** Z-scores of the $-\log(\text{P-value})$ of the GO terms along the pseudotime trajectory of stage 1 cells. **(e)** Logarithm of the P-value of GO terms searched from peaks in accessions ac1~ac7, which are the marker accessions of cluster DP. A1 and DP. A3/4/5 of thymocytes.