

METHOD

Open Access

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



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

\* Correspondence: [qukun@ustc.edu.cn](mailto:qukun@ustc.edu.cn)

<sup>†</sup>Bin Li and Young Li are co-first authors.

<sup>1</sup>Department of Oncology, The First Affiliated Hospital of USTC, Division of Life Sciences and Medicine, Division of Molecular Medicine, Hefei National Laboratory for Physical Sciences at Microscale, University of Science and Technology of China, Hefei 230001, Anhui, China

<sup>3</sup>CAS Center for Excellence in Molecular Cell Sciences, The CAS Key Laboratory of Innate Immunity and Chronic Disease, University of Science and Technology of China, Hefei 230027, Anhui, China

Full list of author information is available at the end of the article

## 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 cell by groups of accessible regions with synergistic signal patterns termed “accessons”. This python-based package greatly improves the accuracy of unsupervised single-cell clustering for many public datasets. It also predicts gene expression, identifies enriched motifs, discovers super-enhancers, and projects pseudotime trajectories. APEC is available at <https://github.com/QuKunLab/APEC>.

**Keywords:** scATAC-seq, Cell clustering, Accesson, Regulome, Pseudotime trajectory

## Background

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 processes [1], such as hematopoietic stem cell (HSC) differentiation [2], embryonic development [3], neuronal activity and regeneration [4, 5], tumor cell metastasis [6], and patient responses to anticancer drug treatment [7]. 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 (sciATAC-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 Seurat [13], SC3 [14], SIMLR [15], and SCANPY [16]. However, most of these algorithms are not directly compatible with a single-cell ATAC-seq dataset, for



© The Author(s). 2020 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

which the signal matrix is much sparser. To characterize scATAC-seq data, the Greenleaf lab developed an algorithm named chromVAR [17], 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], Cicero [18], cisTopic [19], and SnapATAC [20], have also been applied to process single-cell/nucleus ATAC-seq data [10, 12].

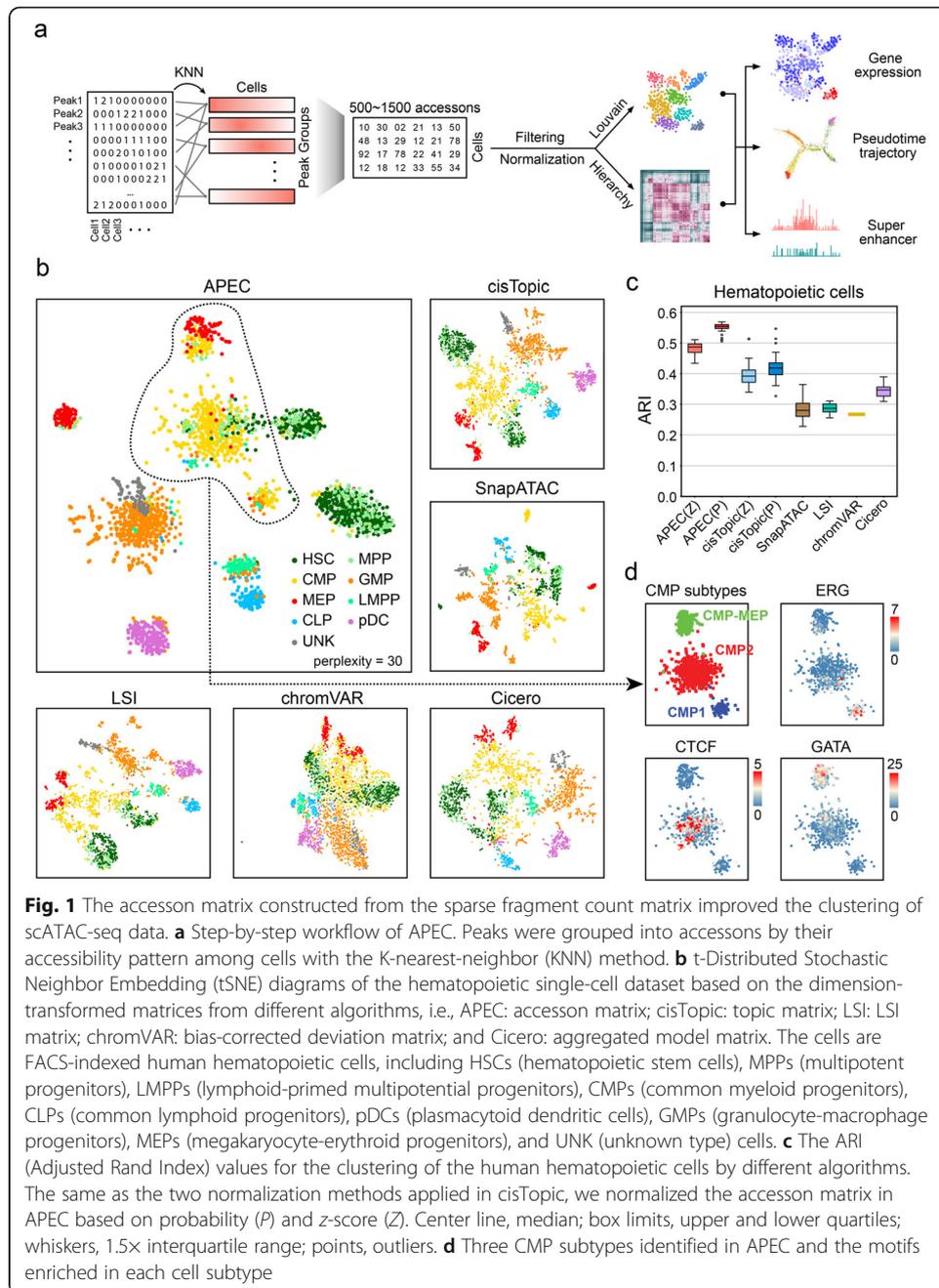
However, great challenges still remain for current algorithms to provide comprehensive analysis on single-cell chromatin accessibility profiles. For instance, (1) how to improve the accuracy of single-cell clustering on cells with minor difference, (2) how to better illuminate the cellular developmental processes based on scATAC-seq data, (3) how to integrate RNA/ATAC profiles and predict the RNA expressions from scATAC-seq data and further reveal the biological insights of the system, and (4) how to manipulate massive sequencing data with high speed, scalability, and stability. Therefore, a fast, scalable, and stable algorithm is needed to precisely categorize cell subgroups, to predict gene expressions and developmental trajectories, and thereby to provide a deeper mechanistic understanding of single-cell epigenetic heterogeneity and regulation.

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, 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 an efficient, accurate, robust, and rapid cell clustering from single-cell ATAC-seq profiles. APEC was also integrated into a head-to-toe program package that has been made available on GitHub (<https://github.com/QuKunLab/APEC>).

## Results

### Accession-based algorithm improves single-cell clustering

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) differentiation [21]. Since the cell sorting strategy in this study was proofed at both bulk and single-cell level [22, 23], we used it as a gold standard. Compared to other state-of-the-art single-cell chromatin accessibility analysis methods, such as cisTopic [19], SnapATAC [20], LSI [11, 12], chromVAR [17], and Cicero [18], this new accession-based algorithm can clearly cluster cells into their corresponding identities according to the Adjusted Rand Index (ARI) (Fig. 1b, c). Here, we adopted the Louvain clustering algorithm for all the methods and assumed that the number of cell types was unknown to ensure that the clustering analysis was unsupervised. To calculate ARI values, we ignored all “unknown” cells and sampled the tunable parameter (e.g., accession number in APEC, and random seed in cisTopic) multiple times for each tool (see the “Methods” section). On average, 67% of cells were correctly classified by APEC with  $ARI = 0.483/0.552$  (using matrices normalized by z-score or



**Fig. 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-neighbor (KNN) method. **b** t-Distributed Stochastic Neighbor Embedding (tSNE) diagrams of the hematopoietic single-cell 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; and 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 UNK (unknown type) cells. **c** The ARI (Adjusted Rand Index) values for the clustering of the human hematopoietic cells by different algorithms. The same as the two normalization methods applied in cisTopic, we normalized the accession matrix in APEC based on probability (P) and z-score (Z). Center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers. **d** Three CMP subtypes identified in APEC and the motifs enriched in each cell subtype

probability, respectively), while cisTopic was the second most accurate method to predict cell identities (average ARI = 0.392/0.418) (Additional file 2: Table S1).

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 Erg [24]; CMP2 cells are enriched with CTCF motif, suggesting that these cells are at the fate decision stage with CTCF-associated chromatin remodeling [25]; CMP-MEP cells are considered as MEP-committed CMPs and are strongly enriched with crucial regulators for MEP differentiation, such as GATA1 [26]. However, these 3 CMP sub-clusters were not as clearly distinguishable and cells were more scattered on the t-

Distributed Stochastic Neighbor Embedding (tSNE) maps generated by other methods (Additional file 1: Figure S1). 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.

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) patients [17]. We found that both APEC and cisTopic were tied for best to classify these cells (Additional file 1: Figure S2a). Interestingly, APEC, cisTopic, and LSI were all capable of almost perfectly separating the three distinct cell types (LMPPs, monocytes, and HL60), with ARI = 1.000/0.988, 0.987/0.987, and 0.969, respectively. However, in terms of clustering the four similar cell types from AML patients, APEC (average ARI = 0.575/0.564) outperformed other tools (Additional file 1: Figure S2b), 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 tunable parameters to ensure the reliability of their predictions (Additional file 1: Figure S2c, Additional file 3: Table S2).

Low sequencing depth and dropouts in single-cell sequencing profiles can severely obscure important cell-cell relationships [27, 28]. To further test the robustness of APEC on datasets with low quality, we first simulated datasets of low sequencing depth and high noise level by randomly extracting reads from the original mapped profile, or setting elements from the original fragment count matrix to zero. We then calculated and compared the ARI values from each method on those down-sampled datasets and found that APEC still exhibited an ARI > 0.7 at sequencing depths as low as 50%, or noise level as high as 40%, confirming its reliability to cluster cells (Additional file 1: Figure S2d&e).

Unlike chromVAR, APEC aggregates the minor differences between similar cells into accessions, which is not necessarily of the same motif. For example, APEC identified prominent super-enhancers around the genes *N4BPI* [29] and *GPHN* [30] in the LSC cells from AML patient 1 (P1-LSC) but not the other cell types (Additional file 1: Figure S3a & b). These two loci were also confirmed as super-enhancers by ROSE [31] (the top 2 candidates in Additional file 4: Table S3). 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 (Additional file 1: Figure S3c & d).

In contrast to Cicero, which aggregates peaks based on their cis-co-accessibilities networks (CCAN) within a certain range of genomic distance [18], APEC combines synergistic peaks genome-wide. Take the human hematopoietic cell dataset as an example, 600 accessions were built from the 54,212 peaks, and each accession contained ~ 40 peaks (median number) which were distributed to ~ 30 CCANs, compare with ~ 5 peaks in each CCAN which were distributed to ~ 5 accessions (Additional file 1: Figure S4a-c). The average distance between peaks in the 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 was 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 (Additional file 1: Figure S4d). Therefore, APEC and Cicero are two completely different approaches.

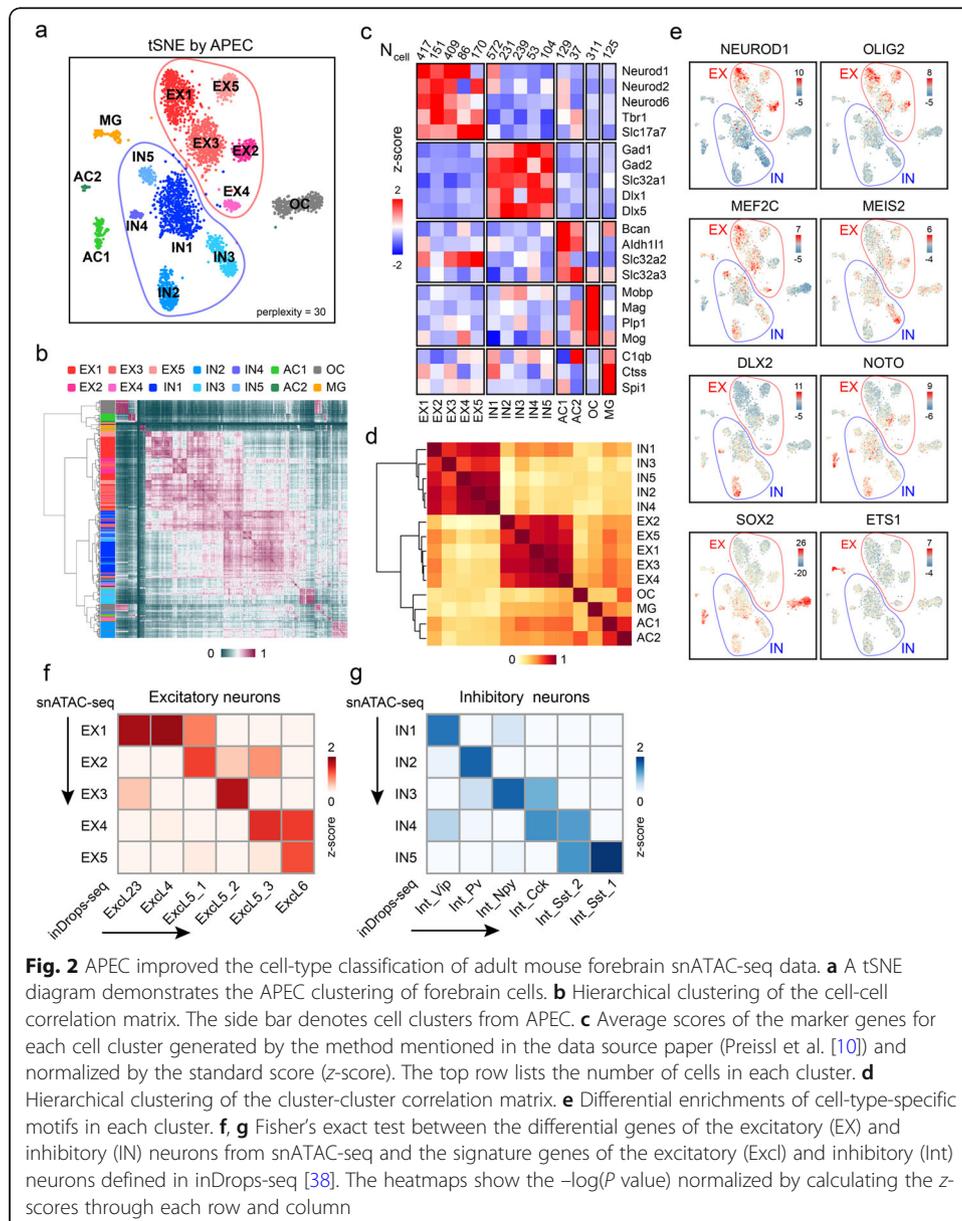
Furthermore, Buenrostro et al. showed that the covariation of the accessible sites across all the cells may reflect the spatial distance between the corresponding peaks [8]. By integrating the chromatin conformation profiles from Hi-C experiments with the scATAC-seq profile for the same cells, we found that peaks in the same accession are spatially much closer to each other than randomly selected peaks (Additional file 1: Figure S5a,  $P$  value  $< 10^{-7}$ ), suggesting that they may belong to the same topologically associated domains (TADs) (Additional file 1: Figure S5b). This discovery was confirmed on datasets from both the GM12878 and K562 cell lines (Additional file 1: Figure S5a & c). Since CTCF is a well-known architectural protein that mediates long-range chromatin looping and bridges genome topology and function [25, 32–35], if APEC does group peaks that are spatially adjacent, the CTCF motif should be enriched in peaks from the same accession. As expected, we found that the CTCF/CTCF motifs were the top two enriched motifs in accessions with more than 500 peaks (Additional file 1: Figure S5d & e), which strongly suggested that accessions may have captured valid long-range relations.

In cisTopic, hundred thousands of peaks were aggregated into dozens of representative “topics” (usually 10~40) that preserve the chromatin accessibility information of each cell via a Bayesian topic modeling method. Cells were then clustered based on their distance in the “topic” space. Although cisTopic is usually more than 10 times slower than other tools (Additional file 1: Figure S6a & b), it is one of the best tools for cell-type clustering. We then compared APEC with cisTopic on the peaks in “accessions” versus those in “topics” on the GM12878 dataset and found that peaks in the same “accession” were spatially closer to each other than the representative peaks in the same “topic” (Additional file 1: Figure S5a). Motif enrichment analysis also suggested that peaks in “accessions” were spatially more coherent than those in “topics”, since the CTCF and CTCFL motifs were top enriched in “accessions” but not in “topics” (Additional file 1: Figure S5d & f).

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 dataset [36], which contains 81,173 high-quality cells. Taking account of all the 436,206 peaks, it took 310 min to cluster 80,000 cells with 1 CPU thread using APEC (Additional file 1: Figure S6a). We also randomly select 100,000 peaks from the entire dataset to test the computer time spent by these tools (Additional file 1: Figure S6b). 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 (Additional file 1: Figure S6c-e). In terms of speed and scalability, LSI and SnapATAC exhibited better performance than the other tools. But APEC is capable with acceptable performance.

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

To evaluate the compatibility and performance of APEC with other single-cell chromatin accessibility detection techniques, such as snATAC-seq [10], transcript-indexed scATAC-seq [37], and sciATAC-seq [11], APEC was also tested with the datasets generated by those experiments. For example, APEC discovered 14 cell subpopulations in adult mouse forebrain snATAC-seq data [10], including four clusters of excitatory neurons (EX1–5), five groups of inhibitory neurons (IN1–5), astroglia cells (AC1&2), oligodendrocyte cells (OC), and microglial cells (MG; Fig. 2a, b). To quantify gene expression level, we defined a gene score as the average signal of the peaks close to its TSS region (Fig. 2c; see the “Methods” section). With that, we identified 5 excitatory subpopulations and 5 distinct inhibitory subpopulations, and all the cell groups were clearly distinguished from each other by hierarchical clustering (Fig. 2d). In contrast,

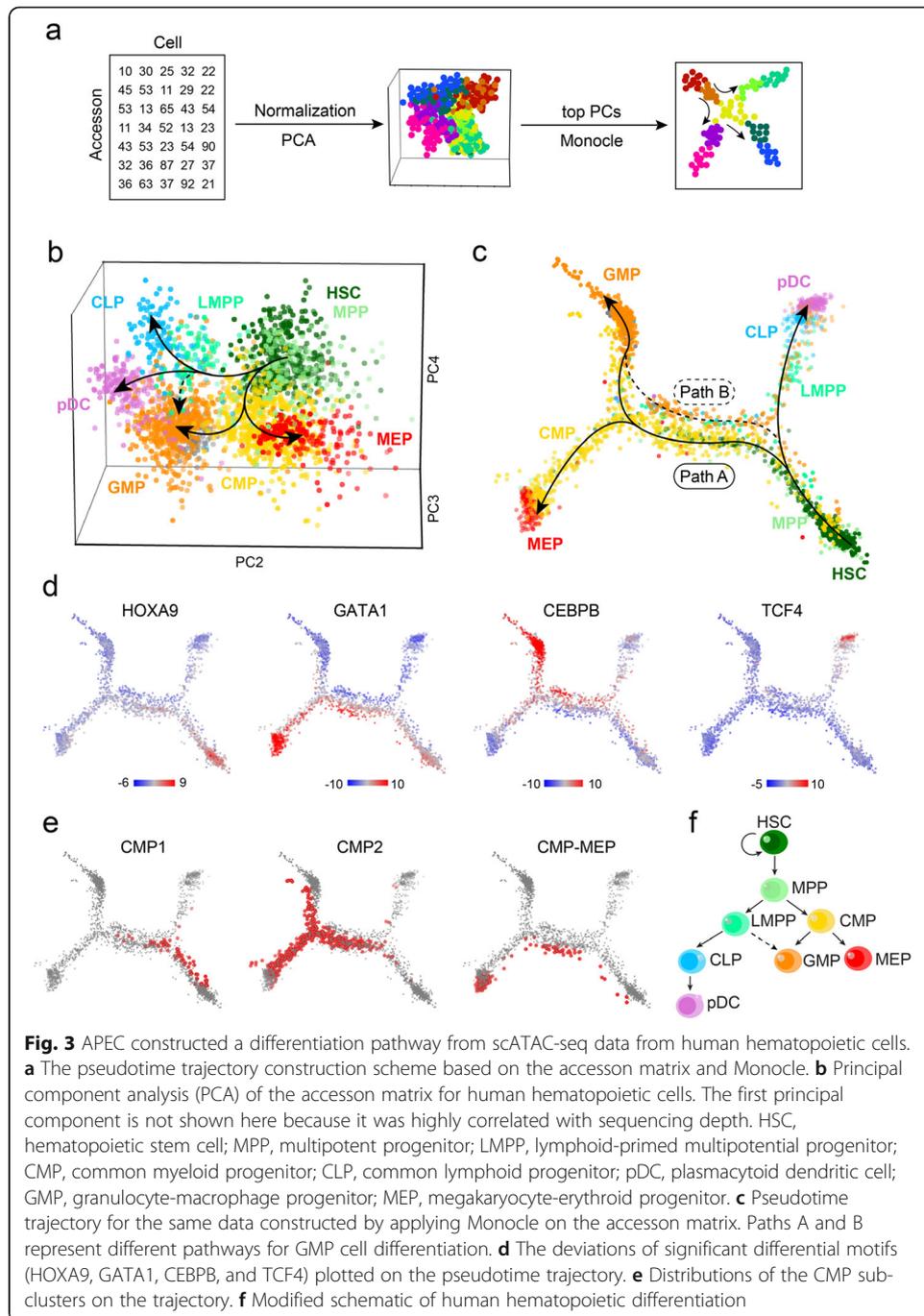


more than 29.7% (946 out of 3034) of cells were unable to be correctly assigned into any subpopulation of interest in the previous study [10]. Besides, scores for several marker genes such as *Neurod6* and *Aldh1l1* were not available in Cicero and cisTopic, respectively, making it difficult to identify the corresponding cell clusters (i.e., c0~c2 in Additional file 1: Figure S7a & b). The SnapATAC algorithms however mixed-clustered the AC1/2 with excitatory neurons in the correlation matrix (Additional file 1: Figure S7c). On the other hand, the motif enrichment analysis module in APEC identified cell-type-specific regulators that are also consistent with previous publications [10]. For example, the NEUROD1 and OLIG2 motifs were generally enriched on excitatory clusters (EX1\2\3\5); the MEF2C motif was more enriched on EX1/3/4/5; the motifs of MEIS2 and DLX2 were differentially enriched on different subtypes of inhibitory neurons (IN3 and IN2/4, 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 can distinguish and categorize single cells with great sensitivity and reliability.

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 multiple excitatory and inhibitory neuronal subtypes in the mouse visual cortex using single-cell inDrops sequencing [38] and provided top 20 signature genes that distinguished these cell subtypes. However, due to the sparseness of snATAC-seq matrix, scores of many signature genes were not strong enough to distinguish the cell sub-clusters. To overcome this, we developed a gene set overlap algorithm to associate cell clusters from scATAC-seq and scRNA-seq profiles (see the “Methods” section). We found that sub-cluster EX1~5 and IN1~5 in snATAC-seq can nicely correspond to the neuron subtypes classified by Hrvatin et al. (Fig. 2f, g). These results highlight the potential advantages of the accession-based approach for the integrative analysis of scATAC-seq and scRNA-seq data.

### **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 [39, 40]. For instance, Monocle [40, 41] constructs the minimum spanning tree, and Wishbone [42] and SPRING [43] construct the nearest neighbor graph from single-cell transcriptome profiles. These tools have been widely used to depict neurogenesis [44], hematopoiesis [45, 46], and reprogramming [47]. APEC integrates the Monocle algorithm into the accession-based method and enables pseudotime prediction from scATAC-seq data [21] and was applied to investigate HSC differentiation lineages (Fig. 3a). Principal component analysis (PCA) of the accession matrix revealed multiple stages of the lineage during HSC differentiation (Fig. 3b) and was consistent with previous publications [2, 21]. 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 commitment [48]. Notably, pDCs from the bone marrow are CD34<sup>+</sup> (Additional file 1: Figure S8a), 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 reports [49].

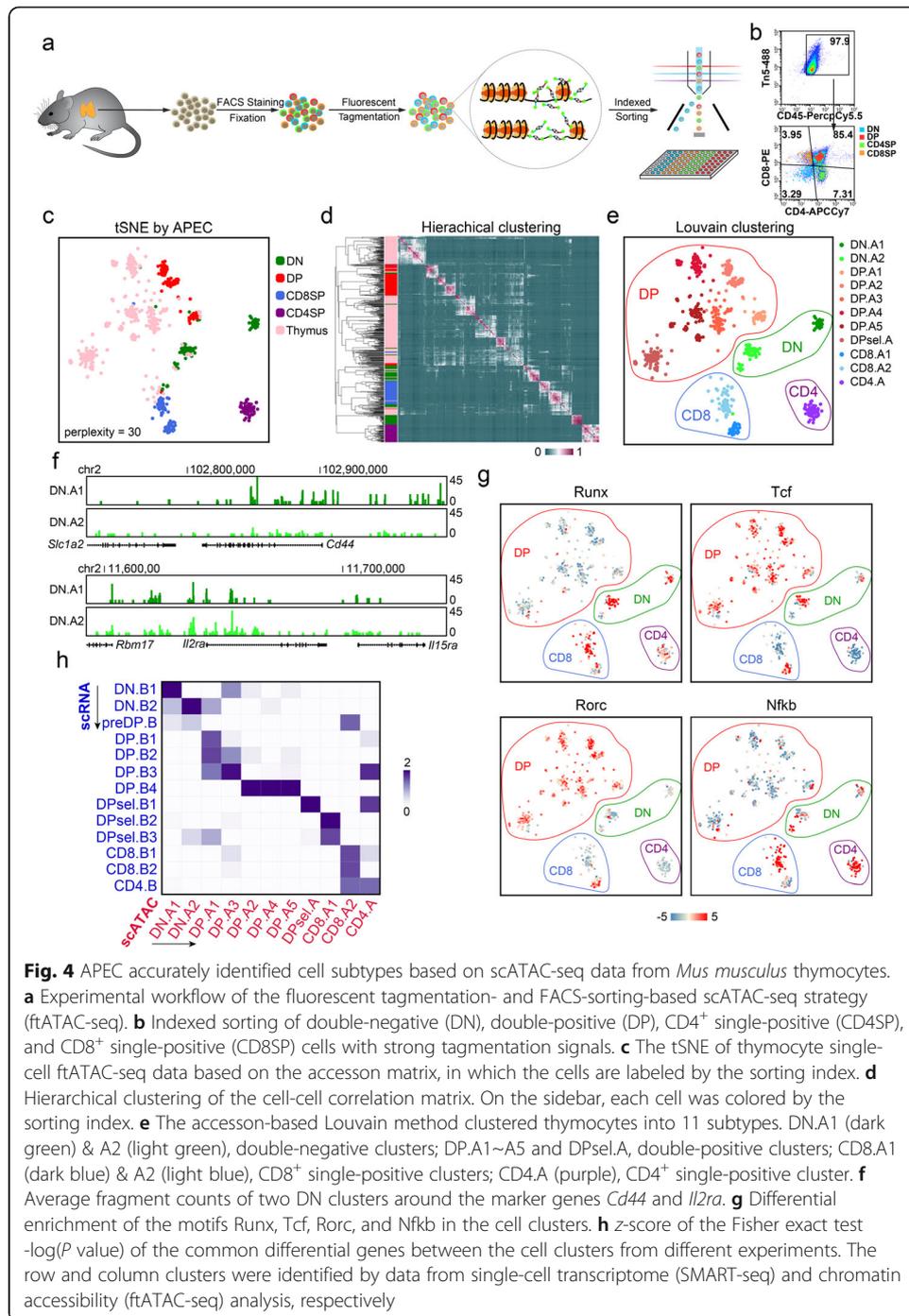
Furthermore, APEC incorporated the chromVAR algorithm to determine the regulatory mechanisms during HSC differentiation by evaluating the deviation of each TF

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, respectively [21] (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) [24–26]. Finally, we generated a hematopoiesis tree based on the APEC analysis (Fig. 3f).

We also 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 of the developmental pathways was constructed (Additional file 1: Figure S8b), 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 other methods, including cisTopic, SnapATAC, LSI, Cicero, and chromVAR (Additional file 1: Figure S8c-g); and (3) when we applied other pseudotime trajectory construction methods, such as SPRING [43], after APEC, a similar though less clear cell differentiation diagram was also obtained, suggesting the reliability of our prediction (Additional file 1: Figure S8h).

### 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) stages [50]. 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-seq [51] and Pi-ATAC-seq [52] 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, Additional file 1: Figure S9a). 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 (Additional file 1: Figure S9b-d). Correlation analysis with the published bulk ATAC-seq data of thymocytes [53] indicates that the cells we sorted in ftATAC-seq were correctly labeled (Additional file 1: Figure S9e). 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 130,685 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, and 2 CD8SP, suggesting that extensive

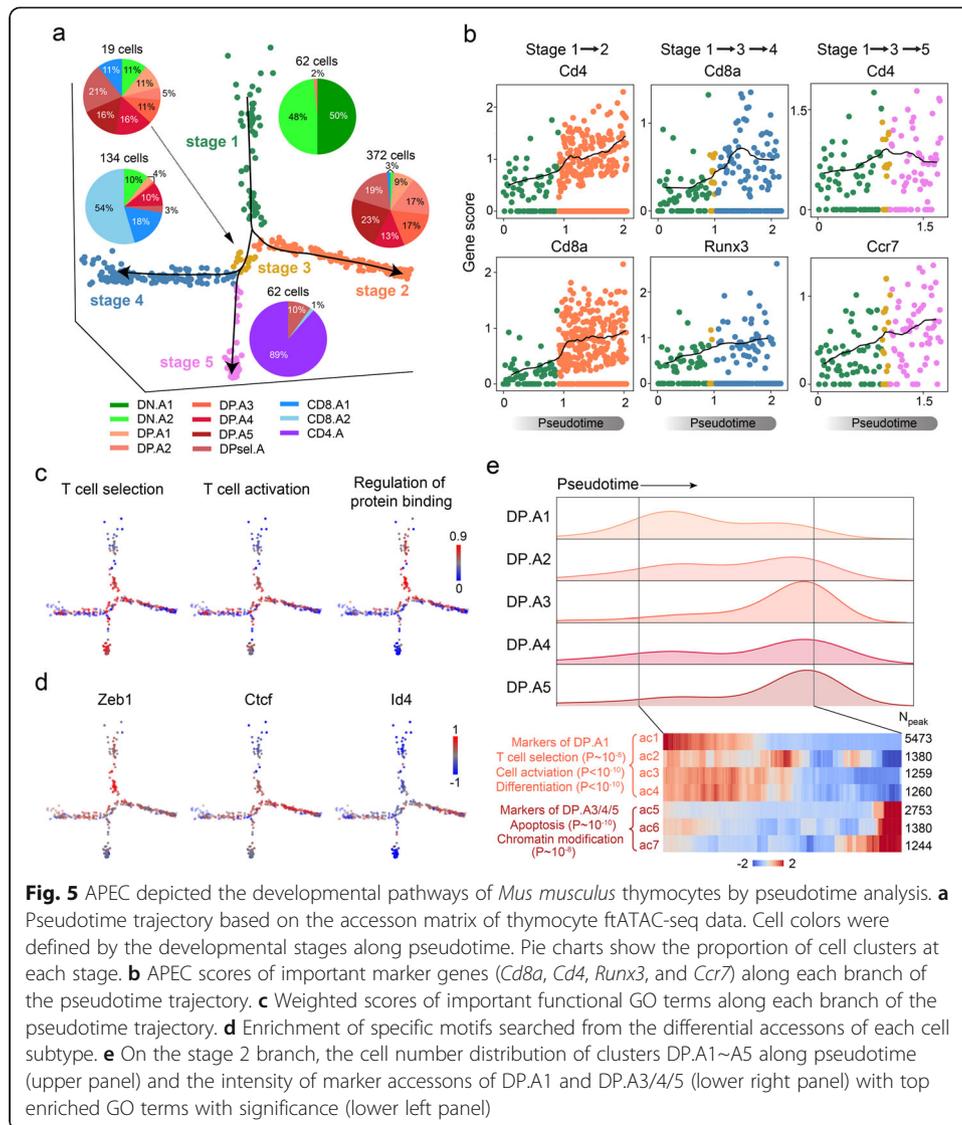
epigenetic heterogeneity exists among cells with the same CD4 and CD8 surface markers (Fig. 4e). For instance, there are four main subtypes of DN cells, according to the expression of the surface markers CD44 and CD25 [54], while two clusters were identified in ftATAC-seq. The accessibility signals around the *Il2ra* (Cd25) and *Cd44* gene loci demonstrated that DN.A1 comprised CD44<sup>+</sup>CD25<sup>-</sup> and CD44<sup>+</sup>CD25<sup>+</sup> DN subtypes (DN1 and DN2), and DN.A2 cells comprised CD44<sup>-</sup>CD25<sup>+</sup> and CD44<sup>-</sup>CD25<sup>-</sup> 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 cells [55], and we observed significant enrichment of the RUNX motif on DN cells and a group of CD8SP cells, similarly the TCF [56, 57], RORC [58], and NFkB [59] 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 (Additional file 1: Figure S10). 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 the “Methods” section). 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 protocol [60]. Unsupervised analysis of gene expression profiles clustered these thymocytes into 13 groups in Seurat [13] (Additional file 1: Figure S11a & b), and each subpopulation was identified based on known feature genes (Additional file 1: Figure S11c & 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 the “Methods” section) 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, Additional file 1: Figure S12a & 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 the last stages of DP cells. Besides Monocle, a



**Fig. 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)

similar developmental pathway can also be constructed by SPRING [43] based on the accession matrix (Additional file 1: Figure S12c). 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 cell selections, including  $\beta$ -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 (Additional file 1: Figure S12d). 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), APEC also identified Zeb1 [61], Ctfc [62], 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 [63, 64]. Associated with the fact that the vast majority of DP thymocytes die because of a failure of positive selection [65], 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, Additional file 1: Figure S12e). 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 [66, 67]. 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 subtypes [68].

## Discussion

Here, we introduced an accession-based algorithm named APEC for single-cell chromatin accessibility analysis, which is generally applicable to profiles generated from a variety of scATAC-seq techniques such as scATAC from the Fluidigm C1 chip [8], snATAC [10], sciATAC [11], Pi-ATAC [52], the 10X platform [69], and in this study the ftATAC. Without relying on any prior information (such as bulk sequencing data or known cell types), 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. 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.

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.

Despite these advantages of APEC, there are still great challenges unaddressed in single-cell chromatin accessibility analysis. For example, the current version of APEC is only capable to anchor scATAC and scRNA data at the group level, but not the single-cell level. APEC predicts the number of cell clusters by the Louvain algorithm, but if there are minor clusters “hidden” in one of the major clusters, one needs to re-perform the clustering analysis within the major cluster to classify the minor clusters. Although APEC is able to identify cell subtypes as rare as 1% of the entire cell population (data not shown), substituting the cell-peak fragment count matrix to the cell-bin count matrix from SnapATAC [20] may further help to identify very rare cell populations. Besides, the current version of APEC does not incorporate the function to correct the batch effect. In general, there are still lots of improvement needed to solve these problems.

## Conclusions

In this study, we presented a new algorithm, APEC, to analyze scATAC-seq profiles based on groups of genomic regions with similar accessibility named accessions. We applied this approach on multiple published datasets and the mouse thymocyte profile generated with ftATAC-seq. Compare to existing methods, APEC effectively improves cell classification and pseudotime construction, making it a powerful tool to study single-cell epigenomic heterogeneity and regulome dynamics.

## Methods

### 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-TCGTCGGCAGCGTCAGATGTGTATAAGAGACAG-3'; and AF488-R2, 5'-AF488-GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAG-3'. Then, 50  $\mu$ 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- $\mu$ m nylon mesh. A total of  $1 \times 10^6$  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% methanol at room temperature for 5 min. After washing twice with  $1 \times$  PBS, the cells were counted again. A total of  $1 \times 10^5$  fixed cells were resuspended in 40  $\mu$ L of  $1 \times$  TD buffer (5 mM Tris-HCl, pH 8.0, 5 mM MgCl<sub>2</sub>, and 10% DMF) containing 0.1% NP-40. Then, 10  $\mu$ L of fluorescent transposomes was added and mixed gently. Fluorescent tagmentation was conducted at 55 °C for 30 min and stopped by adding 200  $\mu$ 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<sup>+</sup>/AF488-Tn5<sup>hi</sup> population were index-sorted into each well of 384-well plates. The 384-well plates used to acquire sorted cells were loaded with 2  $\mu$ 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 tagged cells were incubated at 55 °C for 30 min. Then, 4.2  $\mu$ L of PCR round 1 buffer (1  $\mu$ L of 100  $\mu$ M MgCl<sub>2</sub>, 3  $\mu$ L of 2 $\times$  I-5 PCR mix [MCLAB], and 0.1  $\mu$ L each of 10  $\mu$ M R1 and R2 primers) was 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  $\mu$ L of PCR round 2 buffer (2  $\mu$ L of I-5 PCR Mix, 0.5  $\mu$ L each of Ad1 and barcoded Ad2 primers, and 1  $\mu$ L of ddH<sub>2</sub>O), 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 modifications [70]. Reverse transcription and the template-switch reaction were performed at 50 °C for 1 h with Maxima H Minus Reverse Transcriptase (Thermo Fisher); for library construction, 0.5–1 ng of cDNA was fragmented with 0.05  $\mu$ L of Robust Tn5 transposome in 20  $\mu$ L of TD buffer at 55 °C for 10 min, then purified with 0.8 $\times$  VAHTS DNA Clean Beads (Vazyme Biotech), followed by PCR amplification with Ad1 and barcoded Ad2 primers and purification with 0.6 $\times$  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 datasets 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 GSE74310 [2]; (2) scATAC-seq data for HL-60 cells from GSE65360 [8]; and (3)

scATAC-seq data for hematopoietic development (HSCs, MPPs, CMPs, LMPPs, GMPs, EMPs, CLPs, and pDCs) from GSE96772 [21]. (4) APEC is also compatible with a pre-processed fragment count matrix from the snATAC-seq data for the forebrain of adult mice (p56) from GSE100033 [10]. (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 GSE111586 [36]. (6) The scATAC-seq (GSE65360 [8]) and Hi-C (GSE63525 [71]) 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-pipe [72]. We also implemented the python script in ATAC-pipe [72] 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.

Unlike several previous methods that call peaks from bulk ATAC-seq data or aggregated cell populations according to cell type [20, 21, 73], the APEC pipeline calls peaks from the merged single-cell profiles of all cells using MACS2 to ensure that the entire analysis is unsupervised. We then 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 ( $\pm 2000$  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 ( $\mathbf{C}_{\text{peak}}$ ) 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 datasets. Based on the matrix  $\mathbf{C}_{\text{peak}}$ , 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 datasets, 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. The 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.

- (3) *Cell clustering.* From the 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 Seurat [13] and Scanpy [16], 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 to ensure that the clustering analysis is unsupervised and then performs 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.
- (4) *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 cells [14]. The ARI value can be calculated as follows:

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

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. A higher ARI value indicates more accurate classification of cell types.

- (5) *Characteristics of accession.* The peaks of a same accession can be distant from each other on the genome, and sometimes even on multiple chromosomes. The average number of peaks per accession 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  $\sim 40,000$  and  $150,000$  depending on the total number of cells and the sequencing depth for each cell; thereby, the average number of peaks per accession is around  $\sim 60$ – $250$ . 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. Although default values were chosen to provide better clustering results based on analysis of multiple datasets, users can adjust these parameters as needed.

- (6) *Data format.* APEC uses the MTX file format, which maintains only the non-zero elements in a matrix, to store and process the sparse cell-peak count matrix, and then uses the “Scipy.sparse” library in Python to operate the MTX format matrix for peak grouping and cell clustering. With this approach, APEC is able to process single-cell chromatin accessibility profiles with as many as 80,000 cells and 400,000 peaks within 50 GB RAM. APEC can also access the MTX output files from Cellranger, implying its well compatibility with 10X Genomics datasets.

#### Sampling of accession number

To test if the APEC clustering result is sensitive to the choice of accession numbers, we sampled 100 different accession numbers from 500 to 1500 in steps of 10 and clustered the cells of each dataset 100 times (Fig. 1c and Additional file 1: Figure S2c). APEC generated stable clustering results in terms of the average ARI on these datasets, with a wide range of different accession numbers (Additional file 1: Figure S5e). We used 600 as the default number of accessions 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 cisTopic [74], LSI [11, 12], chromVAR [17], and Cicero [18]. 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 (Fig. 1c and Additional file 1: Figure S2c), 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) *SnapATAC.* As SnapATAC uses a mapping procedure totally different with other tools, we adopt it with default parameters (binsize = 5 k) to build its own fragment count matrix from the merged bam file. We called the function “runJDA” with

- “bin.cov.zscore.lower=-2, bin.cov.zscore.upper=2, pc.num=50, norm.method='normOVE', max.var=5000, do.par=TRUE, ncell.chunk=1000, num.cores=1, seed.use=10, tmp.folder=tempdir()”. We then called the function “runKNN” and sampled the number of principal components from 20 to 40 (step-size = 2) and the number of nearest neighbors from 10 to 20 (step-size = 1). Therefore, we sampled a total of 100 times (10 × 10) to calculate the ARI values. Finally, we called the function “runCluster” with “louvain.lib='leiden', seed.use=10, resolution=1”.
- (3) *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  $\pm 1.5$ . Then, we used the Louvain algorithm to cluster the cells of the LSI-processed matrix.
  - (4) *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.
  - (5) *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 (Additional file 1: Figure S2d). We also random sampled elements from the raw fragment count matrix and set them to zero to simulate noises (Additional file 1: Figure S2e). The manually defined parameters for each method were set to APEC, 600 accessions; cisTopic, random seed 100; SnapATAC, 20 PCs and 15 nearest neighbors; LSI, top 2–6 principle components; Cicero, 10k BPs aggregation distance.

### Gene scores and differential analysis

APEC scores 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 ( $\pm 20,000$  BP by default) 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. After scoring genes, APEC uses the Student's  $t$  test to estimate the significance of each gene between cell clusters and therefore obtained a list of differential genes filtered by  $P$  values and fold changes. We also tested the gene scoring algorithms of other tools such as Cicero, cisTopic, and SnapATAC, and the parameters set in those algorithms were the same as described in the previous section. For scATAC-seq with very sparse signals, APEC first searched for differential accessions between cell clusters and extracted all

the peaks in the differential accessions. APEC then defines genes close to these peaks ( $\pm 20,000$  BP around TSS) as the differential genes.

#### 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 the differential/non-differential 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 (Fig. 4h), if the number of consensus differential genes in both clusters  $a$  and  $b$  is  $G_{11}$ , and the number of differential genes in either cluster  $a$  (or cluster  $b$ ) is  $G_{12}$  (or  $G_{21}$ ), and the number of all the other genes is  $G_{22}$ , then the  $2 \times 2$  matrix  $\mathbf{G}$  can be directly used for Fisher's exact test to evaluate the  $P$  value  $A_{ab}$  between clusters  $a$  and  $b$ . After constructing a matrix  $\mathbf{A}$  filled with  $-\log(A_{ab})$  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. The same algorithm was also used to show the correspondence between the cell sub-clusters (EX1~5 and IN1~5) generated from snATAC-seq data and the cell subtypes identified from inDrops-seq data (Fig. 2f, g).

#### 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. 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 belongs to the same accession. It also requires that the percentage of the putative peaks in one super-enhancer is one of the highest among all genomic areas ( $P$  value  $< 0.01$ ). 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. ROSE [31] was used to confirm the super-enhancers called from APEC, with command line "python ROSE\_main.py -g hg19 -i top\_filtered\_peaks.gff -r P1-LSC.bam -s 12500 -t 2500 -o P1-LSC-SuperEnhancer" applied to the merged bam files of all the P1-LSC cells.

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

To test whether peaks in the same accession are closer in space, we integrated the Hi-C [71] data (GSE63525) and scATAC-seq [8] data (GSE65360) on GM12878 and K562 cells. The spatial correlation of different windows, both intra- and inter-chromosomal, can be directly extracted by Juicer [75]. Pearson's correlation matrix of the intra-chromosomal or inter-chromosomal windows can be calculated from the corresponding observed/expected matrix. We constructed 600 accessions by grouping peaks in the GM12878 and K562 scATAC-seq data in APEC, and removed the accessions 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 accession, termed

“Accession” correlations. For comparison, we also shuffled all peaks in different accessions to make fake accessions and re-collected the Hi-C correlations between windows that contained peaks in each fake accession, termed “Shuffled” correlations. Meanwhile, we also collected Hi-C correlations between windows that contained no peaks, termed “Non-accession” correlations. To calculate the Hi-C correlation within each topic, we used the default parameters of cisTopic to generate the topics, and then set the cutoff of the region score to 0.99 to filter out representative 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 [40, 76]. APEC reduced the dimension of the accession count matrix  $\mathbf{M}$  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 accession 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  $\mathbf{B}$  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  $\mathbf{B}'$ . 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) SnapATAC. We run PCA analysis on the normalized fragment count matrix generated by SnapATAC, and subjected the top 5 PCs to Monocle to build the trajectory.
- (4) 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. As the dimensions had been reduced by the LSI method, we skipped the PCA analysis.
- (5) 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.
- (6) 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, SPRING [43], to the accession count matrix  $\mathbf{M}$  from APEC to reconstruct the pathways for the hematopoietic differentiation dataset and thymocyte developmental dataset. We performed PCA analysis of the accession matrix  $\mathbf{M}$  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 50,000 to reduce computer time, but we recommend using all peaks if the users want to obtain better cell clusters. (1) For the dataset from hematopoietic cells, the  $-\log(Q)$  value threshold of high-quality peaks was set to 35 to retain 54,212 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  $-\log(Q)$  value was set to 5 to retain 38,683 high-quality peaks for subsequent processing. When LMPPs, HL60, and monocytes were added to this dataset with the AML cells, the threshold of  $-\log(Q)$  value was set to 8 to retain 42,139 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 130,685 peaks called by MACS2 were reserved for the fragment count matrix ( $Q$  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 mapping [77], and used RSEM [78] 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 analysis [79].

#### 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 website [80]. 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 suite [81] 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) [81] of each motif in each accession was used to construct a (motifs)  $\times$  (accessions) matrix  $\mathbf{F}$ . The enrichment of each motif on each cell was evaluated by the product of the matrix  $\mathbf{F}$  and the accession count matrix  $\mathbf{M}$ , i.e.,

$$\text{Motif}_{ij} = \sum_k F_{ik} \cdot M_{kj}$$

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.

### Supplementary information

Supplementary information accompanies this paper at <https://doi.org/10.1186/s13059-020-02034-y>.

**Additional file 1: Figure S1.** The 3 subtypes of CMP cells in the tSNE maps generated by other methods. **Figure S2.** 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: accession 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)** Box-plots showing the ARI values for the clustering of the blast and LSC cells from two AML patients. We sampled different tunable parameters for different algorithms. APEC: the accession number; cisTopic: the random seed; SnapATAC: the number of principal components and the number of nearest neighbors; LSI: the number of top SVD components; Cicero: the peak aggregation distance; 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. **(d)** 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. **(e)** The average ARI values of the noised data sampled from the fragment count matrix of the same dataset used in (d). The X-axis represents the percentage of noised elements in the matrix. Shaded error bar: 95% confidence interval. **Figure S3.** 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). **Figure S4.** 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 distances 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)** The distribution of the number of CCANs of peaks from the same accession (left), and the distribution of the number of accessions of peaks from the same CCAN (right). **(d)** Site links discovered by APEC and Cicero. **Figure S5. (a)** Box plots presenting the average spatial distance of peaks in the same accession or topic versus randomly shuffled peaks, and non-accessible genomic regions in the GM12878 cells. Spatial distance was estimated from chromosome conformation

capture (Hi-C) technology. Left panel: Hi-C correlation of intra-chromosomal windows; right panel: Hi-C correlation of inter-chromosomal windows. **(b)** The Hi-C profile of genomic regions between chr1:500,000-21,500,000 in GM12878 cells. 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 positions but close in space. **(c)** Box plots presenting the average spatial distance between peaks in the same accession versus randomly shuffled peaks and non-accessible genomic regions in K562 cells. **(d, e)** Top enriched motifs in the accessions with more than 500 peaks, in GM12878 (d) and K562 (e) cells. **(f)** Top enriched motifs of peaks in topics in GM12878 cells. **Figure S6. (a, b)** The computing time required for different algorithms to cluster cell numbers from 10,000 to 80,000 with all peaks (a) and 100,000 peaks (b). The data were sampled from the single-cell atlas of in vivo mammalian chromatin accessibility. CisTopic was performed using 8 CPU threads and all the other tools with 1 CPU thread. **(c-e)** The ARI values of the clustering results that used different numbers of accessions (c), nearest neighbors (d), and principle components (e). The dataset includes the cells from two AML patients and three cell lines. Default values are noted in red. **Figure S7. (a)** The clustering and cell-type classification of the mouse forebrain dataset by Cicero. Upper panel: cell clusters obtained by Cicero, illustrated in the tSNE diagram. Middle panel: the z-scores of the average gene scores of cell clusters, obtained by Cicero. Lower panel: the hierarchical clustering of the Pearson correlations between cell clusters identified by Cicero. **(b, c)** The clustering and cell-type classification of the same dataset by cisTopic and SnapATAC respectively, as in (a). **Figure S8. (a)** UCSC genome browser track diagram of the normalized fragment count around gene *CD34* for each hematopoietic cell type. **(b-g)** The pseudotime trajectories constructed by the combination of Monocle and the raw peak count matrix, the topic matrix from cisTopic, the normalized count matrix from SnapATAC, the LSI matrix, the aggregated model matrix from Cicero, and the bias corrected deviation matrix from chromVAR, respectively. **(h)** The pseudotime trajectory constructed by the combination of SPRING and the accession matrix from APEC. **Figure S9. (a)** Gating strategy of the mouse thymocytes in ftATAC-seq. **(b-d)** Quality control diagrams for the mouse thymocyte data, including reads numbers and percentage of fragments in peaks for each cell (b), average count of scATAC-seq insertions around TSS regions (c), and statistical distribution of fragment lengths (d). **(e)** The z-score of correlation between the cell types from ftATAC-seq and bulk ATAC-seq data. **Figure S10.** Selected significant motifs enriched in different thymocyte subtypes obtained by the APEC algorithm. **Figure S11.** 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. **Figure S12.** 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  $-\log(P\text{-value})$  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.

**Additional file 2.** Table S1.

**Additional file 3.** Table S2.

**Additional file 4.** Table S3.

**Additional file 5.** Review history.

### Acknowledgements

We thank the Howard Chang lab at Stanford University and CAS Interdisciplinary Innovation Team for the helpful discussion. We thank the USTC supercomputing center and the School of Life Science Bioinformatics Center for providing supercomputing resources for this project.

### Peer review information

Yixin Yao was the primary editor of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

### Review history

The review history is available as Additional file 5.

### 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 the ftATAC-seq technique and performed all scATAC-seq and scRNA-seq experiments with helps from LZ. KL analyzed the scRNA-seq data. JL helped to revise the manuscript. BL, YL, and KQ wrote the manuscript with inputs from all other authors. The authors read and approved the final manuscript.

### Funding

This work was supported by the National Key R&D Program of China (2017YFA0102900 to K.Q.) and by the National Natural Science Foundation of China grants (91940306, 81788101, 91640113, 31970858, and 31771428 to K.Q. and 81871479 to J.L.). It was also supported by the Fundamental Research Funds for the Central Universities WK2070000158 (to K.Q.) and the Anhui Provincial Natural Science Foundation grant BJ2070000097 (to B.L.) and 1908085QH326 (to Y.L.).

### Availability of data and materials

Mouse thymocytes *ft*ATAC-seq data can be obtained from the Gene Expression Omnibus (GEO) with the accession number GSE149480 [82] and is available via <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE149480>. Other published datasets used in this study are also available from NIH GEO with accession numbers GSE74310 [2], GSE65360 [8], GSE96772 [21], GSE100033 [10], GSE111586 [36], and GSE63525 [71]. APEC pipeline can be downloaded from the GitHub repository (<https://github.com/QuKunLab/APEC>) [83] or Zenodo with the access code DOI: <https://doi.org/10.5281/zenodo.3765969> [84].

### Ethics approval and consent to participate

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.

### Competing interests

J.F. is affiliated with HanGen Biotech (chief executive officer).

### Author details

<sup>1</sup>Department of Oncology, The First Affiliated Hospital of USTC, Division of Life Sciences and Medicine, Division of Molecular Medicine, Hefei National Laboratory for Physical Sciences at Microscale, University of Science and Technology of China, Hefei 230001, Anhui, China. <sup>2</sup>HanGene Biotech, Xiaoshan Innovation Polis, Hangzhou 310000, Zhejiang, China. <sup>3</sup>CAS Center for Excellence in Molecular Cell Sciences, The CAS Key Laboratory of Innate Immunity and Chronic Disease, University of Science and Technology of China, Hefei 230027, Anhui, China.

Received: 1 December 2019 Accepted: 30 April 2020

Published online: 12 May 2020

### References

- Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nat Methods*. 2013;10:1213–8.
- Corces MR, Buenrostro JD, Wu B, Greenside PG, Chan SM, Koenig JL, Snyder MP, Pritchard JK, Kundaje A, Greenleaf WJ, et al. Lineage-specific and single-cell chromatin accessibility charts human hematopoiesis and leukemia evolution. *Nat Genet*. 2016;48:1193–203.
- Wu J, Huang B, Chen H, Yin Q, Liu Y, Xiang Y, Zhang B, Liu B, Wang Q, Xia W, et al. The landscape of accessible chromatin in mammalian preimplantation embryos. *Nature*. 2016;534:652–7.
- Jorstad NL, Wilken MS, Grimes WN, Wohl SG, VandenBosch LS, Yoshimatsu T, Wong RO, Rieke F, Reh TA. Stimulation of functional neuronal regeneration from Muller glia in adult mice. *Nature*. 2017;548:103–7.
- Su Y, Shin J, Zhong C, Wang S, Roychowdhury P, Lim J, Kim D, Ming GL, Song H. Neuronal activity modifies the chromatin accessibility landscape in the adult brain. *Nat Neurosci*. 2017;20:476–83.
- Denny SK, Yang R, Chuang CH, Brady JJ, Lim JS, Gruner BM, Chiou SH, Schep AN, Baral J, Hamard C, et al. Nfib promotes metastasis through a widespread increase in chromatin accessibility. *Cell*. 2016;166:328–42.
- Qu K, Zaba LC, Satpathy AT, Giresi PG, Li R, Jin Y, Armstrong R, Jin C, Schmitt N, Rahbar Z, et al. Chromatin accessibility landscape of cutaneous T cell lymphoma and dynamic response to HDAC inhibitors. *Cancer Cell*. 2017;32:27–41 e24.
- Buenrostro JD, Wu B, Litzenburger UM, Ruff D, Gonzales ML, Snyder MP, Chang HY, Greenleaf WJ. Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature*. 2015;523:486–90.
- Lake BB, Chen S, Sos BC, Fan J, Kaeser GE, Yung YC, Duong TE, Gao D, Chun J, Kharchenko PV, Zhang K. Integrative single-cell analysis of transcriptional and epigenetic states in the human adult brain. *Nat Biotechnol*. 2018;36:70–80.
- Preissl S, Fang R, Huang H, Zhao Y, Raviram R, Gorkin DU, Zhang Y, Sos BC, Afzal V, Dickel DE, et al. Single-nucleus analysis of accessible chromatin in developing mouse forebrain reveals cell-type-specific transcriptional regulation. *Nat Neurosci*. 2018;21:432–9.
- Cusanovich DA, Daza R, Adey A, Pliner HA, Christiansen L, Gunderson KL, Steemers FJ, Trapnell C, Shendure J. Multiplex single cell profiling of chromatin accessibility by combinatorial cellular indexing. *Science*. 2015;348:910–4.
- Cusanovich DA, Reddington JP, Garfield DA, Daza RM, Aghamirzaie D, Marco-Ferreres R, Pliner HA, Christiansen L, Qiu X, Steemers FJ, et al. The cis-regulatory dynamics of embryonic development at single-cell resolution. *Nature*. 2018;555:538–42.
- Butler A, Hoffman P, Smibert P, Papalexi E, Satija R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat Biotechnol*. 2018;36:411–20.
- Kiselev VY, Kirschner K, Schaub MT, Andrews T, Yiu A, Chandra T, Natarajan KN, Reik W, Barahona M, Green AR, Hemberg M. SC3: consensus clustering of single-cell RNA-seq data. *Nat Methods*. 2017;14:483–6.
- Wang B, Zhu J, Pierson E, Ramazzotti D, Batzoglu S. Visualization and analysis of single-cell RNA-seq data by kernel-based similarity learning. *Nat Methods*. 2017;14:414–6.
- Wolf FA, Angerer P, Theis FJ. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol*. 2018;19:15.
- Schep AN, Wu B, Buenrostro JD, Greenleaf WJ. chromVAR: inferring transcription-factor-associated accessibility from single-cell epigenomic data. *Nat Methods*. 2017;14:975–8.
- Pliner HA, Packer JS, McFaline-Figueroa JL, Cusanovich DA, Daza RM, Aghamirzaie D, Srivatsan S, Qiu X, Jackson D, Minkina A, et al. Cicero predicts cis-regulatory DNA interactions from single-cell chromatin accessibility data. *Mol Cell*. 2018;71:858–71 e858.
- Bravo Gonzalez-Blas C, Minnoye L, Papanokrati D, Aibar S, Hulselmans G, Christiaens V, Davie K, Wouters J, Aerts S. cisTopic: cis-regulatory topic modeling on single-cell ATAC-seq data. *Nat Methods*. 2019;16:397–400.
- Fang R, Preissl S, Hou X, Lucero J, Wang X, Motamedi A, Shiau AK, Mukamel EA, Zhang Y, Behrens MM, Ecker J, Ren B. Fast and Accurate Clustering of Single Cell Epigenomes Reveals Cis-Regulatory Elements in Rare Cell Types. *bioRxiv* 2019; 615179; <https://doi.org/10.1101/615179>.

21. Buenrostro JD, Corces MR, Lareau CA, Wu B, Schep AN, Aryee MJ, Majeti R, Chang HY, Greenleaf WJ. Integrated single-cell analysis maps the continuous regulatory landscape of human hematopoietic differentiation. *Cell*. 2018;173:1535–48 e1516.
22. Doulatov S, Notta F, Eppert K, Nguyen LT, Ohashi PS, Dick JE. Revised map of the human progenitor hierarchy shows the origin of macrophages and dendritic cells in early lymphoid development. *Nat Immunol*. 2010;11:585–U552.
23. Karamitros D, Stoilova B, Aboukhalil Z, Hamey F, Reinisch A, Samitsch M, Quek L, Otto G, Repapi E, Doondea J, et al. Single-cell analysis reveals the continuum of human lympho-myeloid progenitor cells. *Nature Immunol*. 2018;19:85.
24. Ng AP, Loughran SJ, Metcalf D, Hyland CD, de Graaf CA, Hu YF, Smyth GK, Hilton DJ, Kile BT, Alexander WS. Erg is required for self-renewal of hematopoietic stem cells during stress hematopoiesis in mice. *Blood*. 2011;118:2454–61.
25. Ong CT, Corces VG. CTCF: an architectural protein bridging genome topology and function. *Nat Rev Genet*. 2014;15:234–46.
26. Dore LC, Crispino JD. Transcription factor networks in erythroid cell and megakaryocyte development. *Blood*. 2011;118:231–9.
27. Huang M, Wang J, Torre E, Dueck H, Shaffer S, Bonasio R, Murray JI, Raj A, Li M, Zhang NR. SAVER: gene expression recovery for single-cell RNA sequencing. *Nat Methods*. 2018;15:539–42.
28. van Dijk D, Sharma R, Nainys J, Yin K, Kathail P, Carr AJ, Burdziak C, Moon KR, Chaffer CL, Pattabiraman D, et al. Recovering gene interactions from single-cell data using data diffusion. *Cell*. 2018;174:716.
29. Oberst A, Malatesta M, Aqeilan RI, Rossi M, Salomoni P, Murillas R, Sharma P, Kuehn MR, Oren M, Croce CM, et al. The Nedd4-binding partner 1 (N4BP1) protein is an inhibitor of the E3 ligase Itch. *Proc Natl Acad Sci U S A*. 2007;104:11280–5.
30. Eguchi M, Eguchi-Ishimae M, Seto M, Morishita K, Suzuki K, Ueda R, Ueda K, Kamada N, Greaves M. GPHN, a novel partner gene fused to MLL in a leukemia with t(11;14)(q23;q24). *Genes Chromosomes Cancer*. 2001;32:212–21.
31. Whyte WA, Orlando DA, Hnisz D, Abraham BJ, Lin CY, Kagey MH, Rahl PB, Lee TI, Young RA. Master transcription factors and mediator establish super-enhancers at key cell identity genes. *Cell*. 2013;153:307–19.
32. Splinter E, Heath H, Kooren J, Palstra RJ, Klous P, Grosveld F, Galjart N, de Laat W. CTCF mediates long-range chromatin looping and local histone modification in the beta-globin locus. *Genes Dev*. 2006;20:2349–54.
33. Handoko L, Xu H, Li GL, Ngan CY, Chew E, Schnapp M, Lee CWH, Ye CP, Ping JH, Mulawadi F, et al. CTCF-mediated functional chromatin interactome in pluripotent cells. *Nat Genet*. 2011;43:630–U198.
34. Narendra V, Rocha PP, An DS, Raviram R, Skok JA, Mazzoni EO, Reinberg D. CTCF establishes discrete functional chromatin domains at the Hox clusters during differentiation. *Science*. 2015;347:1017–21.
35. Tang ZH, Luo OJ, Li XW, Zheng MZ, Zhu JJ, Szalaj P, Trzaskoma P, Magalska A, Wlodarczyk J, Rusczycki B, et al. CTCF-mediated human 3D genome architecture reveals chromatin topology for transcription. *Cell*. 2015;163:1611–27.
36. Cusanovich DA, Hill AJ, Aghamirzaie D, Daza RM, Pliner HA, Berletch JB, Filippova GN, Huang XF, Christiansen L, DeWitt WS, et al. A single-cell atlas of in vivo mammalian chromatin accessibility. *Cell*. 2018;174:1309.
37. Satpathy AT, Saligrama N, Buenrostro JD, Wei Y, Wu B, Rubin AJ, Granja JM, Lareau CA, Li R, Qi Y, et al. Transcript-indexed ATAC-seq for precision immune profiling. *Nat Med*. 2018;24:580–90.
38. Hrvatin S, Hochbaum DR, Nagy MA, Cicconet M, Robertson K, Cheadle L, Zilionis R, Ratner A, Borges-Monroy R, Klein AM, et al. Single-cell analysis of experience-dependent transcriptomic states in the mouse visual cortex. *Nat Neurosci*. 2018;21:120–9.
39. Bendall SC, Davis KL, Amir el AD, Tadmor MD, Simonds EF, Chen TJ, Shenfeld DK, Nolan GP, Pe'er D. Single-cell trajectory detection uncovers progression and regulatory coordination in human B cell development. *Cell* 2014, 157:714–725.
40. Trapnell C, Cacchiarelli D, Grimsby J, Pokharel P, Li S, Morse M, Lennon NJ, Livak KJ, Mikkelsen TS, Rinn JL. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol*. 2014;32:381–6.
41. Qiu X, Mao Q, Tang Y, Wang L, Chawla R, Pliner HA, Trapnell C. Reversed graph embedding resolves complex single-cell trajectories. *Nat Methods*. 2017;14:979–82.
42. Setty M, Tadmor MD, Reich-Zeliger S, Angel O, Salame TM, Kathail P, Choi K, Bendall S, Friedman N, Pe'er D. Wishbone identifies bifurcating developmental trajectories from single-cell data. *Nat Biotechnol*. 2016;34:637–45.
43. Weinreb C, Wolock S, Klein AM. SPRING: a kinetic interface for visualizing high dimensional single-cell expression data. *Bioinformatics*. 2018;34:1246–8.
44. Habib N, Li YQ, Heidenreich M, Swiech L, Avraham-Davidi I, Trombetta JJ, Hession C, Zhang F, Regev A. Div-Seq: single-nucleus RNA-Seq reveals dynamics of rare adult newborn neurons. *Science*. 2016;353:925–8.
45. Olsson A, Venkatasubramanian M, Chaudhri VK, Aronow BJ, Salomonis N, Singh H, Grimes HL. Single-cell analysis of mixed-lineage states leading to a binary cell fate choice. *Nature*. 2016;537:698.
46. Zhou F, Li XL, Wang WL, Zhu P, Zhou J, He WY, Ding M, Xiong FY, Zheng XN, Li Z, et al. Tracing haematopoietic stem cell formation at single-cell resolution. *Nature*. 2016;533:487.
47. Treutlein B, Lee QY, Camp JG, Mall M, Koh W, Shariati SAM, Sim S, Neff NF, Skotheim JM, Wernig M, Quake SR. Dissecting direct reprogramming from fibroblast to neuron using single-cell RNA-seq. *Nature*. 2016;534:391.
48. Adolfsson J, Mansson R, Buza-Vidas N, Hultquist A, Liuba K, Jensen CT, Bryder D, Yang L, Borge OJ, Thoren LA, et al. Identification of Flt3+ lympho-myeloid stem cells lacking erythro-megakaryocytic potential a revised road map for adult blood lineage commitment. *Cell*. 2005;121:295–306.
49. Chistiakov DA, Orekhov AN, Sobenin IA, Bobryshev YV. Plasmacytoid dendritic cells: development, functions, and role in atherosclerotic inflammation. *Front Physiol*. 2014;5:279.
50. Germain RN. T-cell development and the CD4-CD8 lineage decision. *Nat Rev Immunol*. 2002;2:309–22.
51. Chen X, Shen Y, Draper W, Buenrostro JD, Litzgenberger U, Cho SW, Satpathy AT, Carter AC, Ghosh RP, East-Seletsky A, et al. ATAC-seq reveals the accessible genome by transposase-mediated imaging and sequencing. *Nat Methods*. 2016;13:1013–20.
52. Chen X, Litzgenberger UM, Wei Y, Schep AN, LaGory EL, Choudhry H, Giaccia AJ, Greenleaf WJ, Chang HY. Joint single-cell DNA accessibility and protein epitope profiling reveals environmental regulation of epigenomic heterogeneity. *Nat Commun*. 2018;9:4590.
53. Yoshida H, Lareau CA, Ramirez RN, Rose SA, Maier B, Wroblewska A, Desland F, Chudnovskiy A, Mortha A, Dominguez C, et al. The cis-regulatory atlas of the mouse immune system. *Cell*. 2019;176:897–912 e820.

54. Godfrey DI, 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.* 1993;150:4244–52.
55. Taniuchi I, Osato M, Egawa T, Sunshine MJ, Bae SC, Komori T, Ito Y, Littman DR. Differential requirements for Runx proteins in CD4 repression and epigenetic silencing during T lymphocyte development. *Cell.* 2002;111:621–33.
56. Ioannidis V, Beermann F, Clevers H, Held W. The beta-catenin–TCF-1 pathway ensures CD4(+)CD8(+) thymocyte survival. *Nat Immunol.* 2001;2:691–7.
57. Yu S, Zhou X, Steinke FC, Liu C, Chen SC, Zagorodna O, Jing X, Yokota Y, Meyerholz DK, Mullighan CG, et al. The TCF-1 and LEF-1 transcription factors have cooperative and opposing roles in T cell development and malignancy. *Immunity.* 2012;37:813–26.
58. Sun Z, Unutmaz D, Zou YR, Sunshine MJ, Pierani A, Brenner-Morton S, Mebius RE, Littman DR. Requirement for RORgamma in thymocyte survival and lymphoid organ development. *Science.* 2000;288:2369–73.
59. Gerondakis S, Fulford TS, Messina NL, Grumont RJ. NF-kappaB control of T cell development. *Nat Immunol.* 2014;15:15–25.
60. Picelli S, Bjorklund AK, Faridani OR, Sagasser S, Winberg G, Sandberg R. Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat Methods.* 2013;10:1096–8.
61. Higashi Y, Moribe H, Takagi T, Sekido R, Kawakami K, Kikutani H, Kondoh H. Impairment of T cell development in deltaEF1 mutant mice. *J Exp Med.* 1997;185:1467–79.
62. Heath H, Ribeiro de Almeida C, Sleutels F, Dingjan G, van de Nobelen S, Jonkers I, Ling KW, Gribnau J, Renkawitz R, Grosveld F, et al. CTCF regulates cell cycle progression of alphabeta T cells in the thymus. *EMBO J.* 2008;27:2839–50.
63. Andres-Barquin PJ, Hernandez MC, Israel MA. Id4 expression induces apoptosis in astrocytic cultures and is down-regulated by activation of the cAMP-dependent signal transduction pathway. *Exp Cell Res.* 1999;247:347–55.
64. Carey JP, Knowell AE, Chinaranagari S, Chaudhary J. Id4 promotes senescence and sensitivity to doxorubicin-induced apoptosis in DU145 prostate cancer cells. *Anticancer Res.* 2013;33:4271–8.
65. Surh CD, Sprent J. T-cell apoptosis detected in situ during positive and negative selection in the thymus. *Nature.* 1994;372:100–3.
66. Huesmann M, Scott B, Kiselow P, von Boehmer H. Kinetics and efficacy of positive selection in the thymus of normal and T cell receptor transgenic mice. *Cell.* 1991;66:533–40.
67. 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.* 1991;173:323–32.
68. Chow SC, Snowden R, Orrenius S, Cohen GM. Susceptibility of different subsets of immature thymocytes to apoptosis. *FEBS Lett.* 1997;408:141–6.
69. Granja JM, Klemm S, McGinnis LM, Kathiria AS, Mezger A, Corces MR, Parks B, Gars E, Liedtke M, Zheng GXY, et al. Single-cell multiomic analysis identifies regulatory programs in mixed-phenotype acute leukemia. *Nature Biotechnology.* 2019;37:1458.
70. Picelli S, Faridani OR, Bjorklund AK, Winberg G, Sagasser S, Sandberg R. Full-length RNA-seq from single cells using Smart-seq2. *Nat Protoc.* 2014;9:171–81.
71. Rao SS, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, Robinson JT, Sanborn AL, Machol I, Omer AD, Lander ES, Aiden EL. A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell.* 2014;159:1665–80.
72. Zuo Z, Jin Y, Zhang W, Lu Y, Li B, Qu K. ATAC-pipe: general analysis of genome-wide chromatin accessibility. *Briefings Bioinformatics* 2018:bby056-bby056.
73. Chen H, Lareau C, Andreani T, Vinyard ME, Garcia SP, Clement K, Andrade-Navarro MA, Buenrostro JD, Pinello L. Assessment of computational methods for the analysis of single-cell ATAC-seq data. *Genome Biol.* 2019;20:241.
74. Bravo Gonzalez-Blas C, Minnoye L, Papasokrati D, Aibar S, Hulselmans G, Christiaens V, Davie K, Wouters J, Aerts S. Cis-topic modelling of single cell epigenomes. *bioRxiv* 2018;370346; <https://doi.org/10.1101/370346>.
75. Durand NC, Shamim MS, Machol I, Rao SS, Huntley MH, Lander ES, Aiden EL. Juicer provides a one-click system for analyzing loop-resolution Hi-C experiments. *Cell Syst.* 2016;3:95–8.
76. Qiu X, Hill A, Packer J, Lin D, Ma YA, Trapnell C. Single-cell mRNA quantification and differential analysis with Census. *Nat Methods.* 2017;14:309–15.
77. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras TR. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics.* 2013;29:15–21.
78. Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics.* 2011;12:323.
79. Macosko EZ, Basu A, Satija R, Nemes J, Shekhar K, Goldman M, Tirosh I, Bialas AR, Kamitaki N, Martersteck EM, et al. Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. *Cell.* 2015;161:1202–14.
80. McLean CY, Bristol D, Hiller M, Clarke SL, Schaar BT, Lowe CB, Wenger AM, Bejerano G. GREAT improves functional interpretation of cis-regulatory regions. *Nat Biotechnol.* 2010;28:495–501.
81. Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, Ren J, Li WW, Noble WS. MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res.* 2009;37:W202–8.
82. Li B, Li Y, Qu K. Single-cell ftATACseq data of mouse thymocytes. *Gene Expression Omnibus*: <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE149480>. Accessed 28 Apr 2020.
83. Li B, Qu K. APEC: single cell epigenomic clustering based on accessibility pattern. *GitHub*: <https://github.com/QuKunLab/APEC>. Accessed 23 Apr 2020.
84. Li B, Qu K. APEC: Accessibility Pattern based Epigenomic Clustering algorithm for single cell ATAC-seq data. *Zenodo*: <https://doi.org/10.5281/zenodo.3765969>. Accessed 23 Apr 2020.

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.