# 1 Single-cell ATAC-seq clustering and differential
# 2 analysis by convolution-based approach

3 **Li Lin[1*], Liye Zhang[1*]**

4 [1]School of Life Science and Technology, ShanghaiTech University, Shanghai, China.

5 [*]Correspondence should be addressed to L.L. (linli@shanghaitech.edu.cn) or L.Z.

6 (zhangly@shanghaitech.edu.cn).

7

# 8 Abstract

9 Single-cell ATAC-seq is a powerful tool to interrogate the epigenetic heterogeneity of cells. Here,

10 we present a novel method to calculate the pairwise similarities between single cells by directly

11 comparing their Tn5 insertion profiles instead of the binary accessibility matrix using a

12 convolution-based approach. We demonstrate that our method retains the biological

13 heterogeneity of single cells and removes undesirable batch effects, which leads to more

14 accurate results on downstream analyses such as dimension reduction and clustering. Based on

15 the similarity matrix learned from epiConv, we develop an algorithm to infer differentially

16 accessible peaks directly from heterogeneous cell population to overcome the limitations of

17 conventional differential analysis through two-group comparisons.

18

# Introduction

19

20    The expression of genes is regulated by a series of transcription factors (TFs) that bind to the

21    regulatory elements of the genome. As the accessible chromatin covers more than 90% TF

22    binding regions, many techniques, such as Assay for Transposase-Accessible Chromatin using

23    sequencing (ATAC-seq), have been developed to detect the accessible states of chromatin[1, 2].

24    Recent technical advancements in ATAC-seq have made it possible to profile the chromatin states

25    of single cells at a high-throughput manner[3-5]. However, both data processing and interpretation

26    of single-cell ATAC-seq (scATAC-seq) data is more challenging than single-cell RNA-seq (scRNA-seq)

27    data owing to low DNA copy number and complexity of chromatin states[1].

28    Up to now, most methods cluster single cells based on a peak by cell matrix (e.g. Buenrostro

29    et al. 2015[6]). Unlike well-annotated RNA transcripts in the genome, the exact locus of regulatory

30    elements is largely uncharacterized and must be learned from the data itself. However, learning

31    cell type specific regulatory elements from cell mixtures is problematic[7]. Moreover, given that

32    there are no golden rules to define functional elements across the genome, the strategies to

33    perform such task varied considerably in different studies[6, 8], and its effect on downstream

34    analyses is largely unknown.

35    Detecting differentially expressed genes (or differentially accessible peaks for ATAC-seq, we

36    call them DE peaks below) is another important task in single cell analysis. In a conventional

37    pipeline, cells are first grouped into several clusters and subsequent differential analysis is

38    performed by comparison between clusters. Thus, the resolution settings (e.g. number of clusters)

39    may have strong effects on the identification of genes or locus accounting for the heterogeneity

40    of cell population. Recently one method incorporated pseudotime as one predictor into the

41    regression model to infer DE peaks, instead of performing two-group comparisons[9]. But it

42    required cells to be properly embedded into one dimensional space (e.g. pseudotime through

43    differentiation process), which greatly limits its application in complex cell population. Moreover,

44    cells still need to be clustered into small groups (50~100 cells). Such processing step overcomes

45    the sparsity of scATAC-seq data but reduces the sample size. In scRNA-seq, an alternative

46    approach is to find highly variable genes instead of differentially expressed genes, which does not

47    require the clustering of cell population to be defined. But this strategy cannot be applied to

48    scATAC-seq as the chromatin state is always binarized. Despite that, several state-of-the-art tools

49    designed for scATAC-seq merge individual peaks into meta features (regulomes, topics, principal

50    components, k-mers, etc.) to overcome the sparsity of data[3, 10, 11]. Subsequent differential

51    analysis is performed on meta features instead of individual peaks. Such strategy may help reveal

52    the epigenetic programs that governs the cell identities but lacks sufficient resolution for the

53    dynamic change of individual peaks.

54        Here, we introduce a novel tool, named epiConv, for scATAC-seq analysis. EpiConv addresses

55    two important questions in scATAC-seq analysis, cell clustering and differential analysis. Unlike

56    most of existing methods, epiConv learns the similarities (or distances) between single cells from

57    their raw Tn5 insertion profiles by a convolution-based approach, instead of a binary accessibility

58    matrix. We demonstrate that epiConv retains biological heterogeneity of single cells and removes

59    unwanted variations derived from multiple batches or sample preparing protocols. Utilizing the

60    similarities learned by epiConv, we also develop an algorithm to infer DE peaks among single cells

61    that can be directly applied to cell mixtures without resolving the intra population structure.

62

# 63 Results

**64 Infer the similarity from Tn5 insertion profiles**

65     First, we give an overview of the algorithm that calculates the similarity between cells from

66 their Tn5 insertion profiles (**Fig. 1**). Given two cells, A with m insertions and B with n insertions in

67 one genomic region, we collapse the insertions into a continuous distribution across the genome

68 by Gaussian smoothing as follows:

69
$$f_{Ai}(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu_{Ai})^2}{2\sigma^2}\right), \; f_A(x) = \sum_i^m f_{Ai}(x)$$

70
$$f_{Bj}(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu_{Bj})^2}{2\sigma^2}\right), \; f_B(x) = \sum_j^n f_{Bj}(x)$$

71     Where $\mu_{Ai}$ is the locus of insertion *i* in cell A, $\mu_{Bj}$ is the locus of insertion *j* in cell B, $f_A(x)$

72 and $f_B(x)$ give the overall chromatin states of cell A and cell B in the given region. The similarity

73 between A and B over the given region ($S_{AB}$) is calculated by the convolution of $f_A(x)$ and $f_B(x)$

74 and can be solved analytically as follows:

$$s_{AB} = \int f_A(x) f_B(x) dx = C \cdot \sum_{i,j} \exp\left(-\frac{\left(\mu_{Ai} - \mu_{Bj}\right)^2}{4\sigma^2}\right)$$

75 Where C is an σ dependent constant. In this study, parameter σ is set to 100 bp. To save running

76 time, long distance (> 4σ) is treated as infinity. Through weighted aggregation of the similarities

77 from all informative regions across the genome and proper normalization with respect to

78 sequencing depth, we can obtain the normalized similarity score between any two cells.

79 Subsequent analyses such as dimension reduction or clustering can be performed on the

80 similarity matrix. We also develop a simplified version of epiConv (epiConv-simp), which can be

81 applied to binary accessibility matrix like existing methods. The simplified version does not

82 perform as well as the full version but always generates similar results and runs much faster. In

83 the benchmarking below, we show the results from both full and simplified versions. Other

84    details of epiConv are provided in Methods section.

85

## EpiConv outperforms other methods in cell lines data

87      We evaluated the performance of epiConv on several datasets and compared it with

88    cisTopic[11], Latent Semantic Indexing (LSI)[3] and SnapATAC, which show better performance than

89    other methods in one recently published benchmarking study[7]. We first applied epiConv to the

90    data from Buenrostro et al. 2015[6]. Specifically, we mixed the data of four cell lines from

91    hematopoietic lineages (K562, GM12878, HL-60 and TF-1) together and tested whether epiConv

92    could cluster single cells correctly based on their biological identities. Given the apparent

93    difference among cell lines, each method performed well in clustering single cells from the same

94    cell line together (**Fig. 2**). However, we found that LSI could not clearly segregate drug-treated

95    and untreated K562 cells. CisTopic segregated treated and untreated K562 cells into two clusters

96    but cells treated by different drugs were still mixed together. Only epiConv-full and SnapATAC

97    grouped K562 cells treated by different drugs into distinct clusters, while epiConv-full showed

98    higher resolution than SnapATAC, yielding the best results. Notably, untreated K562 cells from

99    three replicates were grouped into one cluster without obvious batch effects. Thus, the

100   segregation of cells treated by different drugs was more likely to be attributed to their biological

101   variations rather than batch effects. EpiConv-simp suggested one extra cluster with mixed cell

102   types, performing worse than other methods (**Fig. 2b**). These results highlighted the superiority

103   of directly comparing the Tn5 insertions when performing clustering on highly similar cells

104   (SnapATAC divided the genome into equal-length bins instead of peak calling, which could also be

105   considered as direct comparison on Tn5 insertions but with decreased resolution than

106    epiConv).However, we found that the worse performance of epiConv-simp was partially due to

107    improper denoising method. With an alternative denoising method, epiConv-simp provided good

108    results but still with lower resolution than the full version (top-right in **Fig. 2b**, see

109    **Supplementary Note 1** for the details of alternative denoising method). These results suggested

110    that the matrix still contained the variations derived from drug treatments but they could be

111    easily overwhelmed by noise.

112

### EpiConv removes batch effects in scATAC-seq data

114    Next, we applied epiConv to the data generated by droplet-based protocol from Satpathy et

115    al. 2019[4]. The authors reported detectable batch effects from LSI method that confounded

116    downstream analyses. Here we asked whether epiConv could perform better. We tested the

117    performance of epiConv on two datasets, one dataset containing cells from two batches of

118    unsorted peripheral blood mononuclear cells (PBMCs), two batches of sorted CD4+CD45RA+

119    naïve CD4 T cells and two batches of sorted CD4+CD45RA- memory CD4 T cells (PBMC dataset),

120    and the other dataset containing two batches of sorted CD34+ hematopoietic progenitors (CD34+

121    dataset). Based on our preliminary analyses, epiConv still suffered from batch effects but was less

122    sensitive compared to other methods (data not shown). So, we developed a simple method to

123    remove detectable effects. Although there are many methods to remove batch effects for

124    scRNA-seq data, few studies examined their performance on scATAC-seq data. So, we just

125    compared the results of epiConv to other methods without any batch correction.

126    In PBMC dataset, the majority of cells from two replicates of memory CD4 T cells were

127    clustered into one tightly related group by epiConv and were close to a small fraction of unsorted

128    PBMCs. Two replicates of naive CD4 T cells also showed similar results. Other unsorted PBMCs

129    formed several groups without strong batch effects (**Fig. 3a**). On the contrary, cells were mostly

130    clustered by batches for cisTopic, LSI and SnapATAC (**Fig. 3b, Fig. S1a,b**). These results

131    demonstrated that epiConv successfully removed batch effects. To verify whether epiConv

132    clustered single cells based on their biological identities, we marked single cells according to their

133    annotations from Satpathy et al. 2019[4]. The results of epiConv were also largely consistent with

134    the annotations and revealed all major lineages of PBMCs (T cells, NK cells, B cells and Monocytes)

135    and several subpopulations of T cells (**Fig. S2a-d**). In CD34+ dataset, epiConv still performed

136    better in removing batch effects compared to cisTopic, LSI and SnapATAC (**Fig. 3c,d, Fig. S1c,d**).

137    Based on the annotations from Satpathy et al. 2019[4], the results of epiConv were also consistent

138    with our knowledge on hematopoietic differentiation (**Fig. S2e-h**). Moreover, only epiConv and

139    cisTopic clearly revealed the trajectory of hematopoietic differentiation in unsupervised manner,

140    while the results of epiConv were with higher resolution and less noise than cisTopic.

141        To demonstrate that the power of epiConv was not restricted to specific cell lineages or

142    sample-preparing protocols, we combined scATAC-seq data of adult mouse brain from three

143    experimental protocols, mouse cortex from 10x Genomics, whole mouse brain from droplet

144    single-cell assay for transposase-accessible chromatin using sequencing (dscATAC-seq)[5] and sci-

145    protocols for chromatin accessibility (sci-ATAC-seq)[8]. The dataset contained single cells from 5

146    batches, one from 10x Genomics, two from dscATAC-seq and two from sci-ATAC-seq. Consistent

147    with previous results, epiConv performed better than cisTopic, LSI and SnapATAC in removing

148    batch effects (**Fig. 3e,f, Fig. S1e,f**) and agreed with the annotations from Cusanovich et al. 2018[8]

149    and Lareau et al. 2019[5] by clustering cells with the same identity together (**Fig. S2i,m**). Although

150    cisTopic suffered from batch effects, it largely agreed with the annotations from original articles

151    within each batch (**Fig. S2j,n**). However, LSI and snapATAC performed worse when comparing

152    them with the annotations from original articles (**Fig. S2k,l,o,p**). Although we lacked direct

153    evidence to evaluate which method performed best in clustering cells according to their cell

154    identities, the results of epiConv and cisTopic largely agreed with each other and could be

155    supported by the annotations from original article. Besides that, only epiConv was capable of

156    clustering cells in a batch-independent manner. Finally, we compared the results between full and

157    simplified versions of epiConv. Simplified version was highly consistent with full version on the

158    three datasets described above and also performed better than other methods (**Fig. S3**).

159

160    **EpiConv is scalable with large datasets**

161        As the full version of epiConv do pairwise comparisons between single cells, the step of

162    insertions counting is slower than other methods but can be split into small jobs and run in

163    parallel. Based on our tests, it requires 75 CPU hours for 50 million fragments from 5,000 cells

164    (after removing low quality cells and fragments outside informative regions) and 2,400 CPU hours

165    for 270 million fragments from 20,000 cells. The simplified version runs much faster and can be

166    applied to large datasets. Based on our tests, the simplified version requires 17 hours and 520 GB

167    RAM for the Mouse Cell Atlas dataset[8] (81,173 cells and 436,206 peaks) with single thread, faster

168    than cisTopic (48 hours) but slower than LSI (1 hour). SnapATAC failed to run on the full dataset of

169    Mouse Cell Atlas dataset in the step of calculating Jaccard distance due to the memory limitation

170    for single object in R (This error may depend on the system as Chen et al.[7] reported that

171    SnapATAC could run on the full dataset. Actually, we also encountered the same error for epiConv

172    but we modified our scripts to avoid it). The results of Mouse Cell Atlas dataset by epiConv-simp

173    also largely agreed with the annotations from Cusanovich et al. 2018[8] (**Fig. S4**).

174       Notably, a large proportion of cells were marked as unknown in the Mouse Cell Atlas dataset

175    (**Fig. 4a-d**). In the results of cisTopic , LSI and SnapATAC (we randomly sampled 25% cells in

176    Mouse Cell Altas dataset for SnapATAC), these cells formed a large cluster of their own, showed

177    close relationships with several clusters with known identities but did not overlap with them (**Fig.**

178    **4a-c**). However, unknown cells did not form a single cluster but were mixed with other known cell

179    types in the results of epiConv-simp (mainly associated with 6 clusters with more than 10% cells

180    marked as unknown, **Fig. 4d**). This might suggest a large improvement of epiConv over existing

181    methods. In order to validate our findings, we aggregated the cells with known and unknown cell

182    identities respectively for each cluster. Then we calculated the spearman correlation between the

183    12 aggregated samples over a set of highly accessible peaks (accessible in at least 1% cells from

184    these 6 clusters). We found that all unknown samples showed highest correlations with

185    corresponding known samples within the same clusters (**Fig. 4e**). Individual unknown samples

186    showed low correlations with each other, suggesting that epiConv successfully "demultiplexed"

187    unknown cells by their biological identities. By these results, we confirmed that epiConv showed

188    significant improvements over existing methods on the Mouse Cell Atlas dataset. Combined with

189    other results of epiConv-simp mentioned above, we concluded that in most cases epiConv-simp

190    also proved to be a reliable tool for the investigation of large datasets.

191

192    **EpiConv detects differentially accessible peaks in cell mixtures**

193       In the section below, we aim to develop an algorithm to infer DE peaks directly from cell

194    mixtures. Our algorithm compares the number of accessible cells among each cell's neighbors

195    with the proportion of accessible cells in cell mixture for each peak and turns the binary

196    chromatin states into normalized z-scores, which show the enrichment of accessible cells among

197    neighbors (we call it z-scores below). If the number of cells showing high z-scores for one peak

198    exceeds the threshold, we then consider the peak to be differentially accessible. Details of our

199    algorithm can be found in Methods.

200        In order to test whether the algorithm could detect DE peaks in cell mixture, we first applied

201    our method to one dataset of myoblast differentiation[9]. We found that although epiConv could

202    reconstruct the differentiation process of myoblasts, where cells were roughly ordered by

203    harvesting times (**Fig. 5a,b**), it was difficult to cluster cells due to the continuous differentiation

204    process. Using our algorithm, we detected 37,107 peaks to be differentially accessible during the

205    differentiation process. To show the dynamics of DE peaks, we plotted heatmap of z-scores,

206    where cells and DE peaks were embedded into one-dimensional (1D) space based on the

207    similarity matrix and the spearman correlation of z-scores between peaks (**Fig. 5b**). The results

208    showed approximately half peaks to be more accessible in the early stage of differentiation and

209    others to be more accessible in the later stage. The dynamic changes of z-scores along

210    differentiation was consistent with merged scATAC-seq profiles by harvesting times,

211    demonstrating the reliability of our algorithm (**Fig. 5c**).

212        Next, we want to test the sensitivity of our algorithm. We applied our algorithm to the

213    HSC-MPP-LMPP cluster in the CD34+ dataset. We chose the HSC-MPP-LMPP cluster because up

214    to now, few methods could distinguish MPPs from HSCs in scATAC-seq data duo to the high

215    similarity between them (see the benchmarking study of Chen et al.[7]). Given epiConv already

216     removed most of batch effects, we could include cells from both replicates of CD34+ dataset to

217     increase the statistical power. Through our algorithm, we detected 27,612 DE peaks within the

218     HSC-MPP-LMPP cluster. The dynamic changes of z-scores were highly consistent with the bulk

219     ATAC-seq profiles of FACS-sorted HSCs, MPPs and LMPPs[12](**Fig. 5d,e**). All DE peaks were properly

220     ordered through the 1D embedding and agreed with their accessibility dynamics in both

221     single-cell and bulk samples, suggesting that the co-accessible pattern between peaks could be

222     revealed by z-scores. (**Fig. 5d,e**). Moreover, our results also showed gradual gain or loss of

223     accessibility in a wide range of peaks in HSCs, revealing the continuous cell state transition within

224     HSCs. Although we lacked direct evidence to evaluate whether epiConv clustered HSCs and MPPs

225     into two groups, we could still extract HSC and MPP unique signatures through DE analysis. As

226     scaled heatmap could not reveal the fold change of peaks, we also examined the log2 Fold

227     Change between MPP, HSC and LMPP bulk samples for all detected DE peaks (**Fig. 5f**). Most peaks

228     showed strong difference between MPP/LMPP or HSC/LMPP bulk samples, while MPP or HSC

229     unique peaks just showed weak difference between MPP/HSC bulk samples. As is shown by many

230     single-cell studies, FACS-sorted cells may still be the mixtures of similar cell types. We thought

231     that this could partially explain the weak difference between MPP/HSC bulk samples.

232     Unexpectedly, we also found that LMPPs could be further divided into three groups based on

233     their unique signatures and bulk LMPPs seemed to be the mixture of these three groups. By

234     comparing the z-scores of single cells with bulk samples, we found that they might represent

235     different stages of LMPPs during differentiation (early undifferentiated stage, later stage to GMP

236     and later stage to CLP, see heatmap in the right in **Fig. 5e**). These results demonstrated that

237     inferring DE peaks directly from cell mixtures helped reveal the intra-population structure and

238    intermediate cell states in a signature-driven manner instead of statistical ways.

239        We also applied our algorithm to all cells in CD34+ dataset to test the scalability of our

240    algorithm. Similar with previous results on HSC-MPP-LMPP cluster, we also found a series of

241    peaks that gradually gained or lost accessibility through differentiation (e.g. in MDPs, **Fig. S5a**).

242    The z-scores did not fully capture the chromatin states of bulk samples for a few peaks (**Fig. S5b**).

243    We found that it could be explained by the difference between single-cell and bulk samples (data

244    not shown), probably because there might be some batch effects between them. In fact, it was

245    not difficult to infer DE peaks from distinct clusters. But we demonstrated that our algorithm

246    could also perform such task like conventional methods.

247        We found that sometimes the z-scores did not agree with the binary accessibility profiles of

248    single cells. It was because z-scores were normalized by the library size of single cells. However,

249    we thought that the library size of single cells in droplet-based protocols could reveal the

250    difference of global chromatin states between different cell types. By comparing the chromatin

251    states of neighbors with the background, our algorithm already removed the variation of library

252    size for individual cells. So, we designed another normalization strategy, where the scaling factors

253    of all cells were set to 1. We tried this normalization strategy in cells from replicate 1, where all

254    cells were processed in parallel during experiment and their library size could reflect the global

255    chromatin states. The z-scores were consistent with binary accessibility profiles under this

256    normalization strategy (**Fig. S5c,d**) but did not agree with corresponding bulk samples (**Fig. S5e**).

257

# **Discussions**
258

259        In this study, we developed a novel method to directly compare the Tn5 insertions between

260     single cells and compared it with three existing methods, cisTopic, LSI and SnapATAC. Results

261     demonstrated that our method had several advantages over existing methods. The most

262     significant difference between our algorithm and others is that we calculated the distance

263     between single cells using a convolution-based approach instead of commonly used

264     Euclidean-distance. Although the Jaccard similarity used by SnapATAC is similar to epiConv

265     (Assuming two binary vector A and B, Jaccard similarity is calculated by $\frac{|A \cap B|}{|A \cup B|}$, while epiConv uses

266     $\frac{|A \cap B|}{|A| \cdot |B|}$. Moreover, epiConv assigns weights to different loci), the distance is calculated by

267     Euclidean-distance on principal components. Interestingly, we also found a way to make epiConv

268     mimic the behavior of other methods, making it easy to compare the difference between two

269     forms of distance. Given the similarity matrix $S$ before denoising step, we used Eigen value

270     decomposition to obtain a series of latent features from $S$. Given $Q^T S Q = \Lambda$, where $Q$ is the matrix

271     containing Eigen vectors of $S$ and $\Lambda$ is the diagonal matrix containing Eigen values of $S$, the

272     columns of $Q \Lambda^{\frac{1}{2}}$ can be considered as latent features. Here we used top 50 latent features. By

273     calculating the Euclidean-distance on these latent features, the behavior of epiConv was highly

274     similar with existing methods (we showed the results of PBMC dataset and Mouse Cell Atlas

275     dataset in **Fig. S6**). In PBMC dataset, epiConv became sensitive to batch effects and the batch

276     effects could not be removed by our algorithm (compare **Fig. S6a** with **Fig. 3a,b** and **Fig. S1a,b**).

277     However, within each batch, major cell types could also be distinguished like other methods (**Fig.**

278     **S6b**). In Mouse Cell Atlas dataset, epiConv clustered "unknown" cells into single cluster like other

279     methods (compare **Fig. S6c** with **Fig. 4a-d**). These results clearly demonstrated that different

280     denoising process could have significant effects on our understanding of the cell heterogeneity

281     even when the raw data is identical. We hypothesized that methods trying to capture latent

282   features may suffer from common biases. By using convolution-based approach to define the

283   similarities between cells, epiConv provides a new angle of view in the analysis of sparse

284   epigenetic data. Moreover, epiConv also provides DE analysis in single-cell resolution and in

285   unbiased manner, while no existing methods could perform such task. Thus, we believe that

286   epiConv will have wide applications and improve our understanding on the epigenetic dynamics

287   of single cells.

288

## Methods

290   **Informative region calling for epiConv.** EpiConv takes processed fragments as input file. To call

291   informative regions for epiConv, we first extended Tn5 insertions from both directions using the

292   pileup command in MACS2[13] (-B --extsize 100). Then, we sorted all sites of the genome by their

293   density in decreasing order and selected regions with cumulative density less than 70% of total

294   insertions. These regions were extended from both directions by 100 bp and merged together if

295   having any overlap. Tn5 insertions overlapping with these informative regions (~70% of total

296   reads) were used for downstream analysis. We used such strategy instead of MACS2 because the

297   proportion of reads used in downstream analyses could be easily specified through the threshold

298   of cumulative density. Moreover, this strategy can always obtain some peaks, while MACS2 may

299   fail when the number of cells is low (e.g. < 200, reported by Satpathy et al. 2019[4]). The threshold

300   of cumulative density is determined by the distribution of insertion length. Based on our

301   preliminary analysis, fragments spanning one or more nucleosomes are nosier than fragments

302   from nucleosome-free regions. Thus, the threshold should be close to the proportion of

303   fragments from nucleosome-free regions. For the myoblast and mouse brain datasets, we set the

304    threshold to 50% as they had higher proportion of fragments spanning one or more nucleosomes

305    (data not shown). The major purpose of informative region calling is to calculate the weights for

306    different genomic regions (see below). Additionally, it could remove some background noise.

307    Although it is possible to compare the Tn5 insertions of the whole genome, which might help

308    detect rare cell types, we find that it just increases running time but does not improve the

309    results.

310    **epiConv algorithm.** In the results section, we described the algorithm to calculate the similarity

311    between two cells over one region. Here assume that we have N cells and K regions, with the

312    similarities between any two cells $i$ and $j$ over region $k$ ($s_{ijk}$) being known. First, we weight each

313    region as follows:

$$freq_k = \sqrt{\frac{2}{N(N-1)} \sum\nolimits_{ij} s_{ijk}}$$

$$w_k = log10(1 + freq_k^{-1})$$

314    The form of weight is similar to that used in LSI but the frequency is replaced by a

315    pseudo-frequency estimated from our convolution-based approach. We use such form of weight

316    to increase the contribution of low-density regions to the similarity score. The similarity between

317    cell $i$ and $j$ is calculated using a bootstrap approach. Assuming we perform L replicates (L = 30 in

318    this study) and in each replicate we randomly sample some regions (12.5% of total informative

319    regions in this study). The similarity of $s_{ij}$ is calculated as follows:

$$s_{ij} = \frac{\sum_l log10\left(\sum_{k \in rep_l} s_{ijk} \cdot w_k^2\right)}{L} - log10(lib_i \cdot lib_j)$$

320    where $lib_i$ and $lib_j$ is the library size of cell $i$ and $j$. We normalize the aggregated similarity by

321    $lib_i \cdot lib_j$ because $\sum_{k \in rep_l} s_{ijk} \cdot w_k^2$ can be considered as the sum of $lib_i \cdot lib_j$ random

322    variables with identical distribution given the analytical form of similarity described above.

323    Averaging the similarities from replicates helps reduce the noise compared to simple aggregation

324    of similarities from all regions. But for deep sequencing data, we find that simple aggregation

325    also generates similar results (data not shown).

326        In the simplified version, matrix is first binarized and TF-IDF transformed like LSI[3] (In

327    epiConv-simp, normalization with respect to sequencing depth and peak weighting are identical

328    as LSI). Given TF-IDF matrix $M$ and L bootstrap matrices $M_{rep_l}$ by randomly sampling peaks from

329    $M$, the similarity matrix $S$ can be calculated as follows:

$$S = \frac{\sum_l log10\left(M_{rep_l}{}^T \cdot M_{rep_l}\right)}{L}$$

330    where $M_{rep_l}{}^T \cdot M_{rep_l}$ is the matrix product. Unlike LSI implemented in Cusanovich et al. 2015[3]

331    and Cusanovich et al. 2018[8], we do not filter any peaks. By adopting the formula above, the

332    distance between two insertions $\mu_{Ai} - \mu_{Bj}$ is considered as zero if they are in the same peak or

333    infinite otherwise. Further steps are identical for full and simplified versions.

334        Next, we denoise the similarities between cells by borrowing the information from their

335    neighbors. The denoised similarities are calculated by the number of shared nearest neighbors

336    between two cells. The number of nearest neighbors for each cell is set to 50 in this study. If the

337    dataset contains cells from multiple batches, we force cells to select equal number of nearest

338    neighbors from each batch to remove batch effects. The distance matrix D is calculated by

339    $D = -S_{denoise}$. Although the batch removal strategy can be applied to the similarity or distance

340    matrix generated by various methods, we find that it only works well with epiConv. As mentioned

341    above, it is because that epiConv is less sensitive to batch effects even without any correction.

342        The denoising method above changes the unit of similarity matrix (from continuous values

343    to integer values). Occasionally we find that it may make the results worse (see the results of

344     epiConv-simp for cell lines data in **Fig. 2b**). We also developed an alternative denoising method

345     that keeps the unit of similarity matrix unchanged (**Supplementary Note 1**). Generally, it is noisier

346     than the first method and cannot remove batch effects. But it may perform better when the first

347     method fails (see top-right in **Fig. 2b**).

348     **Pre-processing of ATAC-seq data.** We took the processed fragment file or peak by cell matrix as

349     inputs if available. For the unprocessed data from Buenrostro et al. 2015[6] and bulk samples from

350     Corces et al. 2016[12], we aligned raw reads to the hg19 genome using Bowtie2[14] (-X 2000

351     --no-mixed --no-discordant) and removed reads with mapping quality <10 and duplicates using

352     Picard tools. The start and end of the fragments were adjusted (+5 for forward strand and −4 for

353     reverse strand). We called peaks using MACS2[13] (--nomodel --nolambda --keep-dup all --shift -200

354     --extsize 400) and generated the count matrix by counting the number of Tn5 insertions falling in

355     peaks.

356     For the mouse brain dataset, we randomly sampled 2,000 cells from Channel 1 and Channel

357     2 in Lareau et al. 2019 (dscATAC-seq)[5], 1,000 cells from the mouse cortex data from 10x

358     Genomics and 2,000 cells from two replicates of whole mouse brain in Cusanovich et al. 2018

359     (sciATAC-seq)[8]. The dataset contains 5,000 cells in total. Data from Cusanovich et al. 2018 were

360     converted from mm9 to mm10 using liftOver[15]. Data from 10x Genomics and Cusanovich et al.

361     2018 were re-counted against the peaks called by Lareau et al. 2019 for data integration.

362     For the myoblast dataset, few outlier cells that did not cluster together with the majority of

363     cells were excluded in differential analysis (**Fig. 5a**).

364     **Implement of cisTopic, LSI and SnapATAC**. In cisTopic, the number of topics is set to 20, 30, 40

365     and 50 and automatically decided by cisTopic. For the analysis of cell lines data from Buenrostro

366   et al. 2015[6], in order to explore whether increased number of topics could provide higher

367   resolution for K562 cells, we increase the number of topics from 20 to 100 with a step of 10 but

368   the optimal number of topics is still decided by cisTopic. In LSI, we use the scripts from

369   Cusanovich et al. 2018[8], filter out peaks with frequency < 0.01 and use the top 50 components of

370   singular value decomposition for dimension reduction. In SnapATAC, the bin size was set to 5000.

371   We fixed the number of principal components used for dimension reduction to 30 instead of

372   manually examining the distribution of each component to avoid ambiguity.

373   **Differential analysis algorithm.** The input data is a binarized peak by cell matrix and a similarity

374   matrix between cells. Here we use the peak by cell matrix from previous steps. The similarity

375   matrix is calculated by $S_{denoise} + S/100$ (The similarity matrix is mainly determined by $S_{denoise}$.

376   When two cells have equal number of common neighbors to another cell, the similarities are

377   further determined by the original similarity matrix). For each single cell, we define $k$ cells with

378   highest similarities as its neighbors (including itself). Then for each peak, we test whether it is

379   more likely to be accessible in the cell's neighbors. This problem can be resolved using

380   hypergeometric test, with cells accessible as black balls, cells inaccessible as white balls. The

381   sampling times ($\hat{k}$, the adjusted number of neighbors) is calculated by the total scaling factors of

382   all neighbors divided by the average scaling factors of all cells. The scaling factors of cells were

383   equal to their library sizes or set to 1 for all cells (DE analysis on replicate 1 of CD34+ dataset, see

384   Results). The z-scores are calculated by the number of cells accessible among neighbors and

385   z-normalized by corresponding mean and variance of the null distribution.

386       In differential analyses in this study, the number of neighbors $k$ is set to 5% of total cells. The

387   number of neighbors $k$ defines the size of potential clusters, which serves similar function as the

388 number of clusters in conventional pipeline. However, the results demonstrated that our

389 algorithm with fixed $k$ could still detect DE peaks in clusters with a wide range of size. Here, $k$ is

390 set to 5% in order to make our algorithm more sensitive to DE peaks of small clusters. After

391 obtaining the z-scores, we select peaks with z-score > 2 in at least 10% cells as DE peaks. When

392 we applied our algorithm to all CD34+ cells (whole dataset or replicate 1), we select peaks with

393 z-score > 2 in at least 30% cells as we only want to detect DE peaks between major clusters and

394 the criterion of 10% cells suggested most peaks to be differentially accessible, which was

395 reasonable but not desired.

396  In fact, it is not straightforward to choose a proper threshold for z-score. We find that peaks

397 that do not satisfy the threshold described above may also show weak DE pattern. Here, we use

398 the threshold of 10% cells with z-score >2 because selected peaks can be easily validated by bulk

399 samples. For general purpose, users can set the threshold manually to obtain appropriate

400 number of DE peaks.

401 **Dimension reduction.** We perform dimension reduction of single cells using the uniform

402 manifold projection (UMAP) algorithm[16] by feeding umap with the distance matrix learned by

403 epiConv, cisTopic, LSI and SnapATAC using default settings. The number of reduced components

404 was set to 1 for heatmaps and 2 for scatterplot of cells. We also embed DE peaks into 1D space

405 by feeding umap with the distance matrix that is calculated by one minus spearman correlation

406 of z-scores between peaks.

407 **Bulk sample processing.** For bulk samples of hematopoietic cells from Corces et al. 2016[12], we

408 count the Tn5 insertions against the peaks called from Satpathy et al. 2019[5], normalize the counts

409 by library size and average the normalized counts across all replicates for each cell type. For the

410     myoblast dataset, we de-multiplex the reads, count the Tn5 insertions and normalize the counts

411     by harvesting times.

412     **Data availability.** The cell lines data of Buenrostro et al. 2015[6] is obtained from Gene Expression

413     Omnibus (GEO) accession GSE65360. The data of Satpathy et al. 2019[4] is obtained from GEO

414     accession GSE129785. The data of Lareau et al. 2019[5] is obtained from GEO accession GSE123581.

415     The data of Cusanovich et al. 2018[8] is obtained from Mouse Cell Atlas

416     (http://atlas.gs.washington.edu/mouse-atac/). The data of adult mouse cortex is obtained from

417     10X Genomics website

418     (https://support.10xgenomics.com/single-cell-atac/datasets/1.2.0/atac_v1_adult_brain_fresh_5k

419     ). Myoblasts data[9] is obtained from GEO accession GSE109828. EpiConv is available at Github

420     (https://github.com/LiLin-biosoft/epiConv).

421

422     **Acknowledgements**

423     This project was funded by the National Key Research and Development Program of China

424     (2018YFC1004602), National Natural Science Foundation of China (NSF 31871332) and a startup

425     fund to L.Z. from ShanghaiTech University. We would like to thank Xiaojing Zhao for testing the

426     reproducibility of the study. We would like to thank Yingdong Zhang on his technical support on

427     the HPC platform of ShanghaiTech University.

428

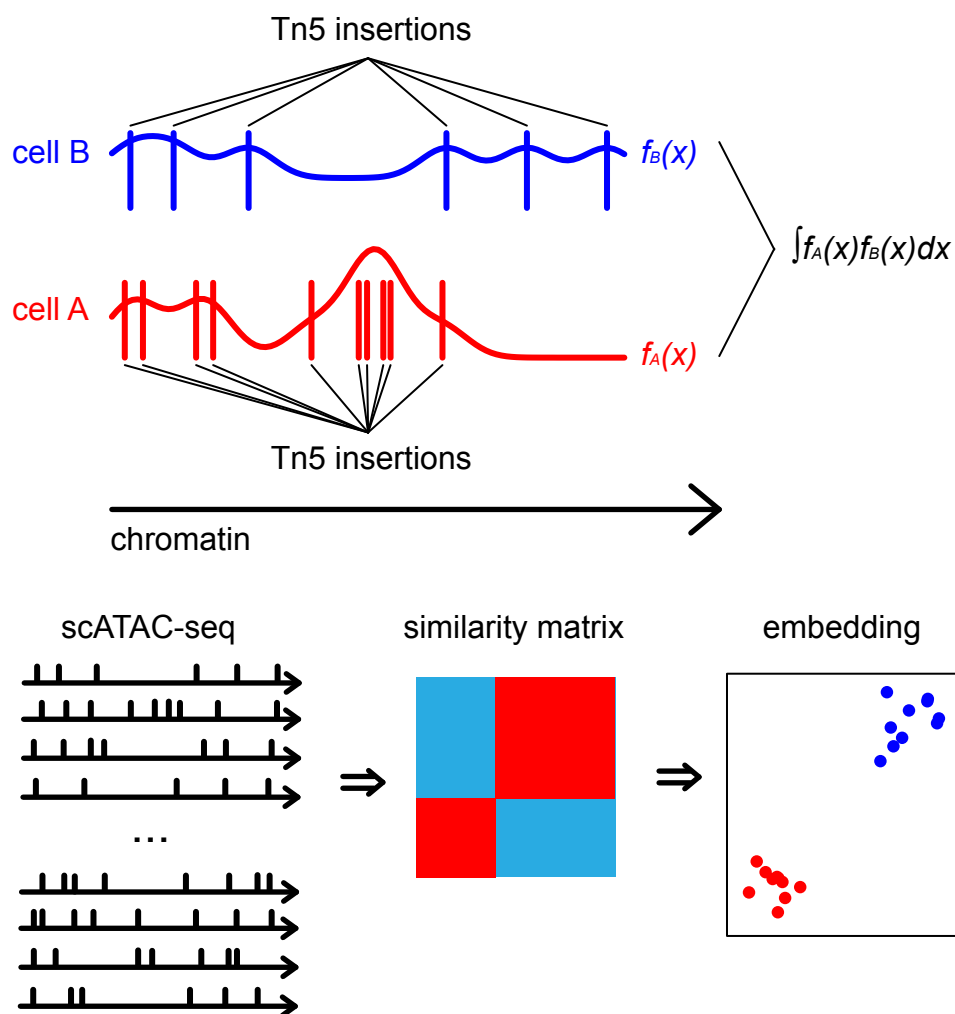429     **Author contributions**

430     L.L. conceived the study, developed the methods and performed the analyses. L.L. and L.Z. wrote

431     the manuscript. L.Z supervised the study.
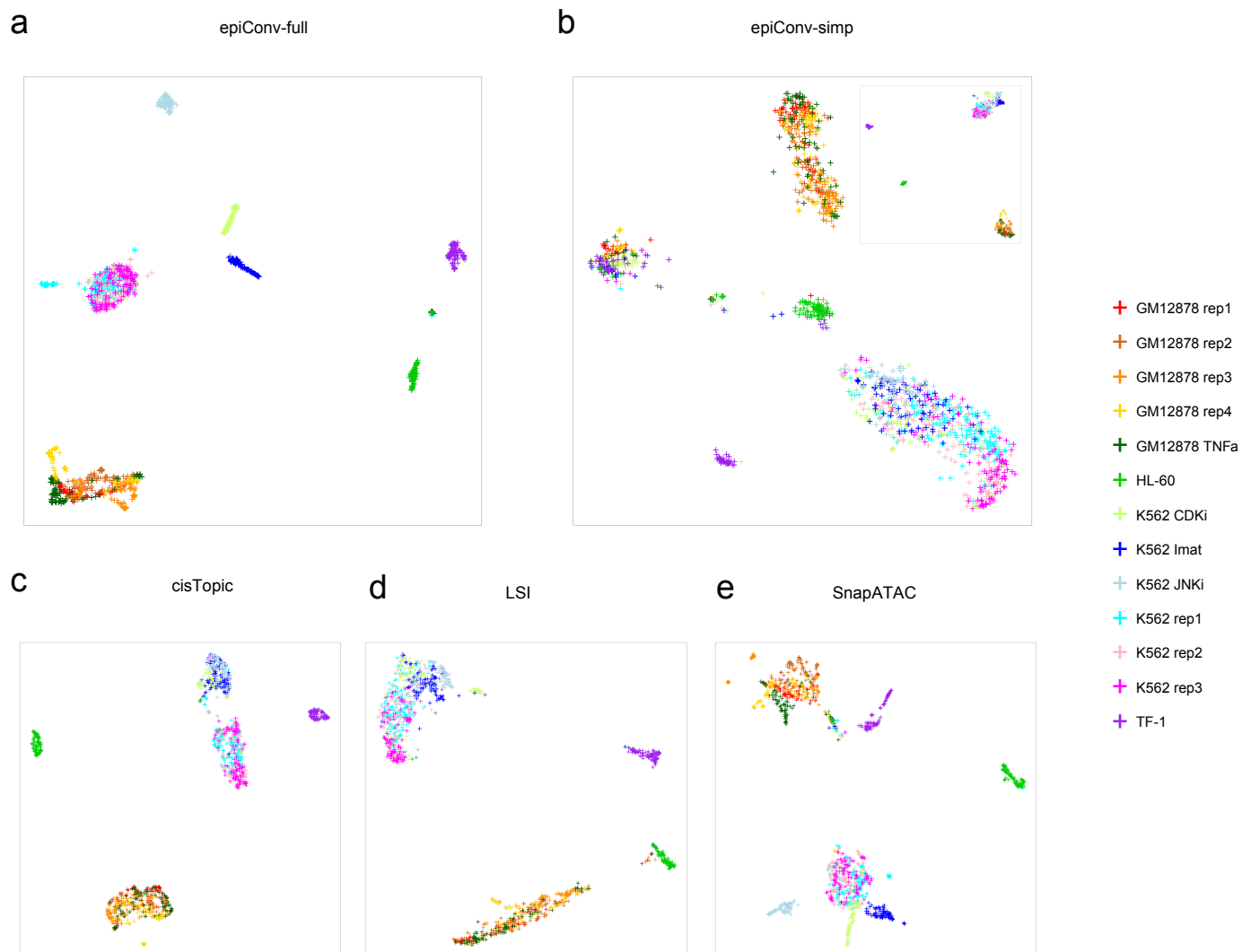
432  **Competing interests**

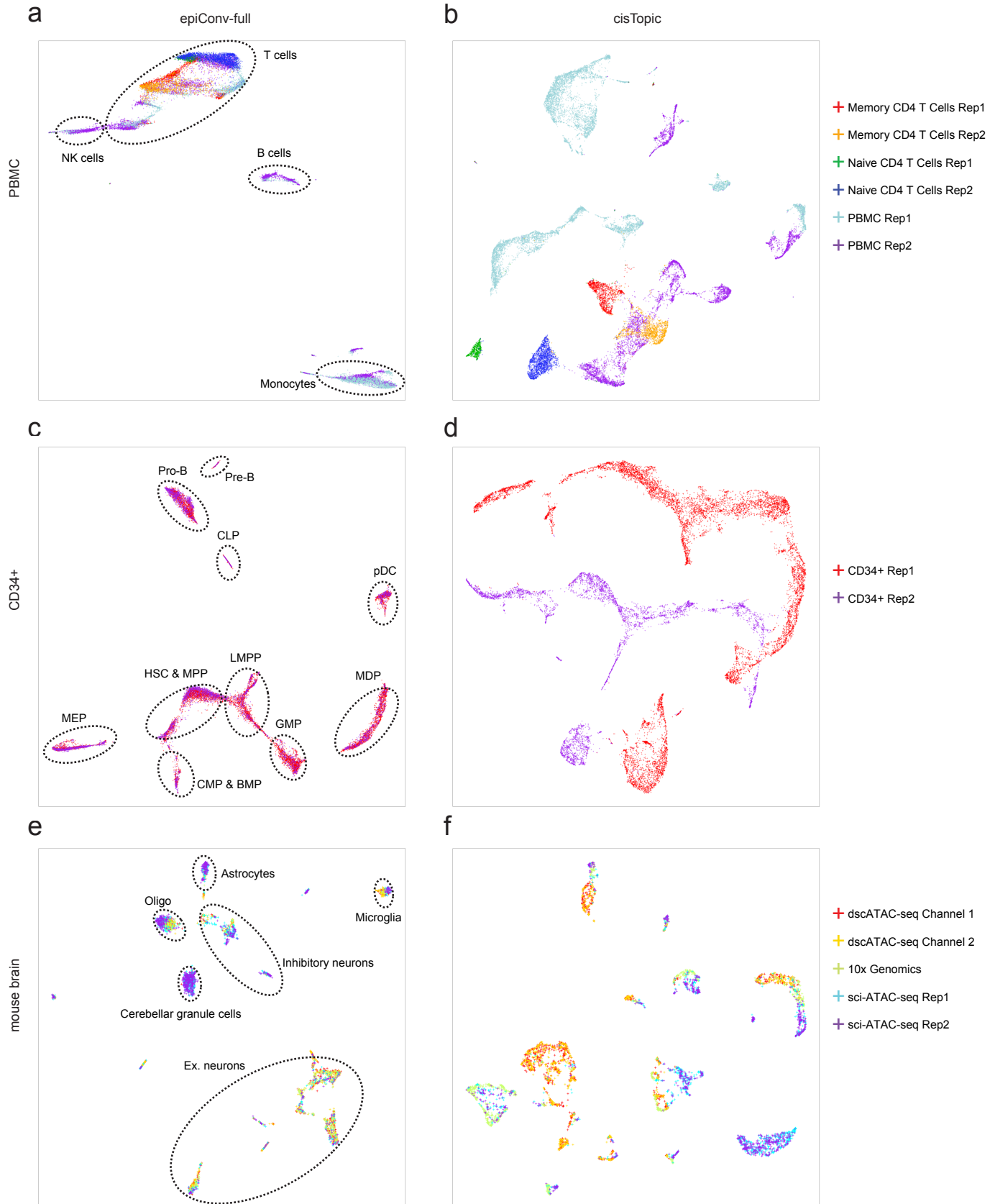433  The authors declare no competing interests.


434

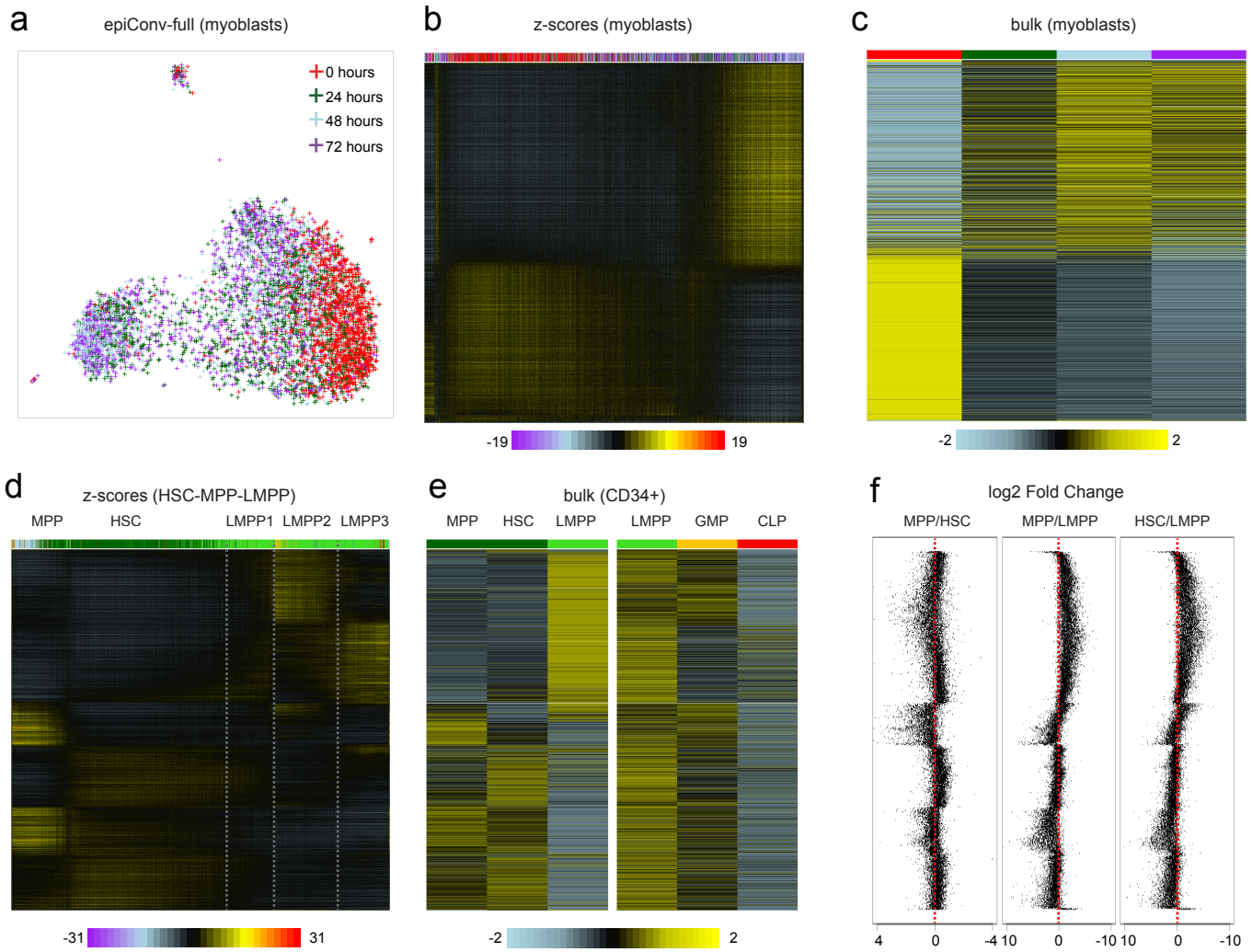435     **Figure 1.** An overview of the epiConv algorithm.

436

a      epiConv-full

b      epiConv-simp

c      cisTopic

d      LSI

e      SnapATAC



+ GM12878 rep1
+ GM12878 rep2
+ GM12878 rep3
+ GM12878 rep4
+ GM12878 TNFa
+ HL-60
+ K562 CDKi
+ K562 Imat
+ K562 JNKi
+ K562 rep1
+ K562 rep2
+ K562 rep3
+ TF-1

437     **Figure 2.** EpiConv performs better than other methods on cell lines data. (**a**) Embedding by

438     epiConv-full. (**b**) Embedding by epiConv-simp. (**c**) Embedding by cisTopic. (**d**) Embedding by LSI (**e**)

439     Embedding by snapATAC.

440

441 **Figure 3.** EpiConv removes batch effects. (**a,b**) Embeddings of PBMC dataset by epiConv-full and

442 cisTopic. (**c,d**) Embeddings of CD34+ dataset by epiConv-full and cisTopic. (**e,f**) Embeddings of the

443 integration of mouse brain data from dscATAC-seq, 10x Genomics and sci-ATAC-seq by

444 epiConv-full and cisTopic. Major cell types in the embeddings of epiConv are circled. Embeddings

445 by LSI and SnapATAC can be found in **Fig. S1**. HSC, hematopoietic stem cells; MPP, multipotent

446 progenitors; LMPP, lymphoid-primed multipotent progenitors; CMP, common myeloid progenitors;

447 BMP, basophil-mast cell progenitors; GMP, granulocyte-macrophage progenitors; MDP,

448 monocyte-dendritic cell progenitors; pDC, plasmacytoid dendritic cells; MEP,

449 megakaryocyte-erythroid progenitors; CLP, common lymphoid progenitors; Oligo,

450 oligodendrocytes; Ex. neurons, excitatory neurons.

451

452     **Figure 4.** EpiConv reveals the identities of unknown cells in Mouse Cell Atlas dataset. (**a**)

453     Embedding by cisTopic. (**b**) Embedding by LSI. (**c**) Embedding by SnapATAC. (**d**) Embedding by

454     epiConv-simp. In (**a-d**), unknown cells and cells showing close relationships with them are colored

455     according to the annotations from Cusanovich et al. 2018. Other irrelevant cells are colored in

456     grey. Six major clusters in (**d**) that contain high proportion of unknown cells are circled. (**e**)

457     Spearman correlations between aggregated samples with known and unknown identities from 6

458     major clusters marked in (**d**). Labels of unknown samples are colored in red. Numbers in the

459     diagonal elements show the correlations between unknown samples and corresponding known

460     samples. Endo I, endothelial I cells; Endo II, endothelial II cells; Ex. neurons, excitatory neurons;

461     HSPC, hematopoietic progenitors; Oligo, oligodendrocytes.

462

a  epiConv-full (myoblasts)

+ 0 hours
+ 24 hours
+ 48 hours
+ 72 hours

b  z-scores (myoblasts)

-19          19

c  bulk (myoblasts)

-2          2

d  z-scores (HSC-MPP-LMPP)

MPP    HSC    LMPP1  LMPP2  LMPP3

-31          31

e  bulk (CD34+)

MPP  HSC  LMPP    LMPP  GMP  CLP

-2          2

f  log2 Fold Change

MPP/HSC    MPP/LMPP    HSC/LMPP

4    0    -4  10    0    -10  10    0    -10

463    **Figure 5.** EpiConv detects differentially accessible peaks in cell mixtures. (**a**) Embedding of

464    myoblast single cells by epiConv-full. (**b**) Accessibility z-scores of myoblast single cells inferred by

465    epiConv. (**c**) Accessibility profiles of aggregated myoblast samples by harvesting times. Cells or

466    aggregated samples in (**b,c**) are colored by harvesting times according to (**a**). (**d**) Accessibility

467    z-scores of HSC-MPP-LMPP single cells inferred by epiConv. (**e**) Accessibility profiles of HSC, MPP

468    and LMPP bulk samples and LMPP, GMP and CLP bulk samples. (**f**) Log2 Fold Change of peaks

469    between HSC, MPP and LMPP bulk samples. Points represent the corresponding peaks in (**d,e**)
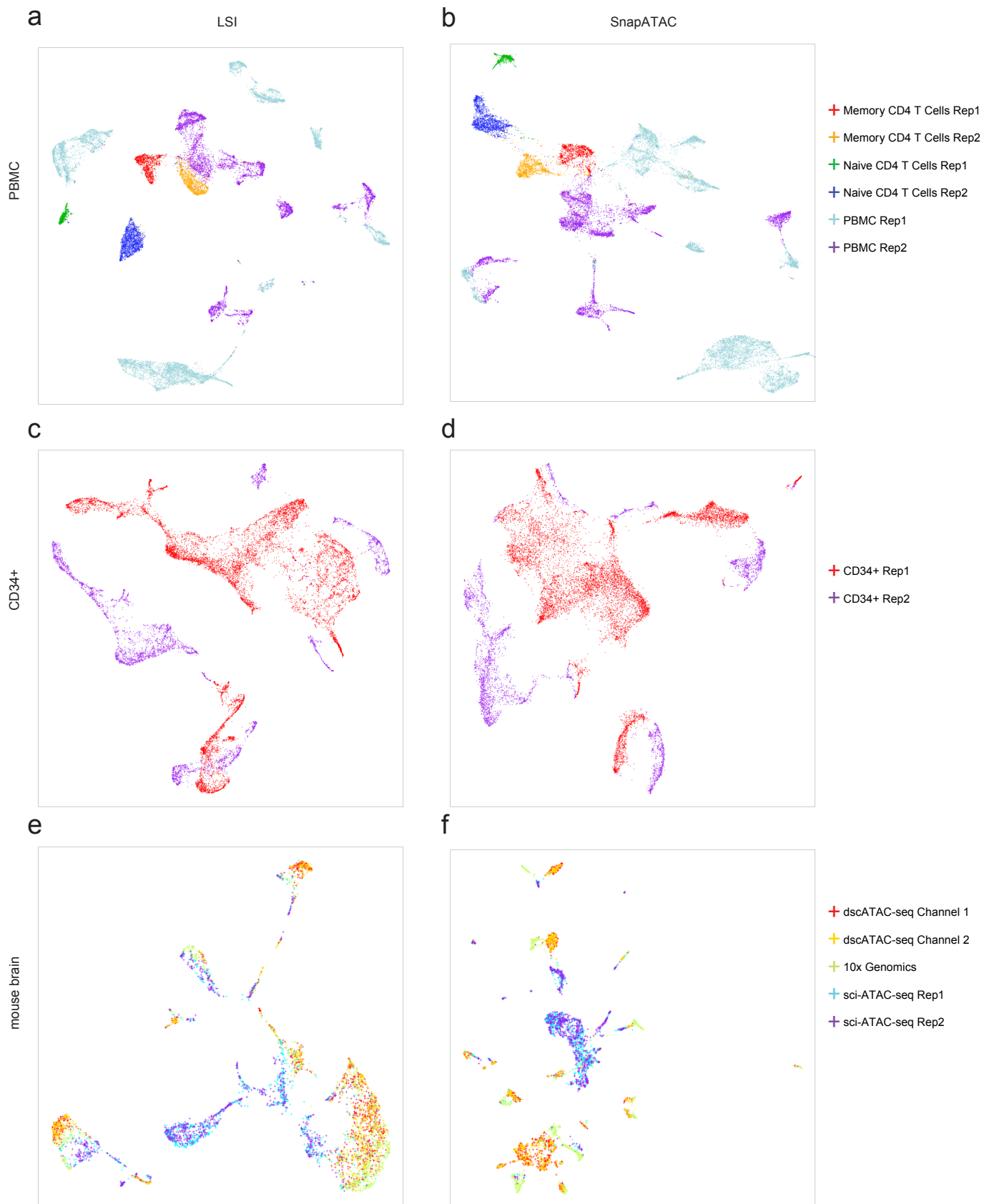
470

# Reference:

1.  Buenrostro, J.D., Giresi, P.G., Zaba, L.C., Chang, H.Y. & Greenleaf, W.J. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nat Methods* **10**, 1213-1218 (2013).

2.  Klemm, S.L., Shipony, Z. & Greenleaf, W.J. Chromatin accessibility and the regulatory epigenome. *Nat Rev Genet* **20**, 207-220 (2019).

3.  Cusanovich, D.A. et al. Multiplex single cell profiling of chromatin accessibility by combinatorial cellular indexing. *Science* **348**, 910-914 (2015).

4.  Satpathy, A.T. et al. Massively parallel single-cell chromatin landscapes of human immune cell development and intratumoral T cell exhaustion. *Nat Biotechnol* **37**, 925-936 (2019).

5.  Lareau, C.A. et al. Droplet-based combinatorial indexing for massive-scale single-cell chromatin accessibility. *Nat Biotechnol* **37**, 916-924 (2019).

6.  Buenrostro, J.D. et al. Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature* **523**, 486-490 (2015).

7.  Chen, H. et al. Assessment of computational methods for the analysis of single-cell ATAC-seq data. *Genome Biol* **20**, 241 (2019).

8.  Cusanovich, D.A. et al. A Single-Cell Atlas of In Vivo Mammalian Chromatin Accessibility. *Cell* **174**, 1309-1324 e1318 (2018).

9.  Pliner, H.A. et al. Cicero Predicts cis-Regulatory DNA Interactions from Single-Cell Chromatin Accessibility Data. *Mol Cell* **71**, 858-871 e858 (2018).

10. Schep, A.N., Wu, B., Buenrostro, J.D. & Greenleaf, W.J. chromVAR: inferring transcription-factor-associated accessibility from single-cell epigenomic data. *Nat Methods* **14**, 975-978 (2017).

11. Bravo Gonzalez-Blas, C. et al. cisTopic: cis-regulatory topic modeling on single-cell ATAC-seq data. *Nat Methods* **16**, 397-400 (2019).

12. Corces, M.R. et al. Lineage-specific and single-cell chromatin accessibility charts human hematopoiesis and leukemia evolution. *Nat Genet* **48**, 1193-1203 (2016).

13. Zhang, Y. et al. Model-based analysis of ChIP-Seq (MACS). *Genome Biol* **9**, R137 (2008).

14. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**, 357-359 (2012).

15. Kent, W.J. et al. The human genome browser at UCSC. *Genome Res* **12**, 996-1006 (2002).

16. McInnes, L., Healy, J., Saul, N. & Großberger, L. UMAP: Uniform Manifold Approximation and Projection. *Journal of Open Source Software* **3** (2018).
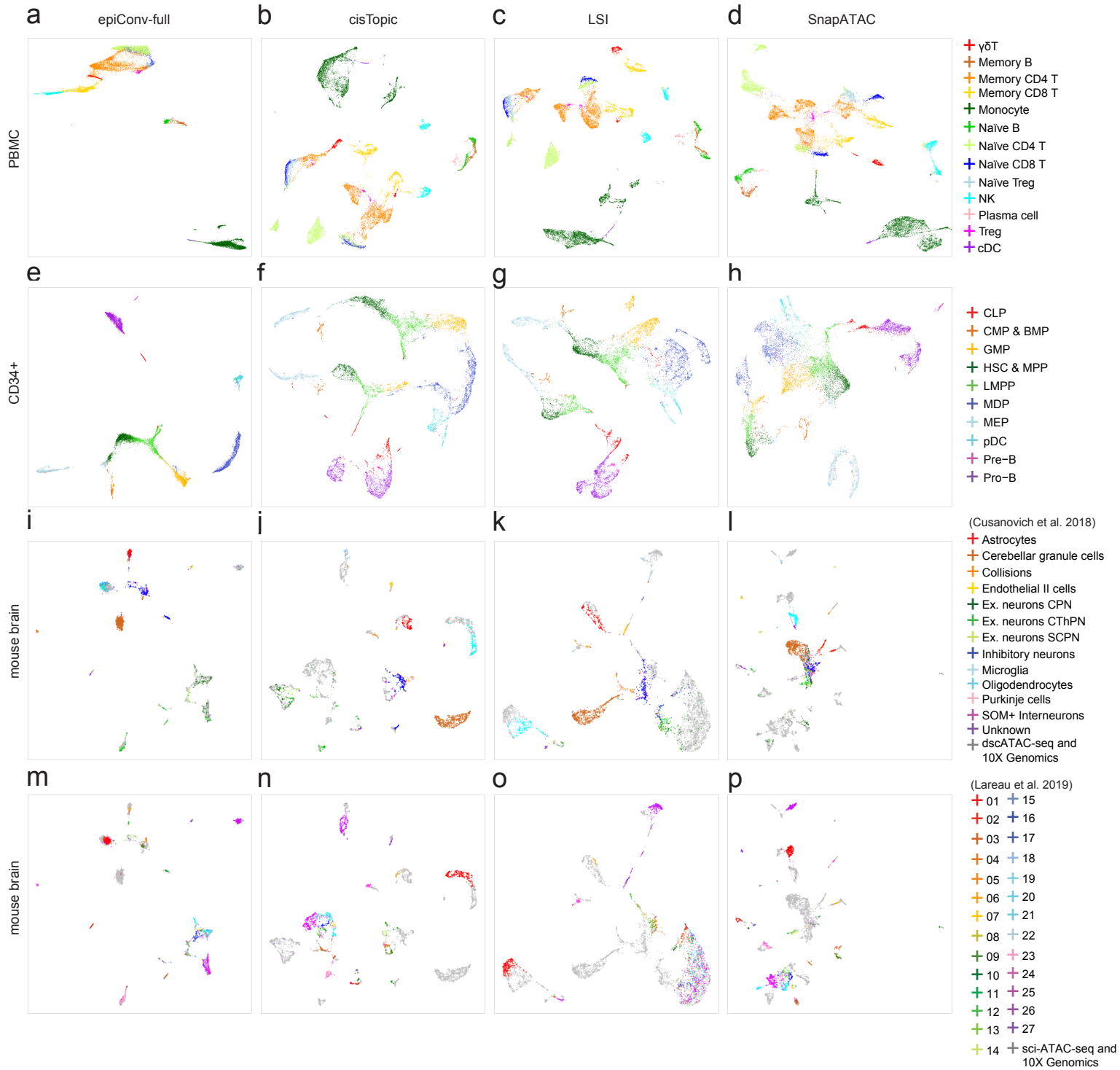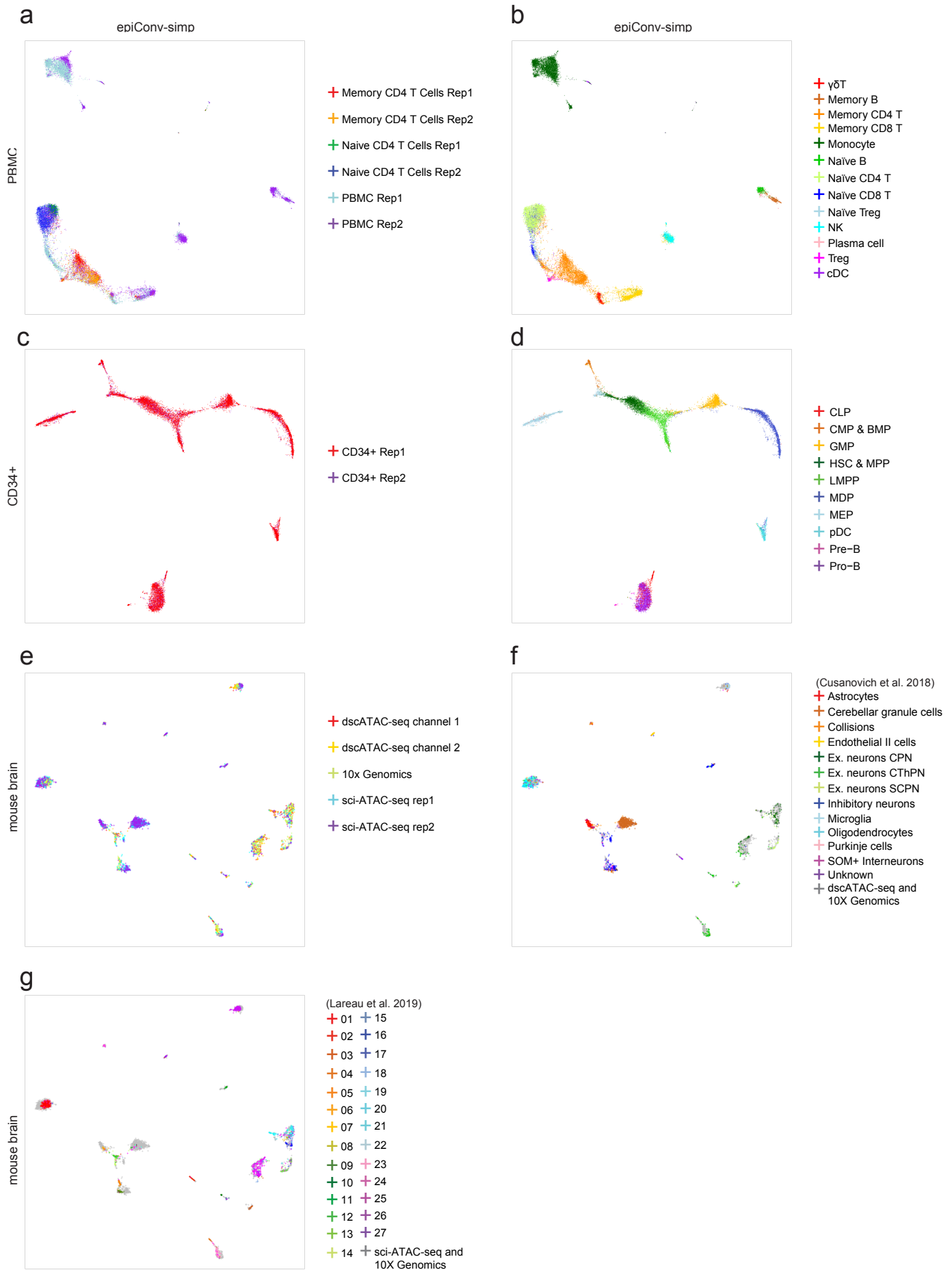
505 **Supplementary materials**
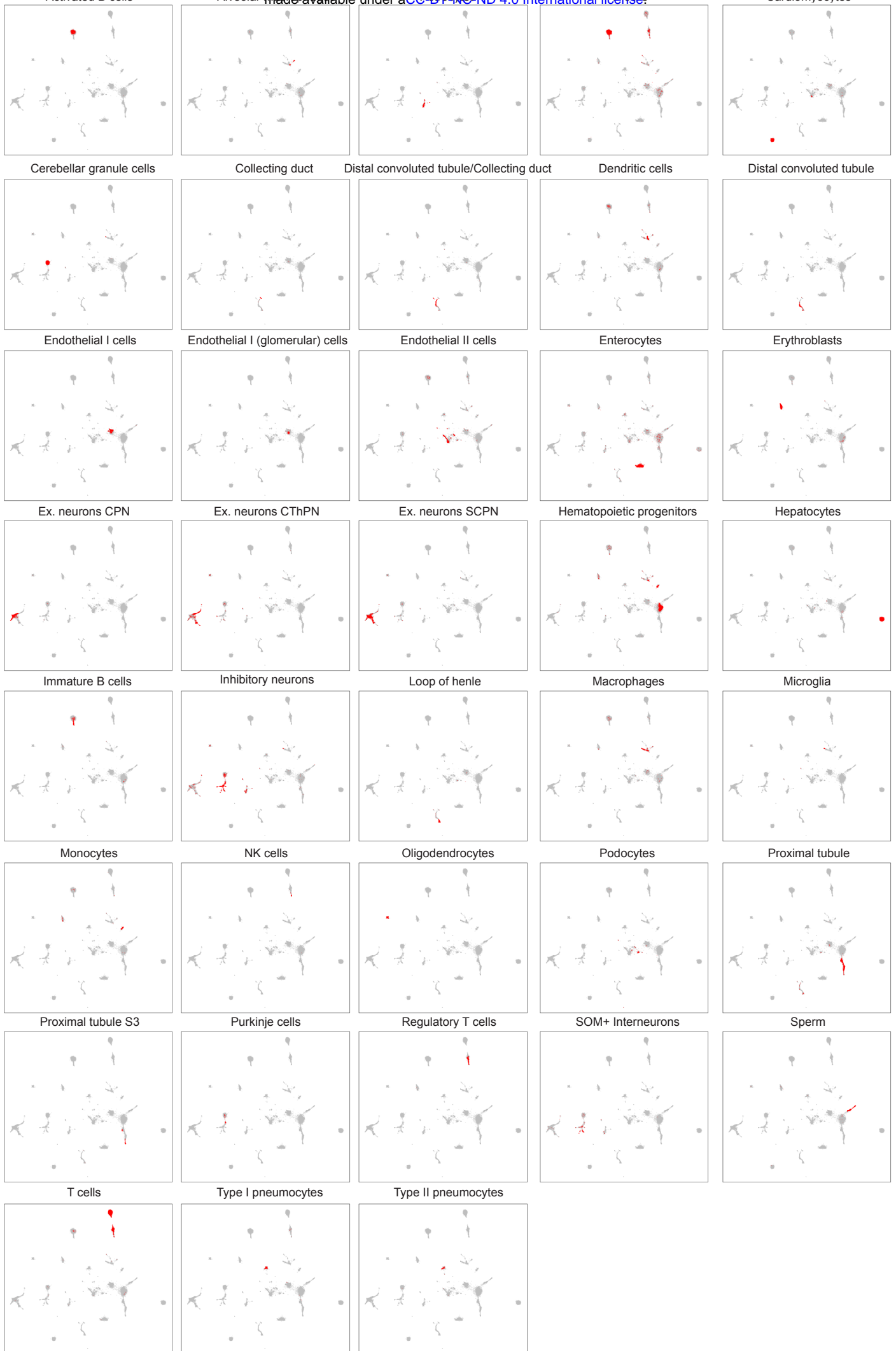
506

507     **Figure S1.** Embeddings of PBMC (**a,b**), CD34+ (**c,d**) and mouse brain (**e,f**) datasets by LSI and

508     SnapATAC .

509

510    **Figure S2.** Comparison of embeddings by epiConv-full, cisTopic, LSI and SnapATAC with cell

511    annotations from original articles. (**a-d**) Embeddings of PBMC dataset by epiConv-full, cisTopic,

512    LSI and SnapATAC, colored by annotations from Satpathy et al. 2019. (**e-h**) Embeddings of CD34+

513    dataset by epiConv-full, cisTopic, LSI and SnapATAC, colored by annotations from Satpathy et al.

514    2019. (**i-p**) Embeddings of mouse brain dataset by epiConv-full, cisTopic, LSI and SnapATAC,

515    colored by annotations from Cusanovich et al. 2018 (**i-l**) and Lareau et al. 2019 (**m-p**).
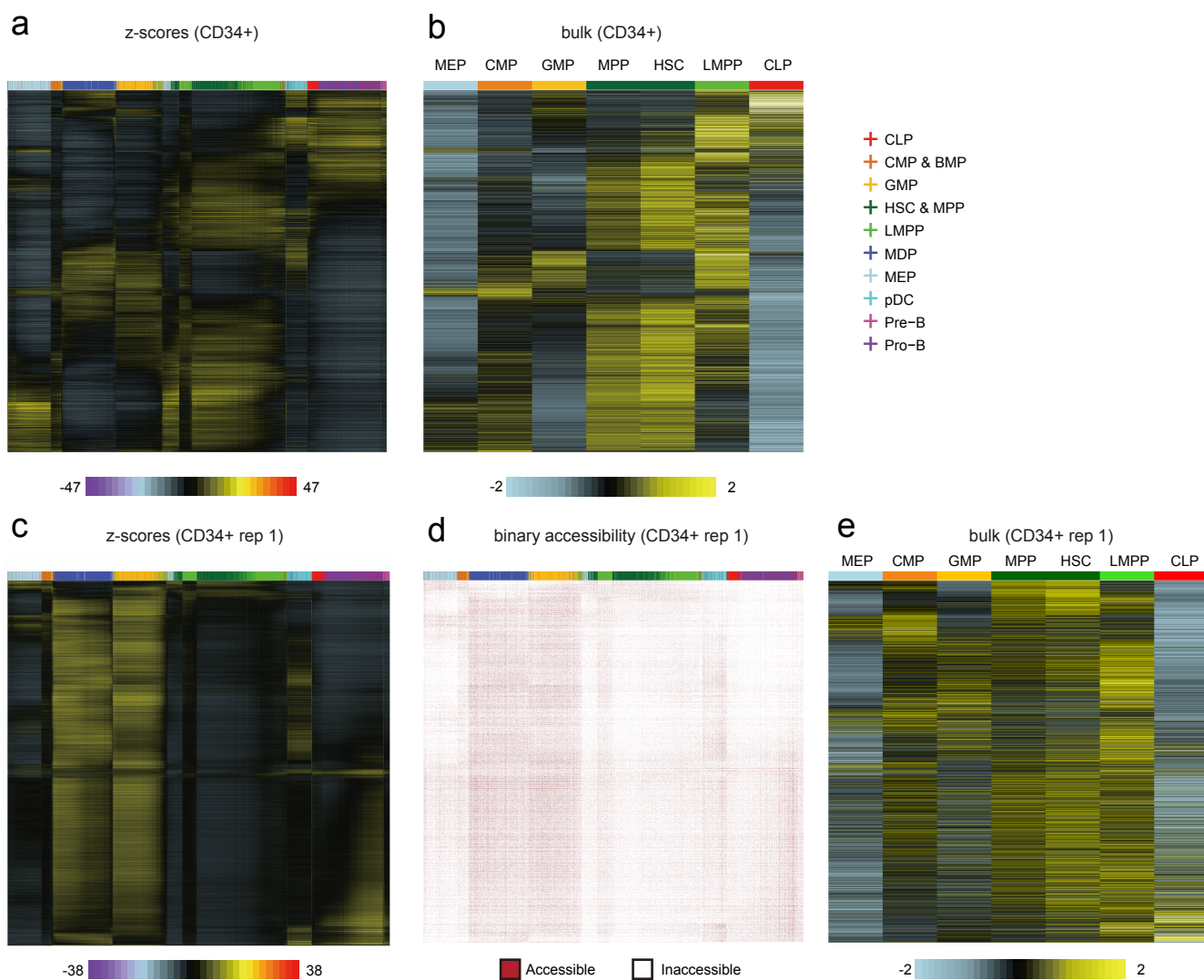
516

517     **Figure S3.** Embeddings of PBMC, CD34+ and mouse brain datasets by epiConv-simp. (**a,b**)

518     Embeddings of PBMC dataset by epiConv-simp, colored by batch (**a**) and annotations from

519     Satpathy et al. 2019 (**b**). (**c,d**) Embeddings of CD34+ dataset by epiConv-simp, colored by batch (**c**)

520     and annotations from Satpathy et al. 2019 (**d**). (**e-g**) Embeddings of mouse brain dataset by

521     epiConv-simp, colored by batch (**e**), annotations from Cusanovich et al. 2018 (**f**) and annotations

522     from Lareau et al. 2019 (**g**).

523

Activated B cells | Alveolar macrophages | Astrocytes | B cells | Cardiomyocytes
Cerebellar granule cells | Collecting duct | Distal convoluted tubule/Collecting duct | Dendritic cells | Distal convoluted tubule
Endothelial I cells | Endothelial I (glomerular) cells | Endothelial II cells | Enterocytes | Erythroblasts
Ex. neurons CPN | Ex. neurons CThPN | Ex. neurons SCPN | Hematopoietic progenitors | Hepatocytes
Immature B cells | Inhibitory neurons | Loop of henle | Macrophages | Microglia
Monocytes | NK cells | Oligodendrocytes | Podocytes | Proximal tubule
Proximal tubule S3 | Purkinje cells | Regulatory T cells | SOM+ Interneurons | Sperm
T cells | Type I pneumocytes | Type II pneumocytes
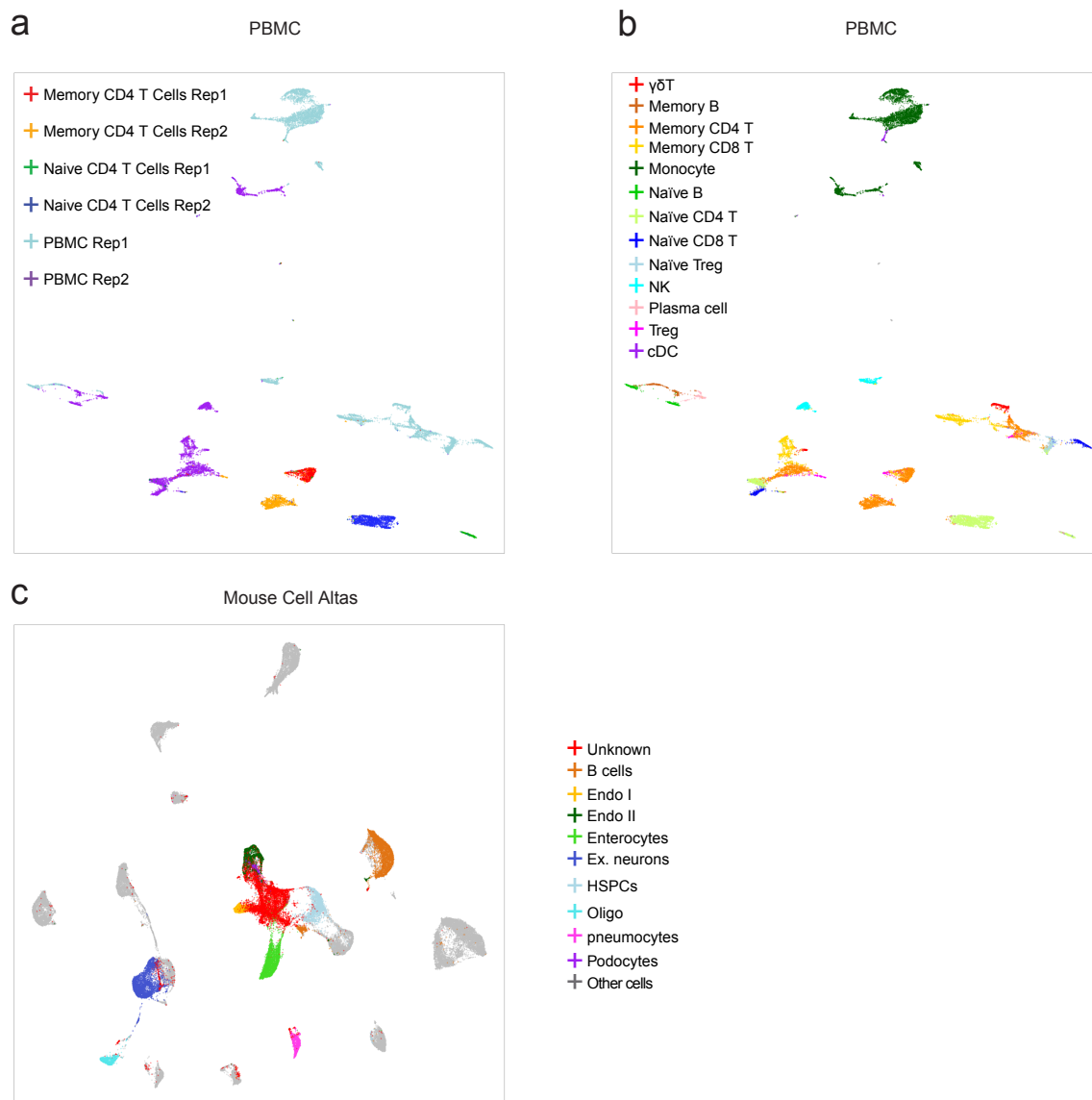
524    **Figure S4.** Embeddings of Mouse Cell Atlas dataset by epiConv-simp. The corresponding cell types

525    are colored in red and other cells are colored in grey.

526

527      **Figure S5.** Accessibility z-scores of CD34+ dataset. (**a**) Accessibility z-scores of CD34+ single cells

528      (**b**) Accessibility profiles of bulk CD34+ samples. (**c**) Accessibility z-scores of CD34+ single cells in

529      replicate 1. (**d**) Binary accessibility profiles of CD34+ single cells in replicate 1. (**e**) Accessibility

530      profiles of bulk CD34+ samples. DE peaks in (**a,b**) are selected from whole CD34+ dataset by

531      setting the scaling factors of cells to their library sizes and DE peaks in (**c-e**) are selected from

532      replicate 1 in CD34+ dataset by setting the scaling factors of cells to 1.

533

534    **Figure S6.** Embeddings of PBMC and Mouse Cell Atlas dataset by inferring latent features from

535    the similarity matrix of epiConv. (**a**) Embedding of PBMC dataset, colored by batch. (**b**)

536    Embedding of PBMC dataset, colored by annotations from Satpathy et al. 2019. (**c**) Embedding of

537    Mouse Cell Atlas dataset, colored by annotations from Cusanovich et al. 2018.

538

539   **Supplementary Note 1**

540       Here, we describe an alternative denoising method that keeps the unit of similarity matrix

541   unchanged. Given N cells and their similarity matrix S where $s_{ij}$ is the similarity between cell $i$ and

542   $j$, we first transform S to a weight matrix W as follows:

$$w_{ij} = \begin{cases} 10^{s_{ij}} \cdot \log 10(lib_i), & i \in j's \ neighbors \\ 0, & i \notin j's \ neighbors \end{cases}$$

543   Where $j$'s neighbors are the top 20 cells with highest similarities to $j$. For each column $j$, we scale

544   the sum of column (excluding the diagonal elements) to a fraction parameter $\theta$ between 0 and 1

545   and the diagonal elements of W are set to $1 - \theta$. Then the sum of each column is equal to 1. The

546   matrix W defines how to mix the information from the cell itself and its neighbors, where $\theta$

547   proportion of information comes from its neighbors and the weight of each neighbor is

548   determined by its similarity to cell $j$ multiplied by its log10 library size, and $1 - \theta$ proportion of

549   information comes from cell $j$ itself. In this study, we set $\theta$ to 0.25. We create a similarity matrix S'

550   where its elements are equal to S except for the diagonal elements (the similarity of each cell to

551   itself, which is not defined for S). The diagonal element $s'_{jj}$ is set to the 99th percentile of column

552   $j$, which can be used to approximate the similarity of cell $j$ to itself. The denoised similarity matrix

553   $S_{denoise}$ is calculated by matrix product of S' and W as follows:

$$S_{denoise} = \frac{S' \cdot W + (S' \cdot W)^T}{2}$$

554   Given $S' \cdot W$ is not a symmetrical matrix, we average $S' \cdot W$ and $(S' \cdot W)^T$ to obtain the

555   denoised matrix. As a proof of the reliability of our algorithm, the upper triangle and lower

556   triangle of $S' \cdot W$ are always close to each other. The distance matrix D is calculated by

557   $D = -S_{denoise}$. Compared to the denoising method described in Methods, the alternative

558   method denoises the data and largely keeps the information of original matrix (including

559   variations from both batch effects and biological heterogeneity).