

 Open access • Posted Content • DOI:10.1101/2021.09.08.459495

scBasset: Sequence-based modeling of single cell ATAC-seq using convolutional neural networks — [Source link](#)

Han Yuan, David R. Kelley

Published on: 10 Sep 2021 - [bioRxiv](#) (Cold Spring Harbor Laboratory)

Topics: Convolutional neural network

Related papers:

- [SCALE method for single-cell ATAC-seq analysis via latent feature extraction.](#)
- [Simultaneous deep generative modeling and clustering of single cell genomic data](#)
- [Learning latent embedding of multi-modal single cell data and cross-modality relationship simultaneously](#)
- [SAILER: Scalable and Accurate Invariant Representation Learning for Single-Cell ATAC-Seq Processing and Integration](#)
- [coupleCoC+: an information-theoretic co-clustering-based transfer learning framework for the integrative analysis of single-cell genomic data](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/scbasset-sequence-based-modeling-of-single-cell-atac-seq-32599wtelr>

1 scBasset: Sequence-based modeling of single cell 2 ATAC-seq using convolutional neural networks

3 Han Yuan^{1,*} and David R Kelley^{1,*}

4 ¹Calico Life Sciences, South San Francisco, CA 94080, USA

5 *correspondence: yuanh@calicolabs.com, drk@calicolabs.com

6 September 8, 2021

7 **1 Abstract**

8 Single cell ATAC-seq (scATAC) shows great promise for studying cellular het-
9 erogeneity in epigenetic landscapes, but there remain significant challenges in
10 the analysis of scATAC data due to the inherent high dimensionality and spar-
11 sity. Here we introduce scBasset, a sequence-based convolutional neural net-
12 work method to model scATAC data. We show that by leveraging the DNA
13 sequence information underlying accessibility peaks and the expressiveness of
14 a neural network model, scBasset achieves state-of-the-art performance across
15 a variety of tasks on scATAC and single cell multiome datasets, including cell
16 type identification, scATAC profile denoising, data integration across assays,
17 and transcription factor activity inference.

18 **2 Introduction**

19 Single cell ATAC-seq (scATAC) reveals epigenetic landscapes at single cell res-
20 olution (Buenrostro et al., 2018). The assay has been successfully applied to
21 identify cell types and their specific regulatory elements, reveal cellular hetero-
22 geneity, map disease-associated distal elements, and reconstruct differentiation
23 trajectories (Satpathy et al., 2019; Miao et al., 2021; Cusanovich et al., 2018).

24 However, there still exist significant challenges in the analysis of scATAC
25 data, due to the inherent high dimensionality of accessible peaks and sparsity
26 of sequencing reads per cell (Bravo González-Blas et al., 2019; Chen et al.,
27 2019). Multiple approaches have been proposed to address these challenges,
28 which can be broadly categorized into two main classes: sequence-free and
29 sequence-dependent methods. Starting from a sparse peak-by-cell matrix gen-
30 erated through aggregation of reads and peak calling in accessible chromatin,
31 most methods represent these annotated peaks as genomic coordinates and ig-
32 nore the underlying DNA sequence. Principal component analysis (PCA) and

33 latent semantic indexing (LSI) perform a linear transformation of the peak-by-
34 cell matrix to project the cells to a low-dimensional space (Pliner et al., 2018;
35 Cusanovich et al., 2018). SCALE and cisTopic model the generative process of
36 the data distribution using latent dirichlet allocation or a variational autoen-
37 coder (Bravo González-Blas et al., 2019; Xiong et al., 2019). These sequence-free
38 methods are able to detect biologically meaningful covariance to effectively rep-
39 resent and cluster or classify cells. However, they ignore sequence information
40 and rely on post-hoc motif matching tools to relate accessibility to transcription
41 factors (TFs). In contrast, sequence-dependent methods such as chromVAR and
42 BROCKMAN represent peaks by their TF motif or k-mer content and aggregate
43 these features across peaks or other regions of interest to learn cell representa-
44 tions (Schep et al., 2017; de Boer and Regev, 2018). While chromVAR directly
45 associates peaks to TFs, emphasizing interpretability, it tends to perform worse
46 in learning cell representations, potentially due to the loss of information from its
47 simple implicit model relating sequence to accessibility through position weight
48 matrices Chen et al. (2019).

49 Here, we propose a more expressive sequence-dependent model based on
50 deep convolutional neural networks (CNNs). CNNs can predict peaks from
51 bulk chromatin profiling assays more effectively than k-mer or TF motif mod-
52 els, exemplified by DeepSEA and Basset (Kelley et al., 2016; Zhou and Troy-
53anskaya, 2015). These models compute explicit embeddings of the sequences
54 underlying peaks via the convolutional layers and implicit embeddings of the
55 multiple “tasks” (which are sequencing experiments) in parameters of the final
56 linear transformation. We extend the Basset architecture to predict single cell
57 chromatin accessibility from sequences, using a bottleneck layer to learn low-
58 dimensional representations of the single cells. We show that by making use of
59 sequence information in a deep learning framework, we outperform state-of-the-
60 art methods for cell representation learning, single cell accessibility denoising,
61 scATAC integration with scRNA, and transcription factor activity inference.

62 **3 Results**

63 **3.1 scBasset predicts single cell chromatin accessibility on** 64 **held-out peaks**

65 scBasset is a deep CNN to predict chromatin accessibility from sequence. CNNs
66 have demonstrated state-of-the-art performance for predicting epigenetic pro-
67 files in bulk data and have been successfully used for genetic variant effect
68 prediction and TF motif grammar inference (Kelley et al., 2016; Zhou and
69 Troyanskaya, 2015; Kelley et al., 2018; Zhou et al., 2018; Agarwal and Shen-
70dure, 2020; Avsec et al., 2021). Here, we move the focus away from maximizing
71 accuracy on held-out sequences and view the model as a representation learn-
72 ing machine. When trained to achieve multiple tasks, the final layer of these
73 models involves a sequence embedded by the convolutional layers and a linear
74 transformation to predict the data in each separate task. The linear transfor-

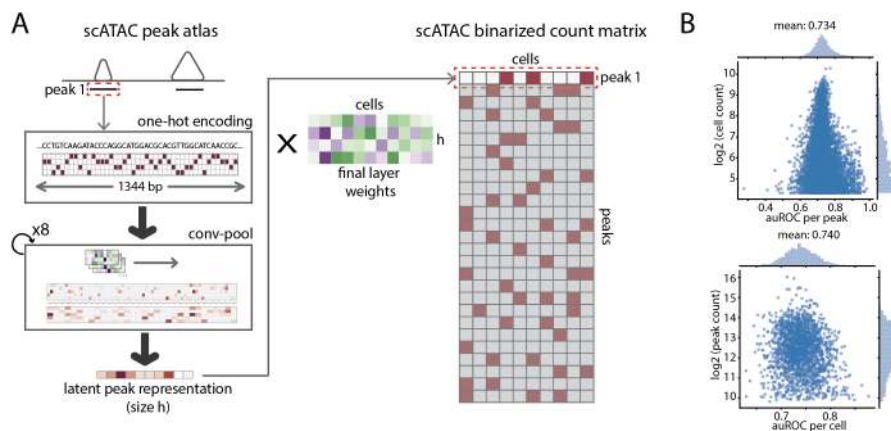


Figure 1: scBasset architecture. A) scBasset is a deep convolutional neural network to predict single cell chromatin accessibility from the DNA sequence underlying peak calls. B) scBasset prediction performance on held-out peaks evaluated by auROC per peak (top) and auROC per cell (bottom) for the Buenrostro2018 dataset.

75 mation matrix comprises a vector representation of each task (here, each single
76 cell), which specifies how to make use of each of the sequence embedding latent
77 variables to predict cell-specific accessibility. In a simple ideal scenario, one can
78 imagine each latent variable representing various regulatory factors such as TF
79 binding or nucleotide composition, and the final transformation specifying how
80 much each cell depends on that factor. We propose that these single cell vec-
81 tors serve as intriguing representations of the cells for downstream tasks such
82 as visualization and clustering.

83 We recommend that users first apply standard processing techniques to bring
84 the raw data to a peak-by-cell binary matrix. scBasset takes as input a 1344
85 bp DNA sequence from each peak’s center and one-hot encodes it as a 4×1344
86 matrix. The input DNA sequence goes through eight convolution blocks, where
87 each block is composed of a 1D convolution, batch normalization, max pooling,
88 and GELU activation layers. Unlike most previous architectures, we follow these
89 by a bottleneck layer of size h intended to learn a low-dimensional representation
90 of the peak via the layer output and the cells via the parameters of the following
91 layer. Finally, a dense linear transformation connects the bottleneck sequence
92 embeddings to predict binary accessibility in each cell (Fig.1a). We apply the
93 standard binary cross-entropy loss function and optimize model parameters with
94 stochastic gradient descent (Methods).

95 To benchmark our approach, we applied scBasset to three public datasets:
96 a scATAC-seq FACS-sorted hematopoietic differentiation dataset (referred to as
97 Buenrostro2018) with 2k cells (Buenrostro et al., 2018), 10x Multiome RNA+ATAC
98 PBMC dataset with 3k cells, and 10x Multiome RNA+ATAC mouse brain

99 dataset with 5k cells. The first dataset provides ground-truth cell type labels from flow cytometry. We consider the multiome datasets to be a valuable resource to validate scATAC methods since they provide independent measurements of gene expression and chromatin accessibility in the same cells.

100
101
102
103 First, we asked how well scBasset can predict accessibility across cells for held out peak sequences to ensure that the model has learned a meaningful relationship between DNA sequence and accessibility using the sparse noisy labels. For held out peaks, we computed the area under the receiver operating characteristic curve (auROC) across peaks for each cell and averaged across cells (referred to as “per peak”). To evaluate cell type specificity, we also computed auROC across cells for each peak and averaged across peaks (referred to as “per cell”). scBasset achieved compelling accuracy levels that indicate successful learning: 0.734 per peak and 0.740 per cell for Buenrostro2018 dataset (Fig.1b), 0.662 per peak and 0.640 per cell for the 10x multiome PBMC, and 0.734 per peak and 0.701 per cell for the 10x multiome mouse brain dataset (Fig.S1). Although these statistics are slightly below the 0.75-0.95 range achieved for bulk DNase samples in the original Basset publication, this is inevitable due to the substantially increased measurement noise due to sparse sequencing for the single cell assay. In support of this claim, we observed that in the 10x multiome PBMC and mouse brain datasets, peaks with very high read coverage are easier to predict (Fig.S1). Given that ubiquitous accessible peaks are known to exist, these peaks are likely truly accessible in all cells and represent a rough upper bound on the achievable accuracy.

122 **3.2 scBasset final layer learns cell representations**

123 We propose that the $h \times \text{cell}$ weight matrix that connects the bottleneck layer to the predictions be used as a low-dimensional representation of the single cells. One requirement for an effective cell representation is removal of the influence of sequencing depth. Thus, we first verified that the intercept vector in the model’s final layer almost perfectly correlates with cell sequencing depth for all datasets (Fig.S2), suggesting that depth has been normalized out from the representations. Next, we compared the cell representations learned by scBasset with other methods both qualitatively and quantitatively. For the Buenrostro2018 dataset, we visualized the cell embeddings in 2D using t-distributed stochastic neighbor embedding (t-SNE) (Fig.2a) and observed differentiation trajectories in the t-SNE space. Compared to other popular methods for scATAC embedding, we observed that chromVAR and PCA have difficulty distinguishing CLP from LMPP, while Cicero, SCALE, cisTopic, and scBasset make the distinction (Fig.S4). Following previous work, we quantified the correctness of cell embeddings by comparing Louvain clustering results with ground-truth cell type labels using the adjusted rank index (ARI) (Chen et al., 2019). scBasset outperforms the other methods according to this metric (Fig.2b,top). Since ARI is sensitive to the hyperparameter choice and stochasticity in the Louvain algorithm, we proposed an alternative method for evaluating cell embeddings. We computed a “label score” by building a nearest neighbor graph based on the cell embed-

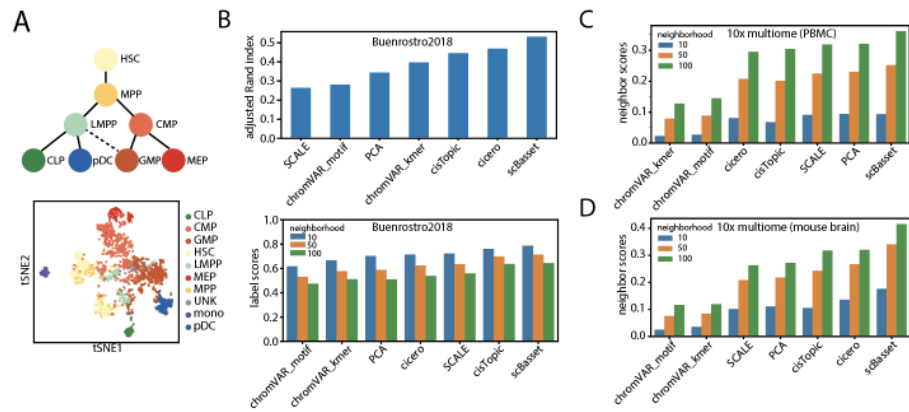


Figure 2: scBasset performance at learning cell representations. A) Top, hematopoietic stem cell differentiation lineage diagram in the Buenrostro2018 study; bottom, t-SNE visualization of cell embeddings learned by scBasset, colored by cell types. B) Top, performance comparison of different cell embedding methods evaluated by adjusted Rand index; bottom, performance comparison of different cell embedding methods evaluated by label score (Methods). C) Performance comparison of different cell embedding methods evaluated by neighbor scores for the 10x multiome PBMC dataset. D) Performance comparison of different cell embedding methods evaluated by neighbor scores for the 10x multiome mouse brain dataset.

143 dings and asked what percentage of each cell’s neighbors share its same label.
 144 For each embedding method, we computed label scores across a range of neigh-
 145 borhoods and observed scBasset consistently outperforms the competitors at
 146 learning cell representations that embed cells of the same type near each other
 147 (Fig.2b,bottom). We also evaluated label scores for each cell type individually
 148 and observed that monocytes are learned best, whereas MPP cells are most
 149 difficult to distinguish (Fig.S3).

150 For the multiome PBMC and mouse brain datasets, we computed an analo-
 151 gue to the label scores for cell embeddings. Since the ground-truth cell types
 152 for the multiome datasets are unknown, we used cluster identifiers from scRNA-
 153 seq Leiden clustering as cell type labels. Again, scBasset outperforms the com-
 154 petitors by this metric across a range of neighborhoods (Fig.S5). For these
 155 multiome datasets, we also computed a “neighbor score”, in which we built
 156 independent nearest neighbor graphs from the scRNA and scATAC and asked
 157 what percentage of each cell’s neighbors are shared between the two graphs.
 158 scBasset outperforms the competitors on both multiome PBMC and multiome
 159 mouse brain datasets when evaluated with neighbor scores across a range of
 160 neighborhoods (Fig.2c,d).

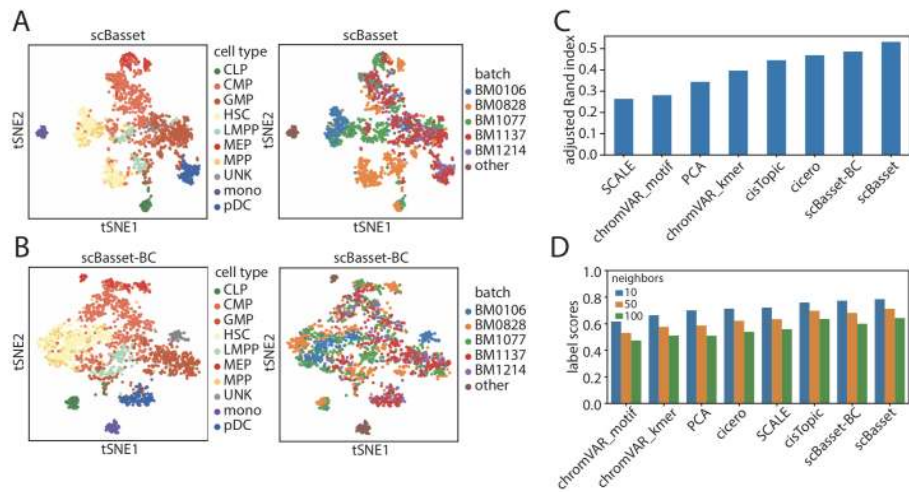


Figure 3: scBasset can be adapted to perform batch correction. A) Cell embeddings learned by scBasset without batch correction, colored by cell type (left) and batch (right). B) Cell embeddings learned by scBasset with batch correction (scBasset-BC), colored by cell type (left) and batch (right). C) Performance comparison of different cell embedding methods to scBasset-BC evaluated by adjusted Rand index. D) Performance comparison of different cell embedding methods to scBasset-BC evaluated by label score.

161 3.3 Batch-conditioned scBasset removes batch effects

162 In the Buenrostro2018 dataset, HSCs cluster into two populations, regardless of
 163 which cell embedding method we apply (Fig.S4). As noted in previous studies,
 164 this is caused by a batch effect due to different donors (Fig.3a) (Buenrostro
 165 et al., 2018; Bravo González-Blas et al., 2019). To correct for this, and batch
 166 effects more generally, we explored modifications to the scBasset architecture.

167 Specifically, after the bottleneck layer, we added a second fully-connected
 168 layer to predict the batch-specific contribution to accessibility (Methods, Fig.S6).
 169 We added the output of the batch layer and cell-specific layer before comput-
 170 ing the final sigmoid. Intuitively, we expect the batch-specific variation will be
 171 captured in this path, whereas the original $h \times$ cell weight matrix will focus on
 172 the remainder of biologically relevant variation.

173 We compared the scBasset cell embedding results before and after batch
 174 correction. We observed an overall mixing of different batches in the t-SNE
 175 space after batch correction. For example, we can see that the two HSC batches
 176 (BM0106 and BM0828) merge into one cluster. In addition, pDC cells from
 177 BM1137 and BM1214 batches previously fell into two distinct sub-clusters, but
 178 are mixed together after batch correction (Fig.3ab). However, we noticed a small
 179 decrease in the cluster evaluation metrics after batch correction. We hypothesize

180 that this is caused by imbalances in cell type distribution from different donors,
181 which are then learned by the batch layer rather than the cell-specific layer.
182 This is also consistent with a recent study’s observation of a trade-off between
183 mixing and cell type separation (Ashuach et al., 2021). Nevertheless, scBasset-
184 BC still outperforms the competitors when evaluated by ARI and is among the
185 top performers when evaluated by label scores (Fig.3cd).

186 As an additional benchmark, we trained scBasset and scBasset-BC on a
187 mixture of PBMC scATAC data from 10x multiome and 10x nextgem chem-
188 istry (Methods). We observed that while there is a strong batch effect between
189 the two chemistries when trained with naive scBasset, scBasset-BC successfully
190 integrated the two datasets (Fig. S6).

191 **3.4 scBasset denoises single cell accessibility profiles**

192 Due to the sparsity of scATAC, the binary accessibility indicator for any given
193 cell and peak contains ample false negatives, such that the data cannot be
194 studied with true single cell resolution and is usually aggregated across cells.
195 However, numerous methods deliver denoised (or imputed) numeric values to
196 represent the accessibility status at every cell/peak combination. scBasset com-
197 utes such values in its sequence-based predictions.

198 From the Buenrostro2018 dataset, we sampled 500 peaks and 200 cells
199 and directly visualized the raw cell-by-peak matrix versus the denoised matrix
200 (Fig.4a). In the raw count matrix, we observed that cells and peaks clustered
201 by sequencing depth, showing no biologically relevant patterns. However, we
202 observed that after scBasset denoising, cells of the same cell type share simi-
203 lar accessibility profiles and hierarchical clustering of cells matched well with
204 ground-truth labels.

205 Several published strategies aggregate scATAC counts in the region around a
206 gene’s transcription start site to estimate its transcription (Granja et al., 2021;
207 Pliner et al., 2018). We propose that effective denoising would improve the corre-
208 lation between these gene accessibility estimates and the gene’s measured RNA
209 expression in multiome experiments. Thus, we computed accessibility scores for
210 each gene by averaging the predicted accessibility values at all promoter peaks
211 before and after denoising (Methods). For both the 10x multiome PBMC and
212 mouse brain datasets, we observed that scBasset denoising improves the con-
213 sistency between gene accessibility and expression ($P < 2.2e-16$, Wilcoxon signed
214 rank test). As one would expect, the improvement is greater for cells with fewer
215 scATAC UMIs (Fig.4b, Fig.S7).

216 Covariance-based methods can also be used to denoise scATAC, and we
217 compared scBasset to SCALE, a sequence-independent method for accessibility
218 denoising using a variational autoencoder. We observed that SCALE gene ac-
219 cessibility scores correlated better than scBasset with gene expression (Fig.S7).
220 Because the two methods take independent approaches (sequence-dependent
221 versus sequence-free), we hypothesized that combining the denoised values from
222 both via a simple average would further improve concordance. Indeed, we ob-
223 served that for both 10x multiome datasets, the combined prediction performs

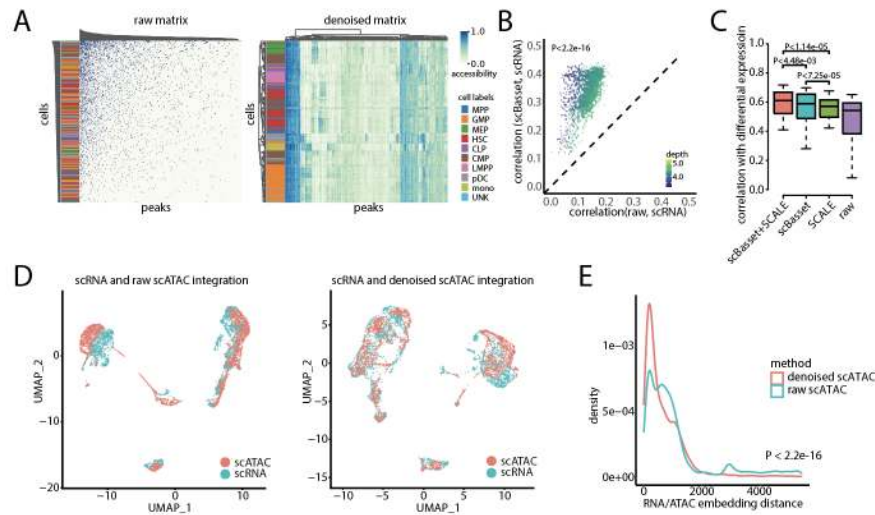


Figure 4: scBasset denoising performance. A) Left, binary count matrix of 200 cells and 500 peaks sampled from Buenrostro2018 dataset, hierarchically clustered by both cells and peaks. Cell type labels annotate the rows. Right, the same matrix and procedure after scBasset denoising. B) Correlation between gene accessibility score and gene expression for each cell before (x-axis) and after denoising (y-axis) for the multiome PBMC dataset. Cells are colored by sequencing depth. C) Comparison of denoising performance on multiome PBMC dataset between raw data, scBasset, SCALE, and scBasset+SCALE combine, evaluated by consistency in differential expression \log_2FC and differential accessibility \log_2FC . We performed Wilcoxon signed rank tests for performance comparisons. D) Left, 10x multiome PBMC RNA (blue) and raw ATAC (red) profile embeddings after integration. Right, 10x multiome PBMC RNA (blue) and denoised ATAC (red) profile embeddings after integration. E) Distribution of the relative distances (Method) between each cell's RNA and ATAC embeddings after integration when using raw ATAC profiles (blue) or denoised ATAC profiles (red). We performed Wilcoxon signed rank test for performance comparison.

224 better than SCALE or scBasset alone when we evaluated consistency with base-
 225 line expression (Fig.S7).

226 Studies have shown that changes in accessibility and expression correlate bet-
 227 ter with each other than their absolute values, and thus would be a more useful
 228 metric for validating accessibility denoising methods (Pliner et al., 2018). We
 229 evaluated scBasset and SCALE accessibility denoising for consistency between
 230 differential expression and differential accessibility. For each cell type cluster as
 231 defined by scRNA in the 10x PBMC dataset, we performed differential expres-
 232 sion and differential accessibility analysis against the rest of the cells. To assess

233 denoising quality, we evaluated the correlation between differential expression
234 log2 fold change (log2FC) and differential accessibility log2FC before and after
235 denoising (Fig.4c).

236 We observed that expression log2FC and accessibility log2FC correlates
237 well even for raw accessibility data ($r=0.47$). Still, consistency is significantly
238 improved after scBasset denoising ($r=0.54$). Interestingly, we observed that
239 even though SCALE correlation exceeded that of scBasset for baseline ac-
240 cessibility/expression, scBasset significantly outperforms SCALE when evalu-
241 ated by differential accessibility/expression ($p<7.25e-05$). We hypothesize that
242 SCALE's reliance on cell-cell covariance encourages cells to be more similar to
243 each other than they actually are and over-smooths (Tjarnberg et al., 2021;
244 Ashuach et al., 2021). scBasset will be less prone to over-smoothing since each
245 peak is considered only through its sequence. As a result, SCALE performs
246 better in denoising baseline accessibility, while scBasset performs better in de-
247 noising differential accessibility, which emphasizes cell identity. As with baseline
248 expression, combining scBasset and SCALE produces greater performance than
249 either method alone (Fig.4c, Fig.S7).

250 Integration of cells independently profiled by scRNA and scATAC into a
251 shared latent space is a key step for many scATAC annotation and analysis
252 methods (Stuart et al., 2019). We hypothesized that scATAC denoising would
253 improve scRNA and scATAC integration performance. In order to evaluate inte-
254 gration performance, we treated the 10x multiome scRNA and scATAC profiles
255 as having originated from two independent experiments. For the 10x multi-
256 ome PBMC dataset, we observed that when we integrated the scRNA profiles
257 with the denoised scATAC profiles, the cells achieve better mixing compared to
258 when we integrated scRNA with raw scATAC profiles (Fig.4d). Quantitatively,
259 we measured the multiome rank distance between the RNA and ATAC em-
260 beddings for each matching cell (Methods). We observed the RNA and ATAC
261 profiles of the same cell are embedded significantly closer to each other when the
262 ATAC profile is denoised compared to the raw ATAC profile (Fig.4e, $P<2.2e-$
263 16). We observed similar results for the 10x multiome mouse brain dataset
264 (Fig.S8).

265 **3.5 scBasset infers transcription factor activity at single** 266 **cell resolution**

267 Transcription factor binding is a major driver of chromatin accessibility (Thur-
268 man et al., 2012). Since scBasset learns to predict accessibility from sequence,
269 we expect the model to capture sequence information predictive of TF binding.
270 To query the single cell TF activity, we leveraged the flexibility of the scBasset
271 model to predict arbitrary sequences. More specifically, we fed synthetic DNA
272 sequences (dinucleotide shuffled peaks) with and without a particular TF motif
273 of interest to a trained scBasset model and evaluated the activity of the motif
274 in each cell based on changes in predicted accessibility (Methods) (Kelley et al.,
275 2016). If a TF is playing an activating role in a particular cell, we expect to see
276 increased accessibility after the TF motif is inserted.

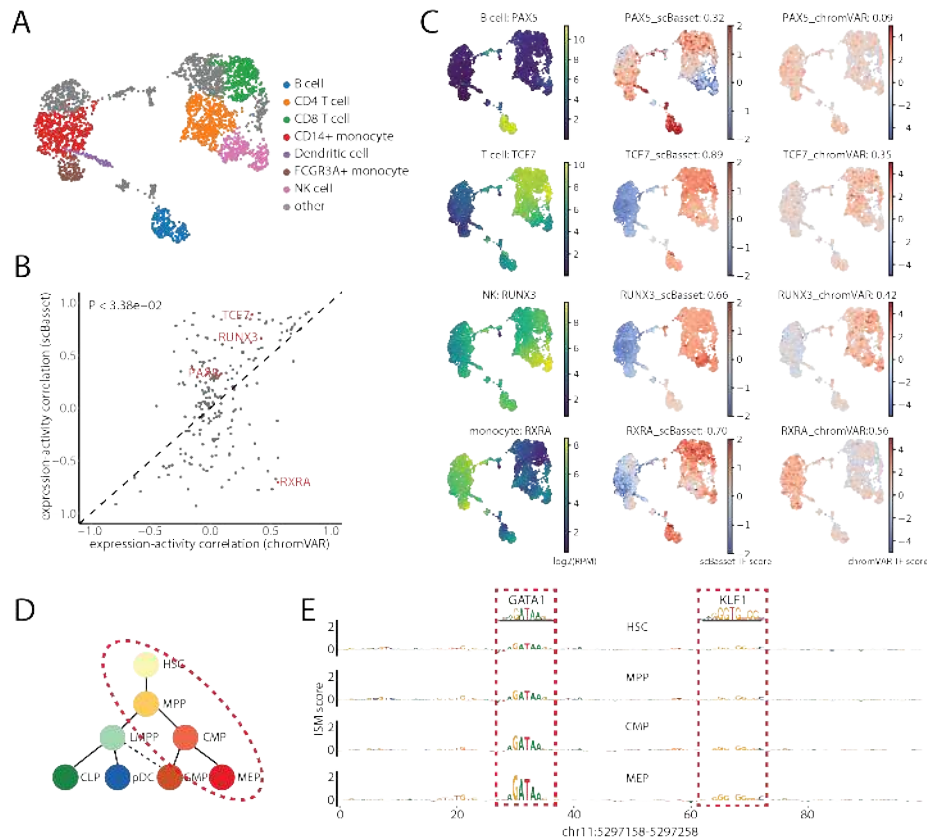


Figure 5: scBasset infers single cell TF activity. A) UMAP showing annotated PBMC cell types. B) Pearson correlation between TF expression and scBasset or chromVAR-predicted TF activity for 203 differentially expressed TFs. The example TFs that we examine in panel C are highlighted in red. C) UMAP visualization of TF expression (left), scBasset TF activity (middle), and chromVAR TF activity (right) for key PBMC regulators. Pearson correlation between inferred TF activity and expression are shown in the title. D) ISM scores for β -globin enhancer at chr11:5297158-5297258 for cells in HSC, MPP, CMP and MEP cell types. Sequences that match GATA1 and KLF1 motifs are highlighted in red boxes.

277 TF regulation in the hematopoietic lineage profiled in the Buenrostro2018
 278 dataset has been studied in detail. We performed motif injection for all 733 hu-
 279 man CIS-BP motifs using the Buenrostro2018-trained model and recapitulated
 280 known trajectories of motif activity. For example, CEBPB, a known regulator
 281 of monocyte development, shows the highest activity in monocytes; GATA1, a
 282 key regulator of the erythroid lineage, is predicted to be most active in MEPs;
 283 HOXA9, a known master regulator of HSC differentiation, has highest predicted

284 activity in HSCs (Fig.S9) (Buenrostro et al., 2018).

285 Previous sequence-based methods such as chromVAR are also able to quan-
286 tify TF motif activity. To comprehensively compare scBasset and chromVAR
287 on this task, we analyzed the 10x PBMC multiome dataset, in which TF ex-
288 pression measured in the RNA can serve as a proxy for its motif's activity. We
289 inferred motif activity for all 733 human CIS-BP motifs using both scBasset
290 and chromVAR. For the 203 TFs that are significantly differentially expressed
291 between cell type clusters, we asked how well the inferred TF activity per cell
292 correlates with its expression. We observed that overall scBasset TF activi-
293 ties correlate significantly better with expression than chromVAR TF activities
294 ($P < 3.38e-02$, Wilcoxon signed rank test) (Fig.5b). This one-sided test is an
295 underestimate of scBasset's performance advantage over chromVAR, since we
296 would expect TF expression and inferred activity to be negatively correlated
297 for repressors. Thus, we evaluated scBasset and chromVAR on activating and
298 repressive TFs separately. For 74 TFs which both methods agreed on a positive
299 TF expression-activity correlation, scBasset predicted TF activities have signif-
300 icantly greater correlation with expression than chromVAR predicted activity
301 ($P < 7.38e-12$, Wilcoxon signed rank test, Fig.S10). For 41 TFs which both meth-
302 ods agreed on a negative TF expression-activity correlation, scBasset predicted
303 TF activities have a significantly lesser correlation (more negative) with ex-
304 pression than chromVAR predicted activity ($P < 1.62e-08$, Wilcoxon signed rank
305 test). This is also true for the 10x multiome mouse brain dataset (Fig.S10).

306 Examining some of the key regulators of PBMC cell types, we observed
307 that scBasset TF activities have better cell type specificity and correlate better
308 with TF expression than chromVAR (Fig.5c). For example, PAX5 is a known
309 master regulator of B cell development (Medvedovic et al., 2011). scBasset pre-
310 dicted B cell specific activity of PAX5, which correlates with PAX5 expression
311 ($r=0.32$), while chromVAR PAX5 activity did not have any cell type speci-
312 ficity or significant PAX5 expression correlation ($r=0.09$). scBasset-predicted
313 activity of the T cell differentiation regulator TCF7 highly correlates with ex-
314 pression ($r=0.89$), while chromVAR TCF7 activity has lesser specificity and
315 expression correlation ($r=0.35$). NK cells have greater expression of RUNX3
316 and scBasset captures this elevated activity in NK cells ($r=0.66$) more effec-
317 tively than chromVAR ($r=0.42$). For monocytes, both scBasset and chromVAR
318 predicted specific activity of CEBPB, with scBasset activity correlating slightly
319 better with expression (0.75 vs. 0.68, Fig.S11). Interestingly, while scRNA-
320 seq suggests monocyte-specific expression of RXRA, scBasset and chromVAR
321 strongly disagree, making opposite predictions for RXRA activity; scBasset pre-
322 dicts RXRA as a repressor ($r=-0.70$) while chromVAR suggests an activating
323 role ($r=0.56$). A literature review revealed stronger evidence that RXRA plays
324 a repressive role in the myeloid lineage through direct DNA binding, which is
325 more consistent with the scBasset prediction (Kiss et al., 2017).

326 Unlike chromVAR, scBasset makes use of an accurate quantitative model
327 that predicts cell type specific accessibility from the DNA nucleotides. Not only
328 are we able to query scBasset for TF activity on a per-cell level, we can also
329 infer TF activity at per-cell per-nucleotide resolution. As a proof of principle,

330 we examined a known enhancer for the β -globin gene that regulates erythroid-
331 specific *beta*-globin expression (Tuan et al., 1985; Li et al., 2002). We performed
332 *in silico* saturation mutagenesis (ISM) for this 100 bp sequence, in which we pre-
333 dicted the change in accessibility in each cell after mutating each position to its
334 three alternative nucleotides. We aggregated to a single score for each position
335 by taking the normalized ISM score for each reference nucleotide (Methods).
336 Fig.5d shows the average ISM score for each cell type in the erythroid lineage.
337 We observed that the most influential nucleotides correspond to GATA1 and
338 KLF1 motifs, which are TFs known to bind to this enhancer region and regu-
339 late β -globin expression (Tallack et al., 2010). Interestingly, GATA1 and KLF1
340 motifs contribute more to the accessibility as the cells differentiate in the ery-
341 throid lineage. In comparison, these two motifs' nucleotides have low scores in
342 cell types outside of the erythroid lineage (Fig.S12). This experiment suggests
343 that scBasset learns the accessibility regulatory grammar at single cell reso-
344 lution and could be used to identify the TFs regulating specific enhancers in
345 individual cells and lineages.

346 4 Discussion

347 In this study we present scBasset, a sequence-based deep learning framework
348 for modeling scATAC data. scBasset is trained to predict individual cell ac-
349 cessibility from the DNA sequence underlying ATAC peaks, learning a vector
350 embedding to represent the single cells in the process. A trained scBasset model
351 can strengthen multiple lines of scATAC analysis, and we demonstrate state-of-
352 the-art performance on several tasks. Clustering the model's cell embeddings
353 achieves greater alignment with ground-truth cell type labels. The model out-
354 puts can be used as denoised accessibility profiles, which improve concordance
355 with RNA measurements. The model learns to recognize TF motifs and their
356 influence on accessibility, and we designed an *in silico* experiment to inject mo-
357 tifs into background sequences to query for TF motif activity in single cells.
358 The model can also be applied to predict the influence of mutations, enabling
359 *in silico* saturation mutagenesis of regulatory sequences of interest at single cell
360 resolution. Compared to previous sequence-based approaches for scATAC anal-
361 ysis such as chromVAR, scBasset achieves better performance at learning cell
362 embeddings and inferring TF activity, because scBasset benefits from a more ex-
363 pressive CNN model that learns more sophisticated sequence features, including
364 non-linear relationships. Compared to previous sequence-free approaches such
365 as cisTopic or SCALE, scBasset achieves better performance in benchmarking
366 tasks and delivers a more interpretable model that can be directly queried for
367 TF activity or identifying regulatory sequences.

368 Sequence-based approaches have several limitations. First, we make use of
369 the reference genome, but many samples will have variant versions, including
370 copy number variations that could lead our models astray. Second, we assume
371 that the regulatory motifs and their interactions generalize across the genome.
372 This assumption may not be entirely true at some genomic loci for which evolu-

373 tion led to bespoke regulatory solutions, such as for X chromosome inactivation
374 in females. However, since scBasset takes a completely independent approach
375 to the covariance-based methods, one can always combine these two types of ap-
376 proaches to further improve their analyses, as we showed for denoising (Fig.4).

377 In addition, we foresee several paths to further improve our method. To
378 improve scBasset memory efficiency in order to scale to extremely large datasets,
379 one could sample mini-batches of both sequences and cells rather than only
380 sequences in our current implementation. Methods such as TF-MoDISco could
381 be applied to scBasset ISM scores for de novo motif discovery (Shrikumar et al.,
382 2018; Avsec et al., 2021). All approaches to scATAC analysis depend on accurate
383 peak calls, and predictive modeling frameworks have been proposed to help
384 identify highly specific regulatory elements (Lal et al., 2021). We expect a
385 neural network model would further improve scATAC peak calling by taking
386 into account sequence information (and accounting for Tn5 transposition bias).
387 Finally, we plan to explore transfer learning approaches in which models are
388 pre-trained on large data compendia before fine-tune training on specific single
389 cell datasets.

390 5 Methods

391 5.1 scATAC-seq preprocessing

392 We downloaded the count matrix and peak atlas files for the Buenrostro2018
393 dataset from GEO (Accession GSE96769) (Buenrostro et al., 2018). Peaks ac-
394 cessible in less than 1% cells were filtered out. The final dataset contains 126,719
395 peaks and 2,034 cells.

396 We downloaded the 10x multiome datasets from 10x Genomics: [https://](https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/2.0.0/pbmc_granulocyte_sorted_3k)
397 [support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/2.0.](https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/2.0.0/pbmc_granulocyte_sorted_3k)
398 [0/pbmc_granulocyte_sorted_3k](https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/2.0.0/pbmc_granulocyte_sorted_3k) for PBMC dataset, and [https://support.](https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/2.0.0/e18_mouse_brain_fresh_5k)
399 [10xgenomics.com/single-cell-multiome-atac-gex/datasets/2.0.0/e18_mouse_](https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/2.0.0/e18_mouse_brain_fresh_5k)
400 [brain_fresh_5k](https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/2.0.0/e18_mouse_brain_fresh_5k) for mouse brain dataset. Genes expressed in less than 5% cells
401 were filtered out. Peaks accessible in less than 5% cells were filtered out.

402 5.2 scRNA-seq preprocessing

403 For the 10x multiome datasets, we processed the expression data with scVI
404 version 0.6.5 with `n.layers=1`, `n.hidden=768`, `latent=64` and a dropout rate of
405 0.2 (Lopez et al., 2018). We trained scVI for 1000 epochs with learning rate of
406 0.001, using the option to reduce the learning rate upon plateau using options
407 `lr_patience` of 20 and `lr_factor` of 0.1. We enabled early stopping when there was
408 no improvement on the ELBO loss for 40 epochs.

409 To generate denoised expression profiles, we used the `get_sample_scale()`
410 function to sample from the generative model 10 times and took the average.
411 We used the learned latent cell representations to build nearest neighbor graphs
412 and perform cell clustering.

413 5.3 Model architecture

414 scBasset is a neural network architecture that predicts binary accessibility vec-
415 tors for each peak based on its DNA sequence. scBasset takes as input a 1344
416 bp DNA sequence from each peak’s center and one-hot encodes it as a 1344×4
417 matrix. The neural network architecture includes the following blocks:

- 418 • 1D convolution layer with 288 filters of size 17×4 , followed by batch nor-
419 malization, Gaussian error linear unit (GELU), and width 3 max pooling
420 layers, which generates a 488×288 output matrix.
- 421 • Convolution tower of 6 convolution blocks each consisting of convolution,
422 batch normalization, max pooling, and GELU layers. The convolution
423 layers have increasing numbers of filters (288, 323, 363, 407, 456, 512) and
424 kernel width 5. The output of the convolution tower is a 7×512 matrix.
- 425 • 1D convolution layer with 256 filters with kernel width 1, followed by batch
426 normalization and GELU, The output is a 7×256 matrix, which is then
427 flattened into a 1×1792 vector.
- 428 • Dense bottleneck layer with 32 units, followed by batch normalization,
429 dropout (rate=0.2), and GELU. The output is a compact peak represen-
430 tation vector of size 1×32 .
- 431 • Final dense layer predicting continuous accessibility logits for the peaks
432 in every cell.
- 433 • (Optional) To perform batch correction, we attach a second parallel dense
434 layer to the bottleneck layer predicting batch-specific accessibility. This
435 batch-specific accessibility is multiplied by the batch-by-cell matrix to
436 compute the batch contribution to accessibility in every cell. This vector is
437 then added to the previous continuous accessibility logits per cell (Fig.S6).
438 L2 regularization can be optionally applied to the cell-embedding path
439 (with hyperparameter λ_1) or the batch-specific path (with hyperparameter
440 λ_2) to tune the contribution of the batch covariate to the predictions.
- 441 • Final sigmoid activation to $[0,1]$ accessibility probability.

442 The total number of trainable parameters in the model is a function of the
443 number of cells in the dataset. Specifically, the model will have $4513960 + 33 \times n_{\text{cells}}$
444 number of trainable parameters.

445 5.4 Training approach

446 We used a binary cross entropy loss and monitored the training area under the
447 receiver operator curve (auROC) after every epoch. We stopped training when
448 the maximum training auROC improved by less than $1e-6$ in 50 epochs. This
449 stopping criteria led to training for around 600 epochs for the Buenrostro2018

450 dataset, 1100 epochs for the 10x multiome PBMC dataset and 1200 epochs for
451 the 10x multiome mouse brain dataset.

452 We focused on training auROC instead of validation auROC for model selec-
453 tion because we observed that the model continues to improve cell embeddings
454 even after the point where the validation auROC has plateaued (Fig.S14). Since
455 our goal in this application is to learn better representations instead of mini-
456 mum generalization loss, we focused on the convergence of the training auROC.
457 In addition, at the bottleneck size of 32, there was only a small drop in gen-
458 eralization performance (validation auROC) when the training auROC reaches
459 convergence (0.734 versus 0.742).

460 We updated model parameters using stochastic gradient descent using the
461 Adam update algorithm. We performed a random search for optimal hyperpa-
462 rameters including: batch size, learning rate, beta1, and beta2 for the Adam
463 optimizer. The best performance was achieved with a batch size of 128, learning
464 rate of 0.01, beta.1 of 0.95, and beta.2 of 0.9995.

465 We focused on the Buenrostro2018 dataset to select the optimal bottleneck
466 layer size. We trained models with bottleneck sizes of 8, 16, 32, 64 and 128 and
467 observed that bottleneck size 32 gives the best performance (Fig.S13).

468 **5.5 Alternative scATAC-seq methods**

469 **5.5.1 PCA**

470 We performed PCA with the scikit-learn python package. We evaluated the per-
471 formance of PCA cell embedding using 10, 20, 30, 40, 50, 60, 70, 80, 90, 100 PCs,
472 with or without the first PC, and reported the model with best performance to
473 compare to scBasset.

474 **5.5.2 cicero**

475 We used Cicero via its R package (Pliner et al., 2018). We ran preprocess_cds()
476 function on the binarized peak by cell matrix with method='LSI', followed by
477 reducedDims() function to learn a vector representation for each cell. PCs whose
478 Pearson correlations with sequencing depth>0.5 are removed.

479 **5.5.3 cisTopic**

480 We used cisTopic via its R package (Bravo González-Blas et al., 2019). We ran
481 runCGSModels() function on the binarized peak by cell matrix with a range of
482 topic numbers (2, 5, 10, 20, 30, 40, 50, 60, 80 and 100) for 200 iterations with
483 burn in periods of 120. For comparison with scBasset, we reported the cisTopic
484 models with the best cell embedding performance.

485 **5.5.4 SCALE**

486 We used SCALE via its command line tool ([https://github.com/jsxlei/](https://github.com/jsxlei/SCALE)
487 SCALE) with parameters -x 0.05 and -min_peaks 500 to filter low quality peaks

488 and cells to avoid exploding gradients (Xiong et al., 2019). We ran SCALE with
489 a range of latent sizes (10, 16, 32, 64) and found that the default latent size
490 of 10 gives the best cell embedding performance. We also added the `-impute`
491 option allowing SCALE to estimate denoised accessibility values.

492 **5.5.5 chromVAR**

493 We used ChromVAR via its R package (Schep et al., 2017). We first created a
494 summarized experiment object from the binary peak by cell matrix, followed by
495 `addGCBias()` using the corresponding genome build. We featurize the sequences
496 into motif space using Jaspas motifs or k-mer space using 6-mers. Next, we
497 computed the deviation z-score matrices for motif and k-mer matches. For
498 each of chromVAR-motif or chromVAR-kmer, we performed PCA on the motif
499 deviation score matrix with 10, 20, 30, 40, 50, 60, 70, 80, 90, 100 PCs and
500 reported the best cell embedding performance to compare to scBasset.

501 When using chromVAR for TF activity inference, we ran chromVAR motif
502 match using CIS-BP motifs instead of the default Jaspas motifs for a fair
503 comparison with scBasset. Then we computed deviation z-scores as previously
504 described.

505 **5.6 Cell embedding evaluation**

506 Adjusted rand index (ARI): We evaluated learned cell embeddings in the Buen-
507 rostro2018 dataset by comparing the clustering to the ground-truth cell type
508 labels. We first built a nearest neighbor graph using scanpy with default
509 `n_neighbors` of 15. Then we followed a previous study to tune for a resolu-
510 tion that outputs 10 clusters (Chen et al., 2019). Finally, we compared the
511 clustering outcome to the ground-truth cell type labels using ARI.

512 Label score: We evaluated the learned cell embeddings using label score for
513 all three datasets. For a given nearest neighbor graph, label score quantifies
514 what percentage of each cell’s neighbors share its same label in a given neigh-
515 borhood. For each cell embedding method, we computed label score across a
516 neighborhood of 10, 50 and 100. Since the ground-truth cell types for the mul-
517 tiome datasets are unknown, we used cluster identifiers from scRNA-seq Leiden
518 clustering as cell type labels.

519 Neighbor score: We evaluated the learned cell embeddings using neighbor
520 score for the 10x multiome datasets. For a 10x multiome dataset, we built in-
521 dependent nearest neighbor graphs from the scRNA (using scVI) and scATAC
522 (using the cell embedding method we want to evaluate) and quantified the per-
523 centage of each cell’s neighbors that are shared between the two graphs across
524 neighborhoods of size 10, 50 and 100.

525 **5.7 Batch correction evaluation**

526 We evaluated scBasset-BC on additional scATAC datasets from mixed PBMC
527 populations from 10x PBMC multiome chemistry (downloaded from <https://>

528 cf.10xgenomics.com/samples/cell-arc/1.0.0/pbmc_granulocyte_sorted_
529 [10k/](https://cf.10xgenomics.com/samples/cell-arc/1.0.0/pbmc_granulocyte_sorted_10k/)) and 10x PBMC nextgem chemistry ([https://cf.10xgenomics.com/samples/](https://cf.10xgenomics.com/samples/cell-atac/2.0.0/atac_pbmc_10k_nextgem/)
530 [cell-atac/2.0.0/atac_pbmc_10k_nextgem/](https://cf.10xgenomics.com/samples/cell-atac/2.0.0/atac_pbmc_10k_nextgem/)). We generated a shared atlas of
531 21,017 peaks from the two datasets by resizing the 10x peak calls from the two
532 datasets to 1000bp and took the intersection. We subsampled 2,000 cells from
533 each dataset and merged them over the shared atlas. We ran scBasset-BC with
534 hyperparameters $\lambda_1=1e-6$ and $\lambda_2=0$.

535 5.8 Denoising evaluation

536 To compute a denoised and normalized accessibility across cells for a query peak
537 with scBasset, we ran a forward pass on the input DNA sequence to compute
538 the latent embedding for the peak. Then we generate the normalized acces-
539 sibility across all cells through dot product of the peak embedding with the
540 weight matrix of the final layer. Notice that since sequencing depth informa-
541 tion is entirely captured by the intercept vector of the final layer, we excluded
542 the intercept term so that scBasset generates denoised profiles normalized for
543 sequencing depth.

544 Our evaluation is based on the hypothesis that effective denoising would
545 improve the correlation between accessibility at genes' promoters and the genes'
546 expression in the multiome measurements (Granja et al., 2021; Pliner et al.,
547 2018). For each gene, we computed a gene accessibility score by averaging
548 the accessibility values for peaks at the gene's promoter (± 2 kb from TSS). We
549 evaluated denoising performance by computing the Pearson correlation between
550 the gene accessibility score and gene expression (after scVI denoising) across all
551 genes for each individual cell.

552 Alternatively, we also evaluated scBasset accessibility denoising for consis-
553 tency between differential expression and differential accessibility. We performed
554 differential gene expression on scVI gene expression for each cell type cluster
555 versus the rest. We also performed differential accessibility analysis on gene
556 accessibility scores for each cell type cluster versus the rest. Then we evaluated
557 performance by computing the Pearson correlation between the gene accessibil-
558 ity score \log_2 FC and gene expression \log_2 FC across all genes for each cell type
559 cluster.

560 5.9 Integration evaluation

561 In order to evaluate integration performance, we treated the 10x multiome
562 scRNA and scATAC profiles as originated from two independent experiments.
563 We summarize the accessibility profile to a gene level by computing gene acces-
564 sibility score as described above and integrated the scRNA and scATAC data by
565 embedding them into a shared space using Seurat `FindTransferAnchors()` and
566 `TransferData()` functions (Stuart et al., 2019).

567 In order to quantify the integration performance, we measured a "multiome
568 rank distance" R_c between the RNA embedding and the ATAC embedding of
569 each cell c . We use R_{rna} to represent the ranking of the Euclidean distance

570 between RNA embedding and ATAC embedding of cell c among all neighbors
571 of c 's RNA embedding, and R_{atac} to represent the ranking of the same distance
572 among all neighbors of c 's ATAC embedding. R_c is computed as the average of
573 R_{rna} and R_{atac} .

574 **5.10 Motif injection**

575 We performed motif injection on scBasset to compute a TF activity score for
576 each TF for each cell. Specifically, we first generated 1000 genomic background
577 sequences by performing dinucleotide shuffling of 1000 randomly sampled peaks
578 from the atlas using `fasta_ushuffle` (Jiang et al., 2008). For each TF in the motif
579 database, we sampled a motif sequence from the position weight matrix (PWM)
580 and inserted into the center of each of the genomic background sequences. We
581 ran forward passes through the model for both the motif-injected sequences
582 and background sequences to predict normalized accessibility across all cells.
583 We took the difference in predicted accessibility between the motif-injected se-
584 quences and background sequences as the motif influence for each sequence. We
585 averaged this influence score across all 1000 sequences for each cell to generate
586 a cell level prediction of raw TF activity. Finally, we z-score normalized the raw
587 TF activities to generate the final TF activity predictions across all cells.

588 We used CIS-BP 1.0 single species DNA database motifs downloaded from
589 <https://meme-suite.org/meme/db/motifs> for our motif analysis (Weirauch
590 et al., 2014).

591 **5.11 *In silico* saturation mutagenesis**

592 We performed *in silico* saturation mutagenesis (ISM) to compute the importance
593 scores of all single nucleotides on a sequence of interest. For each position, we ran
594 three scBasset forward passes, each time mutating the reference nucleotide to an
595 alternative. For each mutation, we compared the accessibility prediction to the
596 prediction with the reference nucleotide to compute the change in accessibility
597 for each cell. We normalized the ISM scores for the four nucleotides at each
598 position such that they sum to zero. We then took the normalized ISM score
599 at the reference nucleotide as the importance score for that position.

600 **6 Code Availability**

601 Code for training and using scBasset model can be found at: [https://github.](https://github.com/calico/scBasset)
602 [com/calico/scBasset](https://github.com/calico/scBasset).

603 **7 Acknowledgements**

604 We thank Vikram Agarwal, Jacob Kimmel and Majed Mohamed for feedback
605 on the manuscript. We thank Sarah Spock for feedback on the code. We also
606 thank Nick Bernstein and Ashlesha Odak for helpful discussions.

⁶⁰⁷ **8 Author Contributions**

⁶⁰⁸ D.R.K. conceived the project. H.Y. and D.R.K. developed the model. H.Y.
⁶⁰⁹ performed the analysis. H.Y. and D.R.K prepared the manuscript.

⁶¹⁰ **9 Competing Interests**

⁶¹¹ H.Y. and D.R.K. are paid employees of Calico Life Sciences.

612 References

- 613 Agarwal, V. and Shendure, J. (2020). Predicting mRNA Abundance Directly
614 from Genomic Sequence Using Deep Convolutional Neural Networks. *Cell*
615 *Reports*.
- 616 Ashuach, T., Reidenbach, D. A., Gayoso, A., and Yosef, N. (2021). PeakVI:
617 A Deep Generative Model for Single Cell Chromatin Accessibility Analysis.
618 *bioRxiv*.
- 619 Avsec, A., Weilert, M., Shrikumar, A., Krueger, S., Alexandari, A., Dalal, K.,
620 Fropf, R., McAnany, C., Gagneur, J., Kundaje, A., and Zeitlinger, J. (2021).
621 Base-resolution models of transcription-factor binding reveal soft motif syn-
622 tax. *Nature Genetics*.
- 623 Bravo González-Blas, C., Minnoye, L., Papasokrati, D., Aibar, S., Hulselmans,
624 G., Christiaens, V., Davie, K., Wouters, J., and Aerts, S. (2019). cisTopic:
625 cis-regulatory topic modeling on single-cell ATAC-seq data. *Nature Methods*.
- 626 Buenrostro, J. D., Corces, M. R., Lareau, C. A., Wu, B., Schep, A. N., Aryee,
627 M. J., Majeti, R., Chang, H. Y., and Greenleaf, W. J. (2018). Integrated
628 Single-Cell Analysis Maps the Continuous Regulatory Landscape of Human
629 Hematopoietic Differentiation. *Cell*.
- 630 Chen, H., Lareau, C., Andreani, T., Vinyard, M. E., Garcia, S. P., Clement,
631 K., Andrade-Navarro, M. A., Buenrostro, J. D., and Pinello, L. (2019). As-
632 sessment of computational methods for the analysis of single-cell ATAC-seq
633 data. *Genome Biology*.
- 634 Cusanovich, D. A., Hill, A. J., Aghamirzaie, D., Daza, R. M., Pliner, H. A.,
635 Berletch, J. B., Filippova, G. N., Huang, X., Christiansen, L., DeWitt, W. S.,
636 Lee, C., Regalado, S. G., Read, D. F., Steemers, F. J., Disteche, C. M., Trap-
637 nell, C., and Shendure, J. (2018). A Single-Cell Atlas of In Vivo Mammalian
638 Chromatin Accessibility. *Cell*.
- 639 de Boer, C. G. and Regev, A. (2018). BROCKMAN: Deciphering variance in
640 epigenomic regulators by k-mer factorization. *BMC Bioinformatics*.
- 641 Granja, J. M., Corces, M. R., Pierce, S. E., Bagdatli, S. T., Choudhry, H.,
642 Chang, H. Y., and Greenleaf, W. J. (2021). ArchR is a scalable software
643 package for integrative single-cell chromatin accessibility analysis. *Nature*
644 *Genetics*.
- 645 Jiang, M., Anderson, J., Gillespie, J., and Mayne, M. (2008). uShuffle: A useful
646 tool for shuffling biological sequences while preserving the k-let counts. *BMC*
647 *Bioinformatics*.
- 648 Kelley, D. R., Reshef, Y. A., Bileschi, M., Belanger, D., McLean, C. Y., and
649 Snoek, J. (2018). Sequential regulatory activity prediction across chromo-
650 somes with convolutional neural networks. *Genome Research*.

- 651 Kelley, D. R., Snoek, J., and Rinn, J. L. (2016). Basset: Learning the regula-
652 tory code of the accessible genome with deep convolutional neural networks.
653 *Genome Research*.
- 654 Kiss, M., Czimmerer, Z., Nagy, G., Bieniasz-Krzywiec, P., Ehling, M., Pap,
655 A., Poliska, S., Boto, P., Tzerpos, P., Horvath, A., Kolostyak, Z., Daniel,
656 B., Szatmari, I., Mazzone, M., and Nagy, L. (2017). Retinoid X receptor
657 suppresses a metastasis-promoting transcriptional program in myeloid cells
658 via a ligand-insensitive mechanism. *Proceedings of the National Academy of
659 Sciences of the United States of America*.
- 660 Lal, A., Chiang, Z. D., Yakovenko, N., Duarte, F. M., Israeli, J., and Buenrostro,
661 J. D. (2021). Deep learning-based enhancement of epigenomics data with
662 AtacWorks. *Nature Communications*.
- 663 Li, Q., Peterson, K. R., Fang, X., and Stamatoyannopoulos, G. (2002). Locus
664 control regions. *Blood*.
- 665 Lopez, R., Regier, J., Cole, M. B., Jordan, M. I., and Yosef, N. (2018). Deep
666 generative modeling for single-cell transcriptomics. *Nature Methods*.
- 667 Medvedovic, J., Ebert, A., Tagoh, H., and Busslinger, M. (2011). Pax5: A
668 Master Regulator of B Cell Development and Leukemogenesis. In *Advances
669 in Immunology*.
- 670 Miao, Z., Balzer, M. S., Ma, Z., Liu, H., Wu, J., Shrestha, R., Aranyi, T., Kwan,
671 A., Kondo, A., Pontoglio, M., Kim, J., Li, M., Kaestner, K. H., and Susztak,
672 K. (2021). Single cell regulatory landscape of the mouse kidney highlights
673 cellular differentiation programs and disease targets. *Nature Communications*.
- 674 Pliner, H. A., Packer, J. S., McFaline-Figueroa, J. L., Cusanovich, D. A., Daza,
675 R. M., Aghamirzaie, D., Srivatsan, S., Qiu, X., Jackson, D., Minkina, A.,
676 Adey, A. C., Steemers, F. J., Shendure, J., and Trapnell, C. (2018). Ci-
677 cero Predicts cis-Regulatory DNA Interactions from Single-Cell Chromatin
678 Accessibility Data. *Molecular Cell*.
- 679 Satpathy, A. T., Granja, J. M., Yost, K. E., Qi, Y., Meschi, F., McDermott,
680 G. P., Olsen, B. N., Mumbach, M. R., Pierce, S. E., Corces, M. R., Shah, P.,
681 Bell, J. C., Jhuttu, D., Nemecek, C. M., Wang, J., Wang, L., Yin, Y., Giresi,
682 P. G., Chang, A. L. S., Zheng, G. X., Greenleaf, W. J., and Chang, H. Y.
683 (2019). Massively parallel single-cell chromatin landscapes of human immune
684 cell development and intratumoral T cell exhaustion. *Nature Biotechnology*.
- 685 Schep, A. N., Wu, B., Buenrostro, J. D., and Greenleaf, W. J. (2017). Chrom-
686 VAR: Inferring transcription-factor-associated accessibility from single-cell
687 epigenomic data. *Nature Methods*.
- 688 Shrikumar, A., Tian, K., Avsec, A., Shcherbina, A., Banerjee, A., Sharmin, M.,
689 Nair, S., and Kundaje, A. (2018). Technical Note on Transcription Factor
690 Motif Discovery from Importance Scores (TF-MoDISco) version 0.5.6.5.

- 691 Stuart, T., Butler, A., Hoffman, P., Hafemeister, C., Papalexi, E., Mauck,
692 W. M., Hao, Y., Stoeckius, M., Smibert, P., and Satija, R. (2019). Com-
693 prehensive Integration of Single-Cell Data. *Cell*.
- 694 Tallack, M. R., Whittington, T., Yuen, W. S., Wainwright, E. N., Keys, J. R.,
695 Gardiner, B. B., Nourbakhsh, E., Cloonan, N., Grimmond, S. M., Bailey,
696 T. L., and Perkins, A. C. (2010). A global role for KLF1 in erythropoiesis
697 revealed by ChIP-seq in primary erythroid cells. *Genome Research*.
- 698 Thurman, R. E., Rynes, E., Humbert, R., Vierstra, J., Maurano, M. T., Hau-
699 gen, E., Sheffield, N. C., Stergachis, A. B., Wang, H., Vernot, B., Garg, K.,
700 John, S., Sandstrom, R., Bates, D., Boatman, L., Canfield, T. K., Diegel,
701 M., Dunn, D., Ebersol, A. K., Frum, T., Giste, E., Johnson, A. K., Johnson,
702 E. M., Kuttyavin, T., Lajoie, B., Lee, B. K., Lee, K., London, D., Lotakis,
703 D., Neph, S., Neri, F., Nguyen, E. D., Qu, H., Reynolds, A. P., Roach, V.,
704 Safi, A., Sanchez, M. E., Sanyal, A., Shafer, A., Simon, J. M., Song, L.,
705 Vong, S., Weaver, M., Yan, Y., Zhang, Z., Zhang, Z., Lenhard, B., Tewari,
706 M., Dorschner, M. O., Hansen, R. S., Navas, P. A., Stamatoyannopoulos, G.,
707 Iyer, V. R., Lieb, J. D., Sunyaev, S. R., Akey, J. M., Sabo, P. J., Kaul, R.,
708 Furey, T. S., Dekker, J., Crawford, G. E., and Stamatoyannopoulos, J. A.
709 (2012). The accessible chromatin landscape of the human genome. *Nature*.
- 710 Tjarnberg, A., Mahmood, O., Jackson, C. A., Saldi, G. A., Cho, K., Christia-
711 en, L. A., and Bonneau, R. A. (2021). Optimal tuning of weighted kNN-
712 And diffusion-based methods for denoising single cell genomics data. *PLoS*
713 *Computational Biology*.
- 714 Tuan, D., Solomon, W., Li, Q., and London, I. M. (1985). The "beta-like-globin"
715 gene domain in human erythroid cells. *Proceedings of the National Academy*
716 *of Sciences of the United States of America*, 82(19):6384–6388.
- 717 Weirauch, M., Yang, A., Albu, M., Cote, A. G., Montenegro-Montero, A.,
718 Drewe, P., Najafabadi, H., Lambert, S., Mann, I., Cook, K., Zheng, H.,
719 Goity, A., vanÁ Bakel, H., Lozano, J.-C., Galli, M., Lewsey, M. G., Huang,
720 E., Mukherjee, T., Chen, X., Reece-Hoyes, J., Govindarajan, S., Shaulsky, G.,
721 Walhout, A., Bouget, F.-Y., Ratsch, G., Larrondo, L., Ecker, J., and Hughes,
722 T. (2014). Determination and Inference of Eukaryotic Transcription Factor
723 Sequence Specificity. *Cell*, 158(6):1431–1443.
- 724 Xiong, L., Xu, K., Tian, K., Shao, Y., Tang, L., Gao, G., Zhang, M., Jiang, T.,
725 and Zhang, Q. C. (2019). SCALE method for single-cell ATAC-seq analysis
726 via latent feature extraction. *Nature Communications*.
- 727 Zhou, J., Theesfeld, C. L., Yao, K., Chen, K. M., Wong, A. K., and Troyanskaya,
728 O. G. (2018). Deep learning sequence-based ab initio prediction of variant
729 effects on expression and disease risk. *Nature Genetics*.
- 730 Zhou, J. and Troyanskaya, O. G. (2015). Predicting effects of noncoding variants
731 with deep learning-based sequence model. *Nature Methods*, 12(10).

⁷³² **10 Supplementary Figures**

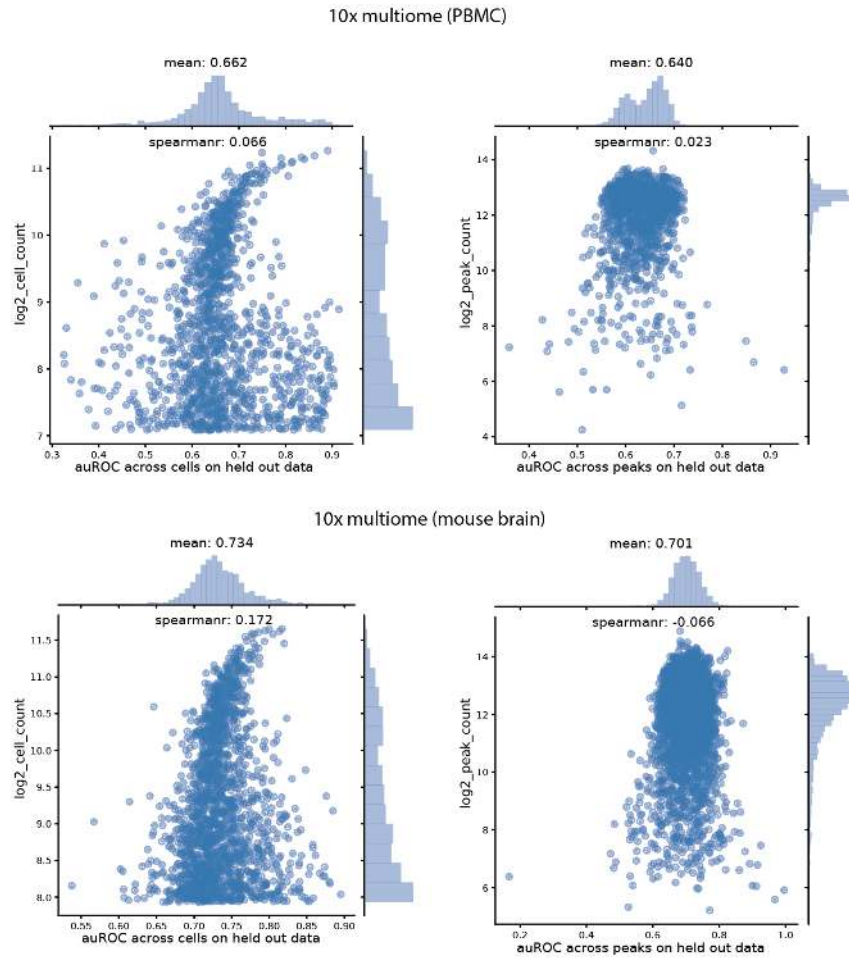


Figure S1: auROC on held-out peaks for 10x multiome PBMC and mouse brain datasets. Top, scBasset prediction performance on held-out peaks evaluated by auROC per peak (left) and by auROC per cell (right) for 10x multiome PBMC dataset. Bottom, scBasset prediction performance on held-out peaks evaluated by auROC per peak (left) and by auROC per cell (right) for 10x multiome mouse brain dataset.

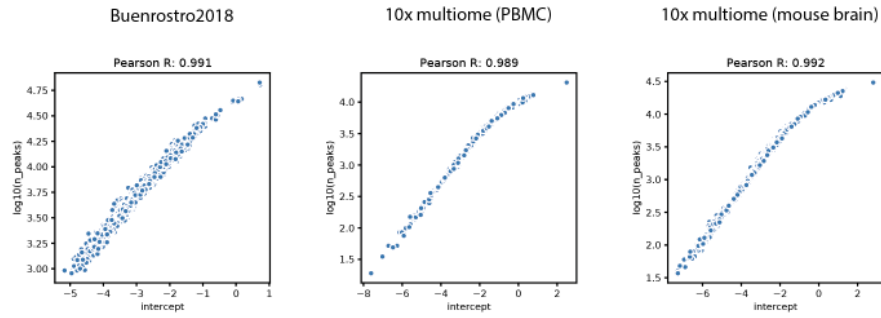


Figure S2: Correlations of final layer intercepts with sequencing depth (\log_{10} UMI) for Buenrostro2018, 10x multiome PBMC and 10x multiome mouse brain datasets (from left to right).

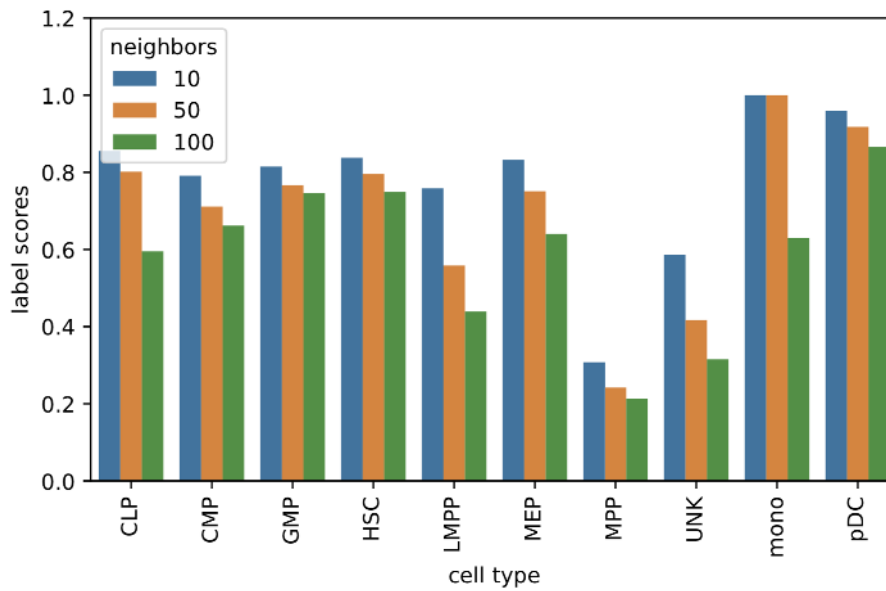


Figure S3: scBasset cell embedding performance as evaluated by label scores for each cell type with a neighborhood of 10, 50 and 100.

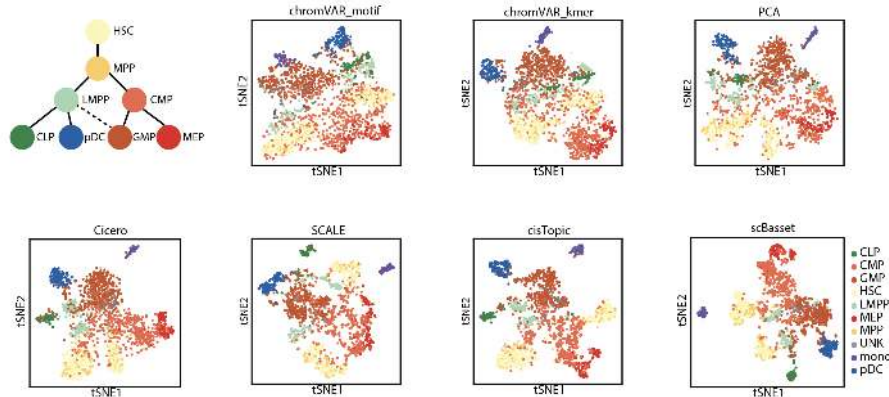


Figure S4: t-SNE visualization of different cell embedding methods, including: chromVAR motif, chromVAR kmer ($k=6$), PCA, cicero (LSI), SCALE, cisTopic and scBasset.

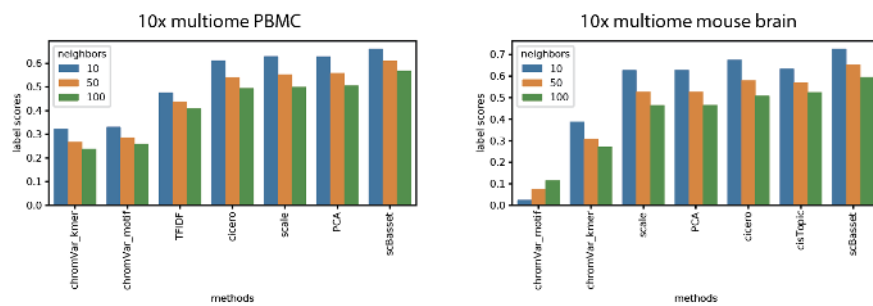


Figure S5: Performance comparison of different cell embedding methods as evaluated by label scores for 10x multiome PBMC (left) and mouse brain (right) datasets.

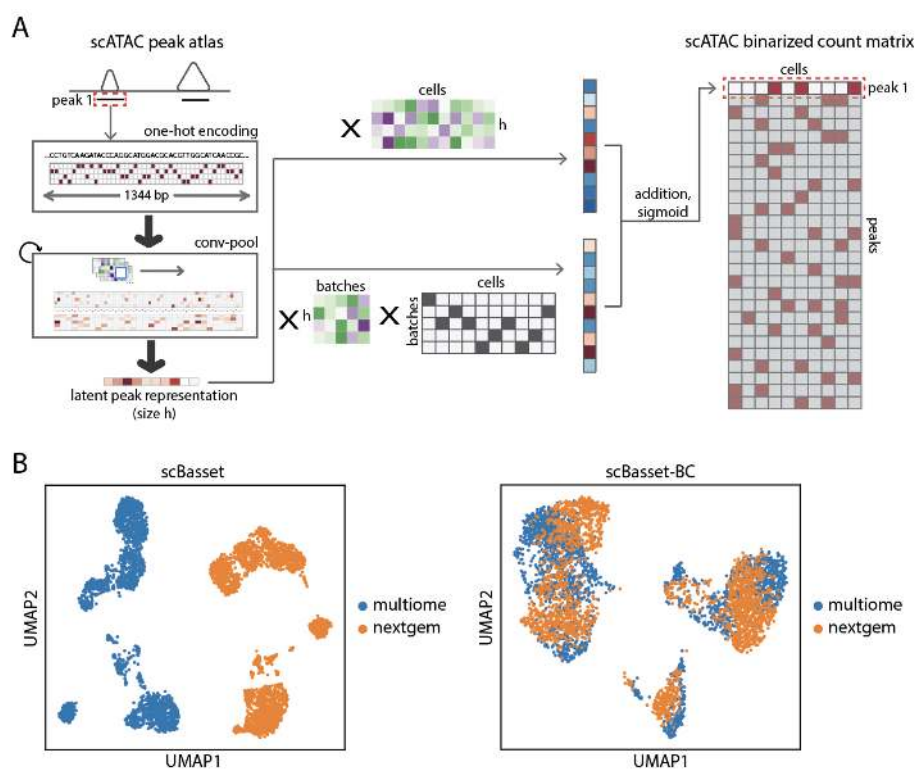


Figure S6: A, Model architecture of scBasset-BC. B) UMAP embeddings of mixed PBMC populations from 10x multiome scATAC and 10x nextgem scATAC chemistries. Left figure shows the embeddings learned by scBasset model. Right figure shows the embeddings learned by scBasset-BC model, where batch is encoded as a covariate.

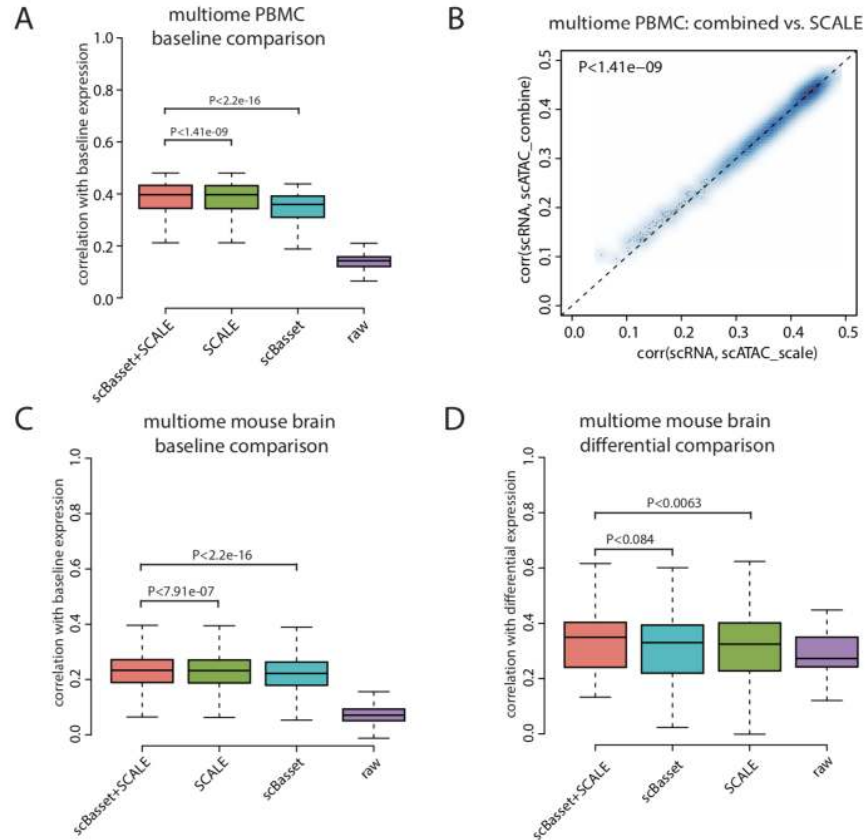


Figure S7: Additional denoising results for 10x multiome datasets. A) Comparison of denoising performance on the multiome PBMC dataset between raw data, scBasset, SCALE, and scBasset+SCALE combined, evaluated by correlation between baseline gene accessibility score and baseline gene expression. B) A scatterplot showing a closer look at the performance comparison between scBasset+SCALE (y-axis) versus SCALE on multiome PBMC dataset, evaluated by correlation between baseline gene accessibility score and baseline gene expression. C) Comparison of denoising performance on the multiome mouse brain dataset between raw data, scBasset, SCALE, and scBasset+SCALE combined, evaluated by correlation between baseline gene accessibility score and baseline gene expression. D) Comparison of denoising performance on multiome mouse brain dataset between raw data, scBasset, SCALE, and scBasset+SCALE combine, evaluated by consistency in differential expression log₂FC and differential accessibility log₂FC. We performed Wilcoxon signed rank tests for performance comparisons.

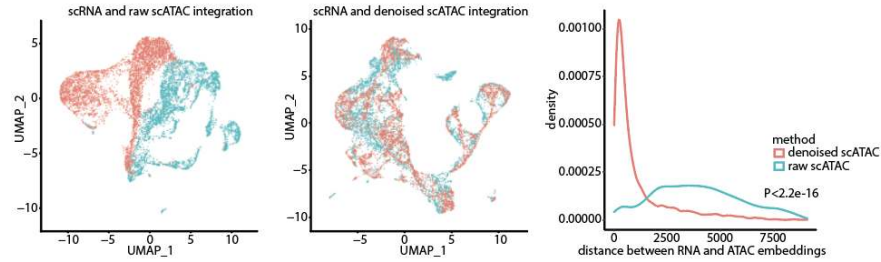


Figure S8: Integration results for the 10x multiome mouse brain dataset. Left, RNA (blue) and raw ATAC (red) profile embeddings after integration. Middle, RNA (blue) and denoised ATAC (red) profile embeddings after integration. Right, distribution of the relative distances (Methods) between each cell's RNA and ATAC embeddings after integration when integrating with raw ATAC profiles (blue) or denoised ATAC profiles (red). We performed Wilcoxon signed rank test for performance comparison.

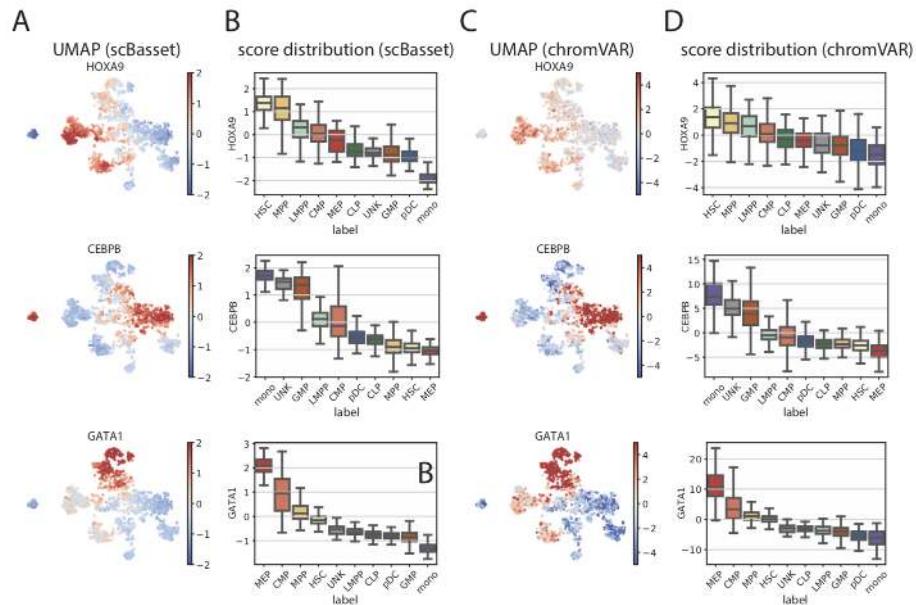


Figure S9: Motif activity inference using scBasset and chromVAR on the Buenrostro 2018 dataset for known regulators. A) UMAPs showing scBasset-predicted TF activity. B) Boxplots showing scBasset-predicted TF activity by cell type. C) UMAPs showing chromVAR-predicted TF activity. B) Boxplots showing chromVAR-predicted TF activity per cell type.

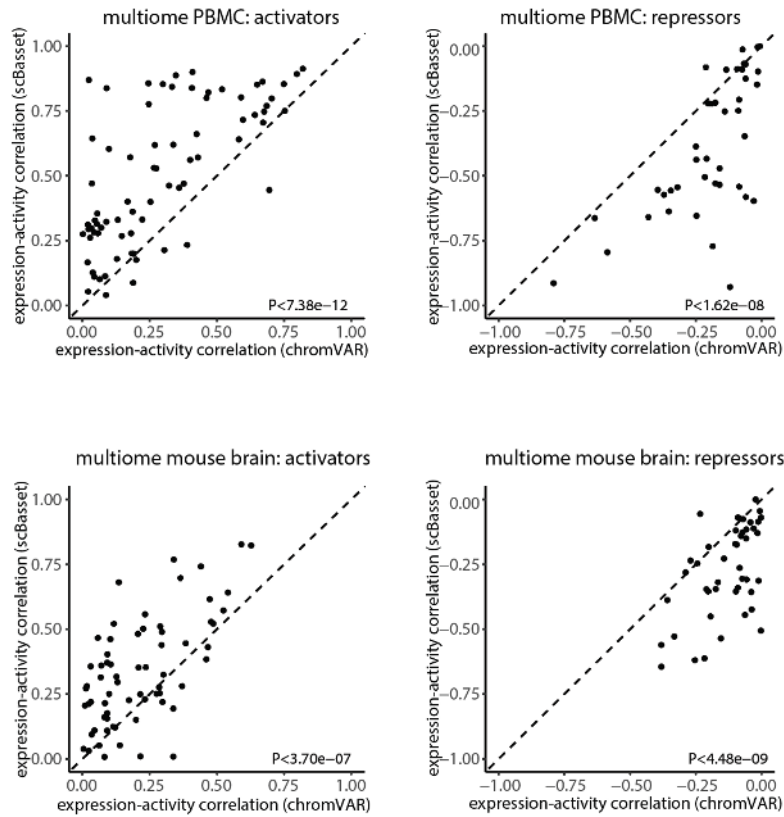


Figure S10: TF expression and TF activity correlation for the 10x multiome datasets. Scatterplots of correlations between chromVAR-inferred activity and expression (x-axis) versus correlations of scBasset-inferred TF activity and expression (y-axis) for activating TFs (left) and repressive TFs (right) in the 10x multiome PBMC (top) and 10x multiome mouse brain (bottom). Activating TFs are TFs which both scBasset and chromVAR agree on a positive correlation between TF expression and activity. Repressive TFs are TFs which both scBasset and chromVAR agree on a negative correlation between TF expression and activity.

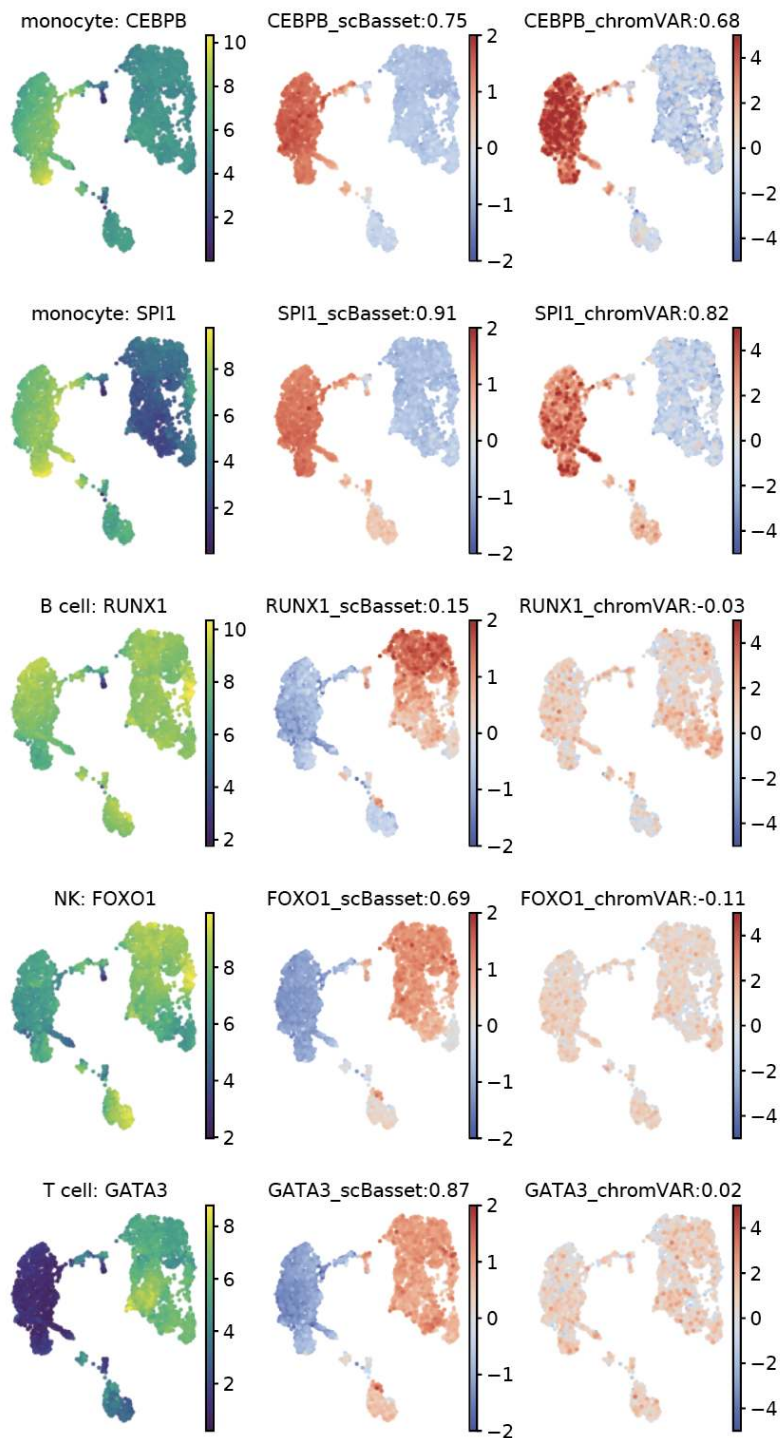


Figure S11: Motif activity inference using scBasset and chromVAR on the 10x multiome PBMC data. UMAP visualization of TF expression (left), scBasset TF activity (middle), and chromVAR TF activity (right) for additional known PBMC regulators. Pearson correlation between inferred TF activity and expression are shown in the titles.

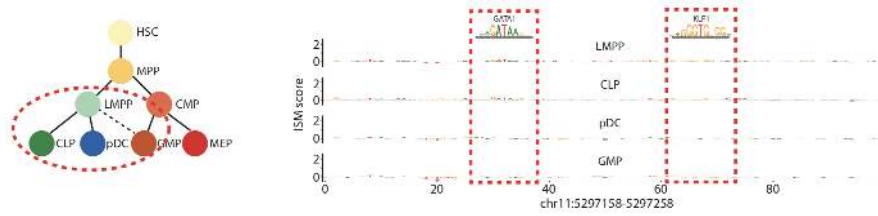


Figure S12: ISM scores for β -globin enhancer at chr11:5297158-5297258 for cells in LMPP, CLP, pDC and GMP cell types. Sequences that match GATA1 and KLF1 motifs are highlighted in red boxes.

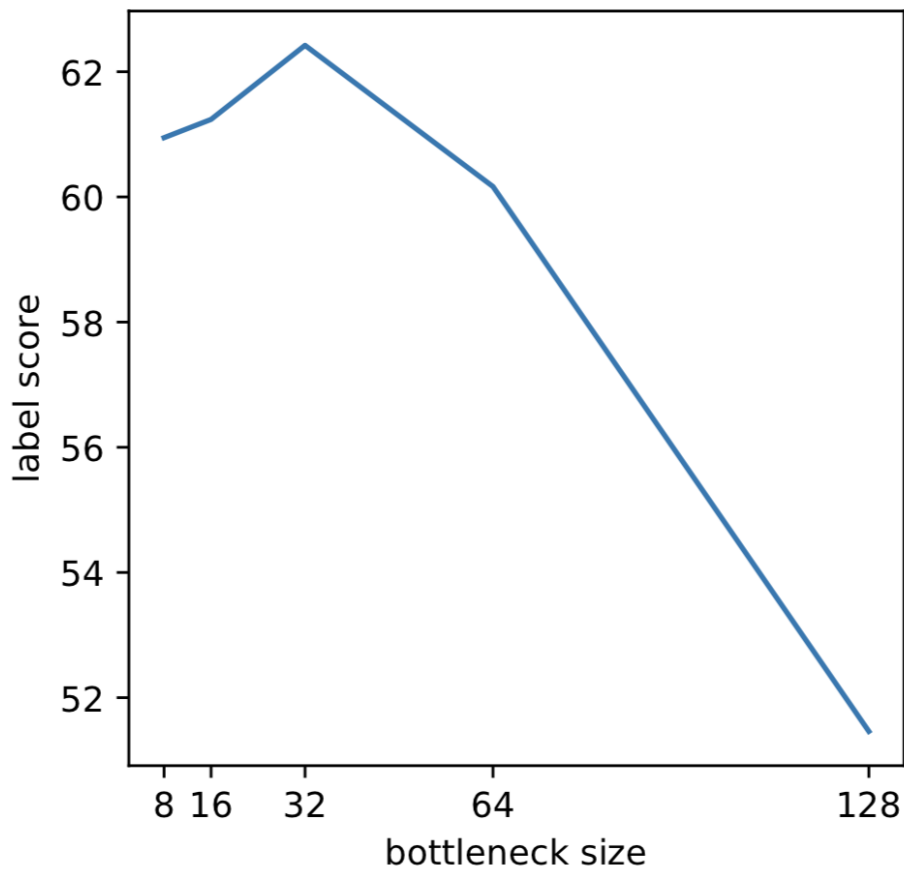


Figure S13: Label scores as a function of scBasset bottleneck layer size in Buenroostro2018 dataset.

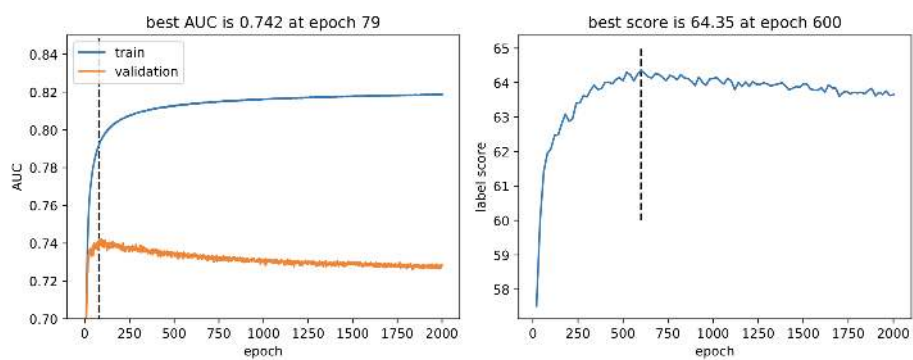


Figure S14: Left, training (blue) and validation auROCs (red) per epoch for the Buenrostro2018 dataset. Right, label scores per epoch.