1    **SIMBA: SIngle-cell eMBedding Along with features**

2

3

4    Huidong Chen[1, 2, 3], Jayoung Ryu[1, 2, 4], Michael E. Vinyard[1, 2, 3, 5], Adam Lerer[6]+, Luca Pinello[1, 2, 3]+

5

6    1. Molecular Pathology Unit, Center for Cancer Research, Massachusetts General Hospital,
7    Charlestown, MA, USA
8    2. Department of Pathology, Harvard Medical School, Boston, MA, USA
9    3. Broad Institute of Harvard and MIT, Cambridge, MA, USA
10   4. Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA
11   5. Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA, USA
12   6. Facebook AI Research

13

14   +Corresponding author

15

16   **Abstract**

17

18   Recent advances in single-cell omics technologies enable the individual and joint profiling of
19   cellular measurements. Currently, most single-cell analysis pipelines are cluster-centric and
20   cannot explicitly model the interactions between different feature types. In addition, single-cell
21   methods are generally designed for a particular task as distinct single-cell problems are
22   formulated differently. To address these current shortcomings, we present *SIMBA*, a graph
23   embedding method that jointly embeds single cells and their defining features, such as genes,
24   chromatin accessible regions, and transcription factor binding sequences into a common latent
25   space. By leveraging the co-embedding of cells and features, SIMBA allows for the study of
26   cellular heterogeneity, clustering-free marker discovery, gene regulation inference, batch effect
27   removal, and omics data integration. SIMBA has been extensively applied to scRNA-seq,
28   scATAC-seq, and dual-omics data. We show that SIMBA provides a single framework that allows
29   diverse single-cell analysis problems to be formulated in a unified way and thus simplifies the
30   development of new analyses and integration of other single-cell modalities. SIMBA is
31   implemented as an efficient, comprehensive, and extensible Python library (https://simba-
32   bio.readthedocs.io) for the analysis of single-cell omics data using graph embedding.

33

34   **Introduction**

35

36   Technology to profile single cells has advanced to several molecular modalities, dramatically
37   advancing our ability to characterize cell states as well as discover key molecular machinery
38   that underlies both development and disease. Individual cells are now measured using multiple
39   molecular modalities, simultaneously. At the same time, single-cell experiments have scaled
40   such that tens of thousands of cells can be routinely profiled. The emergence of single-cell
41   multi-omics technologies allows for the measurements of multiple cellular layers, including
42   genomics, epi-genomics, transcriptomics, and proteomics. Such assays have pioneered an
43   avenue toward a better understanding of the interplay between layers as they jointly define cell

44    states based on diverse genomic and molecular features including genes, regulatory elements,
45    and transcription factors. While single-cell multi-omic assays have quickly evolved towards the
46    incorporation of additional modalities with increasing resolution, harnessing their full potential
47    has posed several significant computational challenges.
48
49    Many single-cell computational methods have been developed for the analysis of one modality
50    (e.g., scRNA-seq or scATAC-seq analysis) [1-4]. Common to these methods is a workflow that
51    includes routine steps such as feature selection, dimension reduction, clustering, and
52    differential feature detection. These "cluster-centric" analysis methods rely on accurately
53    defined clustering solutions to discover meaningful and informative marker features.
54    Unfortunately, clustering solutions may range widely within the space of the user-defined
55    clustering resolution (number of clusters) and the chosen clustering algorithm. These
56    parameters may markedly influence the resulting cluster assignment and clusters may not
57    always correspond to the correct or intended cell populations, thereby leading to inconsistent
58    and potentially misleading biological annotations[5]. Although initial efforts have been made
59    recently to develop clustering-free approaches to discover informative genes, they are
60    specifically designed for extracting gene signatures [6, 7] or identifying perturbations between
61    experimental conditions[8] from scRNA-seq data, and are therefore limited to single-modality
62    and single-task analysis.
63
64    In addition to single-batch/modality analysis, approaches have also been proposed for multi-
65    batch and cross-modality analysis, such as multimodal analysis (distinct cellular parameters are
66    measured in the same cell)[9], batch correction (the same cellular parameter is measure in
67    different batches) [10-12], and integration of multi-omics datasets (distinct cellular parameters are
68    measured in different cells)[11, 12]. These approaches play a critical role in removing batch effects
69    that confound true biological variation, improving the characterization of cell states by
70    leveraging the unique strengths of each assay, and providing insights into the complex
71    mechanisms of gene regulation. However, these tasks are formulated differently from those in
72    single-batch/modality settings and thus require development of new dedicated analysis
73    techniques. Also, while multiple types of cellular features might be present, the relation
74    between features cannot be exploited directly by most current methods. Furthermore, similar
75    to single-batch/modality analysis methods, these methods identify marker features based on
76    groups of cells obtained by clustering and therefore are limited to clustering solutions.
77    Additionally, instead of directly identifying marker features in the integrated space, most batch
78    correction/multi-omics integration methods need to first detect marker features in the
79    uncorrected/unintegrated original space of each batch/modality independently, and then
80    combine them, thus resulting in potentially inconsistent interpretations between
81    batches/modalities.
82
83    To overcome the limitations in both single-batch/modality analysis and multi-batch/cross-
84    modality analysis, we propose SIMBA (**SIngle-cell eMBedding Along with features**), a versatile
85    single-cell embedding method that co-embeds cells and features into a shared latent space, in
86    which various types of tasks can be performed based on the proximity between entities
87    including cells and features such as genes, peaks, and DNA sequences. Unlike existing methods

88      that require featurization of cells, SIMBA directly encodes the cell-feature or feature-feature
89      relations into a large multi-entity graph. For each task, SIMBA constructs a graph, wherein
90      differing entities (i.e., cells and features) are represented as nodes and relations between these
91      entities are encoded as edges. Once the graph is constructed, SIMBA then applies a multi-entity
92      graph embedding algorithm derived from social networking technologies as well as a Softmax-
93      based transformation to embed the nodes/entities of the graph into a common low-
94      dimensional space wherein cells and features can be analyzed based on their distance. Hence
95      SIMBA provides an information-rich embedding space containing cells and all the features,
96      serving as an informative database of entities. Depending on the task, we can define biological
97      queries on the "SIMBA database" by considering neighboring entities of either a cell (or cells) or
98      a feature (or features) at the individual-cell and individual-feature level (**Methods**). For
99      example, the query for a cell's neighboring features can be used to identify marker features
100     (e.g., marker genes or peaks) or to study the interaction between features (e.g., peak-gene)
101     while the query for features' neighboring cells can be used to annotate cells.
102
103     By formulating single-cell analyses as multi-entity graph embedding problems, we show SIMBA
104     can be used to solve popular single-cell tasks in a unified framework that would otherwise
105     require the development of distinct specialized approaches for each task, including: 1)
106     dimensionality reduction techniques for studying cellular states; 2) clustering-free marker
107     detection based on the similarity between single cells and features; 3) Single-cell multimodal
108     analysis and the study of gene regulation; 4) batch correction and omics integration analysis as
109     well as the simultaneous identification of marker features. SIMBA is adapted to these diverse
110     analysis tasks by simply modifying how the input graph is constructed from the relevant single-
111     cell data.
112
113     We extensively tested SIMBA in multiple scRNA-seq, scATAC-seq and dual-omics datasets
114     covering popular single-cell tasks including scRNA-seq analysis, scATAC-seq analysis, multimodal
115     analysis, batch correction, and multi-omics integration. We demonstrate that SIMBA learns the
116     joint low-dimensional representations of both cells and features and thus enables the ability to
117     simultaneously study cellular heterogeneity as well as proximity-based marker feature
118     detection or gene regulation inference in a clustering-free way. We also demonstrate that
119     SIMBA performs better than or comparably to current state-of-the-art methods specifically
120     developed for each task.
121
122     Importantly, we developed a scalable and comprehensive Python package that enables
123     seamless interaction between graph construction, training with PyTorch for graph embedding,
124     and post-training analysis. The SIMBA package not only provides a self-contained framework
125     that covers preprocessing, graph embedding, and visualization, but also is compatible with
126     popular single cell analysis tool Scanpy[2].  SIMBA with detailed documentation and extensive
127     tutorials is available at https://simba-bio.readthedocs.io.
128
129     We believe that SIMBA, as a broadly applicable approach for single cell omics study, not only
130     outperforms current cluster-centric analysis, but also will simplify the burden of developing

131  methods for new single-cell tasks and measurements, while increasing interpretability of
132  cellular mechanisms and functions.
133
## Results
135
### Overview of SIMBA
137  SIMBA is a single-cell embedding method with support for single- or multi- modality analyses
138  that embeds cells and their associated genomic features into a shared latent space, generating
139  interpretable and comparable embeddings of cells and features. It leverages recent graph
140  embedding techniques that have been successful in modeling complex and hierarchical
141  information present in natural languages, social networks, and other domains, as "knowledge
142  graphs". In our case, the graph encodes cells, different components of cellular regulatory
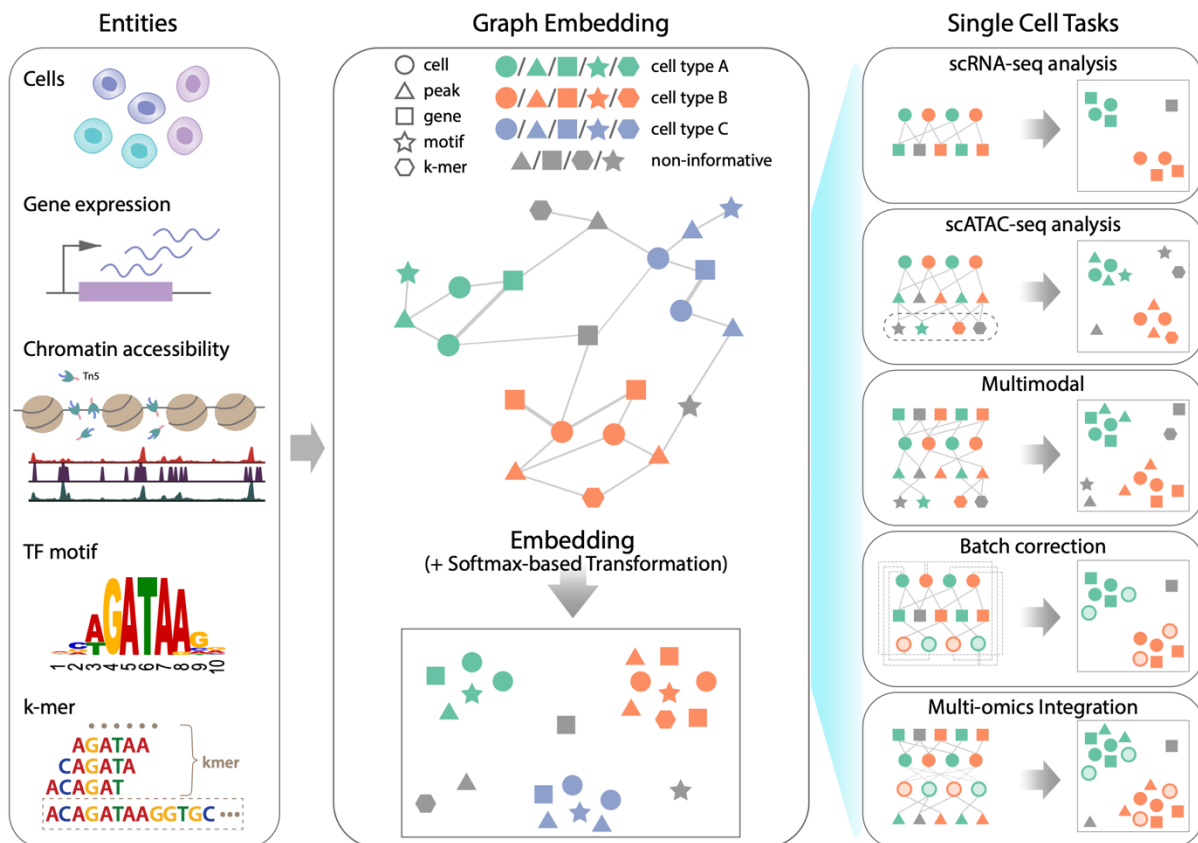143  circuits, and the relations between them.



144

145  **Figure1.** SIMBA framework overview. SIMBA co-embeds cells and various features
146  measured during single-cell experiments into a shared latent space to accomplish both
147  common tasks involved in single-cell data analysis as well as tasks, which remain as open
148  problems in single-cell genomics. **(Left)** Examples of possible biological entities may be
149  encoded by SIMBA including cells, gene expression measurements, chromatin accessible
150  regions, TF motifs, and k-mer sequences found in reads. **(Middle)** SIMBA embedding plot

151 with multiple types of entities into a low-dimensional space. All entities represented as
152 shapes (cell = circle, peak = triangle, gene = square, TF motif = star, k-mer = hexagon) are
153 colored by relevant cell type (green, orange, and blue in this example). Non-informative
154 features are colored dark grey. Within the graph, each entity is a node, and an edge
155 indicates a relation between entities (e.g., a gene is expressed in a cell, a chromatin region
156 is accessible in a cell, or a TF motif/k-mer is present within an open chromatin region,
157 etc.). Once connected in a graph, these entities may be embedded into a shared low-
158 dimensional space, with cell-type specific entities embedded in the same neighborhood
159 and non-informative features embedded elsewhere. **(Right)** Common single-cell analysis
160 tasks that may be accomplished using SIMBA.
161

162 SIMBA first encodes different types of entities such as cells, genes, open chromatin regions
163 (peaks or bins), transcription factor (TF) motifs, and *k*-mers (short sequences of a specific
164 length, *k*), into a single graph (**Fig. 1**, **Methods**) where each node represents an individual entity
165 and edges indicate relations between entities. For example, if a gene is expressed in a cell, an
166 edge is created between the gene and cell. The weight of this edge is determined by the gene
167 expression level.  Similarly, an edge is added between a cell and a chromatin region if the region
168 is open in this cell, or between a chromatin region and a TF motif if the TF motif is found in the
169 region.

170 Once the input graph is constructed, a low-dimensional representation of the graph nodes is
171 then computed using an unsupervised graph embedding method. This graph embedding
172 procedure leverages the PyTorch-BigGraph framework [13], which allows SIMBA to scale to
173 millions of cells (**Methods**). The obtained SIMBA space provides an intuitive way to study gene
174 regulation and the regulatory mechanisms underlying cell differentiation and specification. The
175 resulting joint embedding of cells and features not only reconstructs the heterogeneity of cells
176 but also allows for the discovery of the defining features for each individual cell without relying
177 on a clustering solution, separating cell-type specific features from the non-informative
178 features. In fact, the relationship between cells and features can be explored directly through
179 their mutual proximity in the SIMBA embedding as the distance between embedded nodes
180 reflects their edge probability, which is informative of the potential importance of a feature to a
181 cell and the interplay between features (**Methods**).

182 Therefore, cell-type-specific features such as marker genes, cis-regulatory elements can be
183 discovered without clustering in two different ways. When the labels of cells are known, marker
184 features can be identified as the neighboring features of cells by performing biological queries
185 (**Methods**). When these labels are unknown, marker features can be identified through
186 calculating the imbalance of edge probabilities between a feature and all cells using metrics
187 such as the Gini index (**Methods**).

188 Importantly, graph construction is inherently flexible, enabling SIMBA to be applied to a wide
189 variety of single-cell tasks. In the following sections, we demonstrate the application of SIMBA
190 to several popular single-cell tasks including scRNA-seq, scATAC-seq, multimodal analysis, batch

191  correction and multi-omics integration (**Fig. 1**). Extensions to additional tasks will become
192  readily apparent to the reader and are later discussed.
193
194
195  **SIMBA enables simultaneous learning of cellular heterogeneity and individual-cell-level**
196  **marker genes in scRNA-seq analysis**

197  Single-cell RNA sequencing (scRNA-seq) is the most robust and widely used measurement to
198  profile single cells. **Figure 2a** provides an illustrative overview of the SIMBA graph construction
199  and the resulting low-dimensional embedding matrix of both cells and genes. Here we show
200  how SIMBA enables simultaneous dimensionality reduction and clustering-free marker gene
201  detection in scRNA-seq analysis. We applied SIMBA to a popular PBMCs dataset from 10x
202  Genomics (**Supplementary Table 1**) to illustrate its workflow. After the standard preprocessing
203  steps including normalization and log-transformation, SIMBA discretizes the gene expression
204  matrix into multiple gene expression levels (five levels, by default). The input graph is then
205  constructed wherein two types of nodes –cells and genes are connected by edges that embody
206  the relation between them and are weighted according to the corresponding multiple levels of
207  gene expression. SIMBA then generates embeddings of these nodes through a graph
208  embedding procedure (**Fig. 2a; Methods**).  Depending on the task, we have the full flexibility to
209  visualize either the whole SIMBA embeddings (embeddings of cells and all genes in
210  **Supplementary Fig. 1c**) or the partial SIMBA embeddings (embeddings of cells in **Fig. 2b,** or
211  embeddings of cells and variable genes in **Fig. 2c**, or embeddings of any entities of interest)
212  using visualization tools such as UMAP.

213  When the SIMBA embeddings of cells were visualized, each of the eight cell types, including B
214  cells, megakaryocytes, CD14 monocytes, FCGR3A monocytes, dendritic cells, NK cells, CD4 T,
215  and CD8 T cells, was clearly separated (**Fig. 2b**). When the SIMBA embeddings of both cells and
216  genes were visualized, the co-embedding space showed that SIMBA not only recovered the
217  cellular heterogeneity, but also correctly embedded informative genes close to relevant cell
218  types (**Fig. 2c**). The same set of marker genes used to annotate these cells from Scanpy[2] was
219  highlighted on the UMAP plot. In addition, as a control, we also show the locations of two
220  housekeeping genes *GAPDH* and *B2M*, which would not be expected to associate with any
221  particular cell type. From the UMAP plot, we can see that SIMBA not only was able to embed
222  major-cell-group specific genes to the correct locations (e.g., *IL7R* was embedded into CD4T
223  cells and *MS4A1* was embedded into B cells), but also was robust to rare-cell-group specific
224  genes (e.g., *PPBP* was embedded into megakaryocytes). On the contrary, non-informative or
225  non-cell-type specific genes such as *GAPDH* and *B2M* were embedded in the middle of all cell
226  groups (**Fig. 2c and Supplementary Fig. 1c**).

227  These highlighted genes can be further confirmed with "barcode plot", which visualizes the
228  estimated probability of assigning a feature to a cell by SIMBA based on the recovered edge
229  confidence (**Fig. 2d, Supplementary Fig. 1e, Methods**). An imbalance in probability indicates
230  the association of a gene to a sub-population of cells (often corresponding to known cell-types),
231  whereas a uniform probability distribution indicates a non-cell-type-specific gene. For marker

232  genes (*CST3* for monocytes and dendritic cells, *MS4A1* for B cells, and *NGK7* for NK and CD8T
233  cells), we observed a clear excess in the probability of assigning each gene to their respective
234  cell types. Conversely, for the housekeeping gene *GAPDH*, we observed a more uniform
235  distribution with much lower probability of associating that gene with the top-ranked cells.
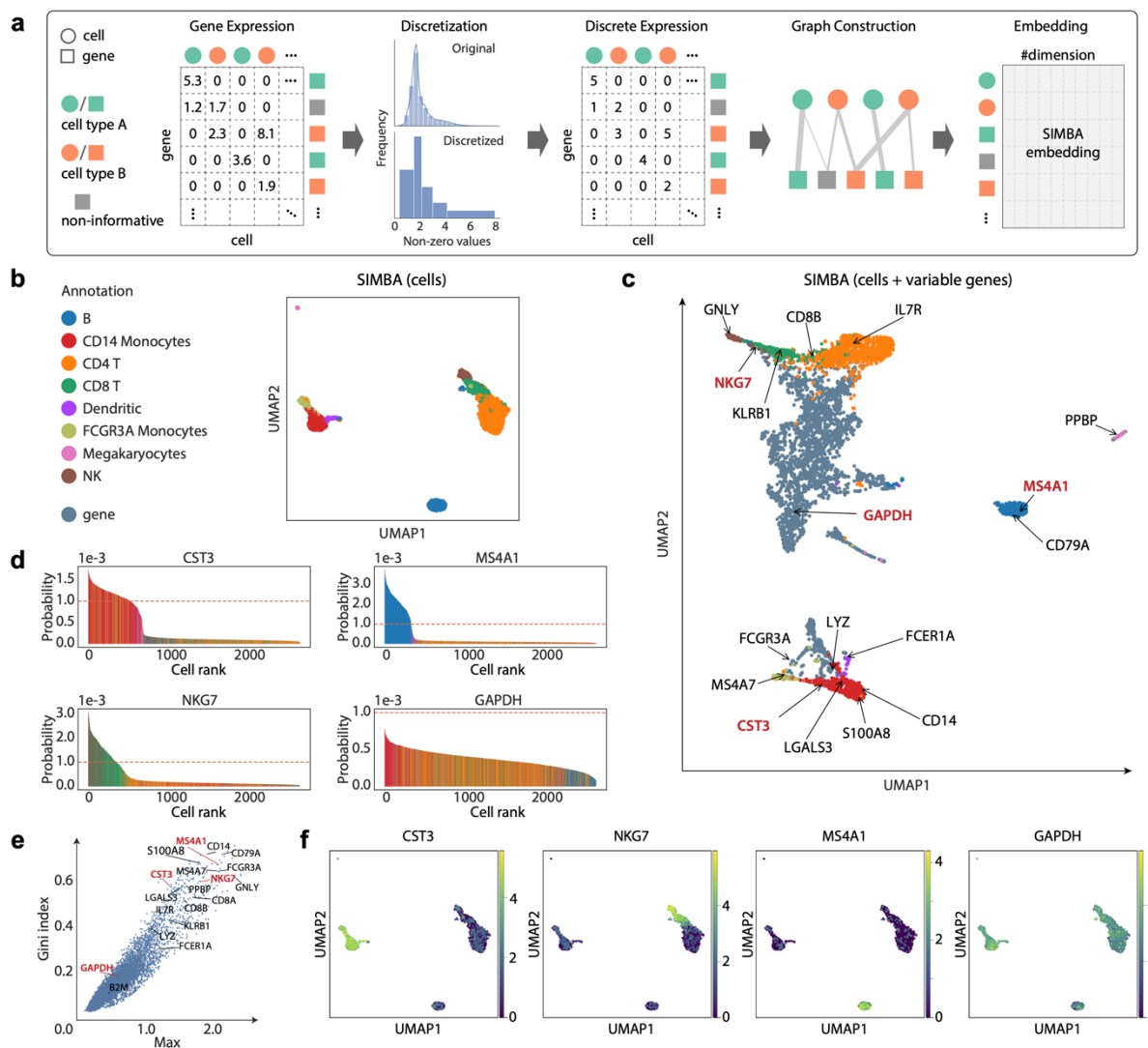


236

**Figure 2.** Single-cell RNA-seq analysis of the 10x PBMCs dataset using SIMBA. **(a)** SIMBA
graph construction and embedding in scRNA-seq analysis. Biological entities including
cells and genes are represented as shapes and colored by relevant cell types (green and
orange). Non-informative genes are colored dark grey. Gene expression measurements
for each cell are organized into a cell-by-gene matrix. These normalized non-negative
observed values undergo discretization into five gene expression levels. Cells and genes
are then assembled into a graph with nodes representing cells and genes, and edges
between them representing different gene expression levels. This graph may then be
embedded into a lower dimensional space resulting in a #entities x #dimension (by default,
50) SIMBA embedding matrix. **(b)** UMAP visualization of SIMBA embeddings of cells

247 colored by cell type. **(c)** UMAP visualization of SIMBA embeddings of cells and variable
248 genes. Cells are colored according to cell type as defined in b. Genes are colored slate
249 blue. Cell-type-specific marker genes and housekeeping genes recovered by Scanpy are
250 indicated with text and arrows. Genes highlighted in red are shown in **d**, **e**, and **f**. **(d)**
251 SIMBA barcode plots of genes *CST3*, *MS4A1*, *NKG7*, and *GAPDH*. The x-axis indicates the
252 ordering of a cell as ranked by the probability for each cell to be associated with a given
253 gene. The y-axis describes the probability. The sum of probability over all cells is equal to
254 1. Each cell is one bar and colored according to cell type as defined in b. **€** SIMBA ranking
255 of genes based on the proposed metrics. All the genes are plotted according to the Gini
256 index against max score. The same set of genes as in **c** are annotated. **(f)** UMAP
257 visualization of SIMBA embeddings of cells colored by gene expression of (left to right):
258 *CST3*, *NKG7*, *MS4A1*, and *GAPDH*.
259

260 SIMBA also provides several quantitative metrics (termed "SIMBA metrics"), including max
261 value, Gini index, standard deviation, and entropy, to assess cell-type specificity of various
262 features without requiring the prior knowledge such as cluster labels, predefined cell types, or
263 known marker genes (**Methods**). As an example, by inspecting the gene metric plot of max
264 value (a measurement of maximum probability, a higher value indicates higher cell-type
265 specificity) vs Gini index (a measurement of imbalance, a higher value indicates higher cell-type
266 specificity), we observed that the marker genes (e.g., *CST3*, *NKG7*, *MS4A1*) fall in the upper right
267 corner, as opposed to housekeeping genes (e.g., *GAPDH*) in the lower left corner (**Fig. 2e**).
268 Similar separation is observed with other metrics (**Supplementary Fig. 1b**). The cell type
269 specificity of the selected marker genes was further confirmed by visualizing their expression
270 pattern on UMAP plots (**Fig. 2f and Supplementary Fig. 1d**), accompanied by SIMBA barcode
271 plots (**Supplementary Fig. 1d**). As a certain feature (e.g., genes) might notably outnumber cells
272 or other features (when multiple types of features are present), SIMBA metrics not only serve
273 as an efficient way of ranking features based on their cell type specificity, but also provides a
274 straightforward way to filter out non-informative (non-cell-type-specific) features so that only
275 the embeddings of cells and informative features will be visualized and the SIMBA space will
276 not be crowded with non-informative features (e.g., house-keeping genes).

277 We next compared the top 600 marker genes identified by SIMBA (based on max value and Gini
278 index) with those identified by the clustering-based statistical-tests method implemented in
279 Scanpy (based on z-score calculated from the two-sided Wilcoxon rank-sum test with a
280 Benjamini-Hochberg p-value correction, one of the statistical tests recommended in Scanpy's
281 tutorial) (**Supplementary Fig. 2a**). Upon comparison, we observed that nearly half of the marker
282 genes discovered by SIMBA overlap with the marker genes identified by Scanpy
283 (**Supplementary Fig. 2a**). However, on inspection of the top non-overlapping marker genes,
284 genes identified by SIMBA are found to be enriched only within certain groups of cells
285 (**Supplementary Figs. 2b and 2c**) while genes identified by Scanpy but not by SIMBA include the
286 housekeeping gene *B2M* and multiple ribosomal protein genes (e.g., *RPS3* and *RPS6*) that are
287 expressed ubiquitously in all cell types (**Supplementary Figs. 2b and 2d**). Furthermore, a
288 combination of different statistical tests proposed in Scanpy is required to recover the genes

289  identified only by SIMBA. For example, IL7R was identified only by using the t-test and FCER1A
290  was identified only by using the Wilcoxon rank-sum test, as also noted in the Scanpy's tutorial,
291  while SIMBA successfully identified both *IL7R* and *FCER1A* as informative genes with a single
292  procedure and without clustering the cells (**Fig. 2e and Supplementary Fig. 1b**). These examples
293  illustrated some limitations of the clustering-based statistical-tests methods.

294  Lastly, we showed that SIMBA does not require variable gene selection, which is an essential
295  step in standard scRNA-seq pipelines such as Seurat or Scanpy. SIMBA produces very similar
296  embeddings for cells with and without variable gene selection (**Fig. 2b** and **Supplementary Fig.
297  2e**), though we observed that variable gene section does improve efficiency of the training
298  procedure.

299  **SIMBA enables simultaneous characterization of cell states and cis-regulatory elements by**
300  **jointly modeling accessible sites and DNA sequences in scATAC-seq analysis**
301
302  As one of the most popular single-cell epigenomic techniques, single-cell assay for transposase-
303  accessible chromatin using sequencing (scATAC-seq) has been widely used to profile regions of
304  open chromatin and identify functional *cis*-regulatory elements such as enhancers and active
305  promoters. In scATAC-seq, cells are characterized by different types of features [14], such as regions
306  of accessible chromatin ("peaks" or "bins") and *cis*-regulatory elements (DNA sequences) within
307  these accessible regions including transcription factor (TF) motifs or *k*-mers.
308
309  Unlike existing methods that can only use peaks/bins or the DNA sequence within them, SIMBA
310  can leverage simultaneously both types of features to learn cell states due to its flexibility in
311  graph construction. Also, as SIMBA encodes cell-feature or feature-feature relations into the
312  graph based on the simple binary presence of a feature, SIMBA does not need additional
313  normalization steps such as term frequency-inverse document frequency (TF-IDF), which is
314  required by most scATAC-seq analyses. When only peaks/bins are used, SIMBA constructs a graph
315  with nodes representing cells and chromatin regions (peaks or bins) and edges indicating the
316  accessibility of the chromatin regions in cells (**Fig. 3a**). When the DNA sequences for chromatin
317  regions are available, SIMBA can also encode DNA sequences including TF motifs and *k*-mers into
318  the graph by adding edges between these entities as nodes and the existing chromatin region
319  nodes. The edges in this case indicate the presence of TF motifs/*k*-mers within these chromatin-
320  accessible regions. Through the embedding procedure, SIMBA generates embeddings of cells
321  along with peaks and DNA sequences (**Methods**). Finally, either the partial SIMBA embeddings
322  (embeddings of cells in **Fig.3b**) or the whole SIMBA embeddings (embeddings of cells and all the
323  features in **Fig.3c**) can be visualized. Therefore, SIMBA enables dimensionality reduction by
324  leveraging both chromatin accessible regions and cis-regulatory sequences. Simultaneously, it
325  highlights the cell-type-specific open chromatin regions and regulatory DNA sequences in a
326  clustering-free way.
327
328  To demonstrate the value of SIMBA embeddings for scATAC-seq analysis, we first applied
329  SIMBA to a scATAC-seq data of 2,034 human hematopoietic cells with FACS-characterized cell
330  types[15](**Supplementary Table 1**). For the SIMBA embeddings of cells alone, as shown in **Fig. 3b**,

331    SIMBA accurately separated cells such that cells belonging to distinct cell types are visually
332    separated. For the SIMBA embeddings of cells together with various types of features, as shown
333    in **Fig. 3c**, SIMBA successfully embedded distinct features from both positional (peaks/bins) as
334    well as sequence-content (TF motifs and k-mers) information together based on their biological
335    relations. Notably, based on SIMBA metrics, these highlighted features that are embedded
336    within each cell type all have high cell-type specificity scores (shown in the upper right part of
337    SIMBA metric plots in **Figure 3d**).
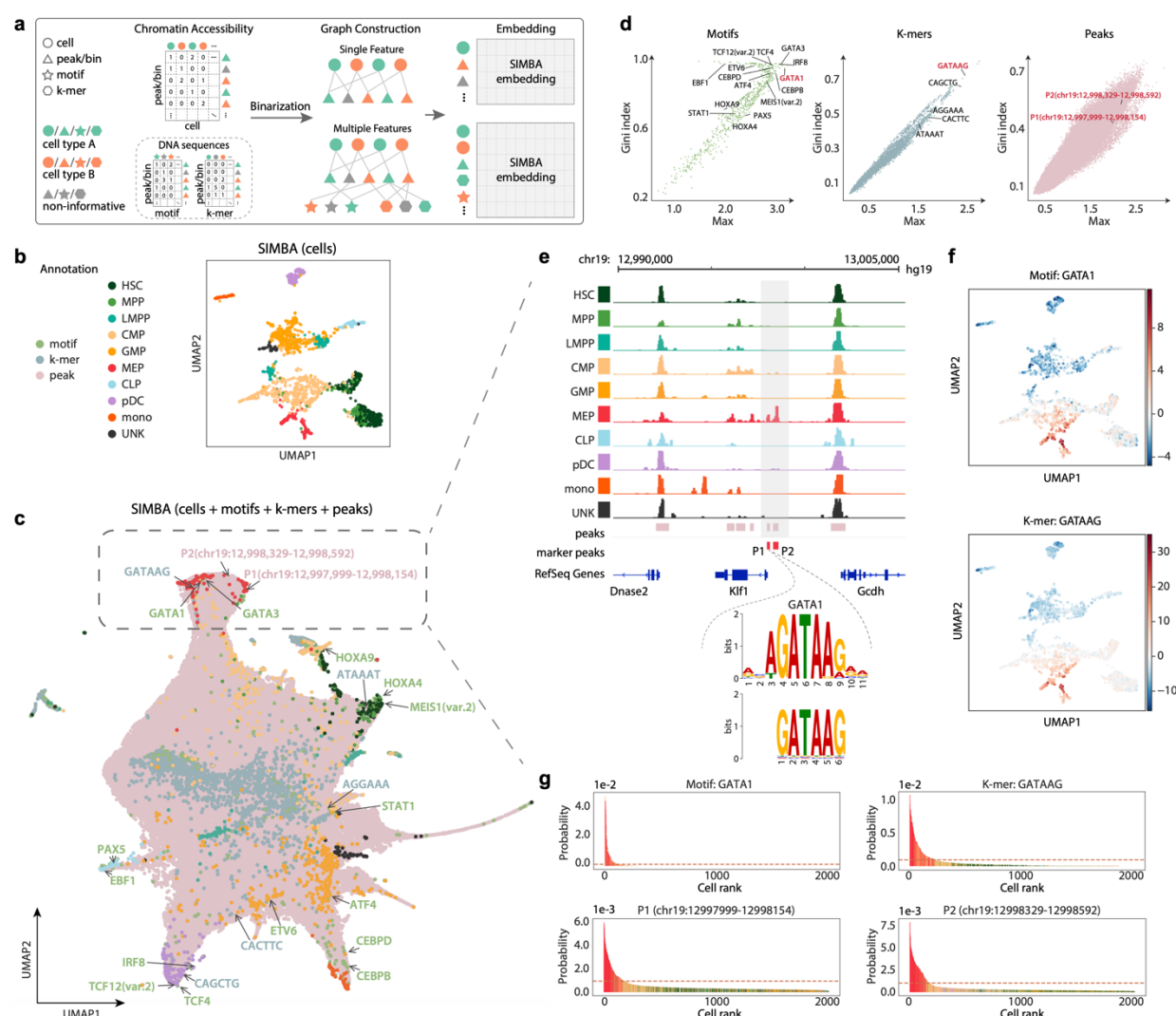
338



339
340

341    **Figure 3.** Single-cell ATAC-seq analysis of the human hematopoiesis dataset using SIMBA.
342    **(a)** SIMBA graph construction and embedding in scATAC-seq analysis. Biological entities
343    including cells, peaks/bins, TF motifs, k-mers are represented as shapes and colored by
344    relevant cell types (green and orange). Non-informative features are colored dark grey.
345    Cells and chromatin accessible features (peaks / bins) are organized into a cell x peaks /
346    bins matrix. When sequence information (TF motif or *k*-mer sequence) within these
347    regions is available, they can be organized into two sub-matrices to associate a TF motif

348     or k-mer sequence with each peak/bin. These constructed feature matrices are then
349     binarized and assembled into a graph. When a single feature (chromatin accessibility) is
350     used, the graph encodes cells and peaks/bins as nodes. When multiple features (both
351     chromatin accessibility and DNA sequences) are used, this graph may then be extended
352     with the addition of TF motifs and k-mer sequences as nodes connected. Finally, SIMBA
353     embeddings of these entities are generated through a graph embedding procedure. **(b)**
354     UMAP visualization of SIMBA embeddings of cells colored by cell type. **(c)** UMAP
355     visualization of SIMBA embeddings of cells and features including TF motifs, k-mers, and
356     peaks. Cells are colored by cell type while motifs, k-mers, and peaks are colored green,
357     blue, and pink, respectively. Cell type specific features that are embedded near their
358     corresponding cell types are indicated as the text labels (colored according to feature type)
359     with arrows. **(d)** SIMBA metric plots of TF motifs, k-mers, and peaks. Cell-type specific
360     features annotated in **(c)** are highlighted. **€** Genomic tracks of aligned scATAC-seq
361     fragments, separated and colored by cell type. Two marker peaks P1 and P2 in red are
362     shown beneath the alignment. Within the peak P1, k-mer GATAAG and its resembling
363     GATA1 motif logo are highlighted. **(f)** UMAP visualization of SIMBA embeddings of cells
364     colored by TF activity scores of the GATA1 motif and k-mer GATAAG enrichment. **(g)**
365     SIMBA barcode plots of the GATA1 motif, k-mer GATAAG, and the two peaks P1 and P2.
366     Cells are colored according to cell type labels described above. Dotted red line indicates
367     the same cutoff used in all four plots.

368

369 Our analysis using SIMBA led to several key findings in human hematopoietic differentiation.

370

371 First, SIMBA identified key master regulators of hematopoiesis. As highlighted in **Fig. 3c**, we
372 observed that motifs of previously reported TFs were embedded near their respective cell types
373 in the UMAP plot. For example, the GATA1 and GATA3 motifs are proximal to megakaryocyte-
374 erythroid progenitor (MEP) cells[16], the PAX5 and EBF1 motifs are near to common lymphoid
375 progenitor (CLP) cells[17], and the CEBPB and CEBPD motifs are proximal to monocyte (mono)
376 population[18].

377

378 Second, SIMBA revealed an unbiased set of DNA sequences, i.e., *k*-mers, that represent
379 important TF binding motifs involved in hematopoiesis. We observed that these *k*-mers were
380 embedded near their resembling TF binding motifs and relevant cell subpopulations (**Fig. 3c** and
381 **3e**, **Supplementary Fig. 3b**), indicating that this methodological framework is capable of *de*
382 *novo* motif discovery. For example, the DNA sequence, CAGCTG is embedded in plasmacytoid
383 dendritic cells (pDCs); this sequence matches the TCF12 binding motif, which controls dendritic
384 cell lineage specification. To further illustrate the interpretability of the SIMBA embeddings of
385 TF motifs and k-mers, we calculated per-cell TF activity scores[19] (high-variance TF motifs/*k*-
386 mers) and visualized them on SIMBA embeddings of cells. As shown in **Figure 3f**, the GATA1 TF
387 motif and k-mer GATAAG that were both embedded in MEP cells by SIMBA, also showed high-
388 level activity in MEP cells. The consistency between SIMBA embedding and TF activity was
389 observed for most of other TF motifs and *k*-mers as well (**Supplementary Fig. 3a, 3b**).

390

391 Third, SIMBA identified differentially accessible chromatin regions that may mediate cell-type
392 specific gene regulation. For example, the two peaks with coordinates chr19:12997999-
393 12998154 (P1) and chr19:12998329-12998592 (P2) that were embedded within MEP cells were
394 almost exclusively observed in MEP cells on KLF1 genome track (**Fig. 3e**). Interestingly, P1,
395 upstream of *KLF1*, contains the *k*-mer GATAAG that matches the GATA1 binding motif, while
396 transcription factor GATA1 is known to regulate the gene *KLF1* and plays a pivotal role in
397 erythroid cell and megakaryocyte development[20]. Therefore, by embedding these MEP-cell-
398 related regulatory elements into the neighborhood of MEP cells, SIMBA demonstrates a novel
399 means of studying the epigenetic landscape of cell differentiation. To further validate the
400 differentially accessible regions identified by SIMBA, we selected 100 peaks at random from
401 each annotated cell type in SIMBA co-embedding space. From the heatmap of chromatin
402 accessibility, we clearly see that the peaks embedded nearby respective types correlate with
403 strong cell-type specificity. This observation is robust to the number of cells within each cell
404 type (**Supplementary Fig. 3c**).

405

406 Available methods for scATAC-seq analysis visualize only cells. While SIMBA diverges from these
407 available workflows, enabling the co-embedding of cells and features, we still qualitatively and
408 quantitatively compared the SIMBA embeddings of cells to state-of-the-art scATAC-seq analysis
409 methods by their ability to distinguish cell types.  Our analyses show that SIMBA overall
410 performs better than the methods evaluated, further demonstrating the wide utility of SIMBA
411 (**Supplementary Figs. 4 and 5; Supplementary Note 1**).

412

413

414 **SIMBA enables simultaneous learning of cellular heterogeneity and gene regulatory circuits**
415 **from integrated analysis of single-cell multimodal data**

416

417 scRNA-seq and scATAC-seq are two of the most widely adopted single-cell sequencing
418 technologies, but they are limited to measuring only a single aspect of cell state at a time. To
419 improve our ability to interrogate cell states, several single-cell dual-omics technologies have
420 been recently developed [21-24] to jointly profile transcriptome and chromatin accessibility within
421 the same individual cells, therefore providing the potential to correlate gene expression with
422 accessible regulatory elements and further delineate the yet elusive principles of gene
423 regulation. This section outlines the SIMBA's ability to simultaneously learn cell heterogeneity
424 as well as gene regulatory circuits from single-cell multiomic data. We applied SIMBA to three
425 recent single-cell dual-omics technologies: SHARE-seq[22], SNARE-seq[21], and a multiome PBMCs
426 dataset from 10x Genomics (**Supplementary Table 1**).

427 **Figure 4a** illustrates the procedure of graph construction and generation of the final SIMBA
428 embedding matrix. Briefly, for scRNA-seq, the gene expression matrix is discretized to generate
429 different levels of gene expression. For scATAC-seq, both the chromatin accessibility matrix and
430 motif/*k*-mer match matrix are binarized. In this graph, there are five entity (node) types,
431 including cells, genes, peaks, motifs, and *k*-mers. For scRNA-seq, an edge indicates whether a
432 gene is expressed in a cell and its weight indicates the gene expression level (five levels, by
433 default). For scATAC-seq, an edge indicates whether a peak is present in a cell or if a TF motif/*k*-

434   mer is present within a peak. Once the graph is constructed, the graph embedding procedure is
435   performed to generate SIMBA embeddings of cells and different types of features. scATAC-seq
436   peaks generally greatly outnumber cells and other features and many of these peaks are non-
437   informative, resulting in them dominating the space if the whole SIMBA embeddings are
438   visualized (**Supplementary Fig. 6a, c**). In such cases, we leverage the flexibility of SIMBA
439   embedding to only visualize the partial SIMBA embeddings to improve the visibility of cells and
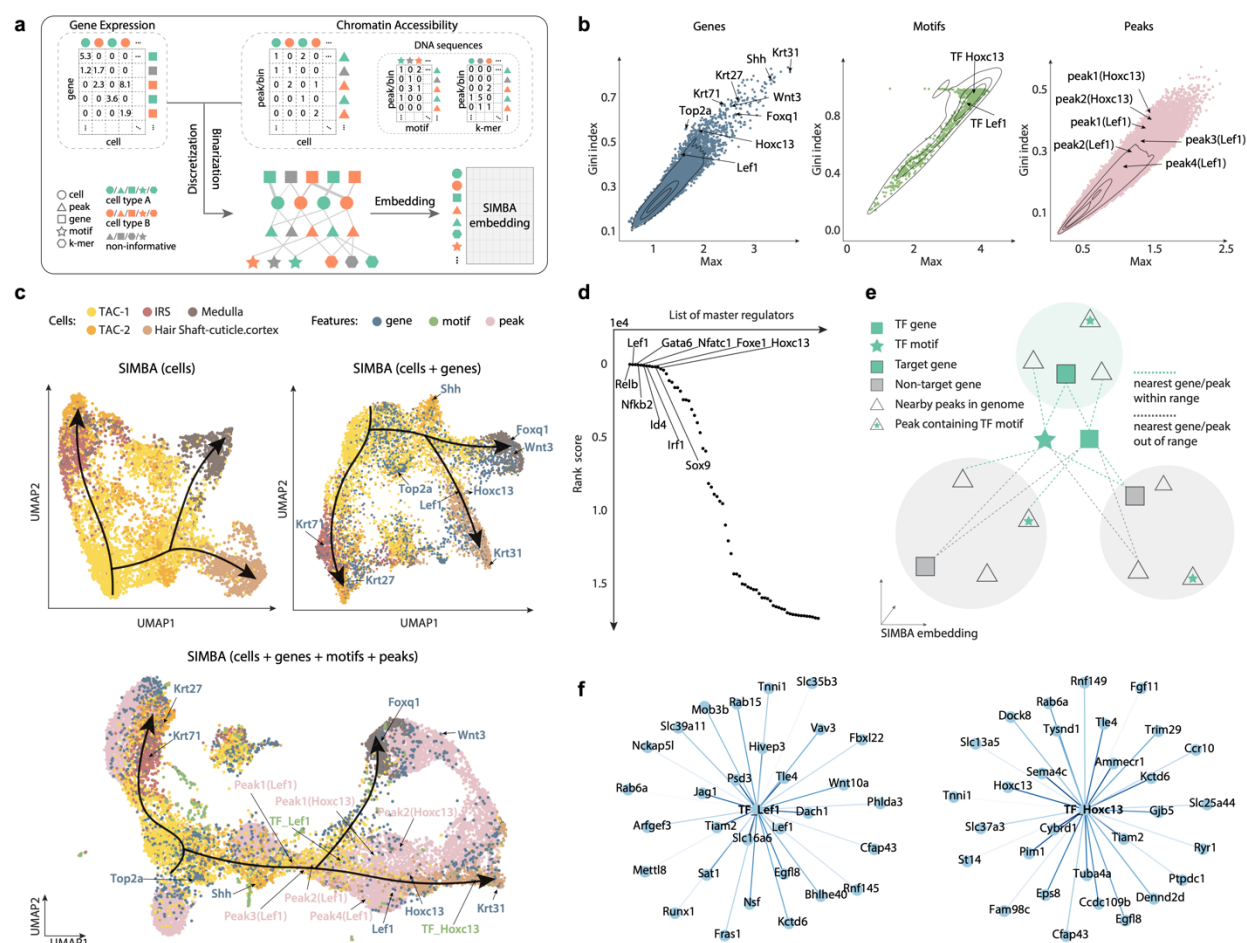440   cell-type-specific features.



441

**Figure 4.** Multimodal analysis of the SHARE-seq hair follicle dataset using SIMBA. **(a)**
SIMBA graph construction and embedding in multimodal analysis. Overview of SIMBA's
approach to multimodal (scRNA-seq + scATAC-seq) data analysis. **(b)** SIMBA metric plots
of genes, TF motifs, and peaks. All these features are plotted according to the Gini index
against max score. Cell-type specific genes, TF motifs, and peaks are highlighted. **(c)** UMAP
visualization of SIMBA embeddings of cells (Top-left), cells and genes (Top-right), and cells
along with genes, TF motifs, and peaks (Bottom). **(d)** Ranked scatter plot of candidate
master regulators as identified by SIMBA. **€** Schematic description of SIMBA's strategy for
identifying target genes given a master regulator. **(f)** Top 30 target genes of transcription
factors Lef1 and Hoxc13 as inferred by SIMBA.

452

453     To demonstrate the usefulness and versatility of the SIMBA embeddings, we analyzed
454     the cell populations undergoing the dynamic process of hair follicle differentiation from
455     mouse skin profiled with SHARE-seq.

456     First, we calculated SIMBA metrics (max values and Gini index scores) to assess the cell-type
457     specificity of different types of features, including genes, TF motifs, and peaks (**Fig. 4b,**
458     **Methods**).  As shown in **Figure 4b**, based on these metrics, we successfully recovered genes
459     associated with hair follicles such as *Lef1* and *Hoxc13*. Similarly, TF motifs and peaks proximal to
460     the genomic loci of these genes also score in the upper right quadrant of the metric plots.
461     SIMBA's cell-type specificity metrics successfully revealed the key genes and regulatory factors
462     important to the hair follicle differentiation process.

463     Next, we visualized and interrogated the SIMBA embeddings of 1) cells; 2) cells and top-ranked
464     genes based on SIMBA metrics; and 3) cells, top-ranked genes and TF motifs based on SIMBA
465     metrics, and the neighboring peaks of these genes and TF motifs by querying the SIMBA space
466     (**Methods**). **Figure 4c** shows the UMAP visualization of the partial SIMBA embeddings of cells
467     and informative features. The UMAP visualization of SIMBA embeddings of cells and the full set
468     of features was also performed (**Supplementary Fig. 6a**).

469     The SIMBA embeddings of cells were able to reveal the three fate decisions from transit-
470     amplifying cells (TACs), including inner root sheath (IRS), medulla, and cuticle/cortex. The
471     SIMBA embeddings of cells and informative features uncovered important genes and regulatory
472     factors along the hair follicle differentiation trajectories. For example, the marker genes *Krt71*,
473     *Krt31*, and *Foxq1* were embedded into their corresponding cell types: IRS, cuticle/cortex, and
474     medulla, respectively. The Lef1 motif was embedded into the beginning of medulla and
475     cuticle/cortex lineages while the Hoxc13 motif was embedded into the late stage of
476     cuticle/cortex differentiation. Peaks near the *Lef1* and *Hoxc13* loci were also embedded into the
477     nearby regions of these genes and motifs, as expected.

478     To show the robustness of SIMBA, we separated the scRNA-seq and scATAC-seq components
479     within the SHARE-seq dataset and performed each respective single-modality analysis. With the
480     consistent embedding results of cells as in multimodal analysis, we further demonstrated that
481     SIMBA embedding procedure is robust to the type and the number of features encoded in the
482     input graph (**Supplementary Fig. 6b,6c**). Each reported marker gene was corroborated using
483     the UMAP plots with cells colored by gene expression as well as using the SIMBA barcode plots.
484     The two aforementioned TF motifs and their respective peak sets were supported by the
485     corresponding SIMBA barcode plots, wherein we observed an imbalanced distribution with high
486     probability towards the correct cell type labels (**Supplementary Fig. 7a-d**).

487     Further, we demonstrated that the SIMBA co-embedding space of cells and features provides
488     the potential to identify master regulators of differentiation and infer their target regulatory
489     genes. To define a master regulator *a priori*, we postulated that both its TF motif and TF gene

490    should be cell-type specific, given that active gene regulation involves both the expression of a

491    TF and accessibility of its binding sites. Thus, the TF motif and TF gene should be embedded

492    closely in the shared latent space. Extending this logic to identify putative master regulators, we

493    assessed the cell-type-specificity of TF motifs and genes based on SIMBA metrics and ranked all

494    potential master regulators based on the distance between the TF motif and the respective TF

495    gene in the shared SIMBA embedding space (**Methods**).  SIMBA successfully identified

496    previously described master regulators such as Lef1, Gata6, Nfatc1, and Hoxc13 as the top

497    master regulators related to lineage commitment in mouse skin (**Fig. 4d, Supplementary Table**

498    **2**). To infer the target genes of a given master regulator, we postulate that in the shared SIMBA

499    embedding space, 1) the target gene is close to both the TF motif and the TF gene; 2) the

500    accessible regions (peaks) near the target gene loci must be close to both the TF motif and the

501    target TF gene. Resting on these assumptions of *cis*-regulatory dynamics, the inference of target

502    genes was performed by calculating the distance between target gene candidates and the

503    respective TF motif and gene. In addition, nearby peaks around the target gene's locus and the

504    presence of TF motif in these nearby peaks are also considered (**Fig. 4e, Methods**). The top 30

505    target genes of TF Lef1 and TF Hoxc13 inferred by SIMBA are shown respectively (**Fig. 4f,**

506    **Supplementary Fig. 7e**). The full list of ranked target genes is provided in **Supplementary Table**

507    **3**. Notably, our approach recovered targets genes that were also reported in the original

508    study[22]. For example, genes *Lef1*, *Jag1*, *Hoxc13*, *Gtf2ird1* are regulated by the TF Lef1, while

509    genes *Cybrd1*, *Hoxc13*, *St14* are regulated by the TF Hoxc13.

510    In addition to SHARE-seq, we also applied SIMBA to another two dual-omics datasets, the

511    mouse cerebral cortex dataset profiled by SNARE-seq[21] (**Supplementary Fig. 8**) and the

512    multiome PBMCs dataset from 10x Genomics (**Supplementary Fig. 9**). By validating the

513    embeddings of cells and features with given cell type labels (**Supplementary Fig. 8a and Fig.9a**),

514    marker genes from the original study (**Supplementary Fig. 8a,b,d and Fig. 9a,b,d**), and

515    differentially accessible chromatin regions (**Supplementary Fig. 8c and Fig. 9c**), we further

516    demonstrate the suitability of SIMBA for multimodal analysis.

517

518    **SIMBA enables simultaneous batch correction and clustering-free marker gene detection**

519

520    Efforts to collect data from single cells has grown to the level of consortia that span multiple

521    institutions with the hopes of finely mapping and characterizing specific tissues. This has

522    brought with it an increased demand for analysis methods that are capable of negating

523    technical covariates inherent to multi-batch data collection, including experimental replicate

524    identity, sample preparation, and sequencing platform. Batch correction that removes the

525    effects of technical covariation while preserving true biological signals is required prior to

526    downstream analysis [25, 26]. Existing methods follow a workflow with four primary steps. The first

527    step is the actual batch correction, which often generates a "batch corrected" latent space. The

528    second step clusters cells in this batch corrected space. Based on the clustering result the third

529    step detects marker genes in the original gene expression space of each batch because the low-

530    dimensional "batch corrected" space is no longer comprised of genes. The fourth step finally

531    combines the marker genes detected from each batch. However, these methods are clustering-

532  dependent and may result in the inconsistent explanation of marker genes as marker genes are
533  detected in each original batch as opposed to the batch-corrected space. Unlike current
534  methods, in addition to embeddings of cells, SIMBA generates comparable embeddings of
535  genes and therefore relieves marker gene discovery from a dependence on the original gene
536  expression space. Thus, SIMBA enables simultaneous batch effect removal and cell-type-
537  specific marker gene detection in the same integrated space without clustering.
538



539
540
541  **Figure 5**. Batch correction analysis of scRNA-seq data using SIMBA. **(a)** SIMBA graph
542  construction and embedding in batch correction analysis. Overview of SIMBA's approach
543  to batch correction across scRNA-seq datasets. Distinct shapes indicate the type of entity
544  (cell or gene). Colors distinguish batches or cell types. **(b)** UMAP visualization of the
545  scRNA-seq human pancreas dataset with five batches of different studies before and after
546  batch correction. Cells are colored by scRNA-seq data source and cell type respectively.
547  Top: UMAP visualization before batch correction; Bottom: UMAP visualization after batch
548  correction with SIMBA; **(c)** UMAP visualization of SIMBA embeddings of cells and genes,
549  with batch effect removed and known marker genes highlighted.

550
551  We first demonstrate that SIMBA readily corrects batch effects and produces joint embeddings
552  of cells and genes across multiple scRNA-seq datasets generated from varying sequencing

553   platforms and cell type compositions. While existing methods for scRNA-seq analysis rely on
554   specialized tools for batch correction, SIMBA works as a stand-alone method obviating the need
555   for prior input data correction when applied to multi-batch scRNA-seq dataset. SIMBA
556   accomplishes batch correction by encoding multiple scRNA-seq datasets into a single graph (**Fig.
557   5a**). Cell nodes across batches are connected to gene nodes through experimentally measured
558   edges as in the previously described scRNA-seq graph construction. Here, the gene nodes are
559   shared between the cell nodes of different batches. In addition to the experimentally measured
560   edges, batch correction is further enhanced through computationally inferred edges drawn
561   between similar cell nodes across datasets using a truncated randomized singular value
562   decomposition (SVD)-based procedure. SIMBA then generates the embeddings of all nodes
563   including cells of each batch and genes from the resulting graph (**Methods**). The SIMBA
564   embeddings of cells naturally represent the batch-corrected space. In addition, the whole
565   SIMBA embeddings of all entities provide the batch-corrected space, in which cells and genes
566   co-exist, and therefore allow for individual-cell-level marker detection by performing biological
567   queries of cells in the SIMBA space (**Methods**). We visualized both SIMBA embeddings of cells
568   (**Fig. 5b**), and the whole SIMBA embeddings of cells and genes (**Fig. 5c**) in UMAP.
569
570   We applied SIMBA to two multi-batch scRNA-seq datasets; a mouse atlas dataset composed of
571   two batches and a human pancreas dataset spanning five batches used in a recent benchmark
572   study[25] (**Supplementary Table 1**). The mouse atlas dataset contains two scRNA-seq datasets
573   with shared cell types from different sequencing platform. The human pancreas dataset
574   contains five samples pooled from five sources using four different sequencing techniques, in
575   which not all cell types are shared across each sample. For both datasets, SIMBA successfully
576   corrected batch effects, evenly mixing batches within annotated cell type clusters, while
577   maintaining the segregation of these clusters in the resulting embedding, indicating
578   preservation of biological signal and elimination of confounding technical covariates (**Fig. 5b,
579   Supplementary Fig. 12b**). It is important to note that the mouse atlas dataset was collected
580   from nine different organ systems, so there exists some expected heterogeneity within cell type
581   labels. Conversely, the human pancreas datasets are curated from a single organ and SIMBA
582   sufficiently separated cell types into transcriptionally distinct, homogeneous cell clusters (**Fig.
583   5b**).
584
585   Through removing batch effects during graph embedding, SIMBA simultaneously identifies cell-
586   type-specific marker genes (**Fig. 5c**). In the absence of the eliminated technical covariation,
587   marker genes are identifiable by performing biological queries for neighboring genes within cell
588   types in the SIMBA embeddings of cells and genes (**Methods**). In the case of unknown cell
589   labels, marker genes can be identified by calculating SIMBA metrics (**Methods**). SIMBA correctly
590   embeds known cell-type-specific marker genes proximal to the correct cell type labels, while
591   non-marker genes were non-proximal to specifically-labelled cells (**Supplementary Fig. 10, 11**).
592   The resulting marker genes recapitulated the clustering-based differential expression (DE)
593   analysis results for each dataset[27-32] (e.g. *Cdh5*, *Tie1*, *Myct1* for endothelial cell and *C1qc*, Fcgr1
594   for macrophage, *S100a8*, *Trem3* for Neutrophil in the mouse atlas dataset and *KIF12* for alpha
595   cell and *KRT19* for ductal cell in the human pancreas dataset) and are shown to be expressed
596   specifically in the queried cell types (**Supplementary Fig. 10, 11**). Taken together, these results

597 distinguish SIMBA from existing batch correction methods that rely on clustering in a batch-
598 corrected space, followed by differential gene expression analysis in the original, uncorrected
599 space of each batch.
600
601 While SIMBA is a versatile graph embedding method that can perform multiple tasks and
602 generate embeddings of both cells and genes, we evaluated the SIMBA embeddings of cells for
603 this task with methods that were specifically designed for batch correction. We considered
604 three widely adopted batch correction methods that demonstrated top-tier performance based
605 on a recent benchmark study[25]: Seurat3, LIGER and Harmony. Our results indicate that SIMBA
606 achieved comparable batch correction performance both qualitatively and quantitatively while
607 enabling simultaneous marker gene detection by providing the additional SIMBA embeddings
608 of genes. (**Supplementary Note 2, Supplementary Figure 12**).
609
610 **SIMBA enables simultaneous multi-omics integration and clustering-free multi-type marker**
611 **feature detection**
612
613 Single-cell assays are now capable of measuring a broad range of cellular modalities and data is
614 being generated that describes cells by varying features sets, which has motivated the need for
615 methods that leverage these features to perform multi-omics integration such that a more
616 comprehensive description of cell state may be learned. This is different from multi-modal
617 analysis because the correspondence between individual cells is unknown. Current multi-omics
618 integration methods follow a similar workflow as the previously described batch correction
619 methods, including: 1) generating a low-dimensional integrated space of cells; 2) clustering cells
620 in the integrated space; 3) detecting marker features in the original feature (e.g., genes, peaks)
621 space of each modality because the low-dimensional integrated space no longer consists of the
622 original features. Unlike existing multi-omics integration methods that cannot directly explore
623 multi-type features in the integrated space and require clustering for identifying marker
624 features, we demonstrate that SIMBA enables simultaneous multi-omics integration and
625 clustering-free detection of distinct marker features, specifically as it is applied to datasets
626 comprised of scRNA-seq and scATAC-seq.
627
628 SIMBA accomplishes this integration by first building one graph for scRNA-seq data and another
629 graph for scATAC-seq data independently as described in previous sections (**Fig. 6a**). To connect
630 these two graphs, SIMBA then calculates gene activity scores by summarizing accessible regions
631 from scATAC-seq data and then infers edges between cells of different assays based on their
632 shared gene expression modules as previously described in the batch correction section. Finally,
633 SIMBA embeds the graph of cells, genes, and peaks into a common, low-dimensional space. The
634 SIMBA embeddings of cells naturally represent the integrated space of multiple modalities.
635 Furthermore, the SIMBA embeddings of all entities provide the integrated space containing cell,
636 genes, and peaks, and therefore enable the individual-cell-level marker detection of multi-type
637 features by performing biological queries of cells in SIMBA space (**Methods**). The SIMBA
638 embeddings of these multi-omics entities can be visualized either partially or as a whole using
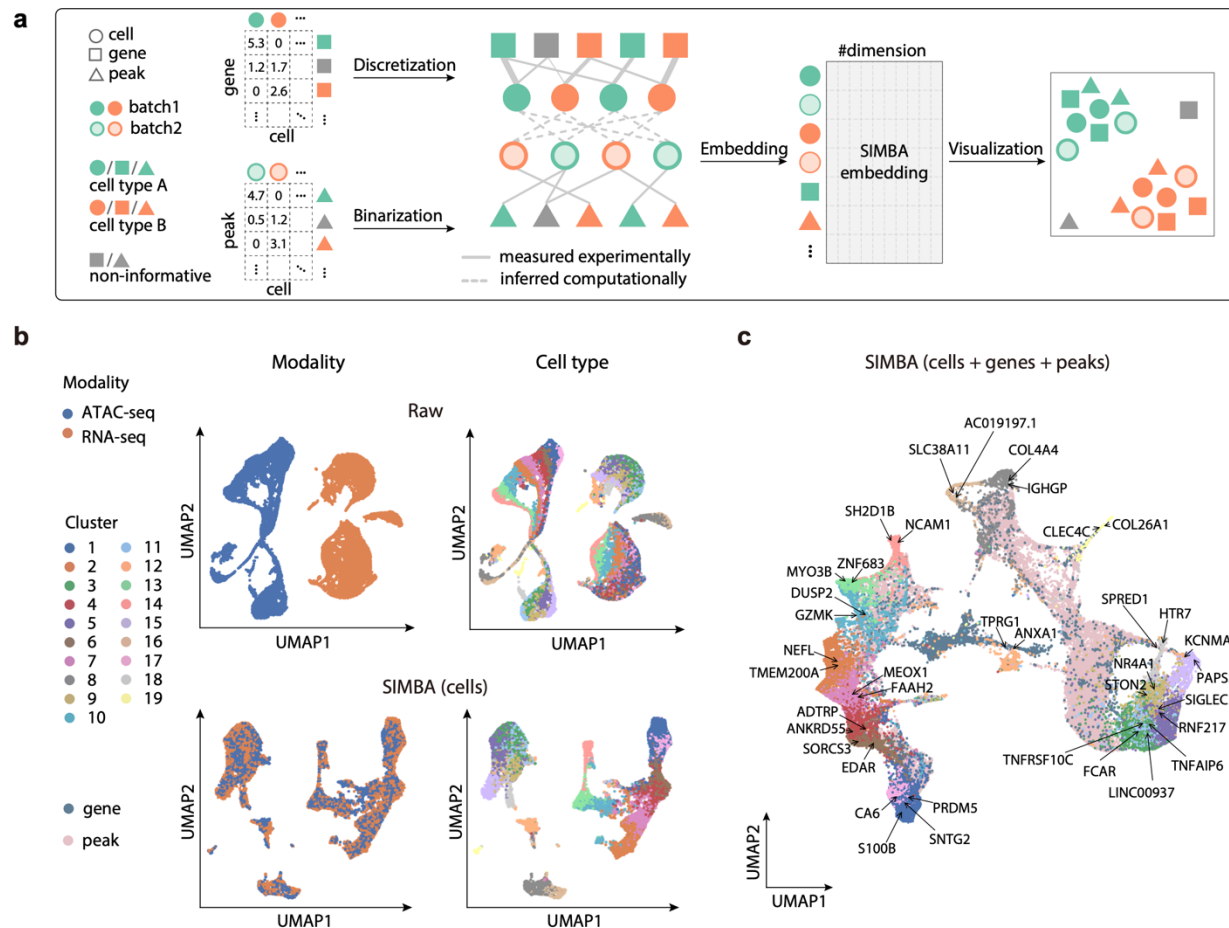639 UMAP or similar visualization tools.
640

**Figure 6**. Multi-omics integration of scRNA-seq + scATAC-seq data using SIMBA. **(a)** SIMBA graph construction and embedding in multi-omics integration. Overview of SIMBA's approach to data integration across scRNA-seq and scATAC-seq. Distinct shapes indicate the type of entity (cell, gene, or peak). Colors distinguish batches or cell types. **(b)** UMAP visualization of the integrated scRNA-seq and scATAC-seq data manually created from the 10x human PBMCs dataset before and after data integration. Cells are colored by single-cell modality and cell type respectively. Top: UMAP visualization before integration; Bottom: UMAP visualization after integration with SIMBA. **(c)** UMAP visualization of SIMBA embeddings of cells, genes, and peaks with two cell modalities integrated and known marker genes highlighted.

To facilitate the evaluation of data integration performance, we created datasets with ground-truth labels by manually splitting the dual-omics datasets into two single-modality datasets (i.e., scRNA-seq and scATAC-seq), in which we know the true matching between cells across the two modalities. We then applied SIMBA to the integration analysis of two case studies where scRNA-seq and scATAC-seq datasets are generated from the SHARE-seq mouse skin dataset and the 10x Genomics multiome human PBMCs dataset, respectively (**Supplementary Table 1**).

660 We first visualized the SIMBA embeddings of cells and observed that SIMBA was able to
661 preserve cellular heterogeneity while evenly mixing the two modalities (**Fig. 6b, Supplementary**
662 **Fig. 15b**). We then visualized the SIMBA embeddings of cells, genes, and top-ranked peaks
663 based on SIMBA metrics and observed that in addition to learning cellular heterogeneity,
664 SIMBA simultaneously identified marker genes and peaks at single-cell resolution. In the co-
665 embedding space, we observed that the neighbor genes of cells (highlighted in UMAP plots),
666 are each exclusively expressed in their corresponding cell types (**Supplementary Figs. 13a-e,**
667 **14a-c,e**). For example, in the SHARE-seq mouse skin dataset, *Foxq1* and *Shh* are located within
668 medulla and TAC-2, respectively; in the 10x PBMCs dataset, *PAPSS2* and *KCNMA1*, which are
669 the marker genes of blood monocytes, are embedded close to each other. Similarly, we
670 observed that the neighbor peaks of cells show a clear cell-type-specific accessibility pattern
671 that is robust to the cluster size of a given cell type (**Supplementary Figs. 13f and 14d).**
672
673 The joint embedding of cells and features produced by SIMBA is fundamentally distinguished
674 from other multi-omics integration methods in that it simultaneously achieves integration as
675 well as marker feature discovery. However, we still sought to compare the SIMBA embeddings
676 of cells with two widely-adopted single-cell multi-omics integration methods, Seurat3 and
677 LIGER, based on their ability to integrate single-cell modalities while persevering cellular
678 heterogeneity (**Supplementary Note 3**). We observed that SIMBA achieved the overall best
679 performance on the mouse skin SHARE-seq dataset and 10x PBMCs multiome dataset.
680

681 ## Discussion

682
683 Multimodal measurements of individual cells offer new and unexplored opportunities to study
684 cell identity as a function of the complex interactions between omic layers. While these
685 datasets offer an exciting potential for discovery, computational analysis methods to fully
686 delineate the cell states and molecular processes across multiple genomic features remain
687 insufficient.
688
689 As presented in this manuscript, SIMBA models cells and measured features as nodes encoded
690 in a graph and employs a scalable and efficient graph embedding procedure to embed nodes of
691 cells and features into a shared latent space. We demonstrate that direct graph representations
692 of single-cell data capture not only the relations between cells and the quantified features of
693 the experiment (e.g., gene expression or chromatin accessibility) but also hierarchical relations
694 between features. An example of such a hierarchical relation is the coordinate-level description
695 of an ATAC-seq peak and the corresponding TF motifs and/or $k$-mer sequences contained
696 within that region. In the resulting joint embedding, proximity-based biological queries can be
697 performed to discover cell-type-specific co-regulatory machinery across modalities. Therefore,
698 SIMBA enables simultaneous learning of cellular heterogeneity and cell-type-specific
699 multimodal features and complements the current gene regulatory network analyses. SIMBA
700 also circumvents the ordinary reliance on cell clustering for cell sub-population feature
701 discovery. We thus avoid user-defined clustering solutions, which may lead to artifactual
702 discovery or false negative results.
703

704     SIMBA has been extensively benchmarked across single-cell modalities and tasks, obtaining
705     better or comparable performance metrics when compared to current state-of-the-art methods
706     developed for the respective task. In contrast to tools developed and optimized for a single,
707     specific task these results suggest a wide applicability of SIMBA's graph-based framework,
708     obviating the need to combine multiple analysis tools.
709

710     Graph embedding methods hold significant promise for the analysis of biological data. Previous
711     applications of graph embedding include functional annotation of genes [33], transcription factor
712     binding to DNA motifs [34] and more recent single-cell RNA-seq analyses [35, 36]. The graph encoding
713     and embedding procedures we have outlined may be potentially improved and extended to
714     better represent biological entities and capture their respective relations.
715

716     Foreseeable extensions of SIMBA may include the analyses of increasingly complex datasets.
717     For example, in the analysis of spatial transcriptomics wherein transcriptomic measurements
718     are mapped to the true cell coordinates within a tissue [37], we can encode the spatial proximity
719     into a SIMBA graph. We also envision extending this framework to data describing 3-D
720     chromatin conformation wherein the interaction between DNA segments can be encoded to
721     represent how regulatory regions are linked to genes[38]. Another potential extension of SIMBA
722     could consider single-cell lineage-tracing datasets[39] wherein both cellular lineage information
723     and gene expression measurements are captured and can be potentially encoded into a SIMBA
724     graph to represent their longitudinal relations. In general, we are interested in the further
725     incorporation of external information and hierarchical relations between features in the graph.
726     We anticipate our comprehensive and extensible SIMBA framework (https://simba-
727     bio.readthedocs.io/) will provide the possibility to leverage *a priori* knowledge for graph
728     embedding and the flexibility to extend to new experimental designs.
729

730     It is likely that multi-omics assays will continue to improve as well as expand in scope. Already,
731     innovation in these data-generating technologies have outpaced the development of
732     corresponding computational frameworks required to gain integrative insights from such rich
733     data. This disparity highlights a need for methods that break through previous limitations and
734     are easily extended to future cell measurements. We believe SIMBA satisfies these conditions
735     as a comprehensive and extensible method for exploring cellular heterogeneity and
736     investigating the regulatory mechanisms that drive cellular diversity while laying a groundwork
737     for the development of new non-cluster-centric analysis methods for single cell omics data.
738
739

740     **Methods**
741

742     **Single-cell data preprocessing**
743

744        a. Single-cell RNA-seq
745        Genes expressed in fewer than three cells were filtered.   Raw counts were library size-
746        normalized and subsequently log-transformed. Optionally, variable gene selection [12] (a
747        python version is implemented in SIMBA that is inspired by Scanpy[2]) may be performed

748    to remove non-informative genes and accelerate the training procedure. Notable
749    differences in the resulting cell embeddings were not observed upon limiting feature
750    input to those identified by variable gene selection but SIMBA embeddings of non-
751    variable genes will not be generated as they are not encoded in the graph.

752

753    b. Single-cell ATAC-seq
754    Peaks present in fewer than three cells were filtered. Optionally, we implemented a
755    scalable truncated-SVD-based procedure to select variable peaks as a preliminary step
756    to additionally filter non-informative peaks and accelerate the training procedure. First
757    the top k principal components (PCs) were selected, with k chosen based on the elbow
758    plot of variance ratio. Then for each of the top $k$ PCs, peaks were automatically selected
759    based on the loadings using a knee point detection algorithm implemented by 'kneed'[40].
760    Finally, peaks selected for each PC were combined and denoted as "variable peaks".
761    Similar to the observation made with scRNA-seq data, the optional step of variable peak
762    selection has a negligible effect on the resulting cell embedding. Despite this minimal
763    impact on the resulting embedding, this feature selection step imparts a significant
764    practical advantage in reducing training procedure time.

765

766    $k$-mer and motif scanning was performed using packages 'Biostrings' and 'motifmatchr'
767    with JASPAR2020[41].  Included in the implementation of SIMBA is a convenient R
768    command line script "scan_for_kmers_motifs.R" , which will convert a list of peaks
769    (formatted in a bed file) to a sparse peaks-by-$k$-mers/motifs matrix, which is stored as
770    an hdf5-formated file.

771

772  **Graph construction (five scenarios)**

773

774    i.  Single-cell RNA-seq analysis
775    The distribution of non-zero values in the normalized gene expression matrix was first
776    approximated using a $k$-means clustering-based procedure. First, the continuous non-
777    zero values were binned into $n$ intervals (by default $n$=5). Bin widths were defined using
778    1-dimensional $k$-means clustering wherein the values in each bin are assigned to the
779    same cluster center. The continuous matrix is then converted into a discrete matrix
780    wherein$1, … , n$ are used to denote $n$ levels of gene expression. Zero values are retained
781    in this matrix.  Then the graph was constructed by encoding two types of entities, cells
782    and genes, as nodes and relations with $n$ different weights between them, i.e., $n$ levels
783    of gene expression, as edges. These $n$ relation weights range from 1.0 to 5.0 with a step
784    size of $5/n$ denoting gene expression levels (lowest: 1.0, highest: 5.0), such that edges
785    corresponding to high expression levels affect embeddings more strongly than those
786    with intermediate or low expression levels. This discretization is implemented in the
787    SIMBA package using the function, "si.tl.discretize()".

788

789   ii.  Single-cell ATAC-seq analysis
790    Peak-by-cell matrices were binarized, with "1" indicating at least one read within a peak
791    and "0" otherwise. The graph was constructed by encoding two types of entities, cells

792     and peaks, as nodes and the relation between them, denoting the presence of a given
793     peak in a cell, as edges. The single relation type was assigned with a weight of 1.0. When
794     the DNA sequence features were available, they were encoded into the graph using $k$-
795     mer and motif sequence entities as nodes. This was performed by first binarizing the
796     peak-by-$k$-mer/motif matrix then constructing an extension to the original peak/cell
797     graph using the peaks, $k$-mers, and motifs as nodes and the presence of these entities
798     within peaks as edges between these additional nodes and the peak nodes. The relation
799     between k-mers and peaks was assigned a weight of 0.02 while the relation between TF
800     motifs was assigned a weight of 0.2. Of note, $k$-mers and motifs may be used
801     independently of each other as node inputs to the graph, depending on the specific
802     analysis task.

803

804   iii.   Multimodal analysis
805     Combination of the above outlined strategies for graph construction of scRNA-seq and
806     scATAC-seq data was used to construct a multi-omics graph.

807

808   iv.   Batch correction
809     A graph for each batch was constructed as described in i).  Edges between cells of
810     different batches were inferred through a procedure based on truncated randomized
811     singular value decomposition (SVD) to link disjoint graphs of different batches. More
812     specifically, in the case of scRNA-seq data, consider two gene expression matrices
813     $X1_{n_1 \times m}$ and $X2_{n_2 \times m}$, where $n_1, n_2$ denotes the number of cells and $m$ denotes the
814     number of the shared features, i.e., variable genes, between datasets. The matrix
815     $X_{n_1 \times n_2}$ was then computed by multiplying $X1$ and $X2$:

816

817 $$X = X1 \times X2^T$$

818

819     Truncated randomized SVD was subsequently performed on $X$:

820

821 $$X \approx U \times \Sigma \times V^T$$

822

823     where $U$ is an $n_1 \times d$ matrix, $\Sigma$ is an $d \times d$ matrix, and $V$ is an $n_2 \times d$ matrix (by
824     default $d = 20$).

825

826     Both $U$ and $V$ were further $L2$ normalized. For each cell in $U$, we searched for $k$
827     nearest neighbors in $V$ and vice versa (by default, $k = 20$). Eventually, only the mutual
828     nearest neighbors between $U$ and $V$ were retained as inferred edges between cells
829     (represented as dashed lines in **Fig. 5a**). The procedure of inferring edges between
830     cells of different batches is implemented in the function "si.tl.infer_edges()" in the
831     SIMBA package.

832

833     For multiple batches, SIMBA can flexibly infer edges between any pair of datasets. In
834     practice, however edges are inferred between the largest dataset(s) or the dataset(s)
835     containing the most complete set of expected cell types and other datasets.

836

837  v. Multi-omics integration

838  scRNA-seq and scATAC-seq graphs were constructed following steps i) and ii),

839  respectively. To infer the edges between cells of scRNA-seq and scATAC-seq, gene

840  activity scores were first calculated for scATAC-seq data[3]. More specifically, for each

841  gene, peaks within 100kb upstream and downstream of the TSS were considered.

842  Peaks overlapping gene body region or within 5kb upstream of gene bodies were

843  given the weight of 1.0. Otherwise, peaks were weighted based on their distances to

844  TSS using the exponential decay function: $e^{\frac{-distance}{5000}}$. Subsequently, the gene score of

845  each gene was computed as a weighted sum of the considered peaks. These gene

846  scores were then scaled to respective gene size. These steps are implemented by the

847  function "si.tl.gene_scores()" in SIMBA. For user convenience, the SIMBA package

848  curates the gene annotations of several commonly used reference genomes, including

849  hg19, hg38, mm9, and mm10. Once gene scores were obtained, the same procedure

850  described in iv) was performed to infer edges between cells profiled by scRNA-seq and

851  scATAC-seq using the function, "si.tl.infer_edges()" in SIMBA.

852

853  The procedure of generating constructed graphs is implemented in the function,

854  "si.tl.gen_graph()" in the SIMBA package.

855

856  **Graph Embeddings with Type Constraints**

857

858  Following the construction of a multi-relational graph between biological entities, we

859  adapted graph embedding techniques from the knowledge graph and recommendation

860  systems literature to construct unsupervised representations for these entities.

861

862  We provide as input a directed graph $G = (V, E)$, where $V$ is a set of entities (vertices)

863  and $E$ is a set of edges, with a generic edge $e = (u, v)$ between a source entity $u$ and

864  destination entity $v$. We further assume that each entity has a distinct known type (e.g.,

865  cell, peak, etc.).

866

867  Graph embedding methods learn a $D$-dimensional embedding vector for each $v \in V$ by

868  optimizing a link prediction objective via stochastic gradient descent, with $D$=50 used

869  for our experiments. We will denote the full embedding matrix as $\theta \in R^{|V| \times D}$ and the

870  embedding for an entity $v$ as $\theta_v$.

871

872  For an edge $e = (u, v)$, we denote $s_e = \theta_u * \theta_v$ as the score for $e$, and optimize a

873  multi-class log loss

874  $$\mathcal{L} = -log \frac{\exp(s_e)}{\sum_{e' \in \mathcal{N}} \exp(s'_e)}$$

875

876   Where $\mathcal{N}$ is a set of "negative sampled" candidate edges generated by corrupting $e$ [42].
877   This log loss objective attempts to maximize the score for all $(u, v) \in E$ and minimize it
878   for $(u, v) \notin E$.
879
880   Negative samples are constructed by replacing either the source or target entity in the
881   target edge $e = (u, v)$ with a randomly sampled entity. However, in graphs like ours
882   where only edges between certain entity types are possible, previous work has shown
883   that it is beneficial to optimize the loss only over candidate edges that satisfy the type
884   constraints[43]. Thus, for e.g., a cell-peak edge we only sample negative candidates
885   between cell and peak entities. This modification is crucial in our setting since most
886   randomly selected edges will be of invalid type (e.g., peak-peak), forcing the
887   embeddings to primarily be optimized for irrelevant tasks (e.g., having low dot product
888   between every pair of peaks).
889
890   Furthermore, it has been frequently observed that in graphs with wide distribution of
891   node degrees, it is advantageous to sample negatives proportional to some function of
892   the node degree to produce more informative embeddings that don't merely capture
893   the degree distribution [13, 44].  For each graph edge in the dataset encountered in a
894   training batch, we produce 100 negatives by corrupting the edge with a source or
895   destination sampled uniformly from the nodes with the correct types for this relation
896   and 100 by corrupting the edge with a source or destination node sampled with
897   probability proportional to its degree[13].
898
899   As with many ML methods, graph embeddings are prone to overfitting in a low-data
900   regime (i.e., low ratio of edges to parameters). We observed overfitting measurable as a
901   gap between training and validation loss on the link prediction task, which we addressed
902   with $L2$ regularization on the embeddings $\theta$,
903

$$\mathcal{L}_{reg} = \mathcal{L} + \lambda \sum_{u \in N} \sum_{d=1}^{D} \theta_{ud}^2.$$

905
906   with $\lambda = wd * wd\_interval$. For weight decay parameter ($wd$), by default it is calculated
907   automatically as $\frac{C}{N_e}$, where $N_e$ is the training sample size (i.e., the total number of edges)
908   and $C$ is a constant. For weight decay interval ($wd\_interval$), we set it to 50 for all
909   experiments.
910
911   We use the PyTorch-BigGraph framework, which provides efficient computation of
912   multi-relation graph embeddings over multiple entity types and can scale to graphs with
913   millions or billions of entities[13]. For 1.3 million cells, the PyTorch-BigGraph training itself
914   takes only ~ 1.5 hours using 12 cores without the requirement of GPU (https://simba-
915   bio.readthedocs.io/en/latest/rna_10x_mouse_brain_1p3M.html).
916

917  The resulting graph embeddings have two desirable properties that we will take
918  advantage of:
919  1.      First-order similarity: for two entity types $T_1$, $T_2$ with a relation between them,
920  edges with high likelihood should have higher dot product; specifically, for any $u \in T_1$,
921  the predicted probability distribution over edges to $T_2$ originating from $u$ is
922  approximated as $\frac{e^{x_u * x_v}}{\sum_{v\prime \in T_2} e^{x_u * x_{v\prime}}}$ .
923  2.      Second-order similarity: within a single entity type, entities that have 'similar
924  contexts', i.e., a similar distribution of edge probabilities, should have similar
925  embeddings. Thus, the embeddings of each entity type provide a low-rank latent space
926  that encodes the similarity of those entities' edge distributions.
927
928  **Evaluation of the model during training**
929
930  During the PyTorch-BigGraph training procedure, a small percent of edges is held out
931  (by default, the evaluation fraction is set to 5%) to monitor overfitting and evaluate the
932  final model. Five metrics are computed on the reserved set of edges, including mean
933  reciprocal rank (MRR, the average of the reciprocal of the ranks of all positives), R1 (the
934  fraction of positives that rank better than all their negatives, i.e., have a rank of 1), R10
935  (the fraction of positives that rank in the top 10 among their negatives), R50 (the
936  fraction of positives that rank in the top 50 among their negatives), and AUC (Area
937  Under the Curve). By default, we show MRR along with training loss and validation loss
938  while other metrics are also available in SIMBA package (**Supplementary Fig. 1a**).  The
939  learning curves for validation loss and these metrics can be used to determine when
940  training has completed. The relative values of training and validation loss along with
941  these evaluation metrics can be used to identify issues with training (underfitting vs
942  overfitting) and tune the hyperparameters weight decay, embedding dimension, and
943  number of training epochs appropriately. For example, in **Supplementary Figure 1a**
944  training can be stopped once the validation loss plateaus. However, for most datasets
945  we find that the default parameters do not need tuning.
946
947  **Softmax transformation**
948
949  PyTorch-BigGraph training provides initial embeddings of all entities (nodes).  However,
950  entities of different types (e.g., cells vs peaks, cells of different batches or modalities)
951  have different edge distributions and thus may lie on different manifolds of the latent
952  space. To make the embeddings of entities of different types comparable, we transform
953  the embeddings of features with the Softmax function by utilizing the first-order
954  similarity between cells (reference) and features (query). In the case of batch correction
955  or multi-omics integration, the Softmax transformation is also performed based on the
956  first-order similarity between cells of different batches or modalities.
957
958  Given the initial embeddings of cells (reference) $(v_{c_1}, \ldots, v_{c_n})$ and features $(v_{f_1}, \ldots, v_{d_m})$,
959  the model-estimated probability of an edge $(c_i, f_j)$ obeys

$$P\left(v_{c_i,f_j}\right) \propto \exp\left(v_{c_i} \cdot v_{f_j}\right)$$

Therefore, if a random edge was sampled from feature $f_j$ to a cell, the model would estimate the distribution over such edges as

$$p_{c_i,f_j} = \frac{\exp\left(v_{c_i} \cdot v_{f_j}\right)}{\sum_{k=1}^{n} \exp\left(v_{c_k} \cdot v_{f_j}\right)}$$

i.e., the Softmax weights between all cells $\{c_i\}$ and the feature $f_j$. We can then compute new embeddings for features as a linear combination of the cell embeddings weighted by the edge probabilities raised to some power.

$$\hat{v}_{f_j} = \frac{\sum_{i=1}^{n} p_{c_i,f_j}^{T^{-1}} v_{c_i}}{\sum_{i=1}^{n} p_{c_i,f_j}^{T^{-1}}}$$

$T$ is a temperature hyperparameter that controls the sharpness of the weighting over cells. At $T = 1$, the cell embeddings are weighted by their estimated edge probabilities; at $T \to 0$, each feature embedding is assigned the cell embedding of its nearest neighbor; at $T \to \infty$, it becomes a discrete uniform distribution, and each query becomes the average of reference embeddings. We set $T = 0.5$ for all the analyses.

These steps are implemented in the function "si.tl.embed()" in the SIMBA package.

**Metrics to assess cell-type specificity**

Four metrics are proposed to assess the cell type specificity of each feature from different aspects, including max value (a higher value indicates higher cell-type specificity), Gini index (a higher value indicates higher cell-type specificity), standard deviation (a higher value indicates higher cell-type specificity), and entropy (a lower value indicates higher cell-type specificity). We observe these four metrics generally give consistent results. For SIMBA metric plot, by default, Gini index is plotted against max value. For feature $f_j$ :

The max value is defined as the average normalized similarity of top $k$ cells (by default, $k$=50). The similarity normalization function is defined as:

$$norm(x_i) = x_i - log\frac{\sum_{j=1}^{n} \exp\left(x_j\right)}{n}$$

Where $i = 1, \dots, n$. $n$ is the number of cells and $x_i$ represents the dot product of $\hat{v}_{f_j}$ and the embedding of cell $i$.

999

1000    The max value is computed as:

1001

1002
$$\max(f_j) = \frac{\sum_{i=1}^{k} norm(x_i)}{k}$$

1003

1004    The Gini index is computed as:

1005

1006
$$\text{gini}(f_j) = \frac{\sum_{i=1}^{n}(2i - n - 1) * p_{c_i, f_j}}{n \sum_{i=1}^{n} p_{c_i, f_j}}$$

1007

1008    The standard deviation is computed as:

1009

1010
$$\text{std}(f_j) = \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (p_{c_i, f_j} - \mu)^2}$$

1011    Where $\mu = \frac{1}{n}\sum_{i=1}^{n} p_{c_i, f_j}$.

1012

1013    Entropy is computed as:

1014
$$\text{entropy}(f_j) = -\sum_{i=1}^{n} p_{c_i, f_j} \log(p_{c_i, f_j})$$

1015

1016    **Queries of entities in SIMBA space**

1017

1018    The informative SIMBA embedding space serves as a database of entities including cells
1019    and features. To query the "SIMBA database" for the neighboring entities of a given cell
1020    or feature, we first build a k-d tree of all entities based on their SIMBA embeddings. We
1021    then search for the nearest neighbors in the tree using Euclidean distance. To do so,
1022    SIMBA query can perform either K-nearest neighbors (KNN) or nearest neighbor search
1023    within a specified radius. SIMBA also provides the option to limit the search to entities
1024    of certain types, which is useful when a certain type of entity significantly outnumbered
1025    others. For example, the K nearest features of a given cell may be all peaks while genes
1026    are the features of interest. In this case, SIMBA allows users to add "filters" to ensure
1027    that nearest neighbor search is performed within the specified types of entities. This
1028    procedure is implemented in the function "st.tl.query()" and its visualization is
1029    implemented in the function "st.pl.query()" in the SIMBA package.

1030

1031    **Identification of master regulators**

1032

1033    To identify master regulators, we take into consideration both the cell type specificity of
1034    each pair of TF motif and TF gene and the distance between them. More specifically, for

1035       each TF motif, first its distances (Euclidean distance by default) to all the genes are
1036       calculated in the SIMBA embedding space. Then the rank of this TF gene among all these
1037       genes is computed. In addition, we also assess the cell type specificity of this pair of TF
1038       motif and TF gene based on SIMBA metrics (by default, max value and Gini index are
1039       used). The same procedure is performed for all TFs. Finally, we identify master
1040       regulators by filtering out TFs with low cell-type specificity and scoring them based on
1041       TF gene rank. This procedure is implemented in the function
1042       "st.tl.find_master_regulators()" in SIMBA package.

1043

1044

1045   **Identification of TF target genes**

1046       Given a master regulator, its target genes are identified by comparing the locations of
1047       the TF gene, TF motif, and the peaks near the genomic loci of candidate target genes in
1048       the SIMBA co-embedding space (**Fig. 4e**). More specifically we first search for $k$ nearest
1049       neighbor genes around the motif (TF motif) and the gene (TF gene) of this master
1050       regulator, respectively ($k$ = 200 by default). The union of these neighbor genes is the
1051       initial set of candidate target genes. These genes are then filtered based on the criterion
1052       that open regions (peaks) within 100kb upstream and downstream of the TSS of a
1053       putative target gene must contain the TF motif.

1054       Next, for each candidate target gene, we compute four types of distances in SIMBA
1055       embedding space: distances between the embeddings of 1) the candidate target gene
1056       and TF gene; 2) the candidate target gene and TF motif; 3) peaks near the genomic locus
1057       of the candidate target gene and TF motif; 4) peaks near the genomic locus of the
1058       candidate target gene and the candidate gene. All the distances (Euclidean distances by
1059       default) are converted to ranks out of all genes or all peaks to make the distances
1060       comparable across different master regulators.

1061       The final list of target genes is decided using the calculated ranks based on two criteria:
1062       1) at least one of the nearest peaks to TF gene or TF motif is within a predetermined
1063       range (top 1,000 by default); 2) the average rank of the candidate target gene is within a
1064       predetermined range (top 5,000 by default). This procedure is implemented in the
1065       function "st.tl. find_target_genes ()" in SIMBA.

1066   **Benchmarking scATAC-seq computational methods**
1067
1068       To compare SIMBA to other scATAC-seq computational methods including SnapATAC [4],
1069       Cusanovich2018[45], and cisTopic[46], we employed the previously developed benchmarking
1070       framework from Chen et al[14](**Supplementary Table 1**). This framework evaluates
1071       different methods based on their ability to distinguish cell types. We applied three
1072       clustering algorithms: k-means clustering, hierarchical clustering, and Louvain on the
1073       feature matrix derived from each method.
1074

1075  For datasets with ground-truth (FACS-sorted labels or known tissue labels), including
1076  simulated bone marrow data, Buenrostro 2018, and sci-ATAC-seq subset, three metrics
1077  including Adjusted Rand Index (ARI), Adjusted Mutual Information (AMI), and
1078  Homogeneity are applied to evaluate the performance. ARI measures the similarity
1079  between two clusters, comparing all pairs of samples assigned to matching or different
1080  clusters in the predicted clustering solution vs the true cluster/cell type label.  AMI
1081  describes an observed frequency of co-occurrence compared to an expected frequency
1082  of co-occurrence between two variables, informing the mutual dependence or strength
1083  of association of these two variables. Homogeneity measures whether a clustering
1084  algorithm preserves cluster assignments towards samples that belong to a single class. A
1085  higher metric value indicates a better clustering solution.

1086

1087  For 10x PBMCs dataset with no ground truth, the Residual Average Gini Index (RAGI)
1088  proposed in the benchmarking study[14] is used as the clustering evaluation metric. RAGI
1089  measures the relative exclusivity of marker genes to their corresponding clusters in
1090  comparison to housekeeping genes, which should demonstrate low specificity to any
1091  given cluster. In brief, the mean Gini Index is computed for both marker genes and
1092  housekeeping genes. The difference between the means is computed to obtain the
1093  average residual specificity (i.e., RAGI) of a clustering solution with respect to marker
1094  genes. A higher RAGI indicates a better separation of biologically distinct clusters.

1095

1096  **Benchmarking single-cell batch correction methods**

1097

1098  The batch correction performance of SIMBA was compared to Seurat3[12], LIGER[11] and
1099  Harmony[10] in two benchmark datasets: the mouse atlas dataset and the human
1100  pancreas dataset (**Supplementary Table 1**). For Seurat3, LIGER and Harmony, the batch
1101  correction was done with the same parameters used in a previous benchmark study[25].

1102

1103  To evaluate the batch integration performance, average Silhouette width (ASW),
1104  adjusted Rand index (ARI), and local inverse Simpson's index (LISI)[10] were calculated for
1105  the batches and cell types using the Euclidean distance as described in a previous
1106  benchmark[25]. To make a fair evaluation, only the cell types that are present in all
1107  batches were considered. We used the same number of dimensions (50) for these
1108  methods and all other parameters were set as in the benchmark.

1109

1110  **Average Silhouette width (ASW)**

1111

1112  Average Silhouette width is the mean value of Silhouette scores calculated from each
1113  cell. Silhouette width measures the relative closeness of cells with the same label
1114  compared to the cells with the different label and ranges from -1 to +1. Silhouette score
1115  for a data point with a label is calculated as

1116

1117
$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

1118
1119 where $a(i)$ is the distance to the closest point with the same label, and $b(i)$ is the
1120 distance to the closest point with different labels. A high Silhouette score means the
1121 point is located more closely with the same label, where a low Silhouette score closer to
1122 -1 means the point is located closer with different labels than that of itself. The ideal
1123 batch correction result will give a low ASW score for batch labels as the point is well
1124 mixed with other batches and a high ASW score for the cell type labels as the cells of the
1125 same cell type should cluster together after the batch correction. The final score is
1126 calculated as the median ASW scores from 20 subsets of randomly sampled 80% cells.
1127
1128 **Average Rand Index (ARI)**
1129
1130 To evaluate the cell type purity, the true cell type labels and the k-means clustering
1131 solution were used to calculate the cell type ARI. To evaluate the batch correction
1132 performance, the true batch labels and the k-means clustering solution were used to
1133 calculate the batch ARI. The final ARI was calculated as the median ARI scores of 20
1134 subsets comprised of randomly sampled 80% cells for batches and cell types,
1135 respectively.  A superior batch correction will have a high cell type ARI (high agreement
1136 between the clustering solution and the true cell type labels), and a low batch ARI ( the
1137 clustering solution is not mainly driven by batches and clusters contain cells with well-
1138 mixed batch labels).
1139
1140 **Local Inverse Simpson's Index (LISI)**
1141
1142 Local Inverse Simpson's Index (LISI) [10] measures the local batch and cell type mixing. For
1143 each data point, it considers the Gaussian kernel weighted distribution of labels in its
1144 neighborhood with a perplexity argument. We set perplexity to 50 40 as in the previous
1145 benchmark study. Using the weighted neighborhood label distribution, the inverse
1146 Simpson's index is calculated as $\frac{1}{\sum_l p(l)}$ where $l$ is the batch or cell type labels and $p(l)$ is
1147 the probability of each label in the local neighborhood obtained with the kernel. For
1148 each cell, the LISI is the expected number of cells to be sampled locally before a cell of
1149 the same label is sampled. A perfect batch correction will have a cell type LISI (cLISI) of 1
1150 and a batch LISI (integration LISI, iLISI) close to the number of batches. The final LISI
1151 score was calculated as the average LISI scores of all cells.
1152
1153 Further details are described in **Supplementary Note 2**.
1154
1155
1156 **Benchmarking single cell multi-omics integration methods**
1157
1158 Two pairs of scRNA-seq and scATAC-seq datasets manually split from the dual-omics
1159 SHARE-seq mouse skin dataset and 10X PBMCs dataset respectively were used for the
1160 modality integration task. For Seurat3 and LIGER, the parameters and preprocessing

were done as described in their documentations. However, for the LIGER analysis of the SHARE-seq mouse skin dataset the parameter 'lambda' was set to 30 and the 'ref_dataset' was set to scATAC-seq to get a better alignment. For the Raw results, the activity matrix of scATAC-seq was constructed using Seurat3 and the first 20 PCs of the scRNA-seq count matrix and the activity matrix were used for the comparison. The integration results generated by each method were evaluated with four metrics— Anchoring distance, anchoring distance rank, Silhouette index, and cluster agreement— as described below.

**Anchoring distance**
The Anchoring distance  was proposed in Dou et al., 2020[47] and is the normalized distance between the matched cells of two modalities (e.g. RNA and ATAC). Here we considered the Euclidean distance and normalized the distance by the mean of the distances calculated between random pairs of cells. The number of pairs randomly sampled was set to 10% of the total number of cells.

**Anchoring distance rank**
Given that the anchoring distance does not account for the local density of cells, we propose a new metric entitled *anchoring distance rank* (ADR). The ADR is based on the normalized rank of the distance between the matched cells of two modalities. For each cell $x_{ij}$ with cell identity i and modality j, the distance between the cell and all the other cells of the other modality j', $d(x_{ij}, x_{kj'}), k = 1, \dots, N$ is calculated, where N is the total number of cells. Then the rank of $r_i = d(x_{ij}, x_{ij'})$ within the calculated distances is normalized by the number of pairs $N - 1$ to obtain the final anchoring rank $m_i = \frac{r_i - 1}{N - 1}$. For each cell, an anchoring rank of 0 indicates an ideal modality integration performance as the matched cells are closest to each other in the embedding.

**Silhouette index**
The silhouette index was calculated as described in 10) based on the cluster assignment wherein each cluster consists of two cells, one cell from a scRNA-seq dataset and one cell from a scATAC-seq dataset.

**Fraction in the same cluster**
Fraction in the same cluster was calculated as the fraction of the matched cells from two modalities in the same cluster. The clusters of cells were generated using Louvain algorithm and the number of clusters is equal to the number of cell types in the dataset.

Further details are described in **Supplementary Note 3**.

## Data availability:

1203    All the datasets used in this study (eight scRNA-seq datasets, four scATAC-seq datasets, and
1204    three dual-omics datasets) are summarized in **Supplementary Table 1**. All these datasets are
1205    curated in the SIMBA package, and they can be easily downloaded and imported directly to
1206    reproduce the analyses presented in this manuscript.
1207

## Code availability:

1209

1210    We provide a comprehensive Python package 'simba' available at
1211    https://anaconda.org/bioconda/simba and https://github.com/pinellolab/simba. All the
1212    proposed procedures are implemented in the "simba" package. 'simba' can be easily installed
1213    with conda "*conda install simba*". We also built a website (https://simba-bio.readthedocs.io),
1214    providing a detailed introduction of the 'simba' software and several SIMBA tutorials for
1215    different types of single-cell analyses presented in this manuscript.
1216

## Acknowledgements

1218

1228

## Author contributions

1230

1231    H.C. and L.P. conceived this project. H.C. developed SIMBA, wrote the SIMBA package, and
1232    performed SIMBA analysis on all datasets. A.L. contributed to the adaption of PyTorch-BigGraph
1233    to single cell analysis. J.R. and H.C. performed the comparison analysis on batch correction and
1234    data integration. M.E.V. and H.C. performed the comparison analysis on scATAC-seq data. L.P.
1235    and A.L. provided guidance and supervised this project. All the authors wrote and approved the
1236    final manuscript.
1237

## Competing interests

1239

1240    The authors declare that they have no competing interests.
1241

1242

1243

1244

## References:

1. Satija, R., Farrell, J.A., Gennert, D., Schier, A.F. & Regev, A. Spatial reconstruction of single-cell gene expression data. *Nat Biotechnol* **33**, 495-502 (2015).

2. Wolf, F.A., Angerer, P. & Theis, F.J. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol* **19**, 15 (2018).

3. Granja, J.M. et al. ArchR is a scalable software package for integrative single-cell chromatin accessibility analysis. *Nat Genet* **53**, 403-411 (2021).

4. Fang, R. et al. Comprehensive analysis of single cell ATAC-seq data with SnapATAC. *Nat Commun* **12**, 1337 (2021).

5. Kiselev, V.Y., Andrews, T.S. & Hemberg, M. Challenges in unsupervised clustering of single-cell RNA-seq data. *Nat Rev Genet* **20**, 273-282 (2019).

6. Cortal, A., Martignetti, L., Six, E. & Rausell, A. Gene signature extraction and cell identity recognition at the single-cell level with Cell-ID. *Nat Biotechnol* (2021).

7. Vandenbon, A. & Diez, D. A clustering-independent method for finding differentially expressed genes in single-cell transcriptome data. *Nat Commun* **11**, 4318 (2020).

8. Dann, E., Henderson, N.C., Teichmann, S.A., Morgan, M.D. & Marioni, J.C. Differential abundance testing on single-cell data using k-nearest neighbor graphs. *Nat Biotechnol* (2021).

9. Hao, Y. et al. Integrated analysis of multimodal single-cell data. *Cell* (2021).

10. Korsunsky, I. et al. Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat Methods* **16**, 1289-1296 (2019).

11. Welch, J.D. et al. Single-Cell Multi-omic Integration Compares and Contrasts Features of Brain Cell Identity. *Cell* **177**, 1873-1887 e1817 (2019).

12. Stuart, T. et al. Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902 e1821 (2019).

13. Lerer, A. et al. Pytorch-biggraph: A large-scale graph embedding system. *arXiv preprint arXiv:1903.12287* (2019).

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

15. Buenrostro, J.D. et al. Integrated Single-Cell Analysis Maps the Continuous Regulatory Landscape of Human Hematopoietic Differentiation. *Cell* **173**, 1535-1548 e1516 (2018).

16. Ferreira, R., Ohneda, K., Yamamoto, M. & Philipsen, S. GATA1 function, a paradigm for transcription factors in hematopoiesis. *Molecular and cellular biology* **25**, 1215-1227 (2005).

17. Tijchon, E., Havinga, J., Van Leeuwen, F. & Scheijen, B. B-lineage transcription factors and cooperating gene lesions required for leukemia development. *Leukemia* **27**, 541-552 (2013).

18. Friedman, A. Transcriptional control of granulocyte and monocyte development. *Oncogene* **26**, 6816-6828 (2007).

19. 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).

1288  20.  Moriguchi, T. & Yamamoto, M. A regulatory network governing Gata1 and Gata2 gene
1289        transcription orchestrates erythroid lineage differentiation. *International journal of*
1290        *hematology* **100**, 417-424 (2014).
1291  21.  Chen, S., Lake, B.B. & Zhang, K. High-throughput sequencing of the transcriptome and
1292        chromatin accessibility in the same cell. *Nat Biotechnol* (2019).
1293  22.  Ma, S. et al. Chromatin Potential Identified by Shared Single-Cell Profiling of RNA and
1294        Chromatin. *Cell* (2020).
1295  23.  Cao, J. et al. Joint profiling of chromatin accessibility and gene expression in thousands
1296        of single cells. *Science* **361**, 1380-1385 (2018).
1297  24.  Zhu, C. et al. An ultra high-throughput method for single-cell joint analysis of open
1298        chromatin and transcriptome. *Nat Struct Mol Biol* (2019).
1299  25.  Tran, H.T.N. et al. A benchmark of batch-effect correction methods for single-cell RNA
1300        sequencing data. *Genome Biol* **21**, 12 (2020).
1301  26.  Leek, J.T. et al. Tackling the widespread and critical impact of batch effects in high-
1302        throughput data. *Nature Reviews Genetics* **11**, 733-739 (2010).
1303  27.  Han, X. et al. Mapping the mouse cell atlas by microwell-seq. *Cell* **172**, 1091-1107. e1017
1304        (2018).
1305  28.  Tabula Muris, C. et al. Single-cell transcriptomics of 20 mouse organs creates a Tabula
1306        Muris. *Nature* **562**, 367-372 (2018).
1307  29.  Baron, M. et al. A single-cell transcriptomic map of the human and mouse pancreas
1308        reveals inter-and intra-cell population structure. *Cell systems* **3**, 346-360. e344 (2016).
1309  30.  Muraro, M.J. et al. A single-cell transcriptome atlas of the human pancreas. *Cell systems*
1310        **3**, 385-394. e383 (2016).
1311  31.  Wang, Y.J. et al. Single-cell transcriptomics of the human endocrine pancreas. *Diabetes*
1312        **65**, 3028-3038 (2016).
1313  32.  Segerstolpe, Å. et al. Single-cell transcriptome profiling of human pancreatic islets in
1314        health and type 2 diabetes. *Cell metabolism* **24**, 593-607 (2016).
1315  33.  Ietswaart, R., Gyori, B.M., Bachman, J.A., Sorger, P.K. & Churchman, L.S. GeneWalk
1316        identifies relevant gene functions for a biological context using network representation
1317        learning. *Genome Biol* **22**, 55 (2021).
1318  34.  Yuan, H., Kshirsagar, M., Zamparo, L., Lu, Y. & Leslie, C.S. BindSpace decodes
1319        transcription factor binding signals by large-scale sequence embedding. *Nat Methods*
1320        (2019).
1321  35.  Li, H., Xiao, X., Wu, X., Ye, L. & Ji, G. scLINE: A multi-network integration framework
1322        based on network embedding for representation of single-cell RNA-seq data. *J Biomed*
1323        *Inform* **122**, 103899 (2021).
1324  36.  Buterez, D., Bica, I., Tariq, I., Andrés-Terré, H. & Liò, P. CELLVGAE: AN UNSUPERVISED
1325        SCRNA-SEQ ANALYSIS WORKFLOW WITH GRAPH ATTENTION NETWORKS. *bioRxiv*
1326        *2020.12.20.423645v1* (2020).
1327  37.  Longo, S.K., Guo, M.G., Ji, A.L. & Khavari, P.A. Integrating single-cell and spatial
1328        transcriptomics to elucidate intercellular tissue dynamics. *Nat Rev Genet* **22**, 627-644
1329        (2021).
1330  38.  Kempfer, R. & Pombo, A. Methods for mapping 3D chromosome architecture. *Nat Rev*
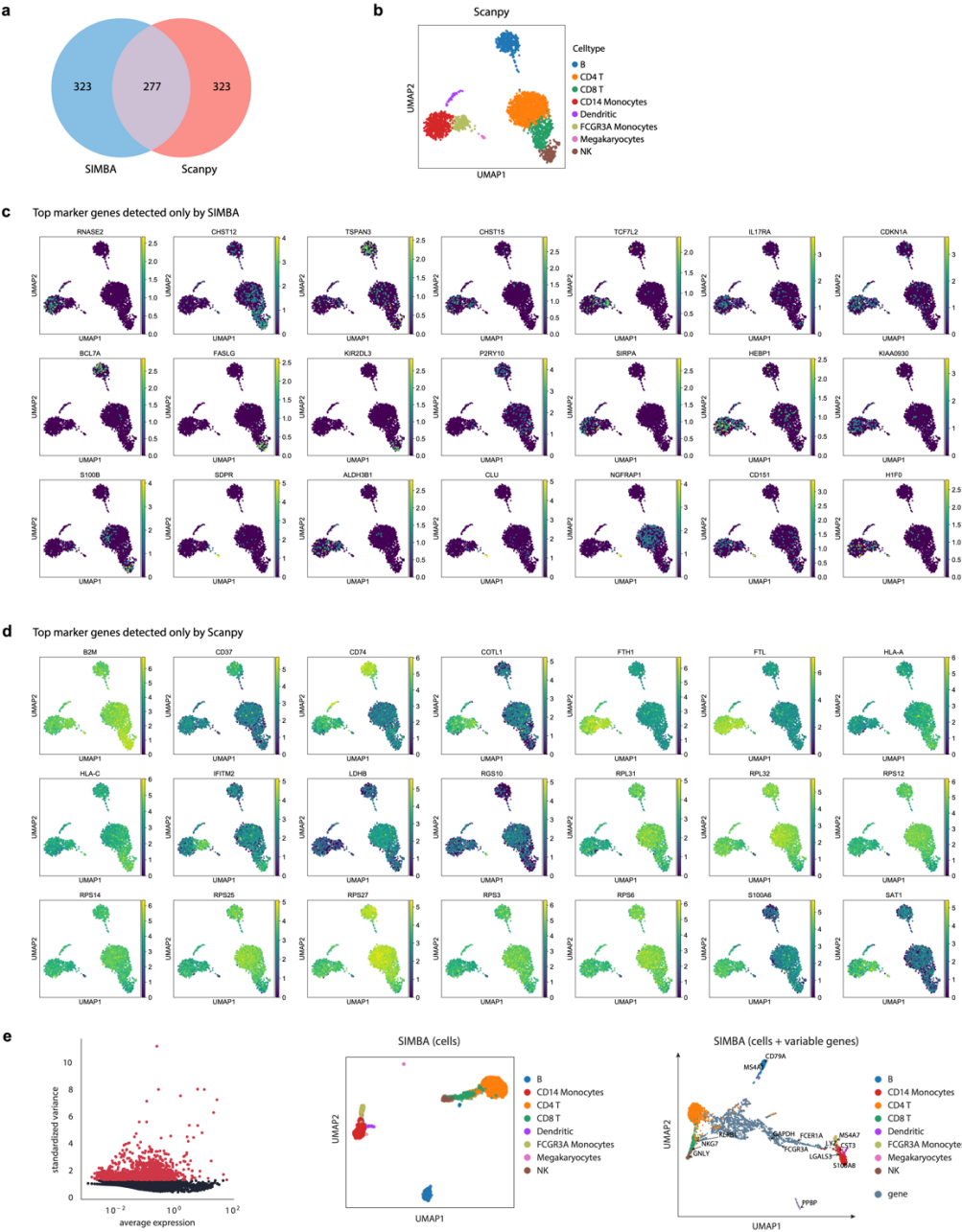1331        *Genet* **21**, 207-226 (2020).

1332 39.  VanHorn, S. & Morris, S.A. Next-Generation Lineage Tracing and Fate Mapping to
1333       Interrogate Development. *Dev Cell* **56**, 7-21 (2021).
1334 40.  Satopaa, V., Albrecht, J., Irwin, D. & Raghavan, B. in 2011 31st international conference
1335       on distributed computing systems workshops 166-171 (IEEE, 2011).
1336 41.  Fornes, O. et al. JASPAR 2020: update of the open-access database of transcription
1337       factor binding profiles. *Nucleic acids research* **48**, D87-D92 (2020).
1338 42.  Kadlec, R., Bajgar, O. & Kleindienst, J. Knowledge base completion: Baselines strike back.
1339       *arXiv preprint arXiv:1705.10744* (2017).
1340 43.  Krompaß, D., Baier, S. & Tresp, V. in International semantic web conference 640-655
1341       (Springer, 2015).
1342 44.  Mikolov, T., Chen, K., Corrado, G. & Dean, J. Efficient estimation of word representations
1343       in vector space. *arXiv preprint arXiv:1301.3781* (2013).
1344 45.  Cusanovich, D.A. et al. The cis-regulatory dynamics of embryonic development at single-
1345       cell resolution. *Nature* **555**, 538-542 (2018).
1346 46.  Bravo Gonzalez-Blas, C. et al. cisTopic: cis-regulatory topic modeling on single-cell ATAC-
1347       seq data. *Nat Methods* **16**, 397-400 (2019).
1348 47.  Dou, J. et al. Unbiased integration of single cell multi-omics data. *bioRxiv*,
1349       2020.2012.2011.422014 (2020).
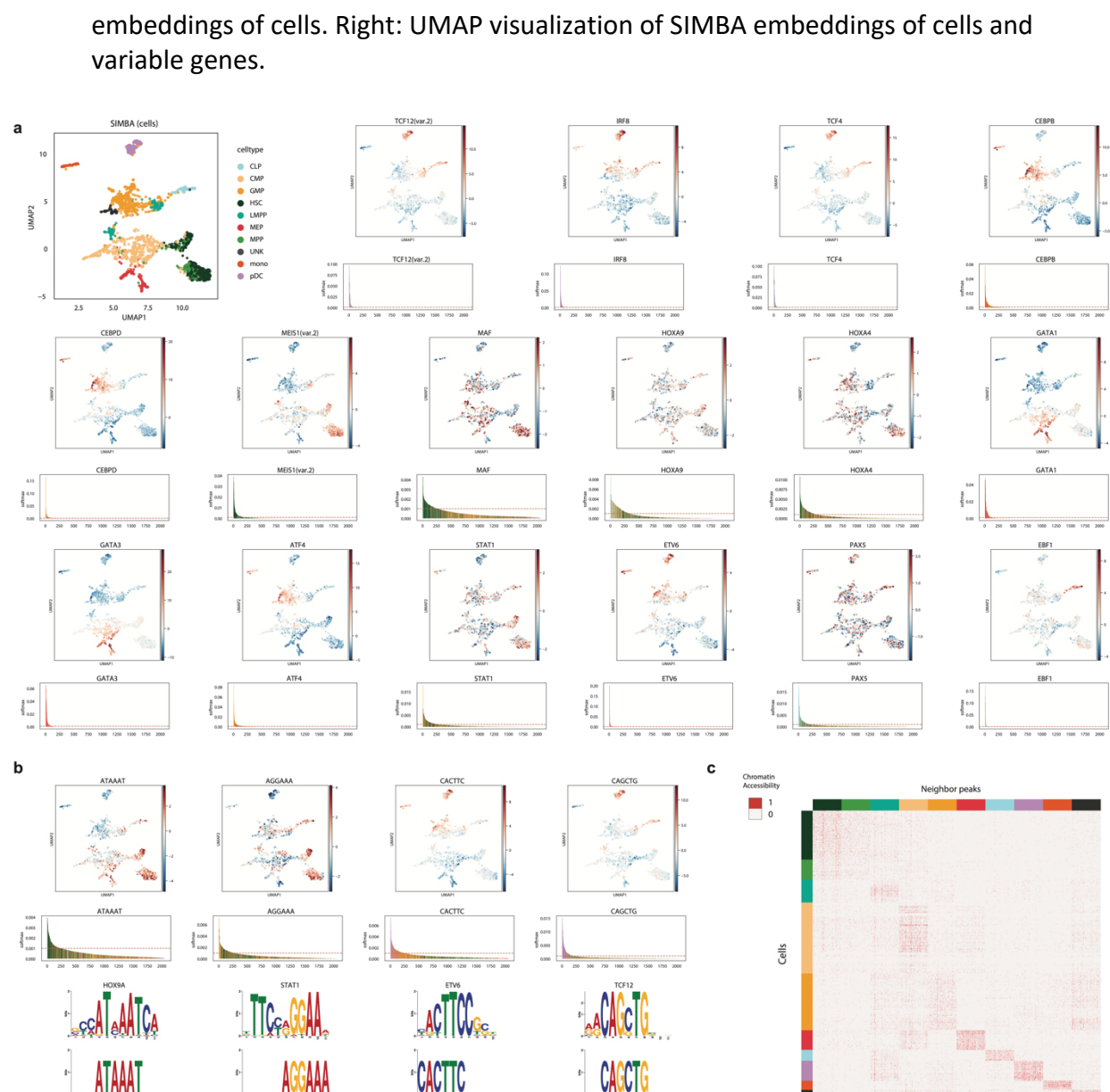1350
1351

## Supplementary Figures

1353



1354
1355

**Supplementary Figure 1**. SIMBA analysis of the scRNA-seq 10x PBMCs dataset.

**a.** Three default metrics used to evaluate SIMBA training procedure, including training loss (top), validation loss (middle), mean reciprocal rank (MRR)

**b.** SIMBA metric plots of genes. All the genes are plotted according to the Gini index against max score, standard deviation (std) against max score, and entropy against max score, respectively. The same set of genes as in Figure 2c are highlighted.

**c.** UMAP visualization of the SIMBA embeddings of cells and the SIMBA embeddings of cells and all genes. Genes highlighted in (b) are also highlighted in the UMAP plot.

**d.** UMAP visualization of the SIMBA embeddings of cells, colored by gene expression of the genes highlighted in (b).

**e.** SIMBA barcode plots of the genes highlighted in (b).

1367

1368



1369

1370 **Supplementary Figure 2**. Comparison of SIMBA with Scanpy on the scRNA-seq 10x PBMCs
1371 dataset.
1372     **a.** Venn diagram of top marker genes identified by SIMBA and Scanpy
1373     **b.** Scanpy-derived UMAP visualization of cells colored by cell type
1374     **c.** Top marker genes detected only by SIMBA. Colored by intensity of gene expression.
1375     **d.** Top marker genes detected only by Scanpy. Colored by intensity of gene expression.
1376     **e.** SIMBA embedding result after implementing variable gene selection. Left: variable gene
1377        selection step implemented in SIMBA. Middle: UMAP visualization of SIMBA

1378    embeddings of cells. Right: UMAP visualization of SIMBA embeddings of cells and
1379    variable genes.
1380



1381
1382
1383    **Supplementary Figure 3**. SIMBA analysis of the *Buenrostro2018* dataset
1384
1385    **a.** UMAP visualization of SIMBA embeddings of cells colored by cell type (top-left), and TF
1386    activity scores of TF motifs calculated with chromVAR, respectively. The SIMBA barcode
1387    plot of each TF motif is shown below the UMAP plot.
1388    **b.** Top: UMAP visualization of SIMBA embeddings of cells colored by TF activity scores of k-
1389    mers calculated with chromVAR. Middle: SIMBA barcode plots of the corresponding k-
1390    mers. Bottom: the matching known motif against the enriched k-mer sequences.
1391    **c.** Heatmap of cells against neighboring peaks of each cell type that are selected in the
1392    SIMBA co-embedding space. Chromatin accessibility is binary and colored accordingly.
1393

1394
1395
1396
1397
1398



1399
1400

**Supplementary Figure 4**. Comparison of SIMBA performance using scATAC-seq peaks and DNA sequence content vs only scATAC-seq peaks. **Top**: UMAP visualization of SIMBA embeddings of cells for each indicated dataset generated from only scATAC-seq peak information. **Bottom**: UMAP visualization of SIMBA embeddings of cells for each indicated dataset generated from scATAC-seq peak information and DNA sequence content information.

1406

**Supplementary Figure 5**. Benchmark of SIMBA against top-performing scATAC-seq analysis methods.

Top: Evaluation of SIMBA and other methods including cisTopic, SnapATAC, *Cusanovich2018* for scATAC-seq analysis using metrics 1) ARI, AMI, and Homogeneity for datasets with ground truth cell type labels and 2) Residual Average Gini Index (RAGI) for the 10x PBMCs dataset without ground truth labels.

Bottom: UMAP visualization of feature matrices produced by each method on each dataset colored by cell type annotation or cluster label.

**Supplementary Figure 6**. SIMBA multimodal analysis of the SHARE-seq hair follicle dataset.

a. SIMBA embedding results when both gene expression and chromatin accessibility are encoded in the graph. Left: UMAP visualization of SIMBA embeddings of cells and genes. Middle: UMAP visualization of SIMBA embeddings of cells along with genes, TF motifs, and k-mers. Right: UMAP visualization of SIMBA embeddings of cells along with genes, peaks, TF motifs, and k-mers.

b. SIMBA embedding results when only gene expression is encoded in the graph. Left: UMAP visualization of SIMBA embeddings of cells. Right: UMAP visualization of SIMBA embeddings of cells and variable genes.

c. SIMBA embedding results when only chromatin accessibility is encoded in the graph. Left: UMAP visualization of SIMBA embeddings of cells. Right: UMAP visualization of SIMBA embeddings of cells and peaks.
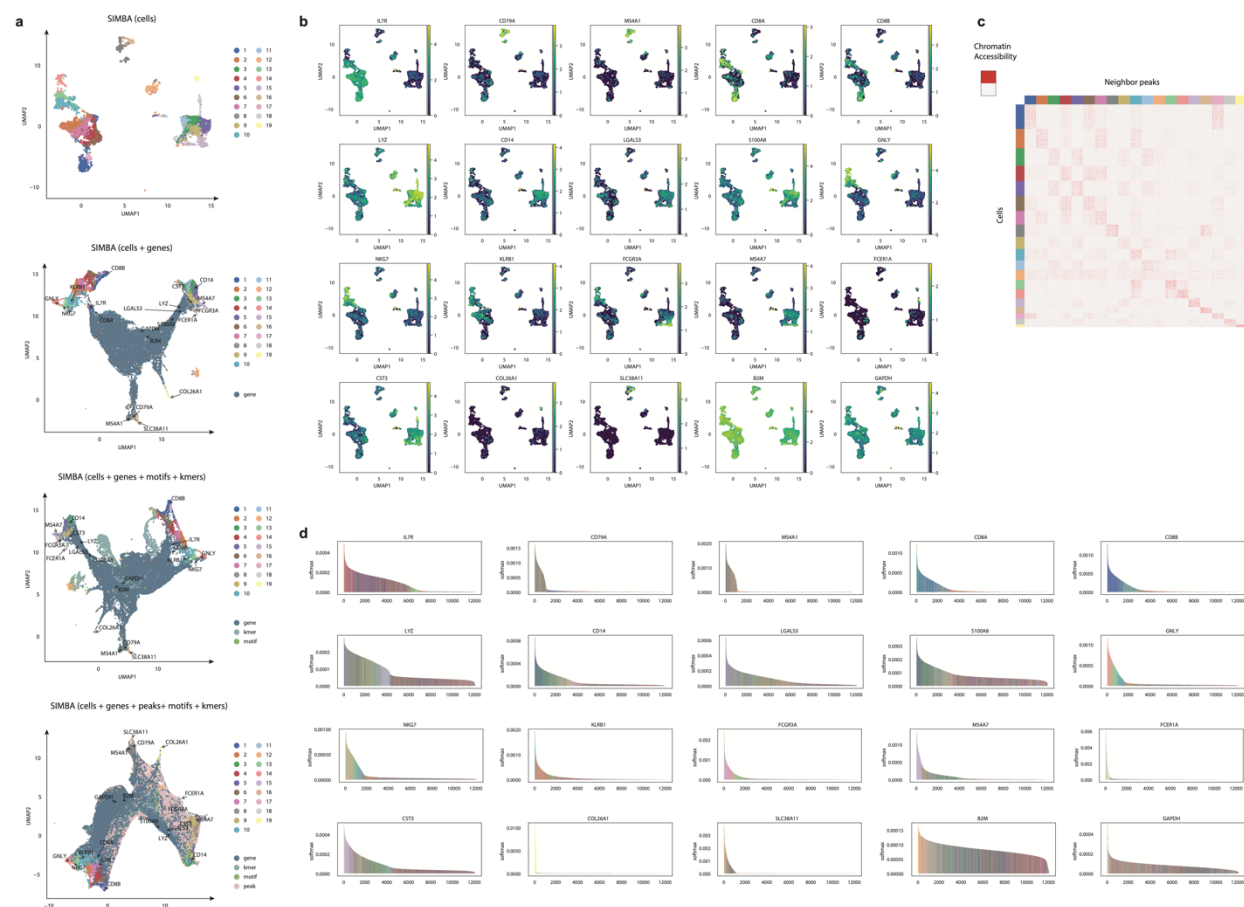
**Supplementary Figure 7**. Cell type specific marker genes and the target genes of master regulators identified by SIMBA in the SHARE-seq hair follicle subset dataset.

    **a.** UMAP visualization of SIMBA embeddings of cells colored by cell type and gene expression intensity.

    **b.** SIMBA barcode plots of each gene plotted above.

    **c.** SIMBA barcode plots of TF motifs *Lef1* and *Hoxc13.*

    **d.** SIMBA barcode plots of peaks near the loci of *Lef1* and *Hoxc13.*

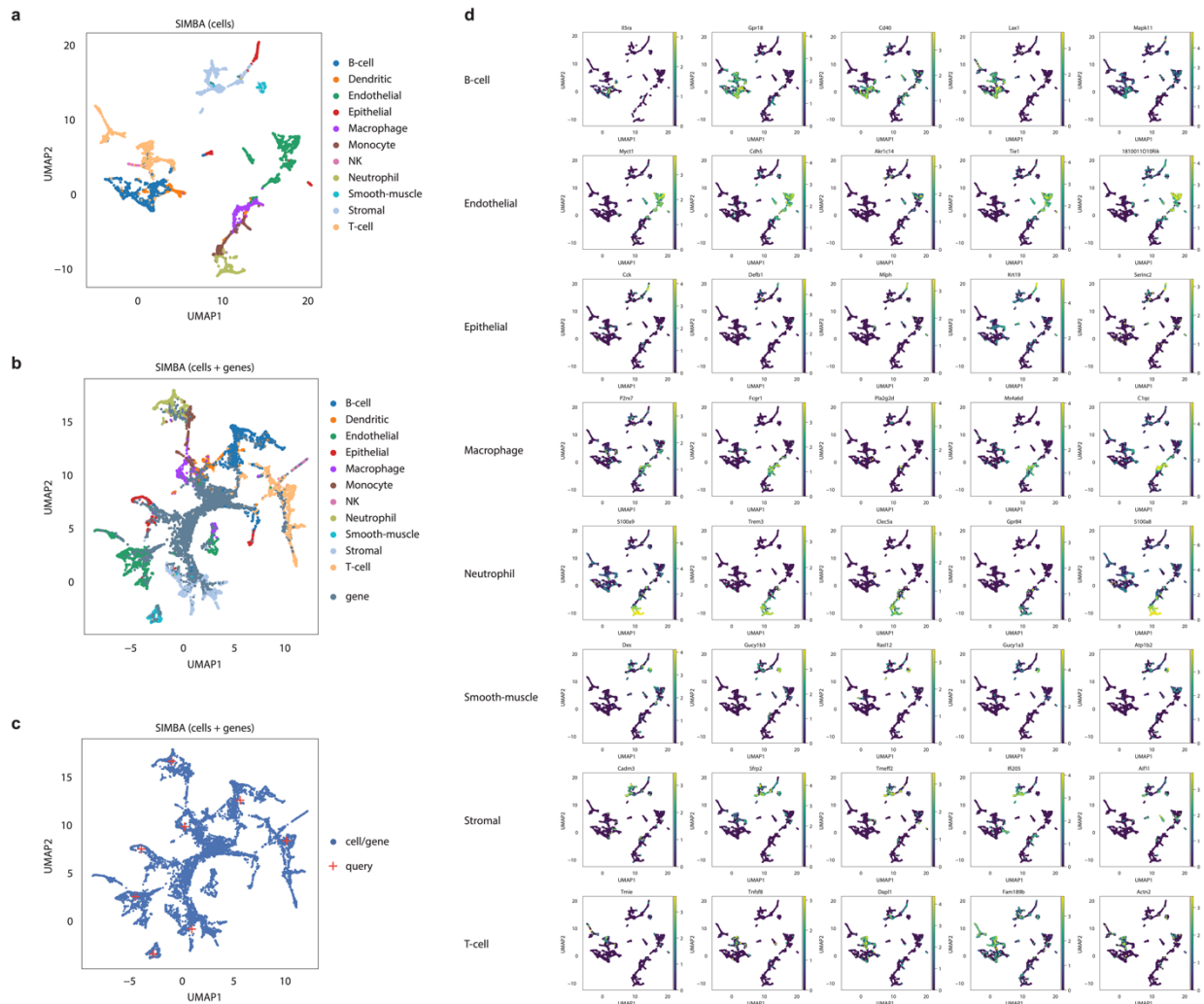    **e.** Top 30 target genes of the master regulators *Relb, Gata6,* and *Nfatc1* as inferred by SIMBA.

1448
1449
1450    **Supplementary Figure 8**. SIMBA multimodal analysis of the SNARE-seq mouse cerebral cortex
1451    dataset.
1452
1453    **a.** From top to bottom: UMAP visualization of SIMBA embeddings of (i) cells (ii) genes
1454         alongside cells (iii) genes, motifs, and k-mers alongside cells (iv) genes, peaks, motifs,
1455         and k-mers alongside cells.
1456    **b.** UMAP visualization of SIMBA embeddings of cells colored by indicated gene expression
1457         intensity.
1458    **c.** Heatmap of cells against neighboring peaks of each cell type that are selected in the
1459         SIMBA co-embedding space. Chromatin accessibility is binary and colored accordingly.
1460    **d.** SIMBA barcode plots of the genes highlighted in (a).
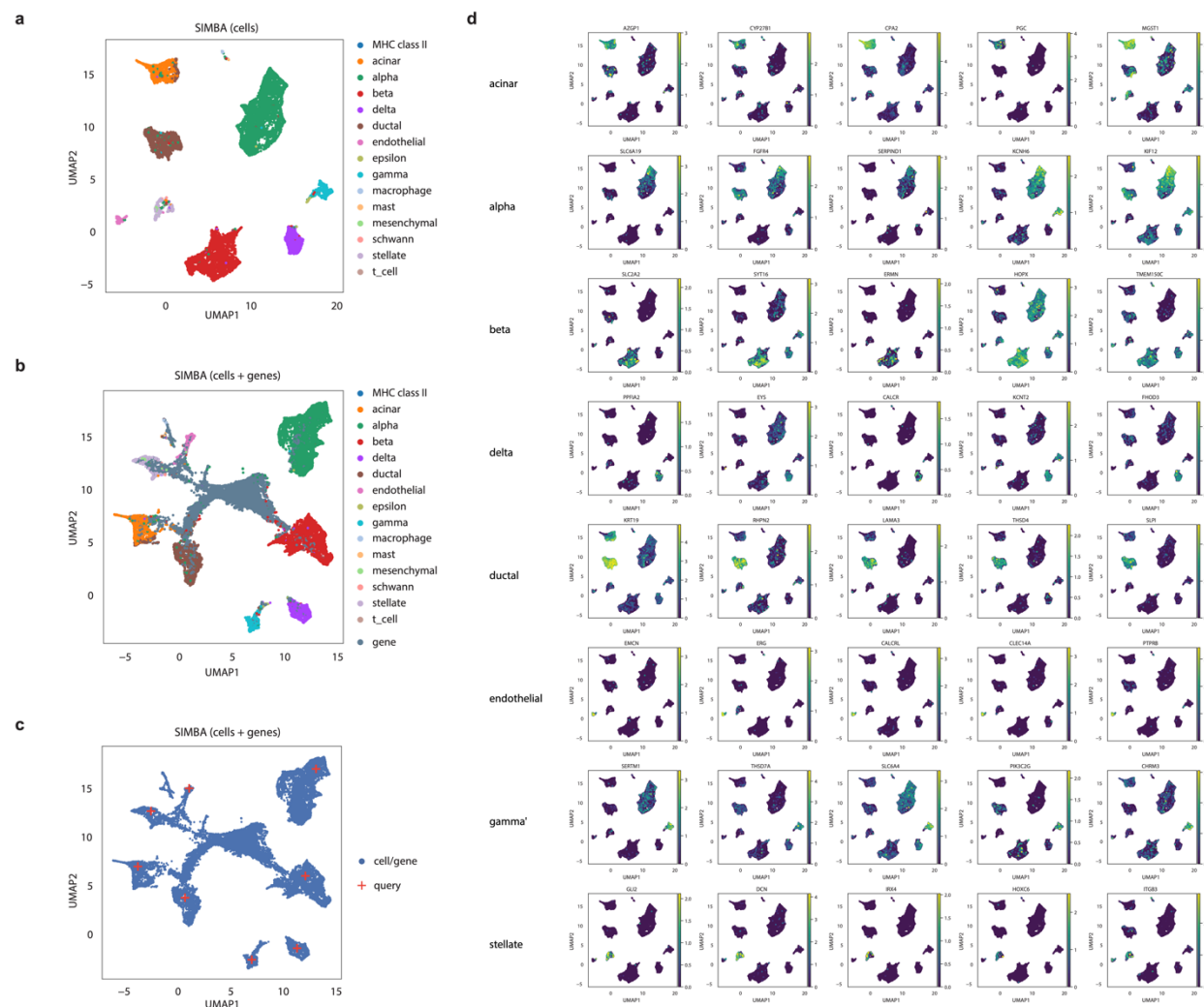1461

**Supplementary Figure 9**. SIMBA multimodal analysis of the 10x multiome PBMCs dataset.

a. From top to bottom: UMAP visualization of SIMBA embeddings of (i) cells (ii) genes alongside cells (iii) genes, motifs, and k-mers alongside cells (iv) genes, peaks, motifs, and k-mers alongside cells.

b. UMAP visualization of SIMBA embeddings of cells colored by indicated gene expression intensity.

c. Heatmap of cells against neighboring peaks of each cluster that are selected in the SIMBA co-embedding space. Chromatin accessibility is binary and colored accordingly.

d. SIMBA barcode plots of the genes highlighted in (a).

**Supplementary Figure 10**. SIMBA-inferred marker genes for the scRNA-seq mouse atlas dataset in batch correction analysis.
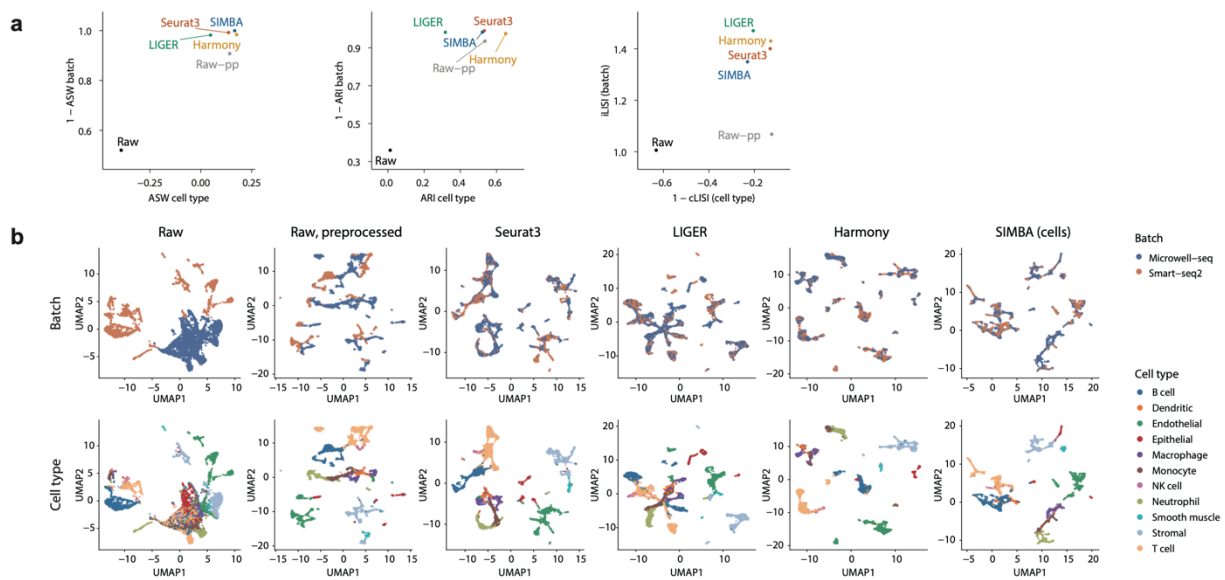
a. UMAP visualization of SIMBA embeddings of cells colored by cell type.
b. UMAP visualization of SIMBA embeddings of cells and genes.
c. UMAP visualization of SIMBA embeddings of cells and genes. Biological "query" points are highlighted with a red "+". Nearby informative genes are colored accordingly.
d. UMAP visualization of SIMBA embeddings of cells colored by indicated gene expression intensity, separated by cell type.
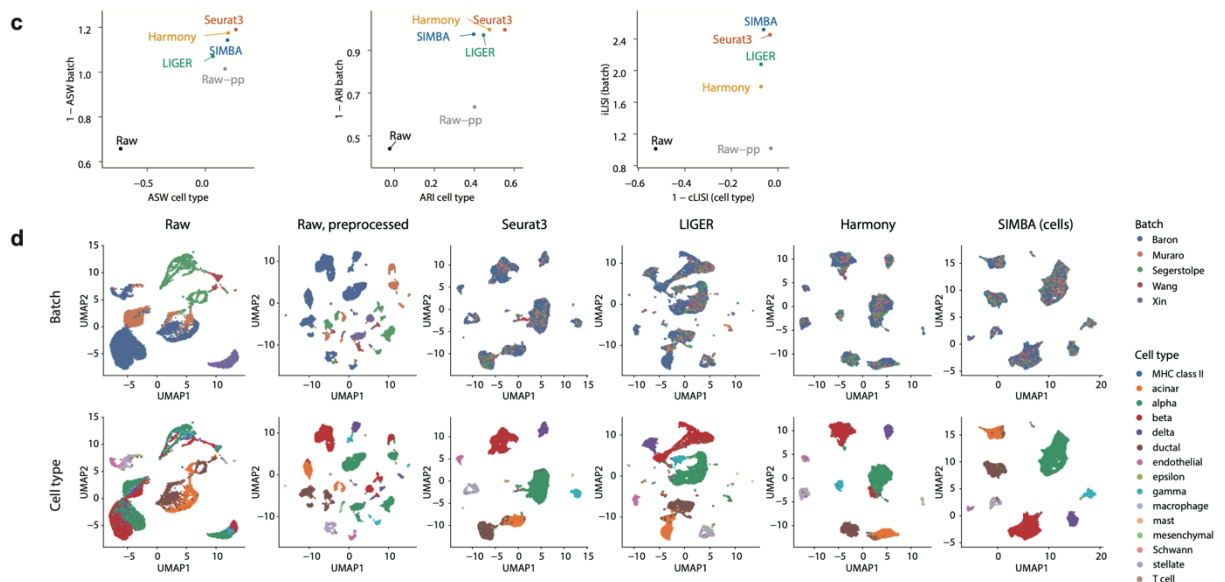
1487
1488
1489 **Supplementary Figure 11**. SIMBA-inferred marker genes for the scRNA-seq human pancreas
1490 dataset in batch correction analysis.
1491
1492     **a.** UMAP visualization of SIMBA embeddings of cells colored by cell type.
1493     **b.** UMAP visualization of SIMBA embeddings of cells and genes.
1494     **c.** UMAP visualization of SIMBA embeddings of cells and genes. Biological "query" points
1495         are highlighted with a red "+". Nearby informative genes are colored accordingly.
1496     **d.** UMAP visualization of SIMBA embeddings of cells colored by indicated gene expression
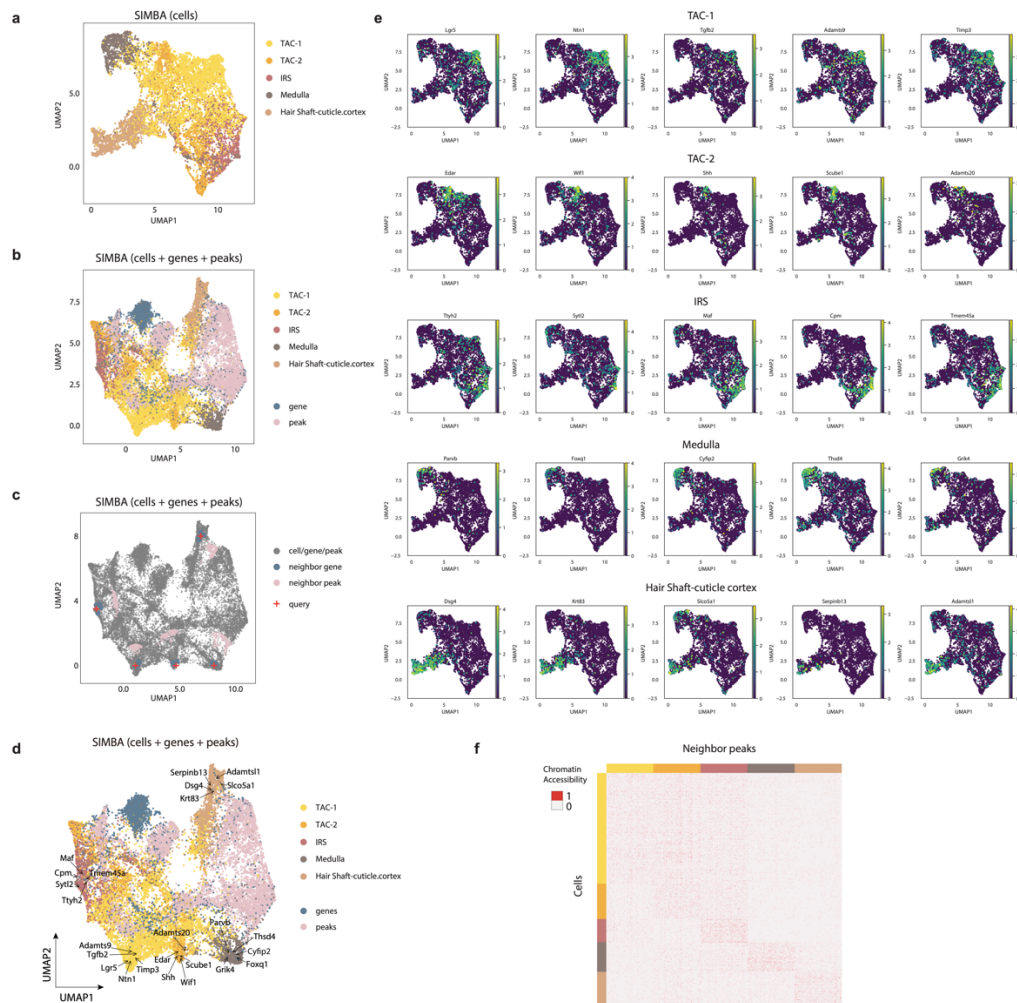1497         intensity, separated by cell type.
1498

**Supplementary Figure 12**. Comparison of SIMBA to other methods for batch correction of the mouse atlas (a-b) and human pancreas scRNA-seq datasets (c-d).

**a, c**. Quantitative comparison of SIMBA with three other batch correction methods including Seurat3, LIGER and Harmony, using, left-to-right: average silhouette width (ASW), adjusted Rand index (ARI), and local inverse Simpson's index (LISI)
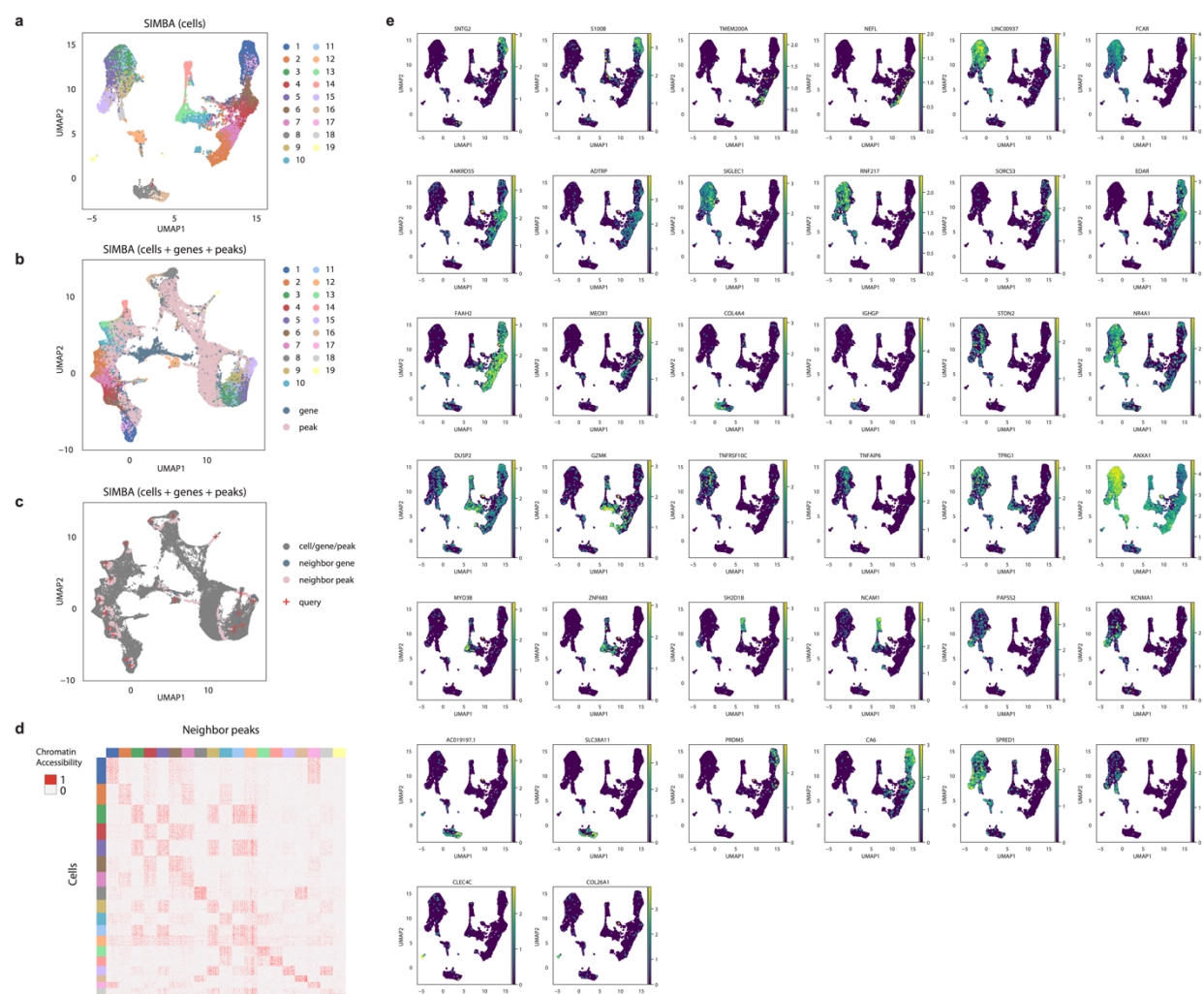
**b, d**. UMAP visualization of raw and preprocessed data alongside the batch corrected results produced by Seurat3, LIGER, Harmony, and SIMBA. Colored by technology (top) and cell type (bottom).

**Supplementary Figure 13**. SIMBA-inferred marker features for the SHARE-seq mouse skin dataset in multi-omics integration analysis.

a. UMAP visualization of SIMBA embeddings of cells with two cellular modalities integrated.

b. UMAP visualization of SIMBA embeddings of cells, genes, and peaks with two cellular modalities integrated.

c. UMAP visualization of SIMBA embeddings of cells, genes, and peaks with two cellular modalities integrated. Biological "query" points are highlighted with a red "+". Nearby informative genes and peaks are colored accordingly.

d. UMAP visualization of SIMBA embeddings of cells, genes, and peaks with two cell modalities integrated and known marker genes highlighted.

e. UMAP visualization of SIMBA embeddings of cells colored by indicated gene expression intensity, separated by cell type.

f. Heatmap of cells against neighboring peaks of each cell type that are selected in the SIMBA co-embedding space. Chromatin accessibility is binary and colored accordingly.

1528



1529
1530
1531 **Supplementary Figure 14**. SIMBA-inferred marker features for the 10x human PBMCs dataset in
1532 multi-omics integration analysis.

1533
1534     **a.** UMAP visualization of SIMBA embeddings of cells with two cellular modalities
1535         integrated.
1536     **b.** UMAP visualization of SIMBA embeddings of cells, genes, and peaks with two cellular
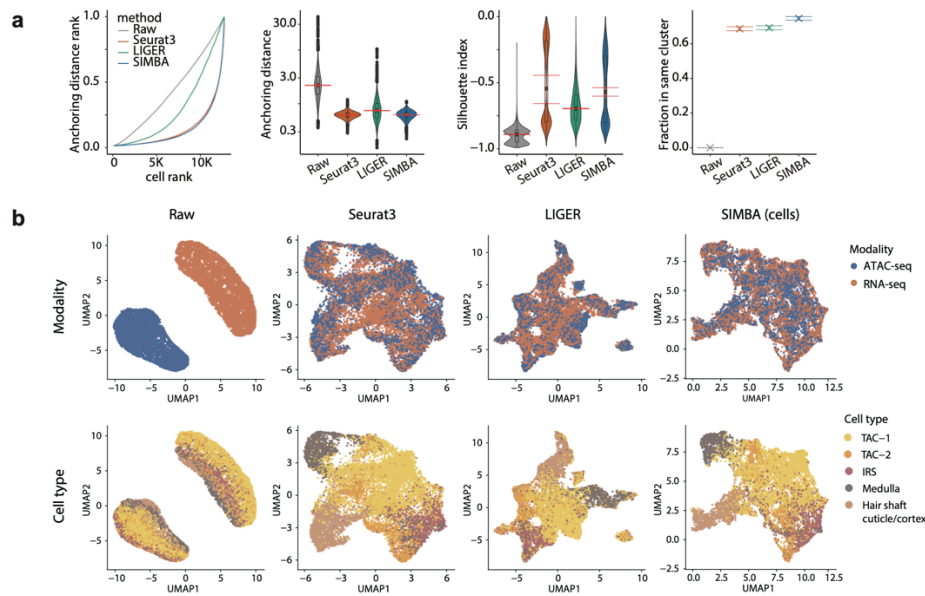1537         modalities integrated.
1538     **c.** UMAP visualization of SIMBA embeddings of cells, genes, and peaks with two cellular
1539         modalities integrated. Biological "query" points are highlighted with a red "+". Nearby
1540         informative genes and peaks are colored accordingly.
1541     **d.** Heatmap of cells against neighboring peaks of each cluster that are selected in the
1542         SIMBA co-embedding space. Chromatin accessibility is binary and colored accordingly.
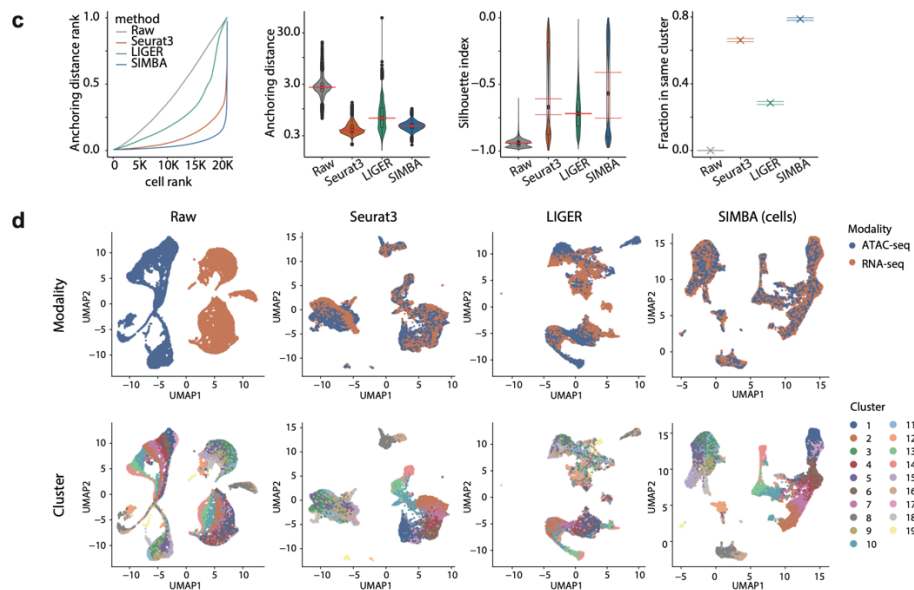1543     **e.** UMAP visualization of SIMBA embeddings of cells colored by indicated gene expression
1544         intensity, separated by cell type.
1545

1546
1547 **Supplementary Figure 15**. Comparison of SIMBA to other methods for multi-omics integration
1548 of the SHARE-seq mouse skin (a-b) and 10x multiome human PBMCs (c-d) datasets.
1549    **a, c**. Quantitative comparison of SIMBA with two other methods including Seurat3,
1550    LIGER for multi-omics integration, using, left-to-right: anchoring distance rank,
1551    anchoring distance, silhouette index, and Fraction in the same cluster.
1552    **b, d**. UMAP visualization of the raw scRNA-seq and scATAC-seq data from the 10x
1553    multiome human PBMCs dataset alongside the integrated results produced by Seurat3,
1554    LIGER, and SIMBA. Colored by data modality (top) and cluster assignment (bottom). The
1555    red intervals of violin plot of Anchoring distance and Silhouette index shows the 95%
1556    of the mean.

1557

1558    **Supplementary Notes**

1559

1560    **Supplementary Note 1: Comparison with scATAC-seq methods**

1561

1562    To assess SIMBA's ability to cluster cell types based on scATAC-seq profiles, we compared SIMBA
1563    with specialized methods specifically designed for this task. We observed that SIMBA yields
1564    consistent embeddings of cells when using either a single feature (peaks) or multiple features
1565    (peaks and DNA sequences from within those peaks) as input to the graph. This comparison was
1566    performed across four scATAC-seq datasets of varying profiling technologies and organisms
1567    (**Supplementary Fig. 4**). Given these differences, to create a fair comparison we used the same
1568    set of features (i.e., peaks) for SIMBA as other methods. SIMBA's performance was compared
1569    against three of the top methods, including SnapATAC[1], *Cusanovich*2018[2], and cisTopic[3]
1570    recommended by our recent benchmark study[4]. This comparison was first made qualitatively
1571    based on UMAP visualization and then quantitatively based on clustering performance. SIMBA
1572    performed as well as or better than each of the methods evaluated. These results comparing
1573    SIMBA to scATAC-seq-specialized methods highlight SIMBA's wide utility for single-cell analyses
1574    (**Supplementary Fig. 5**).

1575

1576    **Supplementary Note 2:  Comparison with batch correction methods**

1577

1578    Multiple methods have now been developed to correct for the technical effects of sample
1579    preparation and data collection in single cells. To assess SIMBA's performance in removing batch
1580    effects, we compared it to Seurat3[5], LIGER[6] and Harmony[7], three top-performing batch
1581    correction methods recommended in a recent benchmark study[8].

1582

1583    Two datasets, including a mouse atlas dataset and a human pancreas dataset (see
1584    **Supplementary Table 1**), were used for the evaluation. The mouse atlas dataset is composed of
1585    two scRNA-seq subsets with shared cell types from different sequencing platforms. The human
1586    pancreas dataset is composed of five samples pooled from five distinct sources using four
1587    different sequencing techniques wherein not all cell types are shared across each sample.

1588

1589    To qualitatively compare these methods, we visualized cells of each dataset before and after
1590    batch-correction in UMAP plots (**Supplementary Fig. 12b,d**). To quantitatively evaluate the
1591    performance of each method, using the benchmarking pipeline laid out in Tran *et al*[8], we
1592    measured the conservation of biological information and batch effect removal based on three
1593    different metrics: average silhouette width (ASW), adjusted Rand index (ARI), and local inverse
1594    Simpson's index (LSI)[7] as in the previously-mentioned benchmark study[8] (**Supplementary Fig.
1595    12a,c**; **Methods**). Each metric measures the relative mixing of class labels, where optimal
1596    performance is associated with maximal mixing in the batch labels and minimal mixing in the
1597    cell type labels.

1598

1599    The "Raw" batch correction results are the first 50 principal components of the horizontally
1600    concatenated gene-by-cell expression count matrix using *stats::prcomp* in R package with
1601    centering and scaling. The "Raw, preprocessed" batch correction used the preprocessed data
1602    with log normalization with scaling factor $10^4$ and selection of 3000 highly variable genes with
1603    Seurat v3 with no restriction on the minimum number of cells and genes.

1604

1605    For batch correction using Seurat v3, default options are used for pancreas dataset whereas for
1606    mouse atlas dataset no cutoff was used for the minimum number of cells and genes as in Tran
1607    *et al.*[8]. The dimension of the batch corrected embedding is set as 50 dimensions following the
1608    default option for *Seurat::RunPCA* and for the consistency with SIMBA.

1609

1610    For batch correction using LIGER, the same arguments are used (lambda = 5, nrep = 3) are used
1611    for *liger::optimizeALS* in Tran *et al*. other than the number of factors k was set as 50 for
1612    consistency with other methods for both datasets.

1613

1614    For batch correction using Harmony, the same arguments are used as in Tran *et al*.[8] other than
1615    the number of dimensions of the output embedding was set to 50 instead of 20. We note that
1616    the output embedding of 20 dimensions would result in the similar result as when used 50
1617    dimensions in these methods.

1618

1619    **Supplementary Note 3**: **Comparison with multi-omics integration methods**

1620

1621    Seurat3 and LIGER are two of the most widely-adopted methods for single-cell data integration.
1622    Here, we demonstrate that SIMBA outperforms these methods on two separate datasets, the
1623    recently published SHARE-seq mouse skin dataset and the similarly recent 10x PBMCs multiome
1624    dataset (**Supplementary Table 1**). We focus on Seurat3 and LIGER as they have explicit
1625    documentation for the task of integrating scRNA-seq and scATAC-seq data.

1626

1627    We first qualitatively evaluated these methods by inspecting UMAP visualization plots. For the
1628    SHARE-seq dataset, we observed that all three methods perform comparably well in mixing
1629    cells of two modalities though LIGER generated particularly small and noisy clusters
1630    (**Supplementary Fig. 15b**). For the 10X PBMCs dataset, SIMBA resulted in the best mixing of
1631    cells belonging to each modality whereas other methods clustered cells separately within the
1632    originating modalities (**Supplementary Fig. 15d**). We next quantitatively assessed the
1633    integration performance of each method using four metrics that measure the distances
1634    between matched cells in the integrated space (**Methods**). In addition to the commonly-used
1635    metrics including anchoring distance, Silhouette index, and Fraction in the same cluster, we
1636    developed an additional metric, *anchoring distance rank* (ADR)*,* which represents the
1637    normalized rank of the distance between matching cells. If two matching cells from scRNA-seq
1638    and scATAC-seq are mutually closest to one another, their ADR will be close to 0 (**Methods**) and
1639    thus a minimized ADR is ideal. Overall SIMBA showed the best performance according to ADR
1640    as well as cluster agreement while showing comparable or better performance according to the
1641    remaining metrics for both datasets (**Supplementary Fig. 15a,c**).

1642

1643

1644 The modality integration procedure for Seurat v3 and LIGER follows the tutorial provided by the
1645 authors (Seurat v3:
1646 https://satijalab.org/seurat/archive/v3.1/atacseq_integration_vignette.html; LIGER:
1647 http://htmlpreview.github.io/?https://github.com/welch-
1648 lab/liger/blob/master/vignettes/Integrating_scRNA_and_scATAC_data.html).

1649

1650 Both Seurat v3 and LIGER formulate the modality integration task between scRNA-seq and
1651 scATAC-seq data as a batch correction task between scRNA-seq and gene activity matrix
1652 constructed from scATAC-seq. In Seurat v3, the gene activity score of a gene is calculated as the
1653 sum of the read counts in the peaks that falls within from 2kb upstream of the TSS to the end of
1654 the gene body.  In LIGER, this score is calculated as the sum of all read counts that falls within
1655 3kb upstream of the TSS to the end of the gene body.

1656

1657 The "Raw" results start from a scRNA-seq count matrix and a gene activity matrix calculated by
1658 Seurat v3. Filtering for the shared genes in both modalities resulted in 16738 genes for the
1659 SHARE-seq mouse skin dataset and 11045 genes for the 10X PBMCs dataset. Gene-by-cell gene
1660 expression matrix and gene activity matrix were horizontally concatenated along matching rows
1661 (genes). The output embedding is the first 20 principal components calculated by the R function
1662 *stats::prcomp*  with centering and scaling.

1663

1664 For the modality integration using Seurat v3, the gene expression count was filtered using the
1665 default parameters *min.cells = 3* and *min.features = 200*. The co-embedding was created as
1666 described in the tutorial of the package using the scRNA-seq. The output embedding consists of
1667 the first 50 principal components, which is the default option of *Seurat::RunPCA*.

1668

1669 For the modality integration using LIGER, the gene expression count and gene activity matrices
1670 were normalized and filtered for the genes that are shared between both matrices. The values
1671 were then scaled according to the tutorial. In applying LIGER to the SHARE-seq mouse skin
1672 dataset, the function, *liger::optimizeALS* was used with the default parameters, k = 20 and
1673 lambda = 5.  The scRNA-seq dataset was indicated as the reference in the function,
1674 *liger::quantile_norm* as described in the documentation. The scRNA-seq and scATAC-seq
1675 modalities of the 10X PBMC multiome dataset were unable to be aligned using the default
1676 parameters. Thus *lambda = 30* and *max.iters = 100* were used for the *liger::optimizeALS*
1677 function and the scATAC-seq dataset was indicated as the reference using the
1678 *liger::quantile_norm* function to ensure a better alignment.

1679

1680 References

1681

1682

1683 1.    Fang, R. et al. Comprehensive analysis of single cell ATAC-seq data with SnapATAC. *Nat*
1684      *Commun* **12**, 1337 (2021).
1685 2.    Cusanovich, D.A. et al. The cis-regulatory dynamics of embryonic development at single-
1686      cell resolution. *Nature* **555**, 538-542 (2018).

1687   3.   Bravo Gonzalez-Blas, C. et al. cisTopic: cis-regulatory topic modeling on single-cell ATAC-
1688        seq data. *Nat Methods* **16**, 397-400 (2019).
1689   4.   Chen, H. et al. Assessment of computational methods for the analysis of single-cell
1690        ATAC-seq data. *Genome Biology* **20**, 241 (2019).
1691   5.   Stuart, T. et al. Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902
1692        e1821 (2019).
1693   6.   Welch, J.D. et al. Single-Cell Multi-omic Integration Compares and Contrasts Features of
1694        Brain Cell Identity. *Cell* **177**, 1873-1887 e1817 (2019).
1695   7.   Korsunsky, I. et al. Fast, sensitive and accurate integration of single-cell data with
1696        Harmony. *Nat Methods* **16**, 1289-1296 (2019).
1697   8.   Tran, H.T.N. et al. A benchmark of batch-effect correction methods for single-cell RNA
1698        sequencing data. *Genome Biol* **21**, 12 (2020).
1699