

scJoint: transfer learning for data integration of single-cell RNA-seq and ATAC-seq

Yingxin Lin^{1,2§}, Tung-Yu Wu^{3§}, Sheng Wan⁴, Jean Y.H. Yang^{1,2}, Wing H.
Wong^{3,5,6*}, and Y. X. Rachel Wang^{1*}

¹School of Mathematics and Statistics, The University of Sydney, NSW, Australia.

²Charles Perkins Centre, The University of Sydney, NSW, Australia.

³Department of statistics, Stanford University, CA, USA.

⁴Institute of Electronics, National Chiao Tung University, Hsinchu, Taiwan.

⁵Department of Biomedical Data Science, Stanford University, CA, USA.

⁶Bio-X Program, Stanford University, CA, USA.

[§]Equal contribution.

*To whom correspondence should be addressed. Email: Y. X. Rachel Wang,
rachel.wang@sydney.edu.au; Wing H. Wong, whwong@stanford.edu.

Abstract

Single-cell multi-omics data continues to grow at an unprecedented pace, and while integrating different modalities holds the promise for better characterization of cell identities, it remains a significant computational challenge. In particular, extreme sparsity is a hallmark in many modalities such as scATAC-seq data and often limits their power in cell type identification. Here we present scJoint, a transfer learning method to integrate heterogeneous collections of scRNA-seq and scATAC-seq data. scJoint uses a neural network to simultaneously train labeled and unlabeled data and embed cells from both modalities in a common lower dimensional space, enabling label transfer and joint visualization in an integrative framework. We demonstrate scJoint consistently provides meaningful joint visualizations and achieves significantly higher label transfer accuracy than existing methods using a complex cell atlas data and a biologically varying

25 multi-modal data. This suggests scJoint is effective in overcoming the heterogeneity in different
26 modalities towards a more comprehensive understanding of cellular phenotypes.

27 **Introduction**

28 Advances in single-cell technologies have enabled comprehensive studies of cell heterogeneity,
29 developmental dynamics, and cell communications across diverse biological systems at an
30 unprecedented resolution. There are a variety of protocols profiling the transcriptomics, as ex-
31 emplified by single-cell RNA-seq (scRNA-seq). In addition, a number of technologies have been
32 developed for other molecular measurements in individual cells towards building a more holistic
33 view of cell functions, including chromatin accessibility, protein abundance, and methylation [1].

34
35 In particular, single-cell ATAC-seq (scATAC-seq) is an epigenomic profiling technique
36 for measuring chromatin accessibility to discover cell type specific regulatory mechanisms
37 [2, 3]. scATAC-seq offers a complementary layer of information to scRNA-seq, and together
38 they provide a more comprehensive molecular profile of individual cells and their identities.
39 However, it has been noted that the extreme sparsity of scATAC-seq data often limits its power
40 in cell type identification [4]. In contrast, large amounts of well-annotated scRNA-seq datasets
41 have been curated as cell atlases [5, 6], motivating us to transfer cell type information from
42 scRNA-seq to scATAC-seq for better classification of cell types in an integrative analysis
43 framework.

44
45 A number of methods exist to denoise, batch correct, and perform integration of single-omics
46 data across multiple experiments for both transcriptomic data [7–12] and scATAC-seq data
47 [13]. However, direct applications of these methods to multi-omics data integration are com-
48 putationally challenging and often suboptimal, since different modalities have vastly different
49 dimensions and sparsity levels. Recently, a growing number of methods have been proposed to
50 address the need for integrative analysis across different modalities. When the data consist of
51 simultaneous multi-modal measurements within the same cell [14, 15], methods like scAI [16]
52 and MOFA+ [17] have been developed based on factor analysis and joint clustering. In general,
53 these paired measurements are technically more challenging and costly to perform. More
54 commonly, different modalities are derived from different cells taken from the same or similar
55 populations. In this setting, most existing methods are broadly based on manifold alignment

56 [18–20] to match the distributions of different modalities globally in a latent space, matrix
57 factorization (Liger [21], coupledNMF [22]), or using correlations to identify nearby cells across
58 modalities (Conos [23], Seurat [24]). While these methods have demonstrated promising results
59 in integrating multiple modalities measured in cells from the same tissue, requiring distributions
60 to match globally in manifold alignment is too restrictive for more complex data compositions
61 as typically seen in cell atlases, where measurements for different modalities are derived from
62 different tissues and cell types. Furthermore, matrix factorization and correlation-based meth-
63 ods designed for unpaired data require a separate feature selection step prior to integration for
64 dimension reduction, and the method’s performance can be sensitive to which genes are selected.

65

66 Here, we present an end-to-end transfer learning method, scJoint, that effectively integrates
67 scRNA-seq and scATAC-seq data using a neural network approach (Figure 1a). Our method is
68 agnostic to the selection of highly variable genes and adds flexibility to the alignment of the
69 two modalities when their cell types do not fully overlap. It is well established that in addition
70 to having high prediction power, the hidden units of neural networks are able to learn implicit
71 representations from the underlying data distribution [25]. Hence, by leveraging information
72 from annotated scRNA-seq datasets, we use the same encoder to simultaneously train the two
73 modalities so that (1) implicit features reflecting the annotations can be learnt by a hidden layer
74 in an embedding space, and (2) unlabeled data from the ATAC domain can be aligned to similar
75 points in the same embedding space. In contrast to methods that need a preliminary dimension
76 reduction step, scJoint contains a *novel loss function* to explicitly incorporate dimension reduc-
77 tion as part of the feature engineering process in transfer learning, allowing the low dimensional
78 features to be updated throughout training and removing the need for selecting highly variable
79 genes. This integrative framework enables scJoint to transfer cell type labels from scRNA-seq to
80 scATAC-seq data and construct a joint embedding for the two modalities. By applying scJoint to
81 integrate two mouse cell atlases (scRNA-seq [5] and scATAC-seq [26]) and a multi-modal data
82 with paired protein measurements (Figure 1b), we demonstrate our method achieves considerably
83 higher label transfer accuracy and integration quality over existing methods.

84 **Results**

85 **scJoint for co-training labeled and unlabeled data**

86 The core of scJoint is a semi-supervised approach to co-train labeled data (scRNA-seq) and
87 unlabeled data (scATAC-seq), where we address the main challenge of aligning these two
88 distinct data modalities via a common lower dimensional space. scJoint consists of three
89 main steps (Figure 1a). Step 1 performs joint dimension reduction and modality alignment
90 in a common embedding space through a novel neural network based dimension reduction
91 (NNDR) loss and a cosine similarity loss respectively. The NNDR loss extracts orthogonal
92 features with maximal variability in a vein similar to PCA, while the cosine similarity loss
93 encourages the neural network to find projections into the embedding space so that majority
94 parts of the two modalities can be aligned. The embedding of scRNA-seq is further guided by
95 a cell type classification loss, forming the semi-supervised part. In Step 2, treating each cell
96 in scATAC-seq data as a query, we identify the k-nearest neighbors (KNN) among scRNA-seq
97 cells by measuring their distances in the common embedding space, and transfer the cell type
98 labels from scRNA-seq to scATAC-seq via majority vote. In Step 3, we further improve the
99 mixing between the two modalities by utilising the transferred labels in a metric learning loss.
100 Joint visualization of the datasets is obtained from the final embedding layer using standard
101 tools including tSNE [27] and UMAP [28]. scJoint requires simple data preprocessing with the
102 input dimension equal to the number of genes in the given datasets after appropriate filtering.
103 Chromatin accessibility in scATAC-seq data is first converted to gene activity scores [29, 30]
104 allowing for the use of a single encoder with weight sharing for both RNA and ATAC.

105

106 We next compared scJoint with methods recently developed and applied to the integration of
107 scRNA-seq and scATAC-seq, including Seurat v3 [24], Conos [23] for label transfer accuracy,
108 and additionally Liger [21] (as a representative matrix factorization method) for evaluating the
109 joint embedding of the two modalities.

110 **scJoint shows accurate and robust performance on large atlas data.**

111 We demonstrate the performance of scJoint in a complex scenario, where the heterogeneity of
112 cell types and tissues in atlas data poses significant challenges to data integration. We applied
113 our method to integrate two mouse cell atlases: the Tabula Muris atlas [5] for scRNA-seq data

114 and the atlas in [26] for scATAC-seq data, containing 73 cell types (96,404 cells from 20 organs,
115 two protocols) and 29 cell types (81,173 cells from 13 tissues) respectively (the latter including
116 a group annotated as “unknown”), of which 19 cell types are common. We focus our initial
117 evaluation on the subset of the atlas data containing 101,692 cells from the 19 overlapping
118 cell types only. Here, we transferred cell type labels from scRNA-seq to scATAC-seq and
119 compared the results with the original labels in [26] for accuracy; these original labels were also
120 used to evaluate the quality of joint visualizations. An inspection of the tSNE plots shows our
121 method effectively mixes the three protocols (FACS, droplet, ATAC) while providing a better
122 grouping of the cells in terms of previously defined cell types than the other methods (Figure
123 2a, Supplementary Figure S1). This observation is confirmed by the quantitative evaluation
124 metrics, with scJoint showing significantly higher cell type silhouette coefficients than all the
125 other methods and similar modality silhouette coefficients as Seurat and Liger. Overall, scJoint
126 has the highest median F1-score of silhouette coefficients, achieving a better trade-off between
127 removing the technological variations in modalities and maintaining the cell type signals (Figure
128 2b, Supplementary Figure S2). In terms of label transfer accuracy, scJoint assigned 84% of the
129 cells to the correct type, 14% and 13% higher than Seurat and Conos (Figure 2d, Supplementary
130 Figure S3).

131

132 To assess the robustness of the label transfer results, we performed a stability analysis on
133 this subset of atlas data by subsampling 80%, 50%, 20% of the cells from scRNA-seq as the
134 training data. Even when only 20% of the cells were used for training, scJoint maintained a
135 high accuracy and small variance (Figure 2c), suggesting that scJoint is potentially applicable to
136 situations where only a subset of the scRNA-seq data is annotated.

137 **Label transfer using highly heterogeneous atlas data refines cell type anno-** 138 **tations in scATAC-seq.**

139 We next performed the more challenging task of integrating the full atlas data. Since the
140 scRNA-seq atlas contains more cell types than the scATAC-seq atlas, we use this application
141 to illustrate how transferred labels can refine and provide new annotations to ATAC cells. To
142 compare with the original labels, tSNE plots were constructed in the same way as [26], using
143 singular value decomposition of the term frequency-inverse document frequency (TF-IDF)
144 transformation of scATAC-seq peak matrix (Figure 3a). We observe that scJoint labels cells

145 close together in this ATAC visualization space in a more consistent way than the other methods.
146 Qualitatively this is supported by scJoint’s higher overall accuracy rate (77% compared with
147 60% for Seurat and 55% for Conos).

148
149 Examining the transferred labels further, we find scJoint labels a group of cells (originally
150 labeled as “unknown” or “endothelials”) as “stromal cells” (4352 cells) and “fibroblasts” (1602
151 cells), which are two cell types not present in the original ATAC labels. These cells show high
152 gene activity scores for *Col1a1*, *Col1a2*, *Dcn* and *Ccdc80*, all of which are markers with high
153 expression levels in stromal cells and fibroblasts but low expression levels in endothelial cells
154 from the scRNA-seq data (Figure 3b). Hence, the new annotations are more consistent with the
155 marker expression levels.

156
157 More interestingly, we note scJoint allows us to annotate 5931 cells labeled as ‘unknown’
158 in [26] with probability score greater than 0.80. These cells are clearly clustered into groups
159 in the tSNE visualization of scJoint’s embedding space (Figure 3c), with the main groups being
160 endothelial cells, stromal cells, neurons and B cells. Using cell type markers identified from the
161 scRNA-seq data, the aggregated gene activity scores of these ATAC cells show clear differential
162 expression patterns (Figure 3d).

163 **scJoint enables accurate integration of single-cell multi-modal data across** 164 **biological conditions.**

165 We demonstrate scJoint is capable of incorporating additional modality information to RNA-seq
166 and ATAC-seq and applicable to experiments with different underlying biological conditions.
167 We consider multi-modal measurements profiling gene expression levels or chromatin acces-
168 sibility simultaneously with surface protein levels, which can be obtained via CITE-seq [31]
169 and ASAP-seq [32]. We analyzed CITE-seq and ASAP-seq data from a T cell stimulation
170 experiment in [32], which sequenced cells with these two technologies in parallel. A total
171 of 18,088 cells were studied under two conditions: one with stimulation of anti-CD3/CD28
172 in the presence of IL-2 for 16 hours and the other without stimulation as control. We first
173 clustered and annotated these cells using CiteFuse [33]. Compared to the cell type labels in the
174 original study, we were able to identify cellular subtypes with CiteFuse, further annotating five
175 subgroups in T cells. Next, we performed integration analysis of CITE-seq and ASAP-seq by

176 concatenating gene expression or gene activity vectors with protein measurements. The analysis
177 was performed in two scenarios: within the stimulated and control condition separately and
178 across the two conditions.

179

180 In both scenarios, scJoint generated a better joint visualization of the two technologies
181 (Figure 4a, Supplementary Figures S4, S5). In particular, in the case where stimulated and
182 control cells are combined, subtypes of T cells (e.g. naive CD8+, effector CD8+, naive CD4+,
183 and effector CD4+) are clearly separated while cells from the two technologies are well mixed
184 (Figure 4a-b). The median cell type silhouette coefficient of scJoint is 0.51, outperforming
185 the other three methods by a large margin (Seurat 0.11, Conos 0.13, and Liger -0.06). With
186 the highest silhouette coefficient F1 scores (median F1 score: 0.59) representing a 16% - 28%
187 improvement over the other methods, scJoint demonstrates the best balance between removing
188 technical variations and preserving biological signals (Figure 4c, Supplementary Figure S6).

189

190 Moreover, scJoint achieves higher accuracy in label transfer under all scenarios (88% in
191 control, 84% in stimulation, and 87% in the combined case), compared with Seurat (80% in
192 control, 79% in stimulation, and 75% combined) and Conos (53% in control, 67% in stimulation,
193 and 56% in combined) (Figure 4d and Supplementary Figure S7). In addition, the transferred
194 labels of scJoint from the two scenarios (control / stimulation alone, and combined) are highly
195 consistent, with 95% of cells having the same annotation, substantially greater than Seurat (84%)
196 and Conos (59%) (Supplementary Figure S8).

197 **Integration of multi-modal data with scJoint captures additional biological** 198 **signals in cell types and conditions**

199 In the combined analysis of stimulation and control, we find that the joint embedding generated
200 by scJoint contains additional information that allows for the identification of a cellular
201 subtype. In the CiteFuse annotation of ASAP-seq data, we labeled one cluster of 142 cells with
202 ambiguous marker expression as “unknown”. Interestingly, in the joint visualization of scJoint,
203 while these “unknown” cells are labeled as “natural killer cells (NK)” by label transfer, they are
204 still clearly separated from the majority of NK cells and form a small cluster together with cells
205 from CITE-seq. We then examined the gene and protein expression levels of NK cell and T cell
206 markers in this subgroup. We find these cells have high expression of CD3 and GNLY at gene

207 level as well as CD3, CD56, CD57, and CD244 at protein level, but low expression of CD8A
208 and CD4. This suggests these cells may be natural killer T cells, a minority of immune cells
209 in PBMC sample (Figure 4e, Supplementary Figure S9) [34]. By contrast, although these cells
210 lack CD8 expression, the other methods are unable to distinguish them from effector CD8+ T
211 cells in their visualizations (Figure 4e, Supplementary Figure S10).

212

213 Lastly, by appropriately aligning the two technologies in the embedding space, scJoint is able
214 to reveal the biological difference between stimulation and control within the same cell type. In
215 the joint visualization of scJoint, three subtypes of T cells (naive CD4+, naive CD8+, effector
216 CD4+) are less well mixed between the two conditions than the other cell types, consistent with
217 the stimulation experiment aiming to activate T cells. In particular, the naive CD4+ T cells show
218 the most notable separation between the two conditions (Figure 4a). We then performed differ-
219 ential expression analysis of the scRNA-seq part of CITE-seq within each cell type across the
220 two conditions using MAST [35]. We find that the naive CD4+ T cells have the largest number
221 of unique differentially expressed genes ($FDR < 0.01$) (Supplementary Figure S11a). Simi-
222 larly, differential proteins analysis of both CITE-seq and ATAC-seq using wilcoxon rank sum
223 test on the log-transformed protein abundances also suggests that naive CD4+ T cells have the
224 most unique differential proteins compared with other cell types ($FDR < 0.01$) (Supplementary
225 Figure S11b-c).

226 **scJoint shows versatile performance on paired measurements of scRNA-seq** 227 **and scATAC-seq.**

228 Although scJoint is designed for integrating unpaired data, it is still directly applicable to paired
229 data. Such an application also enables us to compare its performance with methods that incor-
230 porate pairing information and use the pairing information to validate the label transfer results.
231 We consider the integration of adult mouse cerebral cortex data generated by SNARE-seq [14],
232 a technology that can profile gene expression and chromatin accessibility in the same cell. In
233 addition to Seurat and Liger, we compared scJoint with two other methods designed specifically
234 for paired data, scAI [16] and MOFA+ [17]. In our assessment, all the unpaired methods (scJoint,
235 Seurat, Liger) treat the RNA and ATAC parts of SNARE-seq as two separate datasets, while the
236 paired methods take the pairing information into account. We find that scJoint is able to provide
237 clear groupings of cells according to cellular subtypes (Figure 5a) and achieves comparable or

238 better cell type silhouette coefficients (Figure 5b) than the paired methods. This suggests that
239 scJoint is versatile enough to be applied to paired data, which are becoming increasingly popular.

240

241 Comparing the performance among the unpaired methods, scJoint has the highest medians
242 in cell type silhouette coefficients and F1-scores (Figure 5b, Supplementary Figure S13). For
243 label transfer, scJoint achieves an accuracy rate of 70.9%, retaining better performance than the
244 other two methods (70.1% for Seurat and 49.5% for Conos). Looking closer at the performance
245 in each cell type, scJoint performs the best in 10 out of 22 cell types in terms of F1 scores for
246 classification (Supplementary Figure S14). Together, these results suggest that scJoint performs
247 the best among the unpaired methods and on par with the paired methods, despite treating paired
248 data as separate.

249 Discussion

250 scJoint approaches the integration of scRNA-seq and scATAC-seq as a domain adaptation
251 problem in transfer learning, using the same neural network to co-train labeled data from the
252 source domain (RNA) and unlabeled data from the target domain (ATAC) following a different
253 distribution. scRNA-seq data serve as a natural source domain for transferring information to
254 other modalities due to rapidly growing collections of annotated public data and RNA-focused
255 computational tools that can output accurate classifications [36]. Using mouse cell atlases and
256 multi-modal data with protein measurements, we demonstrate scJoint achieves significantly
257 higher label transfer accuracy and provides better joint visualizations than other methods even
258 when 1) the data is highly complex and heterogeneous and 2) meaningful biological conditions
259 are mixed with technical variations. We have shown that integrative analysis of single-cell
260 multi-omics data by scJoint facilitates re-annotation of cell types in scATAC-seq and discovery
261 of new subtypes not present in training data.

262

263 scJoint provides a concise training framework with one main tuning parameter in the
264 construction of cosine similarity loss. As shown in Supplementary Figure S15a, our results
265 are quite stable with respect to the choice of this parameter. Similar to other methods based
266 on neural networks, the number of hidden nodes in the architecture and other optimization
267 details can be considered tunable as well, although they do not appear to affect our results
268 (Supplementary Figure S15b).

269

270 The superior performance and robustness of scJoint illustrate its utility as a tool to au-
271 tomatically label cells from other modalities given an annotated scRNA-seq database. By
272 embedding all cells in a common lower dimensional space, scJoint assigns a probability score
273 to a cell type prediction by combining the softmax probabilities of its nearest neighbors. As
274 we vary the level of cutoff, the accuracy of scJoint still consistently outperforms the other
275 methods (Supplementary Figure S16). The robustness of scJoint was demonstrated through
276 subsampling experiments, where the stability of our results implies the method can be applied to
277 partially labeled databases. Despite being a semi-supervised method guided by labeled data, the
278 dimension reduction component in our design lends it sufficient flexibility to preserve implicit
279 data signals, including biological variations induced by experimental conditions and additional
280 cellular subtypes. One can conceivably extend scJoint to an unsupervised setting, replacing the
281 softmax prediction layer with a decoder minimizing reconstruction loss.

282

283 Although designed for unpaired data, scJoint is still directly applicable to paired data
284 and generates joint visualizations with cells coherently grouped by cell types. In the current
285 training scheme, the pairing information between RNA and ATAC is only used to validate the
286 label transfer results. We expect that adapting scJoint to take paired vectors during training
287 would enhance its performance on this type of data, and this would be especially useful in the
288 unsupervised setting mentioned above.

289

290 We have focused on scATAC-seq as an example of epigenomic data, but in principle scJoint
291 extends to other modalities such as methylation data, provided the input can be summarized
292 as gene-level scores. While the gene-level summaries are amenable to generalization and
293 widely adopted by unpaired integration methods, this step itself is also a limitation as improper
294 aggregation can incur information loss. Extending scJoint to directly handle epigenomic data at
295 locus level will require designing a separate encoder that is suitable for the high dimensionality
296 and remains easy to train, and we will pursue this for future work.

297

298 In summary, we have developed scJoint as a generalizable transfer learning method for per-
299 forming integrative analysis of single-cell multi-omics data. scJoint was shown to effectively
300 integrate multiple types of measurements from both unpaired or paired profiling, outperforming
301 other methods in label transfer accuracy and providing joint visualizations that remove technical

302 variations while preserving meaningful biological signals. scJoint's ability to integrate multi-
303 omics data by capturing various aspects of cell characteristics unique to different data modalities
304 will facilitate a more comprehensive view of cell functions and cell communications.

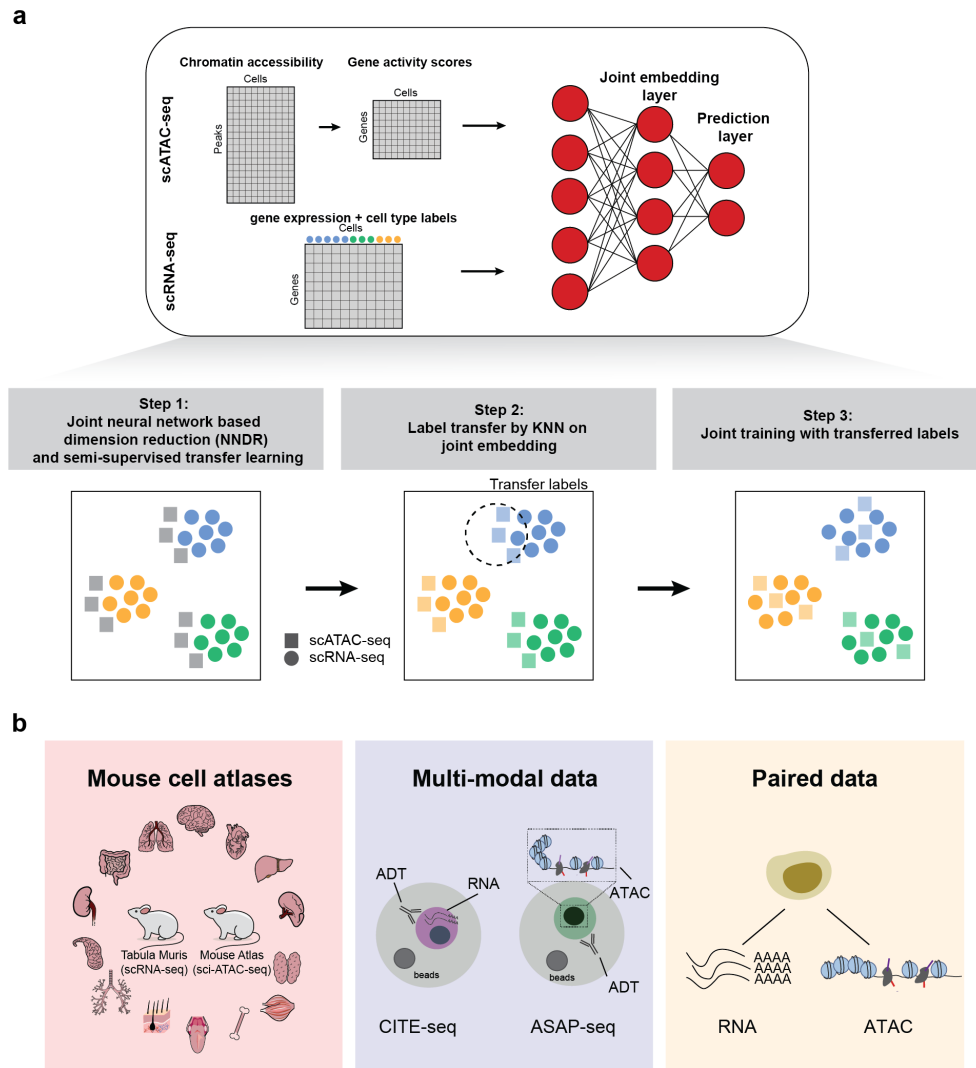


Figure 1: (a) Overview of scJoint. The input of scJoint consists of one (or multiple) gene activity score matrix, calculated from the accessibility peak matrix of scATAC-seq, and one (or multiple) gene expression matrix including cell type labels from scRNA-seq experiments. The method has three main steps: (1) Joint NNDR and semi-supervised transfer learning; (2) Cell type label transfer by k-nearest neighbor in joint embedding space; (3) Joint training with transferred labels. (b) Three data collections used in this study: (1) Mouse cell atlases; (2) Multi-modal data from PBMC; (3) Paired data from adult mouse cerebral cortex data generated by SNARE-seq.

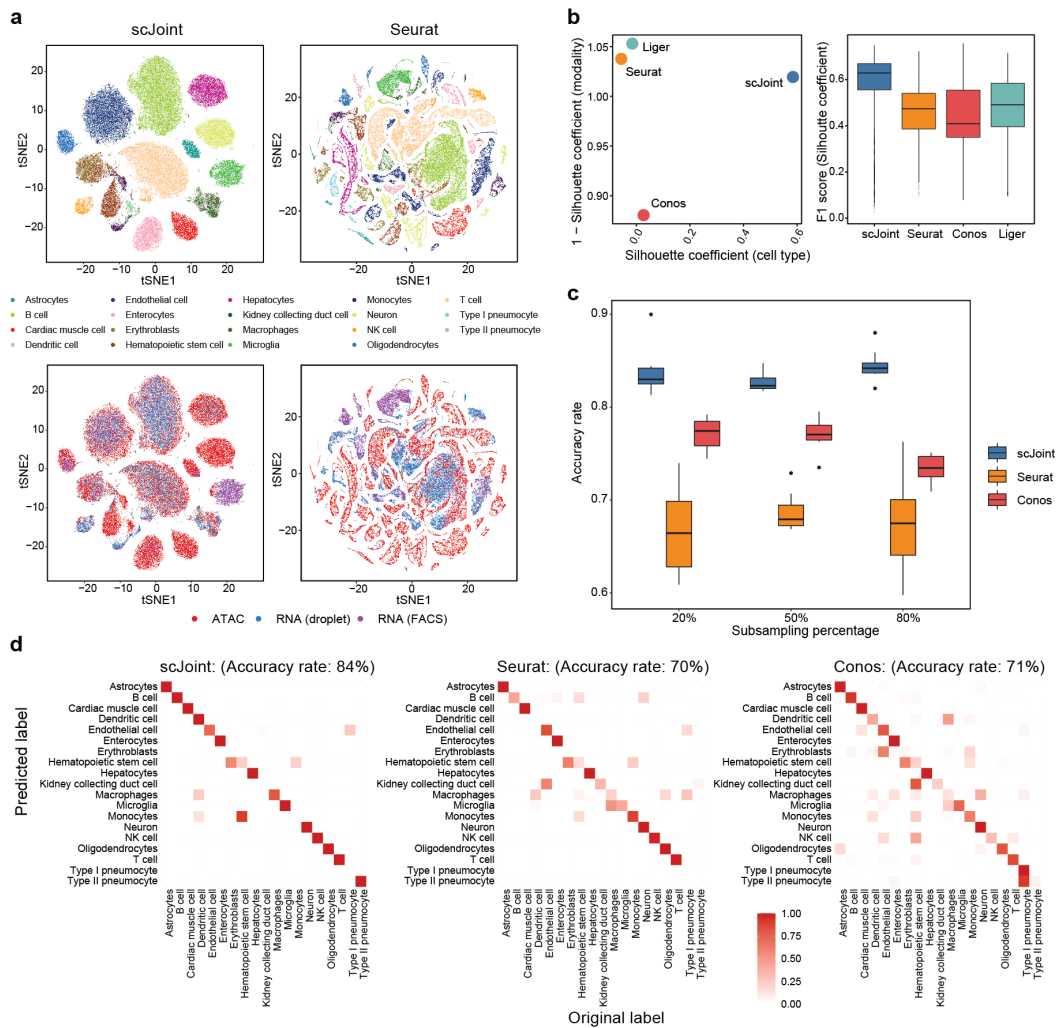


Figure 2: Analysis of mouse cell atlas subset data containing 19 overlapping cell types from RNA and ATAC. (a) tSNE visualization of scJoint (left column) and Seurat (right column), colored by cell types defined in [26] (first row) and three protocols (second row). (b) Scatter plot of mean silhouette coefficients for scJoint, Liger, Seurat, and Conos (left panel), where the x-axis shows the mean cell type silhouette coefficients and the y-axis shows ‘1 - mean modality silhouette coefficients’; ideal outcomes would lie in the top right corner. Boxplots of F1 scores of silhouette coefficients for scJoint, Liger, Seurat, and Conos (right panel). (c) Accuracy rates of scJoint, Seurat and Conos using 20%, 50% and 80% of cells from scRNA-seq data as training data. 10 random subsamplings were performed for each setting to generate the variance. (d) Predicted cell types and their fractions of agreement with the original cell types given in [26] for scJoint (left panel), Seurat (middle panel) and Conos (right panel). Clearer diagonal structure indicates better agreement.

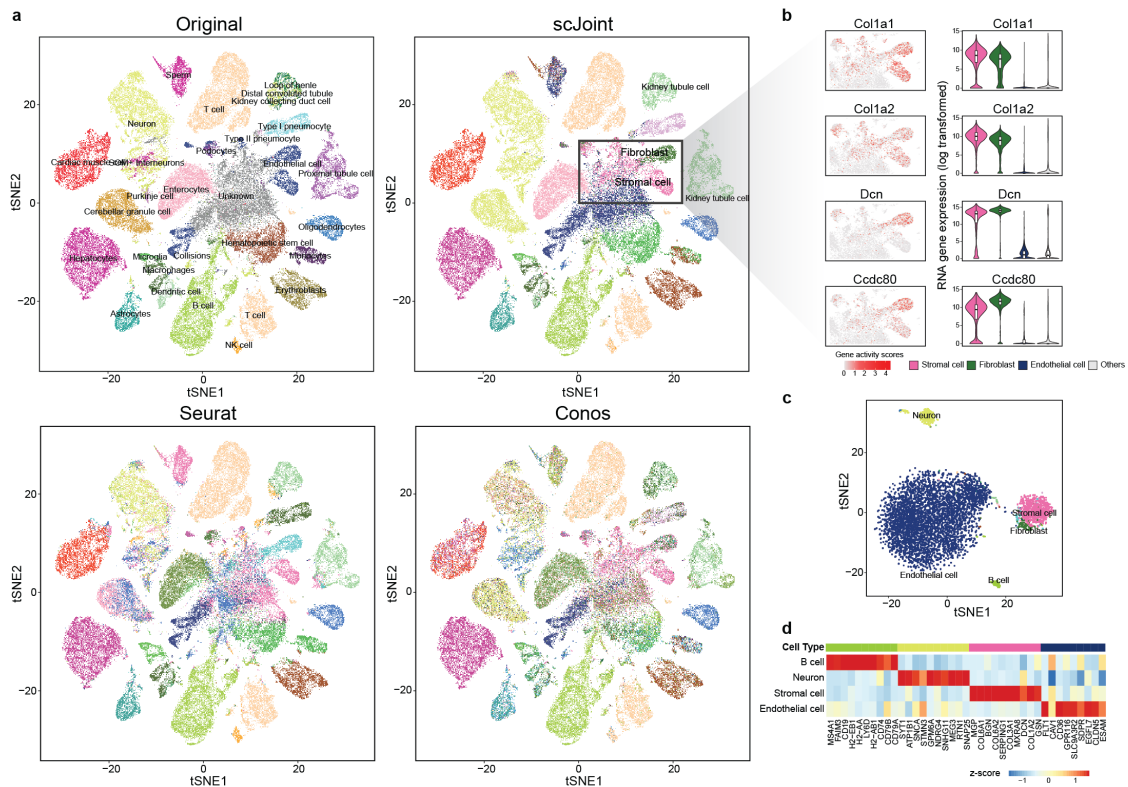


Figure 3: Analysis of mouse cell atlas full data. (a) A 2×2 panel of tSNE plots generated from top 100 dimensions of singular value decomposition of the TF-IDF transformed ATAC-seq data, colored by the original labels (top left), scJoint transferred labels (top right), Seurat transferred labels (bottom left), and Conos transferred labels (bottom right). (b) Marker expressions in stromal cells and fibroblasts: Col1a1, Col1a2, Dcn and Ccdc80. The left column shows the gene activity scores of the markers in ATAC-seq data (4352 stromal cells, and 1602 fibroblasts). The right column shows the log-transformed gene expression of the markers in stromal cells, fibroblasts, endothelial cells versus others; all cells here are taken from the FACS scRNA-seq data. (c) tSNE plot of cells originally labeled as ‘unknown’ and annotated by scJoint with probability scores greater than 0.80, colored by predicted cell types (5931 cells). (d) Heatmap of z-scores of average gene activity scores, calculated from cells aggregated by predicted cell types in ATAC. The rows indicate the top four predicted cell types by size. The columns indicate the top differential expressed genes of the corresponding cell type in RNA.

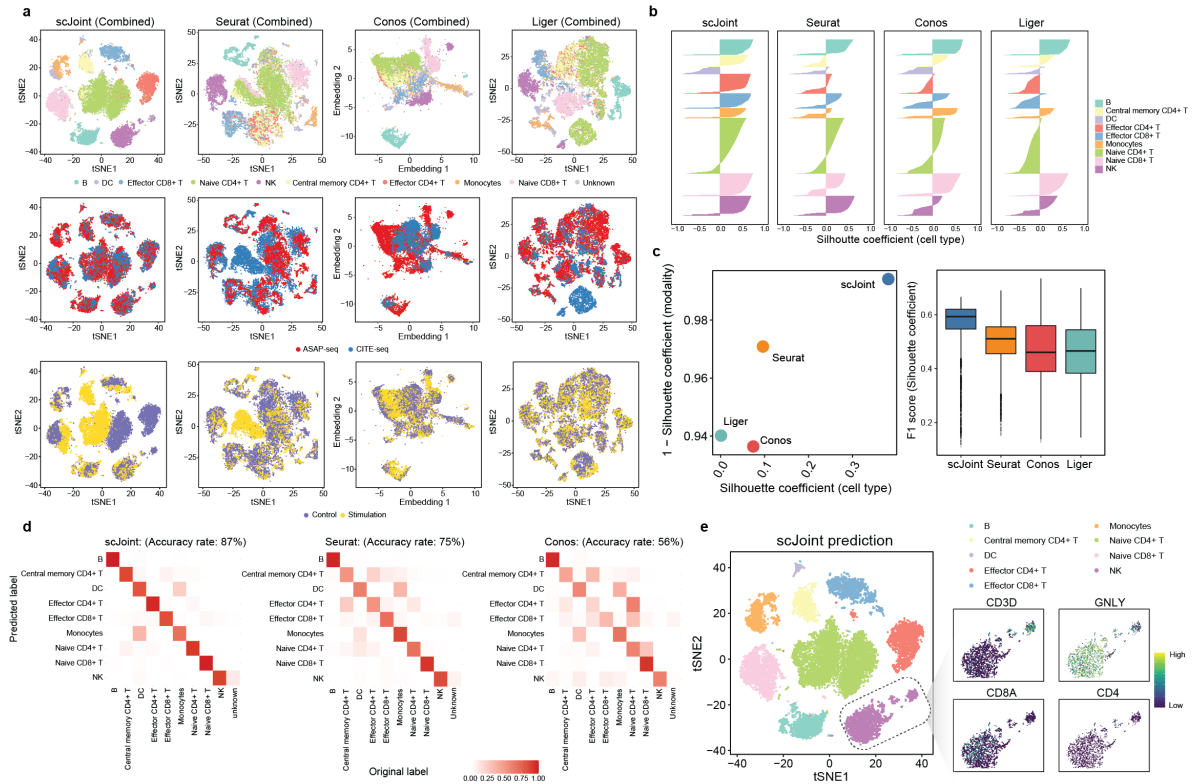


Figure 4: Integration of multi-modal PBMC data across biological conditions. (a) tSNE visualization of scJoint (first column), Seurat (second column), Conos (third column) and Liger (fourth column) of PBMC data generated from CITE-seq and ASAP-seq, colored by cell type obtained from CiteFuse and manual annotations (first row), technology (second row), and biological condition (third row). (b) Barplots of cell type silhouette coefficients for scJoint, Seurat, Conos and Liger for all cells, colored by cell type. Larger values on the x-axis indicate better grouping. (c) Scatter plot of mean silhouette coefficients for scJoint, Seurat, Conos and Liger (left), where the x-axis denotes the mean cell type silhouette coefficients, and the y-axis denotes 1 - mean modality silhouette coefficients; ideal outcomes would lie in the top right corner. Boxplots of F1 scores of silhouette coefficients for scJoint, Liger, Seurat, and Conos (right). (d) Heatmaps comparing the original labels and the transferred labels of scJoint, Seurat and Conos. Clearer diagonal structure indicates better agreement. (e) tSNE visualization of scJoint colored by the predicted cell types with gene expression levels of CD3D, NKG7, CD8A and CD4 in natural killer cells.

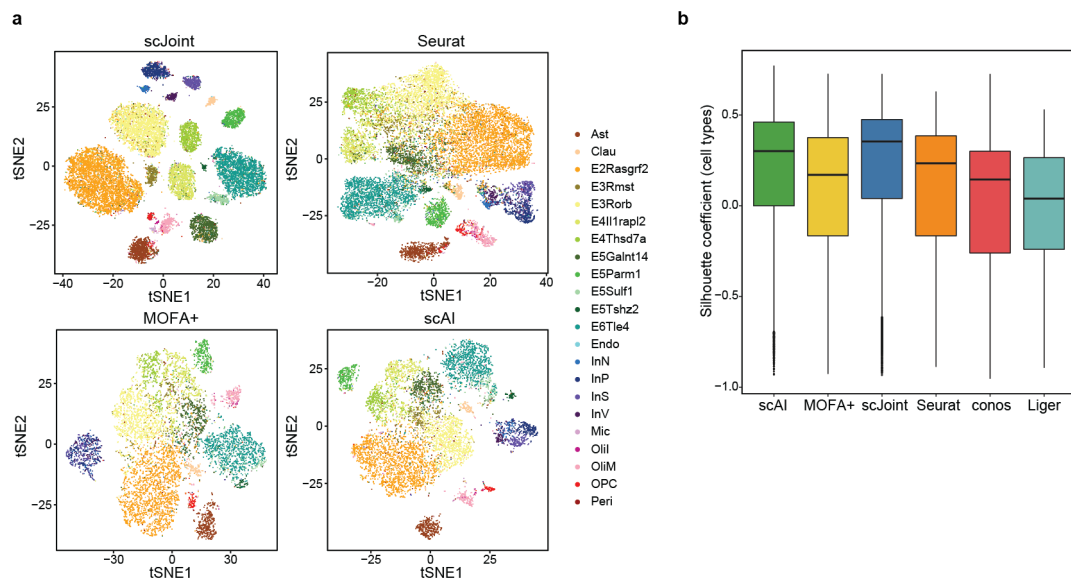


Figure 5: Analysis of paired gene expression and chromatin accessibility data from SNARE-seq. (a) tSNE visualization of SNARE-seq data for scJoint, Seurat, MOFA+ and scAI, colored by cell types given in [14]. All unpaired methods treat the RNA and ATAC parts of SNARE-seq as two separate data. (b) Boxplots of cell type silhouette coefficients for scJoint, Seurat, Conos and Liger, colored by methods.

305 **Methods**

306 **Architecture and training of scJoint**

307 The neural network in scJoint consists of one input layer and two fully connected layers. The
308 input layer has dimension equal to the number of genes common to the expression matrix
309 of scRNA-seq and the gene activity matrix of scATAC-seq, after simple filtering (see Data
310 preprocessing). Now that the two modalities have matching input features, we co-train them
311 using the same encoder which is equivalent to weight sharing. The first fully connected layer has
312 64 neurons with linear activation and serves as the joint low dimensional embedding space that
313 captures aligned features from all cells. visualizations of clustering structure can be obtained
314 by applying tSNE or UMAP to the output of the embedding layer. The second fully connected
315 layer has dimension equal to the number of cell types in scRNA-seq data. Through a softmax
316 transformation, this layer outputs a probability vector for cell type prediction. For cells in
317 scRNA-seq, this layer can be trained in a supervised fashion using the cross entropy loss.

318

319 Given S scRNA-seq experiments with expression matrices and T scATAC-seq experiments
320 with gene activity score matrices, assume suitable intersections have been taken so that all
321 matrices have the same set of genes. Let $\{x_i^{(s)}\}_{i=1}^{N_s}$ be the expression profiles of cells after
322 preprocessing from a scRNA-seq dataset indexed by $s \in \{1, \dots, S\}$, and $\{y_i^{(s)}\}_{i=1}^{N_s}$ be the
323 corresponding cell type annotations. Here each $x_i^{(s)}$ is a G -dimensional vector, where G
324 is the number of genes; $y_i^{(s)} \in \{1, \dots, K\}$, where K is the number of cell types; N_s is
325 the number of cells in experiment s . Similarly, let $\{x_i^{(t)}\}_{i=1}^{N_t}$ be the vectors of gene activity
326 scores after preprocessing from the t -th scATAC-seq dataset with N_t cells ($t \in \{1, \dots, T\}$),
327 whose cell types are unlabeled. The neural network is parametrized by a set of weights and
328 biases, collectively denoted θ . Let $f_{\theta,i}^{(s)} = f(x_i^{(s)}; \theta) \in \mathbb{R}^D$, $D = 64$, be the output of the
329 embedding layer when the input $x_i^{(s)}$ has gone through a transformation of f parametrized by
330 θ . Similarly $g_{\theta,i}^{(s)} = \text{softmax}(h(f(x_i^{(s)}; \theta)))$, where h denotes the output from the prediction
331 layer that goes through the softmax transformation. Thus $g_{\theta,i}^{(s)}$ is a probability vector after the
332 softmax transformation. $f_{\theta,i}^{(t)}$ and $g_{\theta,i}^{(t)}$ are defined in the same way for input $x_i^{(t)}$ from scATAC-seq.

333

334 The training of scJoint consists of three steps.

335 **Step 1: Joint neural network based dimension reduction (NNDR) and semi-supervised**
 336 **transfer learning**

337 We first perform joint dimension reduction and feature alignment by imposing suitable loss
 338 functions on the outputs of the two fully connected layers. A mini-batch \mathcal{B}_0 of data for
 339 training is constructed by sampling equal-sized subsets of cells from each dataset, that is,
 340 $\mathcal{B}_0 = \{\mathcal{B}^{(s)}\}_{s=1}^S \cup \{\mathcal{B}^{(t)}\}_{t=1}^T$, where each subset $\mathcal{B}^{(s)}$ (or $\mathcal{B}^{(t)}$) has B cells.

341 1. *NNDR Loss*. In a spirit similar to PCA, the NNDR loss aims to capture low dimensional,
 342 orthogonal features when projecting each data batch into the embedding space. For now
 343 we omit the dataset-specific superscript with the understanding that this loss function is
 344 applied to each $\mathcal{B}^{(s)}$ and $\mathcal{B}^{(t)}$. Given input vectors $\{x_b\}_{b \in \mathcal{B}}$, define $\bar{f}_{\theta, \cdot} = \frac{1}{B} \sum_{b \in \mathcal{B}} f_{\theta, b} \in$
 345 \mathbb{R}^D , and $\Sigma_{\theta, \cdot}$ as the sample correlation matrix. The NNDR loss is:

$$347 \quad \mathcal{L}_{\text{NNDR}}(\mathcal{B}, \theta) = \left(\frac{1}{BD} \sum_{b \in \mathcal{B}} \sum_{j=1}^D |f_{\theta, b}(j) - \bar{f}_{\theta, \cdot}(j)| \right)^{-1} + \frac{1}{D^2} \sum_{i \neq j} |\Sigma_{\theta, \cdot}(i, j)|$$

$$348 \quad \quad \quad + \frac{1}{BD} \sum_{b \in \mathcal{B}} \sum_{j=1}^D |\bar{f}_{\theta, \cdot}(j)|.$$

350 Note that to minimize this loss, we maximize the variability within each coordinate (inverse
 351 of the first term) and minimize the correlation between all coordinate pairs (the second
 352 term) to achieve orthogonality. The last term tries to fix the means of all coordinates
 353 near zero for model identifiability, preventing θ from drifting to unstable regions of the
 354 parameter space.

355 2. *Cosine similarity loss*. This loss is applied to the embedding layer outputs from $\mathcal{B}^{(t)}$ and
 356 $\mathcal{B}_R = \cup_{s=1}^S \{\mathcal{B}^{(s)}\}$, for every t , and attempts to maximize the similarity between best
 357 aligned ATAC and RNA data pairs. Let p be the fraction of data pairs we expect to have
 358 high cosine similarity scores. Setting $p < 1$ accounts for situations where RNA and ATAC
 359 do not share all their cell types. We set $p = 0.8$ for all the results presented in the paper,
 360 and our results appear to be stable with respect to this parameter (Supplementary Figure
 361 S15a) when the cell types fully overlap. Recall that for a pair of general vectors (u, v) , the
 362 cosine similarity is defined as $\cos(u, v) = \langle u, v \rangle / (\|u\| \|v\|)$. For each $x_b^{(t)}$ with $b \in \mathcal{B}^{(t)}$,
 363 we find the corresponding $i(b) \in \mathcal{B}_R$ with input $x_{i(b)}$ that maximizes $\cos(f_{\theta, b}^{(t)}, f_{\theta, i(b)})$. From
 364 $\mathcal{B}^{(t)}$, we then choose the top p fraction of cells with the highest cosine score and denote the
 365 index set \mathcal{I}_p . (\mathcal{I}_p has size $\lfloor Bp \rfloor$.) The loss is given by

$$\mathcal{L}_{\text{cos}}(\mathcal{B}^{(t)}, \mathcal{B}_R, \theta) = -\frac{1}{[Bp]} \sum_{b \in \mathcal{I}_p} \text{cos}(f_{\theta,b}^{(t)}, f_{\theta,i(b)}).$$

3. *Cross entropy loss.* For every $\mathcal{B}^{(s)}$ with cell type annotations $\{y_b^{(s)}\}_{b \in \mathcal{B}^{(s)}}$, we apply the cross entropy loss to the prediction layer after softmax transformation to supervise the learning of scRNA-seq datasets:

$$\mathcal{L}_{\text{entropy}}(\mathcal{B}^{(s)}, \theta) = -\frac{1}{B} \sum_{b \in \mathcal{B}^{(s)}} \sum_{k=1}^K 1(y_b^{(s)} = k) \log g_{\theta,b}^{(s)}(k),$$

where $1(\cdot)$ is an indicator function.

In Step 1, the final loss function we minimise with respect to θ for a mini-batch \mathcal{B}_0 is

$$\mathcal{L}_1(\mathcal{B}_0, \theta) = \sum_{s=1}^S (\mathcal{L}_{\text{NNDR}}(\mathcal{B}^{(s)}, \theta) + \mathcal{L}_{\text{entropy}}(\mathcal{B}^{(s)}, \theta)) + \sum_{t=1}^T (\mathcal{L}_{\text{NNDR}}(\mathcal{B}^{(t)}, \theta) + \mathcal{L}_{\text{cos}}(\mathcal{B}^{(t)}, \mathcal{B}_R, \theta)).$$

Step 2: Cell type label transfer by KNN in joint embedding space

The output of Step 1 is a joint embedding space that has roughly aligned RNA and ATAC with cells from either modality lying close if they have similar low dimensional representations in this space. Therefore using the embedding vectors for cells in all the datasets and calculating the Euclidean distances, we can determine the KNN among all RNA cells for each cell i in ATAC; denote this set of RNA cells $\mathcal{N}(i)$. The cell type label of i is estimated via majority vote using $\{y_j\}_{j \in \mathcal{N}(i)}$. All the results in the paper were obtained from using 30 nearest neighbors. Let the majority cell type be k^* , then the probability score of cell type prediction for cell i in ATAC is an average of its nearest neighbors in RNA. Since for each $j \in \mathcal{N}(i)$, $g_{\theta,j}$ is already a probability vector after the softmax transformation, we take $p_{\theta,j} = g_{\theta,j}(k^*)$ as the probability score of RNA cell j in the majority class $\mathcal{M}(i) \subset \mathcal{N}(i)$. For other $j \in \mathcal{N}(i) \setminus \mathcal{M}(i)$, we threshold the probability score as 0. Then the probability score of ATAC cell i is calculated as

$$\hat{p}_{\theta,i} = \frac{1}{30} \sum_{j \in \mathcal{M}(i)} p_{\theta,j}.$$

Step 3: Joint training with transferred cell type labels

In the final step of the training, we refine the joint embedding space and improve mixing of cells from the same cell type using the transferred labels from Step 2. We include an additional loss

393 function commonly used in metric learning for enhancing embedded clustering structure given
394 labeled data. The other loss functions and network architecture remain the same as Step 1 with
395 ATAC cells and their transferred labels added to $\mathcal{L}_{\text{entropy}}$.

396

397 For each cell type $k \in \{1, \dots, K\}$, we initialize the class center $c_k \in \mathbb{R}^D$ randomly. We
398 construct mini-batches of cells from all the datasets in the same way as Step 1. Now that all cells
399 have cell type labels (given or transferred), for convenience we will refer to cells in a mini-batch
400 \mathcal{B}_0 without explicitly labeling which dataset they come from. For a given \mathcal{B}_0 , we first update the
401 class centers by taking the average of c_k and $\{f_{\theta,b}\}$ with $b \in \mathcal{B}_0$ and $y_b = k$. Let the updated
402 centers be c'_k . As the number of mini-batches grows, the influence of the initial c_k becomes
403 negligible. The metric learning loss we use is the center loss:

$$404 \quad \mathcal{L}_{\text{center}}(\mathcal{B}_0, \theta) = \frac{1}{|\mathcal{B}_0|K} \sum_{b \in \mathcal{B}_0} \sum_{k=1}^K \|f_{\theta,b} - c'_k\|^2 \mathbf{1}(y_b = k).$$

405

406 The total loss function we minimise in Step 3 is given by

$$407 \quad \mathcal{L}_{\text{scJoint}}(\mathcal{B}_0, \theta) = \mathcal{L}_1(\mathcal{B}_0, \theta) + \mathcal{L}_{\text{center}}(\mathcal{B}_0, \theta).$$

408

409 We perform a final round of majority vote by KNN using distances in the embedding space.
410 If the prediction of any ATAC cell is different from Step 2, we update both its prediction and
411 probability score in the same way as Step 2. Before visualization with tSNE, all embedding
412 vectors are normalized using L_2 norm.

413 **Training details**

414 The batch size B was set to 256 in all cases. The other training details including learning rate
415 and number of training epochs used in each dataset can be found in Table S1. We started all the
416 training with learning rate set to 0.01, since a large learning rate has the benefit of faster training.
417 However, if the values of the loss functions were observed to have too much fluctuation, we
418 decreased the learning rate to 0.001 for more stable training.

419 **Data preprocessing**

420 • *Mouse atlas data.* The processed gene expression matrix and the cell type annota-
421 tion of the Tabula Muris mouse data of scRNA-seq were downloaded from <https://tabula-muris.ds.czbiohub.org/>, which have 41965 cells from protocol
422

423 fluorescence-activated cell sorting (FACS) and 54439 cells from microfluidic droplets
424 (droplet). The quantitative gene activity score matrix and the cell type annotation of Mouse
425 sci-ATAC-seq Atlas were downloaded from <https://atlas.gs.washington.edu/mouse-atac/>, including 81173 cells in total. The number of common genes be-
426 tween two modalities is 15519. We manually checked the cell type annotations from the
427 original studies and re-annotated the labels such that the naming convention is consistent
428 across the datasets. For example, the cell type “Cardiac muscle cell” in the sci-ATAC-seq
429 dataset was changed to “Cardiomyocytes”. We also combined some of the cellular sub-
430 types in the sci-ATAC-seq data to increase the percentage of overlapping labels between
431 two atlases for evaluation. More specifically, we combined “Regulatory T cell” and “T
432 cell” into “T cell”; “Immature B cell”, “Activated B cell” and “B cell” into “B cell”; “Ex-
433 citatory neurons” and “Inhibitory neurons” into “Neuron”.

435 • *SNARE-seq data*. The SNARE-seq data from adult mouse cerebral cortex was downloaded
436 from the National Center for Biotechnology Information (NCBI) Gene Expression Om-
437 nibus (GEO) accession number GSE126074 [14], with both raw gene expression and DNA
438 accessibility measurements available for the same cell. The fastq files were downloaded
439 from the Sequence Read Archive (SRA) for SRP183521. We first derived the fragment
440 files from the fastq files using `sinto fragments` (`sinto v0.7.2`), and then generated
441 the gene activity matrix using Signac (`v1.1.0.9000`) [30]. The cell type information was
442 obtained from the original study [14]. We filtered out the cells that were originally labeled
443 as “Misc” (cells of miscellaneous cluster), resulting in a dataset with 9190 cells and 15725
444 genes for the integrative analysis.

445 • *Multi-modal data (CITE-seq and ASAP-seq PBMC data)*. The ASAP-seq and CITE-seq
446 data were downloaded from GEO accession number GSE156478 [32], which included the
447 fragment files and antibody-derived tags (ADTs) matrices for ASAP-seq, the raw unique
448 molecular identifier (UMI) and ADT matrices for CITE-seq, from both control and stim-
449 ulated conditions. The gene activity matrices for ASAP-seq were generated by Signac.
450 Most of the thresholds we used for quality control metrics were consistent with those in
451 the original paper [32]. The control and stimulated CITE-seq were filtered based on the
452 following criteria: mitochondrial reads greater than 10%; number of expressed genes less
453 than 500; total number of UMI less than 1000; total number of ADTs from the rat iso-
454 type control greater 55 and 65 in the control and stimulated conditions respectively; total

455 number of UMI greater than 12,000 and 20,000 for the control and stimulated conditions
456 respectively; total number of ADTs less than 10,000 and 30,000 for control and stimulated
457 conditions respectively. We further filtered out cells that were classified as doublets in
458 original study. For the ASAP-seq data, we filtered out cells with the number ADTs more
459 than 10,000 and number of peaks more than 100,000. Finally, 4502 cells (control) and
460 5468 cells (stimulated) from ASAP-seq, 4644 cells (control) and 3474 cells (stimulated)
461 from CITE-seq were included in the downstream analysis. The number of common genes
462 across the four matrices is 17441 and the number of common ADTs is 227. We used
463 CiteFuse to integrate the peak matrix or gene expression matrix with their corresponding
464 protein expression and obtain clustering for ASAP-seq and CITE-seq within each condi-
465 tion separately [33]. For ASAP-seq, the similarity matrices of the chromatin accessibility
466 are calculated by applying the Pearson correlation to the TF-IDF transformation of the
467 peak matrix. We then followed the procedure described in [37] to annotate the clusters.

468 For scJoint, all the gene expression matrices and gene activity score matrices were binarized as 0
469 or 1, with 1 representing any non-zero original values, as the final input for training. Binarization
470 scales the two modalities so that their distributions have the same range and reduces the noise
471 level in the data for easier co-training.

472 **Settings used in other methods**

473 For the unpaired data (mouse cell atlases and multi-modal data from CITE-seq and ASAP-seq),
474 we benchmarked the performance of scJoint against three other methods designed for integrating
475 unpaired single-cell multi-modal data: Seurat (v3), Conos and Liger. We compared the label
476 transfer accuracy with Seurat and Conos and the joint visualizations with all three methods.
477 For the paired data (SNARE-seq), we further compared joint visualizations with two methods
478 specifically designed for paired data, scAI and MOFA+. For all the unpaired methods, we used
479 gene activity matrices derived from the above data preprocessing step as input for scATAC-seq.
480 For the two paired methods, we used the peak matrices of scATAC-seq data as input. Detailed
481 settings used in each method are as follows.

- 482 • *Seurat*. R package Seurat v3.2.0 [24] was used for all the datasets. The raw count
483 matrix of scRNA-seq and unnormalized gene activity score matrix of scATAC-seq were
484 used as input, which were then normalized using the `NormalizeData` function in Seu-
485 rat. Noted that for the CITE-seq and ASAP-seq data, the input was a concatenated

486 matrix of log-transformed normalized gene expression data/gene activity score matrix
487 and log-transformed ADTs matrix. Top 2000 most variable genes were selected from
488 scRNA-seq using `FindVariableFeatures` with `vst` as method. To identify the an-
489 chors between scRNA-seq and scATAC-seq data, `FindTransferAnchors` function
490 was used with “cca” as reduction method. The scATAC-seq data was then imputed using
491 `TransferAnchors` function, where the anchors were weighted by latent semantic in-
492 dexing (LSI) reduced dimension of scATAC-seq. Principal component analysis was then
493 performed on the merged matrix of scRNA-seq data and imputed scATAC-seq data. For all
494 the datasets, 30 principal components (PCs) were used for joint visualization with tSNE
495 (function `RunTSNE`).

496 For the mouse cell atlas data, we first integrated the two scRNA-seq datasets (FACS and
497 droplet) using `FindIntegrationAnchors` and `IntegrateData`, and then the inte-
498 grated matrix was scaled using `ScaleData` and used as reference to find transfer anchors.

499 • *Conos*. R package `conos` v1.3.1 [23] was used for all the datasets. Function
500 `basicP2proc` in `pagoda2` package (v0.1.2) was performed to process the raw
501 count matrix of scRNA-seq and unnormalized gene activity score matrix of scATAC-
502 seq. The joint graph was built using `buildGraph` with `k=15`, `k.self=5`, and
503 `k.self.weigh=0.01`, which were set as suggested in the tutorial for integrating
504 RNA and ATAC ([http://pklab.med.harvard.edu/peterk/conos/atac_](http://pklab.med.harvard.edu/peterk/conos/atac_rna/example.html)
505 [rna/example.html](http://pklab.med.harvard.edu/peterk/conos/atac_rna/example.html)). The joint visualization of scRNA-seq and scATAC-seq were
506 generated using `largeVis` by `embedGraph`, which is the default visualization in `Conos`.

507 • *Liger*. R package `liger` v0.5.0 [21] was used for the datasets. The raw count matrix of
508 scRNA-seq and unnormalized gene activity score matrix of scATAC-seq were used as in-
509 put, which were normalized using `normalize` function in `liger`. Highly variable genes
510 were selected using the scRNA-seq. For the mouse cell atlas data, both FACS and droplet
511 scRNA-seq data were used to select features. For all the datasets, number of factors was
512 set to 20 in `optimizeALS`. tSNE was then performed on the normalized cell factors
513 to generate the joint visualization of scRNA-seq and scATAC-seq (function `runTSNE` in
514 `liger`).

515 • *scAI*. R package `scAI` v1.0.0 [16] was used for the integration of SNARE-seq data. The raw
516 count matrix of scRNA-seq and raw peak matrix of scATAC-seq were used as input. We

517 ran scAI using `run_scAI` by setting the rank of the inferred factor set as 20, `do.fast =`
518 `TRUE`, and `nrun = 1`, with other parameters set as default, as suggested in the pipeline in
519 the github repository. tSNE plots were generated using `reducedDims` function in scAI.

520 • *MOFA+*. R package MOFA2 v1.0 [17] was used for the integration of SNARE-seq data.
521 Following the suggested integration tutorial for SNARE-seq in the github repository, we
522 first selected top 2500 most variable genes using `FindVariableFeatures` in Seurat
523 package with `vst` as method and top 5000 most variable ATAC peaks with `disp` as
524 method. By subsetting the counts matrix of scRNA-seq and peak matrix of scATAC-seq
525 with the selected features, we ran MOFA+ by setting the number of factors as 10, with
526 other parameters set as default. tSNE plots were generated using `run_tsne` function in
527 MOFA2.

528 Evaluation metrics

529 Joint embedding evaluation - Silhouette coefficients

530 To evaluate whether the joint embeddings from different methods show clustering structure
531 reflecting biological signals or technical variations, we calculated the silhouette coefficient
532 for each cell by considering two different groupings: (1) grouping based on the modalities
533 (scRNA-seq or scATAC-seq), called the modality silhouette coefficient ($s_{modality}$); (2) grouping
534 based on known cell types, called the cell type silhouette coefficient ($s_{cellTypes}$). Note that for
535 the atlas data, we consider FACS and droplet in scRNA-seq as two distinct technologies and the
536 modality silhouette coefficient has three groups (FACS, droplet, ATAC) in the calculation. For
537 SNARE-seq, the paired methods (scAI and MOFA+) have no modality silhouette coefficients
538 since each cell has one paired profile of RNA and ATAC. An ideal joint visualization should have
539 low modality silhouette coefficients, suggesting the removal of the technical effect, and large
540 cell type silhouette coefficients, indicating the cells are grouped by cell types. The euclidean
541 distance for all methods except Conos is obtained from the tSNE embedding. For Conos, the
542 distance is obtained from the `largeVis` embedding, which is the method's default output.

543

544 We then summarize the two silhouette coefficients by calculating an F1-score as follows:

$$545 F1_{sil} = \frac{2 \cdot (1 - s'_{modality}) \cdot s'_{cellTypes}}{1 - s'_{modality} + s'_{cellTypes}},$$

546 where $s' = (s + 1)/2$. A higher F1 score indicates better performance in the alignment of the
547 modalities as well as the preservation of biological signals.

548 **Accuracy evaluation of transferred labels**

549 We evaluated the accuracy of label transfer from two aspects: (1) Overall accuracy rate; (2)
550 Cell type classification F1-score. The overall accuracy rate was computed only accounting for
551 the common cell types between scRNA-seq and scATAC-seq data. The cell type classification
552 F1-score is the harmonic mean of precision and recall of each cell type.

553 **Software availability**

554 scJoint was implemented using PyTorch (version 1.0.0) with code available at [https://](https://github.com/SydneyBioX/scJoint)
555 github.com/SydneyBioX/scJoint.

556 **Acknowledgments**

557 The authors gratefully acknowledge the following funding sources: Research Training Program
558 Tuition Fee Offset and Stipend Scholarship and Chen Family Research Scholarship to Y.L.; Aus-
559 tralian Research Council Discovery Project grant (DP170100654) to J.Y.H.Y.; Australian Re-
560 search Council DECRA Fellowship (DE180101252) to Y.X.R.W; NIH grants R01 HG010359
561 and P50 HG007735 to W.H.W.

562 **Author contributions**

563 T.W., W.H.W. and Y.X.R.W. conceived and designed this project; Y.L., T.W. and S.W. per-
564 formed data preprocessing, model development, and evaluation of results; J.Y.H.Y., W.H.W. and
565 Y.X.R.W. supervised the execution; Y.L., J.Y.H.Y., W.H.W. and Y.X.R.W. wrote the manuscript.
566 All authors read and approved the manuscript.

567 **Conflict of interest**

568 The authors declare that they have no conflict of interest.

569 **References**

- 570 1. Stuart, T. & Satija, R. Integrative single-cell analysis. *Nature Reviews Genetics* **20**, 257–
571 272 (2019).
- 572 2. Berger, S. L. The complex language of chromatin regulation during transcription. *Nature*
573 **447**, 407–412 (2007).
- 574 3. Klemm, S. L., Shipony, Z. & Greenleaf, W. J. Chromatin accessibility and the regulatory
575 epigenome. *Nature Reviews Genetics* **20**, 207–220 (2019).
- 576 4. Pott, S. & Lieb, J. D. Single-cell ATAC-seq: strength in numbers. *Genome Biology* **16**, 172
577 (2015).
- 578 5. Schaum, N. *et al.* Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris:
579 The Tabula Muris Consortium. *Nature* **562**, 367 (2018).
- 580 6. Regev, A. *et al.* Science forum: the human cell atlas. *Elife* **6**, e27041 (2017).
- 581 7. Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for
582 single-cell transcriptomics. *Nature methods* **15**, 1053–1058 (2018).
- 583 8. Wang, J. *et al.* Data denoising with transfer learning in single-cell transcriptomics. *Nature*
584 *methods* **16**, 875–878 (2019).
- 585 9. Lin, Y. *et al.* scMerge leverages factor analysis, stable expression, and pseudoreplication
586 to merge multiple single-cell RNA-seq datasets. *Proceedings of the National Academy of*
587 *Sciences* **116**, 9775–9784 (2019).
- 588 10. Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony.
589 *Nature methods*, 1–8 (2019).
- 590 11. Wang, T. *et al.* BERMUDA: a novel deep transfer learning method for single-cell RNA
591 sequencing batch correction reveals hidden high-resolution cellular subtypes. *Genome bi-*
592 *ology* **20**, 1–15 (2019).
- 593 12. Amodio, M. *et al.* Exploring single-cell data with deep multitasking neural networks. *Na-*
594 *ture methods*, 1–7 (2019).
- 595 13. Xiong, L. *et al.* SCALE method for single-cell ATAC-seq analysis via latent feature extrac-
596 tion. *Nature communications* **10**, 1–10 (2019).
- 597 14. Chen, S., Lake, B. B. & Zhang, K. High-throughput sequencing of the transcriptome and
598 chromatin accessibility in the same cell. *Nature biotechnology* **37**, 1452–1457 (2019).

- 599 15. Cao, J. *et al.* Joint profiling of chromatin accessibility and gene expression in thousands of
600 single cells. *Science* **361**, 1380–1385 (2018).
- 601 16. Jin, S., Zhang, L. & Nie, Q. scAI: an unsupervised approach for the integrative analysis
602 of parallel single-cell transcriptomic and epigenomic profiles. *Genome biology* **21**, 1–19
603 (2020).
- 604 17. Argelaguet, R. *et al.* MOFA+: a statistical framework for comprehensive integration of
605 multi-modal single-cell data. *Genome Biology* **21**, 1–17 (2020).
- 606 18. Welch, J. D., Hartemink, A. J. & Prins, J. F. MATCHER: manifold alignment reveals cor-
607 respondence between single cell transcriptome and epigenome dynamics. *Genome biology*
608 **18**, 1–19 (2017).
- 609 19. Amodio, M. & Krishnaswamy, S. MAGAN: Aligning biological manifolds. *arXiv preprint*
610 *arXiv:1803.00385* (2018).
- 611 20. Liu, J., Huang, Y., Singh, R., Vert, J.-P. & Noble, W. S. Jointly embedding multiple single-
612 cell omics measurements. *BioRxiv*, 644310 (2019).
- 613 21. Welch, J. D. *et al.* Single-cell multi-omic integration compares and contrasts features of
614 brain cell identity. *Cell* **177**, 1873–1887 (2019).
- 615 22. Duren, Z. *et al.* Integrative analysis of single-cell genomics data by coupled nonnegative
616 matrix factorizations. *Proceedings of the National Academy of Sciences* **115**, 7723–7728
617 (2018).
- 618 23. Barkas, N. *et al.* Joint analysis of heterogeneous single-cell RNA-seq dataset collections.
619 *Nature methods* **16**, 695–698 (2019).
- 620 24. Stuart, T. *et al.* Comprehensive integration of single-cell data. *Cell* **177**, 1888–1902 (2019).
- 621 25. Yosinski, J., Clune, J., Nguyen, A., Fuchs, T. & Lipson, H. Understanding neural networks
622 through deep visualization. *arXiv preprint arXiv:1506.06579* (2015).
- 623 26. Cusanovich, D. A. *et al.* A single-cell atlas of in vivo mammalian chromatin accessibility.
624 *Cell* **174**, 1309–1324 (2018).
- 625 27. Maaten, L. v. d. & Hinton, G. Visualizing data using t-SNE. *Journal of machine learning*
626 *research* **9**, 2579–2605 (2008).
- 627 28. McInnes, L., Healy, J. & Melville, J. Umap: Uniform manifold approximation and projec-
628 tion for dimension reduction. *arXiv preprint arXiv:1802.03426* (2018).

- 629 29. Pliner, H. A. *et al.* Cicero predicts cis-regulatory DNA interactions from single-cell chromatin accessibility data. *Molecular cell* **71**, 858–871 (2018).
630
- 631 30. Stuart, T., Srivastava, A., Lareau, C. & Satija, R. Multimodal single-cell chromatin analysis with Signac. *bioRxiv* (2020).
632
- 633 31. Stoeckius, M. *et al.* Simultaneous epitope and transcriptome measurement in single cells. *Nature methods* **14**, 865 (2017).
634
- 635 32. Mimitou, E. P. *et al.* Scalable, multimodal profiling of chromatin accessibility and protein levels in single cells. *bioRxiv* (2020).
636
- 637 33. Kim, H. J., Lin, Y., Geddes, T. A., Yang, J. Y. H. & Yang, P. CiteFuse enables multi-modal analysis of CITE-seq data. *Bioinformatics* **36**, 4137–4143 (2020).
638
- 639 34. Godfrey, D. I., MacDonald, H. R., Kronenberg, M., Smyth, M. J. & Van Kaer, L. NKT cells: what's in a name? *Nature Reviews Immunology* **4**, 231–237 (2004).
640
- 641 35. Finak, G. *et al.* MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome biology* **16**, 1–13 (2015).
642
643
- 644 36. Abdelaal, T. *et al.* A comparison of automatic cell identification methods for single-cell RNA sequencing data. *Genome biology* **20**, 194 (2019).
645
- 646 37. Maecker, H. T., McCoy, J. P. & Nussenblatt, R. Standardizing immunophenotyping for the human immunology project. *Nature Reviews Immunology* **12**, 191–200 (2012).
647