

1 **MultiMAP: Dimensionality Reduction and Integration of Multimodal Data**

2
3 Mika Sarkin Jain^{1,2*}, Krzysztof Polanski², Cecilia Dominguez Conde², Xi Chen^{2,3}, Jongeun
4 Park^{2,4}, Lira Mamanova², Andrew Knights², Rachel A. Botting⁵, Emily Stephenson⁵, Muzlifah
5 Haniffa^{2,5}, Austen Lamacraft¹, Mirjana Efremova^{2,6*}, Sarah A. Teichmann^{1,2*}

6
7 Affiliations:

- 8 1. Theory of Condensed Matter, Dept Physics, Cavendish Laboratory, University of
9 Cambridge, JJ Thomson Ave, Cambridge CB3 0HE, UK
- 10 2. Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10
11 1SA, UK.
- 12 3. Southern University of Science and Technology, 1088 Xueyuan Ave, Nanshan,
13 Shenzhen, Guangdong Province, China, 518055
- 14 4. KAIST, 291 Daehak-ro, Eoeun-dong, Yuseong-gu, Daejeon, South Korea
- 15 5. Biosciences Institute, Newcastle University, Newcastle upon Tyne, NE2 4HH, UK
- 16 6. Barts Cancer Institute, Queen Mary University of London, London, UK

17
18 *Co-corresponding authors: mikasarkinjain@gmail.com, st9@sanger.ac.uk,
19 m.efremova@qmul.ac.uk

20 **Abstract**

21 Multimodal data is rapidly growing in many fields of science and engineering, including single-cell
22 biology. We introduce MultiMAP, an approach for dimensionality reduction and integration of
23 multiple datasets. MultiMAP recovers a single manifold on which all of the data resides and then
24 projects the data into a single low-dimensional space so as to preserve the structure of the
25 manifold. It is based on a framework of Riemannian geometry and algebraic topology, and
26 generalizes the popular UMAP algorithm¹ to the multimodal setting. MultiMAP can be used for
27 visualization of multimodal data, and as an integration approach that enables joint analyses.
28 MultiMAP has several advantages over existing integration strategies for single-cell data,
29 including that MultiMAP can integrate any number of datasets, leverages features that are not
30 present in all datasets (i.e. datasets can be of different dimensionalities), is not restricted to a
31 linear mapping, can control the influence of each dataset on the embedding, and is extremely
32 scalable to large datasets. We apply MultiMAP to the integration of a variety of single-cell
33 transcriptomics, chromatin accessibility, methylation, and spatial data, and show that it
34 outperforms current approaches in preservation of high-dimensional structure, alignment of
35 datasets, visual separation of clusters, transfer learning, and runtime. On a newly generated
36 single-cell Assay for Transposase-Accessible Chromatin using sequencing (scATAC-seq) and
37 single-cell RNA-seq (scRNA-seq) dataset of the human thymus, we use MultiMAP to integrate
38 cells along a temporal trajectory. This enables the quantitative comparison of transcription factor
39 expression and binding site accessibility over the course of T cell differentiation, revealing
40 patterns of transcription factor kinetics.

41

42 Introduction

43

44 Multimodal data is rapidly growing in many fields of science and engineering, including single-cell
45 biology. Emerging single-cell technologies are providing high-resolution measurements of
46 different features of cellular identity, including single-cell assays for gene expression, protein
47 abundance^{2,3}, chromatin accessibility⁴, DNA methylation⁵, and spatial resolution⁶. Large scale
48 collaborations including the Human Cell Atlas international consortium^{7,8} are generating an
49 exponentially increasing amount of data of many biological tissues, using a myriad of
50 technologies. Each technology provides a unique view of cellular biology and has different
51 strengths and weaknesses. Integrating these measurements in the study of a single biological
52 system will open avenues for more comprehensive study of cellular identity, cell-cell interactions,
53 developmental dynamics, and tissue structure⁹.

54

55 The integration of multi-omic data poses several challenges¹⁰. Different omics technologies
56 measure distinct unmatched features with different underlying distributions and properties and
57 hence produce data of different dimensionality. This makes it difficult to place data from different
58 omics in the same feature space. Additionally, omics technologies can also have different noise
59 and batch characteristics which are challenging to identify and correct. Further, as multi-omic data
60 grows along two axes, the number of cells per omic and the number of omics per study, integration
61 strategies need to be extremely scalable.

62

63 Most data integration methods project multiple measurements of information into a common low-
64 dimensional representation to assemble multiple modalities into an integrated embedding space.
65 Recently published methods employ different algorithms to project multiple datasets into an
66 embedding space, including canonical correlation analysis (CCA)¹¹, nonnegative matrix
67 factorization (NMF)¹² or variational autoencoders¹³. In the field of genomics, single-cell
68 transcriptomics, as a well-established method, often serves as a common reference, facilitating
69 the transfer of cell type annotation and data across multiple technologies and modalities. While
70 these methods can be tremendously powerful, they require correspondence between the features
71 profiled across omics technologies. Another limitation of many existing methods is they are
72 challenged by scaling to large datasets.

73

74 Here we introduce a method that overcomes all these limitations: MultiMAP, an approach for the
75 dimensionality reduction and integration of multiple datasets. MultiMAP integrates data by
76 constructing a non-linear manifold on which diverse high-dimensional data reside and then
77 projecting the manifold and data into a shared low-dimensional space. In contrast to other
78 integration strategies for single-cell data, MultiMAP can integrate any number of datasets, is not
79 restricted to a linear mapping, leverages features that are not present in all datasets (*i.e.* datasets
80 can be of different dimensionalities), can control the influence of each dataset on the embedding,
81 and is effortlessly scalable to large datasets. The ability of MultiMAP to integrate datasets of
82 different dimensionalities allows the strategy to leverage information that is not considered by
83 methods that operate in a shared feature space. (e.g. MultiMAP can integrate the 20,000-feature
84 gene space of scRNAseq data together with a 100,000-feature peak space of scATACseq data).

85
86 We apply MultiMAP to challenging synthetic multimodal data, and demonstrate its ability to
87 integrate a wide range of single-cell omics datasets. Finally, we apply the approach to the study
88 of T cell development with new scATACseq data from fetal thymi. We show that MultiMAP can
89 co-embed datasets across different technologies and modalities, while at the same time
90 preserving the structure of the data, even with extensive biological and technical differences. The
91 resulting embedding and shared neighborhood graph (MultiGraph) can be used for simultaneous
92 visualisation and integrative analysis of multiple datasets. With respect to single cell genomics
93 data, this allows for standard analysis on the integrated data, such as cluster label transfer, joint
94 clustering, and trajectory analysis.
95

96 Results

97

98 The MultiMAP Framework

99

100 We introduce MultiMAP, an approach for integration and dimensionality reduction of multimodal
101 data based on a framework of Riemannian geometry and algebraic topology. MultiMAP takes as
102 input any number of datasets of potentially differing dimensions. MultiMAP recovers geodesic
103 distances on a single latent manifold on which all of the data is uniformly distributed. The distances
104 are calculated between data points of the same dataset by normalizing distances with respect to
105 a neighborhood distance specific to the dataset, and between data points of different datasets by
106 normalizing distances between the data in a shared feature space with respect to a neighborhood
107 parameter specific to the shared feature space. These distances are then used to construct a
108 neighborhood graph (MultiGraph) on the manifold. Finally, the data and manifold space are
109 projected into a low-dimensional embedding space by minimizing the cross entropy of the graph
110 in the embedding space with respect to the graph in the manifold space. MultiMAP allows the
111 user to modify the weight of each dataset in the cross entropy loss, allowing the user to modulate
112 the contribution of each dataset to the layout. Integrated analysis can be performed on the
113 embedding or the graph, and the embedding also provides an integrated visualization. The
114 mathematical formulation of MultiMAP is elaborated in Supplementary Methods.
115

116 In order to study MultiMAP in a controlled setting, we first applied it to two synthetic examples of
117 multimodal data (Methods). The first synthetic data consists of points sampled randomly from the
118 canonical 3D “Swiss Roll” surface and the 2D rectangle (Figure 2a). The dataset is considered
119 multimodal data, because samples are drawn from different feature spaces but describe the same
120 rectangular manifold. In addition, we are given the position along the manifold of 1% of the data.
121 This synthetic setting illustrates that MultiMAP can integrate data in a nonlinear fashion and
122 operate on datasets of different dimensionality, because data points along a similar position on
123 the manifold are near each other in the embedding (Figure 2b). The MultiMAP embedding

124 properly unrolls the Swiss Roll dataset, indicating that the projection is nonlinear. The embedding
125 also appears to preserve aspects of both datasets; the data is curved and at the same time
126 unrolled.

127
128 To determine if MultiMAP can effectively leverage features unique to certain datasets, we used
129 the MNIST database¹⁴, where handwritten images were split horizontally with thin overlap (Figure
130 2c; see Methods for details). The two datasets can be considered multimodal because they have
131 different feature spaces but describe the same set of digit images. The thin overlapping region of
132 the two halves is not enough information to create a good embedding of the data (Figure 2c).
133 Many distinct digits are similar in this thin central sliver, and hence they cluster together in the
134 feature space of this sliver. Indeed, in a UMAP projection of the data in the shared feature space
135 of this overlap, the clusters of different digits are not as well separated as in the UMAP projections
136 of each half (Figure 2c).

137
138 A multimodal integration strategy that effectively leverages all features would use the features
139 unique to each half to separate different digits, and the shared space to bring the same digits from
140 each dataset close together (Figure 2d). We show that with MultiMAP the different modalities are
141 well mixed in the embedding space and the digits cluster separately, despite mostly different
142 feature spaces and noise being added to only the second dataset. This indicates that MultiMAP
143 is leveraging the features unique to each dataset and is also robust to datasets with different
144 noise.

145
146 Moreover, MultiMAP has weight parameters ω^v which control the contribution of each dataset \mathbf{X}^v
147 to the final embedding, allowing the user to modulate which dataset has a greater influence on
148 the MultiMAP embedding. When a dataset's weight is larger, its structure has a larger contribution
149 to the MultiMAP embedding. Our results show that when integrating the MNIST data, for different
150 choices of ω^v , the datasets remain well integrated in the embedding space (Extended Data Figure
151 1a,b).

152
153 Finally, to illustrate that our assumption of a shared manifold is robust to variable levels of overlap
154 across datasets, we used MultiMAP to integrate datasets with varying numbers of shared clusters
155 in the MNIST data (Extended Data Figure 2). Our results show that MultiMAP is able to effectively
156 integrate datasets that have only 1 out of 10 clusters shared between them. The transfer accuracy,
157 silhouette score, and structure score of the MultiMAP integration remained largely constant as
158 the number of overlapping clusters is varied, demonstrating that MultiMAP is highly robust to
159 differences in populations between datasets.

160
161

162

163 MultiMAP integration of single-cell transcriptomics and chromatin 164 accessibility

165

166 Having shown that MultiMAP succeeds in integrating synthetic data, we apply the technique to
167 real biological data. Epigenomic regulation underlies gene expression and cellular identity. Hence,
168 integration of single-cell transcriptomics and epigenomics data provides an opportunity to
169 investigate how epigenomic alterations regulate gene expression to determine and maintain cell
170 identity. In addition, effective integration with transcriptomics data can improve the sensitivity and
171 interpretability of the sparse scATAC-seq data.

172

173 To assess MultiMAP's ability to integrate transcriptomic and epigenomic data, we applied the
174 approach to integrate our previously generated high-coverage scATAC-seq data of mouse
175 splenocytes¹⁵ and generated corresponding single-cell transcriptomic profiles of the same tissue.
176 The high coverage of the plate-based scATAC-seq data as well as the published cluster
177 annotations of the subpopulations served as a good ground truth example to validate our method.
178 The analysis of the transcriptomics data revealed similar subpopulations to the published
179 scATAC-seq dataset, in addition to two RNA-specific clusters: a subpopulation of B cells with
180 higher expression of Interferon-Induced (Ifit) genes and a subpopulation of proliferating cells
181 (Extended Data Figure 3a,b).

182

183 MultiMAP effectively integrated the two datasets, using both gene activity scores and the cell type-
184 specific epigenetic information outside of gene bodies. The different modalities are well mixed in
185 the embedding space and cells annotated as the same type are close together, regardless of the
186 modality for different choices of ω^v (Figure 3a, Extended Data Figure 1c,d). Next, we jointly
187 clustered cells from both datasets using the MultiGraph. This produced clusters with markers
188 corresponding to known cell types¹⁵ (Extended Data Figure 3c). The annotations produced by this
189 joint clustering were generally consistent with independent annotations of each dataset (Figure
190 3c). Two of the clusters determined to be proliferating cells and B cells with upregulated Ifit genes
191 were found only in the scRNA-seq data, as expected (Figure 3a, Extended Data Figure 3b). In
192 addition, the integration produced by MultiMAP is robust to different choices of the weight
193 parameters (Extended Data Figure 1c).

194

195 Further, we used the MultiGraph to directly predict the cell types of the scATAC-seq given the cell
196 types of the scRNA-seq. Figure 3d shows the confusion matrix of the predictions, illustrating that
197 cells were generally annotated correctly. This illustrates the ability of MultiMAP to leverage
198 annotation efforts of one omic technology to inform those of another. Interestingly, a small subset
199 of cells from scRNA-seq previously annotated as T cells is now clearly separated on the MultiMAP
200 plot, and clusters close to the B cells (Figure 3a, Extended Data Figure 3). Doublet detection
201 confirmed that this cluster is composed of doublet T/B cells. These doublets are spread
202 throughout the UMAP plot of the scRNA-seq data, but are clearly distinct on the MultiMAP plot

203 (Extended Data Figure 3). This illustrates the power of MultiMAP both as a visualization tool, and
204 to reveal new populations of cells.

205
206 Next, we applied MultiMAP to integration of multiple batches from each data modality, to assess
207 the ability to account for batch effects. For this purpose, we used recently published scRNA-seq
208 and scATAC-seq data of human bone marrow and peripheral blood mononuclear cells¹⁶. This
209 dataset consists of 16 experimental samples, representing different experimental batches.
210 Another challenge is that cells are not in discrete clusters but rather on a continuum. MultiMAP is
211 able to simultaneously correct batch effects and modality differences, integrating all 16 datasets
212 into a consistent embedding (Figure 3e). The different modalities are well mixed in the embedding
213 and cells of the same type are close together, regardless of modality or batch. The cell type
214 annotations of all of the data were taken from the original publication¹⁶, so they provide a good
215 ground truth and independent validation of MultiMAP. Additionally, MultiMAP is able to correct
216 batch effects present in different omics technologies. Applying MultiMAP to just the scRNA-seq
217 data produces embedding that properly integrates cells of the same type regardless of batch, and
218 the same is true when MultiMAP is applied to only the scATAC-seq data (Figure 3f). It is also
219 evident in this figure that clusters with cell types unique to a batch remain unmixed in the
220 embedding. This indicates that MultiMAP is not forcing incompatible data to integrate and
221 demonstrates that MultiMAP can integrate datasets even if they have extensive technical
222 differences.

223 MultiMAP integration of multiple modalities of mouse brain cells

224
225 Recent advances in spatial sequencing technology enable the simultaneous measurement of
226 gene expression and spatial locations of single-cells, facilitating the study of tissue structure⁶.
227 While these technologies provide spatial information, they often measure only a small fraction of
228 the genes measured by scRNA-seq. Integration of spatial measurements and scRNA-seq has the
229 potential to provide spatial context to scRNA-seq data as well as to reveal finer grained biological
230 differences in the spatial data by leveraging the greater number of cells and genes present in
231 scRNA-seq data.

232
233 We applied MultiMAP to the integration of a Drop-seq scRNA-seq data of the mouse frontal
234 cortex¹⁷ and STARmap *in situ* gene expression dataset¹⁸. Despite the differences between the
235 two dataset in the number of measured genes (only 1020 in STARmap) and the number of cells
236 (71640 in Drop-seq versus 2137 in STARmap), our integrated analysis shows that MultiMAP
237 successfully integrates the datasets. Clustering the integrated data using the MultiGraph
238 produced clusters with markers corresponding to known cell types (Figure 4a,b). One of the
239 clusters, the claustrum, was found only in the scRNA-seq data, as expected. Integration with
240 MultiMAP also resulted in improved cluster annotation for both datasets. The excitatory L4
241 neurons were previously only present in the STARMap data, as the motor cortex and prefrontal
242 cortex that are part of the frontal cortex are considered to lack a layer 4 in mice¹⁹. However, after
243 the integration we also identified L4 cells in the scRNA-seq data previously annotated as L5
244 neurons (Figure 4a,c, Extended Data Figure 4). A similar population of pyramidal cells located

245 between layers 3 and 5 were recently identified both with anatomical and single-cell studies^{20,21}.
246 This was confirmed by expression of marker genes associated with L4, including *Cux2* and *Rorb*
247 (Extended Data Figure 4). This illustrates the power of MultiMAP to reveal new cell types.

248
249 MultiMAP also improves visualization of the STARmap data. Before integration with MultiMAP,
250 many of the cell types of the spatial data did not cluster separately and were visually hard to
251 distinguish. In comparison, the MultiMAP embedding of the STARmap data exhibits tighter cell
252 type clusters and increased separation between cell types (Figure 4e). This improvement was
253 measured by the average Silhouette score in the embedding space, which is significantly larger
254 for MultiMAP (Figure 4e).

255
256 Integration with MultiMAP also enabled us to spatially locate all the joint cell types in the
257 STARmap data, allowing study of the spatial structure of the tissue (Figure 4d). The pyramidal
258 neurons localize to layers 2-6 and oligodendrocytes localize to the layer below the cortex,
259 whereas the interneurons do not appear to exhibit spatial organization. These observations are
260 all consistent with the known spatial architecture of the mouse visual cortex¹⁸.

261
262 To investigate the performance of MultiMAP on the integration of more than two modalities, we
263 applied the approach to integrate recently published multi-omics datasets of the mouse primary
264 motor cortex²⁰ consisting of 9 separate datasets, including 7 single-cell or single-nucleus
265 transcriptomics datasets, one single-nucleus chromatin accessibility, and one single-nucleus
266 DNA methylation (snmC-seq). MultiMAP successfully co-embedded more than 600,000 single-
267 cell or -nucleus samples assayed by six molecular modalities and identified the previously
268 published cell subpopulations. The MultiMAP embedding displays good mixing of clusters from
269 different modalities when the clusters correspond to the same cell type. Cell type annotations
270 were taken from the original publication of the data, so they provide a good ground truth and an
271 independent validation of MultiMAP. We further see that cell types that exist in one modality, but
272 not in the others, are not falsely aligned in the embedding space. This indicates that MultiMAP is
273 not forcing incompatible data to integrate.

274
275 Finally, using the integration of scRNA-seq with the STARmap data, as well as the integration of
276 the multi-omics spleen data, we assessed the impact of using only shared vs. all features. We
277 find that using all features greatly improves the integration and results in embeddings that are
278 visually and quantitatively superior, according to four performance metrics (Extended Data Figure
279 5). This illustrates that non-shared features can be extremely helpful, and demonstrates an
280 advantage of MultiMAP over other methods which do not consider non-shared features.

281

282 Benchmarking

283
284 We assessed and benchmarked the performance of MultiMAP against several popular
285 approaches for integrating single-cell multi-omics, including Seurat³¹¹, LIGER¹², Conos²² and
286 GLUER²³.

287

288 These integration approaches differ in key regards, summarized in Figure 5d. We used a diversity
289 of performance metrics to comprehensively compare MultiMAP with other approaches, including
290 transfer accuracy, silhouette score, alignment, preservation of the structure, and runtime. With
291 these metrics, we quantified the separation of the joint clusters, how well mixed the datasets were
292 after integration and how well they preserved the structure in the original datasets to investigate
293 whether the methods integrate populations across datasets without blending distinct populations
294 together.

295

296 To this end, we generated single-nucleus data from human Peripheral Blood Mononuclear Cells
297 (PBMCs) using the Multiome ATAC + RNA kit. We obtained a PBMC atlas of 6,344 nuclei of high-
298 quality ATAC + RNA profiles. We analysed and annotated the RNA and ATAC data separately,
299 revealing all the known major PBMC types: CD14 and CD16 monocytes, cDCs and pDCs, naive
300 and effector CD4 and CD8 T cells, Tregs, MAIT and gamma-delta T cells, NK and ILCs, naive
301 and memory B cells and plasmablasts (Extended Data Figure 6a). Most cell types were well
302 separated in both modalities with the exception of the NK and ILC clusters and the gamma-delta
303 and the CD8 effector T cells that blended together in the ATAC data.

304

305 We used the PBMCs as a gold standard dataset to benchmark MultiMAP against the four other
306 methods. As shown in the co-embedding and the metrics, MultiMAP successfully integrated the
307 cell types across modalities and outperformed other methods (Extended Data Figure 6b,c). The
308 label transfer accuracy was particularly striking, with MultiMAP achieving a much higher score
309 compared to other methods.

310

311 Furthermore, we also benchmarked MultiMAP using a variety of multi-omic data with published
312 cell type annotations, including the transcriptomics and chromatin accessibility spleen data,
313 scRNA-seq and STARmap of the visual cortex, and the multi-omics data of the primary cortex.
314 For all datasets, MultiMAP achieves top or near top performance on all metrics (Figure 5a,b). The
315 embeddings produced by MultiMAP prove superior for transferring cell type annotations between
316 datasets, separating clusters of different cell populations, integrating datasets in a well-mixed
317 manner, and capturing the high-dimensional structure of each dataset. Critically, MultiMAP is
318 faster than all other benchmarked methods, and significantly faster than LIGER and Seurat 3
319 (Figure 5c). Seurat 3 and LIGER were not able to scale to the primary cortex data of 600k,
320 producing out-of-memory errors despite access to 218 GB of RAM.

321

322 Finally, to assess the batch correction performance of MultiMAP, we also applied it on three
323 scRNA-seq studies of the human pancreas²⁴⁻²⁶ that were recently used for comparison of eight
324 batch correction methods²⁷. Even though the main purpose of MultiMAP is the integration of
325 several different omic technologies, MultiMAP outperformed all other well established batch
326 correction methods in the field, demonstrating that MultiMAP can correct batches and integrate
327 multiple omics data simultaneously (Extended Data Figure 7).

328 MultiMAP reveals patterns of T cell maturation along a multi-omic 329 trajectory

330

331 Single-cell transcriptomics has enabled reconstruction of developmental trajectories and the
332 study of dynamic processes such as differentiation and reprogramming. Bulk RNA-seq and
333 ATAC-seq data has further revealed regulatory events driving these processes²⁸. However, joint
334 analysis of single-cell expression and chromatin accessibility profiles along a time course
335 trajectory would allow the study of dynamic chromatin regulation alongside gene expression and
336 elucidate epigenomic drivers of transcriptional change^{29,30}.

337

338 In order to investigate the potential of integrating multi-omic data along a common differentiation
339 trajectory, we focused on T cell development in the thymus. The thymus is an organ essential for
340 the maturation and selection of T cells. Precursor cells migrate from the fetal liver and bone
341 marrow to the thymus where they develop into different types of mature T cells³¹. We recently
342 provided a comprehensive single-cell transcriptomics atlas of the human thymus during
343 development, childhood, and adult life, and computationally predicted the trajectory of T cell
344 development from early progenitors to mature T cells³¹. To expand on this and further investigate
345 the gene regulatory mechanisms driving T cell development, we generated single-cell
346 transcriptomics and chromatin accessibility data from a human fetal thymus sample at 10 weeks
347 of gestation.

348

349 Clustering of the scRNA-seq data revealed cell types identified in our recently published
350 transcriptomic cell atlas of the thymus³¹, including several clusters of T cells across different
351 stages of development, fibroblasts, endothelial cells, erythrocytes, thymic epithelial cells (TECs),
352 NK and ILC3 cells, and macrophages and dendritic cells (Extended Data Figure 8). The sparse
353 scATAC-seq and the continuous nature of cell types along the maturation trajectory made it
354 difficult to cluster the ATAC cells into different T cell types (Extended Data Figure 8). However,
355 the integration with MultiMAP and the joint clusters obtained using the MultiGraph corresponded
356 to the published thymus cell types³¹ (Figure 6 a,b), allowing us to correctly annotate the cell types
357 of the scATAC-seq data.

358

359 We then selected the T cell populations identified from the joint clustering and performed diffusion
360 map pseudotime analysis using the alignment MultiMAP graph. The reconstructed development
361 trajectory showed a continuous differentiation with the same trend as the published study, starting
362 from early double negative (DN) CD4-CD8-, gradually progressing to double positive (DP)
363 CD4+CD8+ T cells, and then differentiating into single positive (SP) mature CD8+ or CD4+ T
364 cells. Hallmark genes of T cell differentiation varied along the inferred pseudotime in a manner
365 consistent with³¹ (Figure 6d), serving as validation of the trajectory inference and the integration
366 produced by MultiMAP.

367

368 To identify transcription factors (TFs) that potentially regulate T cell development, we studied
369 changes in TF expression and TF binding site accessibility along the differentiation trajectory. The
370 top variable TFs/TF binding sites along the trajectory included many TFs that have been

371 previously shown to be involved in T cell differentiation, including GATA3, SPI1, MEF2C, ERG,
372 TCF3, TCF4, TFAP4, MYBL2, STAT1, NR4A2 and others^{28,31,32} (Figure 6e, Extended Data Figure
373 9). The TFs that most varied along the trajectory were found to show changes in motif accessibility
374 at the transition between the late DN and early DP stage of differentiation as shown before³².

375
376 Moreover, our integrated trajectory allowed us to identify TFs where changes in motif accessibility
377 and expression of the TF itself were closely coordinated, for example ZEB1, IRF1, REL, FOS and
378 others, suggesting that these TFs actively regulate their target genes immediately and directly
379 (Figure 6e). In contrast, for TFs such as ETS1, JUN etc., gene expression of the TF significantly
380 precedes the accessibility of the corresponding TF binding sites, suggesting that additional
381 regulatory mechanisms are potentially required for opening of the TF motifs.
382

383 Discussion

384
385 Here we present a novel approach for dimensionality reduction and integration of multimodal data
386 which takes into account the full data sets, even when they have different feature spaces.
387 MultiMAP embeds all datasets into a shared space to preserve both the manifold structure of
388 each dataset independently, as well as in shared feature spaces. This enables both visualization
389 and streamlined downstream analyses. Crucially, our method can incorporate different types of
390 features, such as gene expression and open chromatin peaks or intergenic methylation, and thus
391 takes advantage of the full power of multi-omics data. Ignoring the features unique to one dataset
392 (as in most existing methods), may omit important information, for instance distinguishing features
393 of certain subpopulations of cells and yield an integrated embedding that does not distinctly
394 cluster all subpopulations.

395
396 An additional advantage of MultiMAP is that the influence of each dataset on the shared
397 embedding can be modulated. This is useful when integrating datasets of different qualities, or
398 when aligning a query dataset to a reference dataset. Comparison with existing methods for
399 integration shows that MultiMAP outperforms or has close to best performance in every aspect
400 investigated. MultiMAP is a robust and effective method for dimensionality reduction and
401 integration of multimodal data, and is extremely fast and scalable to massive datasets.
402

403 Using synthetic examples to illustrate the power of the method, we show that MultiMAP leverages
404 the features unique to each dataset to effectively integrate and reduce the dimensionality of the
405 data, and is also robust to data with noise. Throughout our applications of MultiMAP to diverse
406 single-cell multi-omic data, we demonstrate that our method can facilitate integration across
407 transcriptomic, epigenomic, and spatially resolved datasets, and derive biological insights jointly
408 from multi-omic single-cell data. In addition, our method can align datasets across different
409 technologies and modalities even with extensive biological and technical differences. Crucially,
410 we show that MultiMAP is flexible enough to integrate datasets with different clusters and cell
411 populations, illustrating that MultiMAP is applicable even when its central hypothesis is not strictly

412 reflected by the data. The multimodal integration of three or more omics technologies opens many
413 opportunities for the comprehensive study of tissues.

414
415 We note that our method is based on the hypothesis that multi-omics data are uniformly distributed
416 on a latent manifold. A hypothesis of this sort, about the distribution of data in a latent space, is a
417 central feature of many existing integration strategies. For example, CCA-based strategies
418 (including Seurat and Conos) assume that the data reside in a maximally-correlated manner in a
419 latent space which is a linear projection of the original data. MultiMAP, in contrast, does not make
420 as strong an assumption because we do not restrict the latent manifold to a linear projection of
421 the data. While this kind of hypothesis is often realistic for data generated from the same tissue,
422 there may be cases where this is not strictly the case. In practice, we find that MultiMAP can
423 successfully accommodate datasets that depart from this central hypothesis, *i.e* when clusters
424 and cell populations are not shared across all datasets that are being integrated.

425
426 Perhaps the greatest potential lies in applying MultiMAP to datasets beyond those considered
427 here. Integrative analysis with MultiMAP can be used to compare healthy and diseased states,
428 and identify pathologic features, or to uncover cell-type specific responses to perturbations. Other
429 examples include the integration of data across species to enable studying the evolution of cell
430 states and identifying conserved cell types and regulatory programs. Along similar lines, the
431 integration of *in vivo* with *in vitro* models such as organoids will reveal the quality or faithfulness
432 of cells in a dish relative to their native counterparts. Finally, given the rapid development of joint
433 multimodal single cell genomics methods (e.g. CITEseq for protein and RNA, joint snRNA- and
434 ATACseq), it is relevant to point out that MultiMAP can be applied to multi-omic data acquired
435 both from different cells as well as from the same cells.

436
437 In summary, given the broad appeal of dimensionality reduction methods (e.g. PCA, tSNE,
438 UMAP), and the growth of multimodal data in many areas of science and engineering, we
439 anticipate that MultiMAP will find wide and diverse use.

440 441 AUTHOR CONTRIBUTIONS

442
443 M.S.J., M.E. and S.A.T. conceived the study. M.S.J conceived and developed MultiMAP. M.S.J,
444 created the codebase with contributions from M.E., and K.P. C.D.C., J.-E.P., R.A.B, E.S, L.M.,
445 A.E. and X.C. generated the single cell data. M.E. analysed the data and interpreted the results
446 with contributions from M.S.J. and S.A.T. M.S.J., M.E. and S.A.T. wrote the manuscript with
447 contributions from A.L., X.C., and C.D.C. All authors read and accepted the manuscript.

448 DATA AVAILABILITY

449
450 scRNA-seq and scATAC-seq data generated for this publication were being deposited in
451 ArrayExpress: E-MTAB-9769 for scRNA-seq of mouse splenocytes, E-MTAB-9840 and E-MTAB-
452 9828 for scRNA-seq and scATAC-seq of the thymus (username: me5@sanger.ac.uk; password:
453 rmachcde). The Multiome RNA+ATAC PBMC data is in the process of being deposited.

454 CODE AVAILABILITY

455

456 MultiMAP is publicly available at github.com/Teichlab/MultiMAP.

457

458 ACKNOWLEDGEMENTS

459

460 We thank Jana Eliasova for the graphical illustrations. We are grateful to Emma Dann, Natsuhiko
461 Kumasaka and Zhihan Xu for critical feedback on the manuscript. We thank Ruben Chazarra-Gil
462 and Vladimir Yu Kiselev for the comparison of different batch correction methods on the pancreas
463 dataset. M.S.J. was supported by a Gates Cambridge Scholarship. J.-E.P. was supported by
464 EMBO Long-Term and Advanced Fellowships. M.E. is funded by Barts Charity. S.A.T. is funded
465 by Wellcome (WT206194). The study was supported by Wellcome Human Cell Atlas Strategic
466 Science Support (WT211276/Z/18/Z) and the Chan Zuckerberg Initiative (CZF2019-002445).

467

468 COMPETING INTEREST STATEMENT

469

470 S.A.T. has received remunerations for consulting and Scientific Advisory Board work from
471 Genentech, Biogen, Roche and GlaxoSmithKline as well as Foresite Labs over the past three
472 years.

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

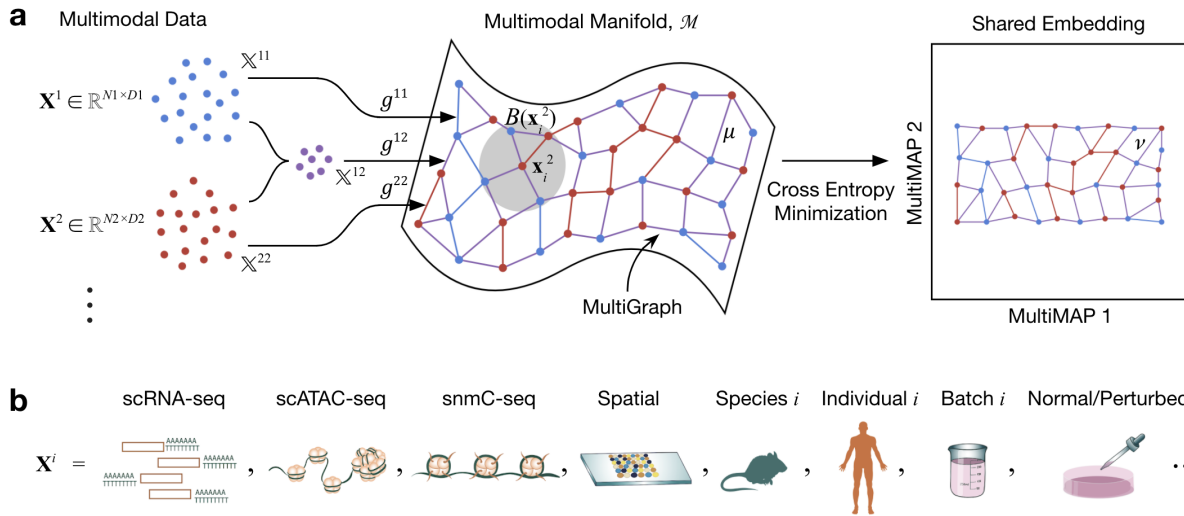
496

497

498 **Figures**

499

500



501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

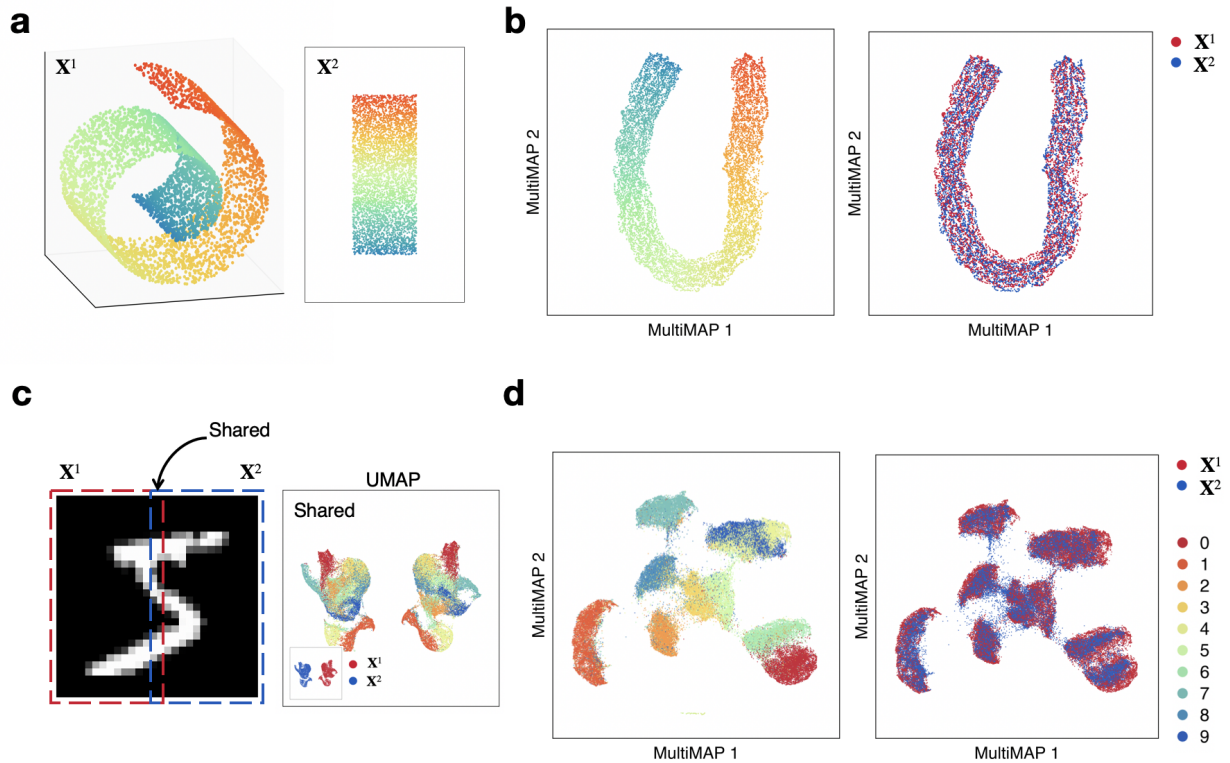
519

520

521

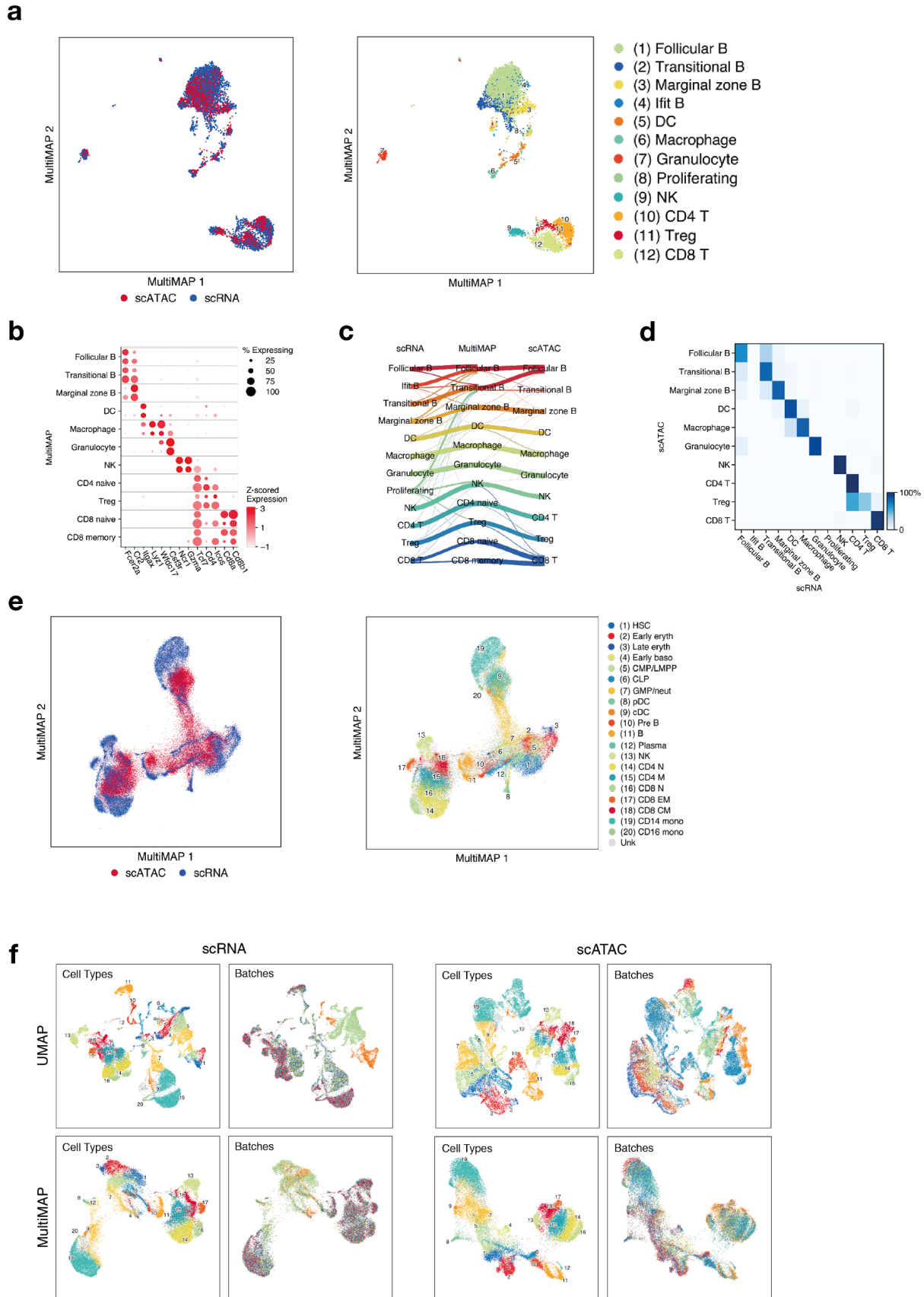
522

Figure 1. Schematic of MultiMAP. **a.** MultiMAP takes any number of datasets, including those of differing dimensions, recovers geodesic distances on a single latent manifold on which all data lie, constructs a neighborhood graph (MultiGraph) on the manifold, and then projects the data into a single low-dimensional embedding. Integrated analysis and visualisation can be performed on the embedding or graph. Variables are discussed in Methods. \mathbf{X}^i is dataset i , \mathbf{x}_i^j is a point in \mathbf{X}^i , M is the shared manifold, $B(\mathbf{x}_i^2)$ is a ball on M centered at \mathbf{x}_i^2 , X^{ij} is the ambient space of M in the coordinate space with data containing points from datasets i and j , g^{ij} is the metric of M in the space X^{ij} , μ is the membership function of the fuzzy simplicial set on the manifold, ν is the membership function of the fuzzy simplicial set in the low-dimensional space. **b.** In the field of cell atlas technologies, encompassing single cell genomics and spatial technologies, MultiMAP can be applied to integrate across different omics modalities, species, individuals, batches, and normal/perturbed states.



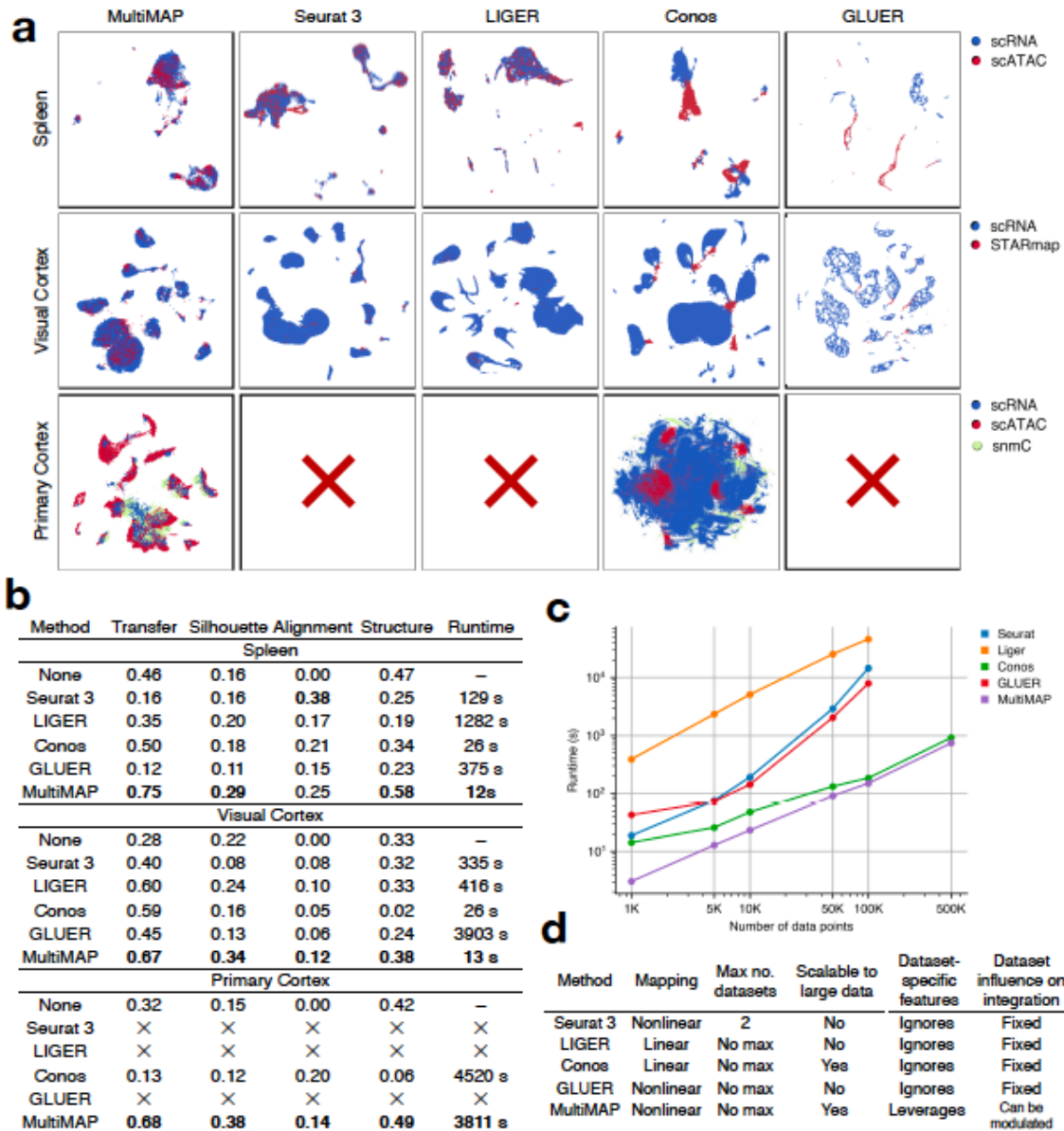
523
524

525 **Figure 2. MultiMAP applied to synthetic data.** **a.** Data sampled from the 3D Swiss Roll (X^1)
526 and a 2D rectangle (X^2). **b.** Shared embedding of both datasets produced by MultiMAP. Color
527 indicates position along the manifold (a,b). **c.** Left (X^1) and right (X^2) halves of MNIST
528 handwritten digit images with a 2 pixel wide shared region. Gaussian noise is added to the left
529 half. UMAP projections of each half and the shared region. **d.** Shared embedding of both MNIST
530 halves (including Gaussian noise introduced for the left half) produced by MultiMAP. Each color
531 is a different handwritten digit (0-9 as shown in the key). This illustrates that MultiMAP leverages
532 both shared and unshared features to integrate multimodal datasets.



534 **Figure 3. MultiMAP integration of single-cell transcriptomics and chromatin accessibility.**
535 **a.** MultiMAP visualization of the integration of published scATAC-seq¹⁵ and newly generated
536 scRNA-seq data of the mouse spleen (n=1), colored by omic technology (left hand panel) and
537 independent cell type annotations of each omic technology (right hand panel). **b.** Dot plot
538 showing the z-score of the mean log-normalised gene expression and gene activity scores of
539 known markers of each identified joint cluster. The top dot of each row shows the cells from the
540 scRNA-seq data, and the bottom dot represents the cells from the scATAC-seq data. **c.**
541 Riverplot showing correspondence between the joint clusters and the independent annotations
542 of the scATACseq and scRNAseq data. **d.** Confusion matrix of label transfer from the
543 scRNAseq to the scATACseq. **e.** MultiMAP visualization of the integration of single-cell
544 transcriptomics and chromatin accessibility of human bone marrow and peripheral blood
545 mononuclear cells¹⁶ colored by omic technology (left hand panel) and by the published cell type
546 annotation (right hand panel). **f.** UMAP (panels in top row) and MultiMAP (panels in bottom row)
547 visualization of the scRNA-seq and scATAC-seq data colored by cluster annotation and batch,
548 showing the effective batch correction of both modalities using MultiMAP.
549

552 **a.** MultiMAP visualization of scRNA-seq¹⁷ (n=2) and spatial STARmap¹⁸ (n=2) data of the
553 mouse brain, colored by omic technology and joint clusters identified with the MultiGraph. **b.** Dot
554 plot showing mean log-normalised gene expression of known markers of each identified joint
555 cluster. The top dot in each row represents cells from the scRNA-seq data, and the bottom dot
556 represents cells from the scATAC-seq data. **c.** Riverplot showing correspondence between the
557 joint clusters, and the independent annotations of the scATACseq and scRNAseq data. **d.**
558 Spatial locations of the STARmap cell, colored by the joint clusters. **e.** UMAP and MultiMAP
559 visualizations of the STARmap dataset. The silhouette score as employed here quantifies the
560 separation of clusters, and the higher value for MultiMAP shows the better cluster separation as
561 compared to UMAP. **f.** MultiMAP visualization of the integration of single-cell transcriptomics,
562 chromatin accessibility, and DNA methylation of the mouse primary cortex, colored by omic
563 technology and the published cell type annotation²⁰.
564
565



566

567 **Figure 5. Benchmarking MultiMAP against existing approaches.** a. Embeddings returned by

568 multi-omic integration methods on different datasets. “X” indicates that the method terminated

569 due to an out-of-memory error (218 GB RAM). b. Comparison of each method in terms of

570 transfer learning accuracy (“Transfer”), separation of cell type clusters as quantified by

571 Silhouette coefficient (“Silhouette”), mixing of different datasets as measured by fraction of

572 nearest neighbours that belong to a different dataset (“Alignment”), preservation of high-

573 dimensional structure as measured by the Pearson correlation between distances in the high-

574 and low-dimensional spaces (“Structure”), and runtime. c. Wall-clock time of multi-omic

575 integration methods on different sized datasets. Seurat 3 and LIGER produced out-of-memory

576 errors when run on 500,000 data points (218 GB RAM). To produce these datasets we

577 subsampled the mouse primary cortex scRNA-seq and scATAC-seq data²⁰ using geometric

578 sketching³³. The datasets were subsampled so that there are equal number of cells in the

579 scRNA-seq and scATAC-seq data until 100,000 cells. Since the scATAC-seq data had 81,196
580 cells in total, for the 500,000 cells comparison, we used an scRNA-seq of 418,804 cells. **d.**
581 Comparison of capabilities and properties of each method. “Mapping” refers to the nature of the
582 mapping employed by the method; “Max no. datasets” refers to the upper limit in terms of
583 numbers of datasets accepted by the method; “Scalable to large data” refers to allowing a total
584 of over 500,000 cells; “Data-set specific features” is whether the integration method allows
585 information that is not shared across datasets; and “Dataset influence on integration” is whether
586 the user can modulate the weighting of a given dataset relative to the others during the
587 integration.

588

589

590

591

592

593

594

595

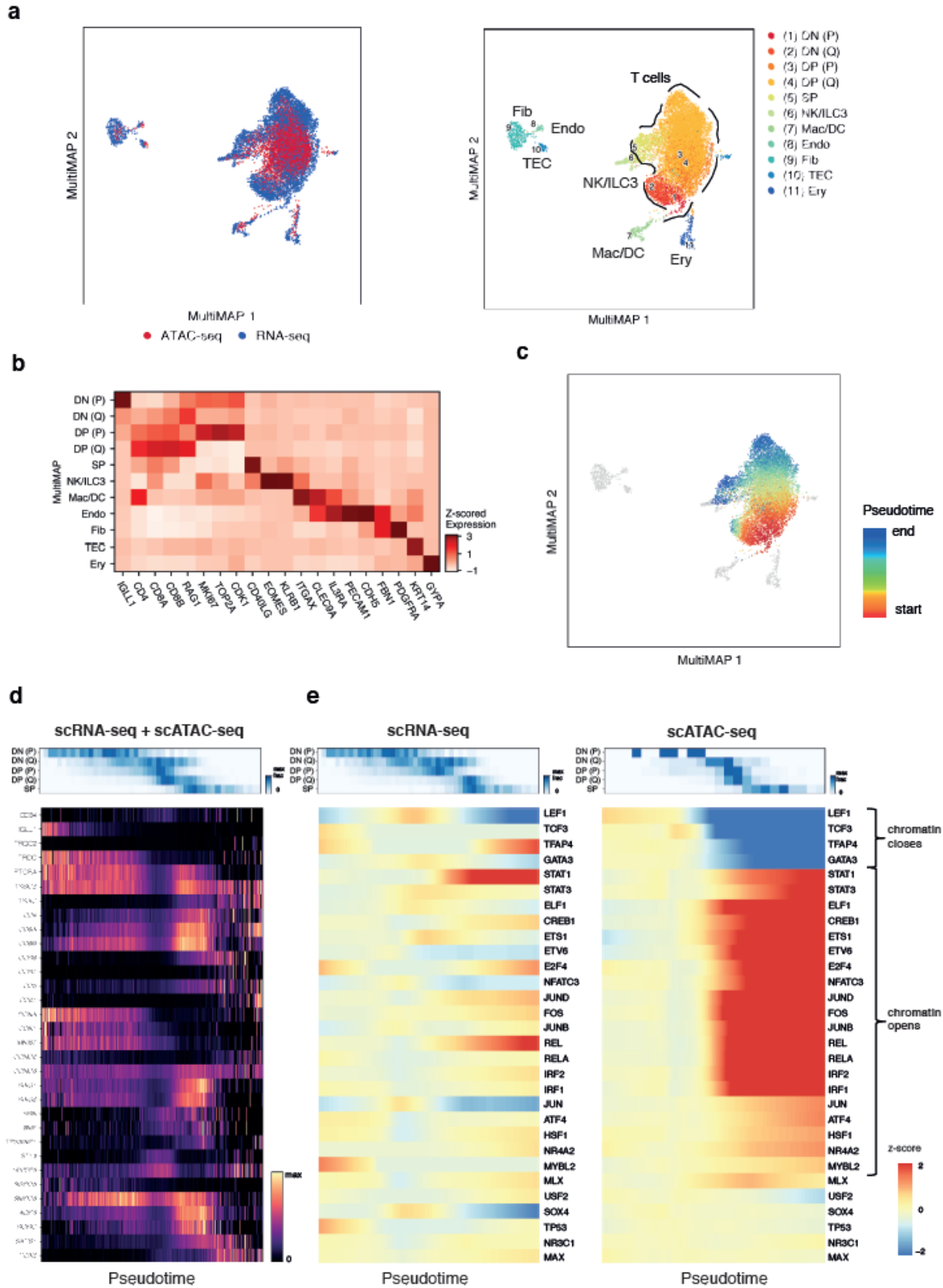
596

597

598

599

600



602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643

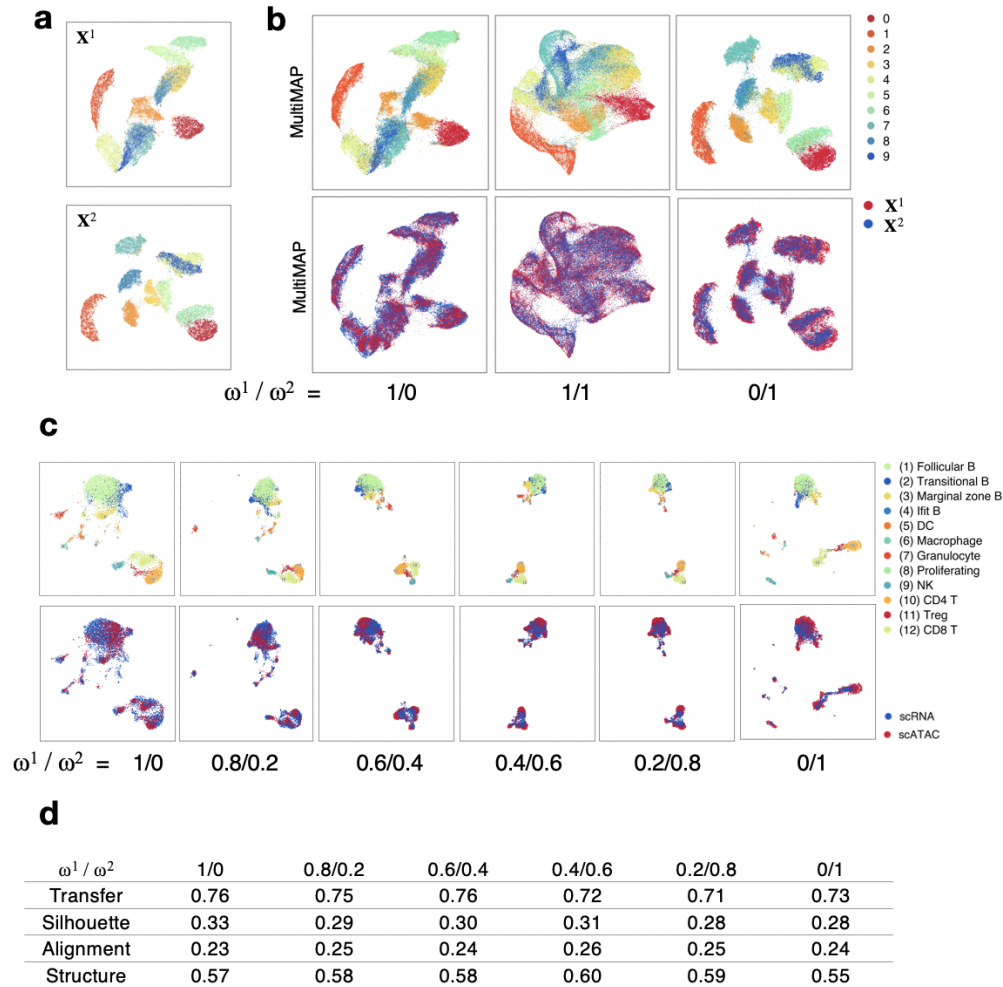
Figure 6. Integration of scRNAseq and scATACseq data of human fetal thymus reveals transcriptional regulatory principles of T cell development

a. MultiMAP visualization of scRNA-seq and scATAC-seq datasets of the human fetal thymus (n=1), colored by modality and joint clusters identified using the MultiGraph. **b.** Heatmap of gene expression and gene activity scores of key markers of the joint clusters identified using the MultiGraph. **c.** Inferred pseudotime using the MultiGraph recovers the T cell differentiation trajectory. Color indicates pseudotime from red (early, beginning) to blue (late, end). **d.** Heatmap of the gene expression and gene activity scores over pseudotime of genes known to be involved in T cell development. **e.** Smoothed heatmaps of the z-score of the gene expression and motif accessibility of the most variable transcription factors over pseudotime. The motif accessibilities of TFs that varied most in time show changes in accessibility at the transition between the late DN and early DP stage of differentiation. This includes TFs such as GATA3, ZEB1 for which the chromatin at the binding sites closes at that transition, and TFs for which the chromatin at the binding sites opens, such as E2F4, ETS1 and others.

644

645 **Supplementary Figures**

646



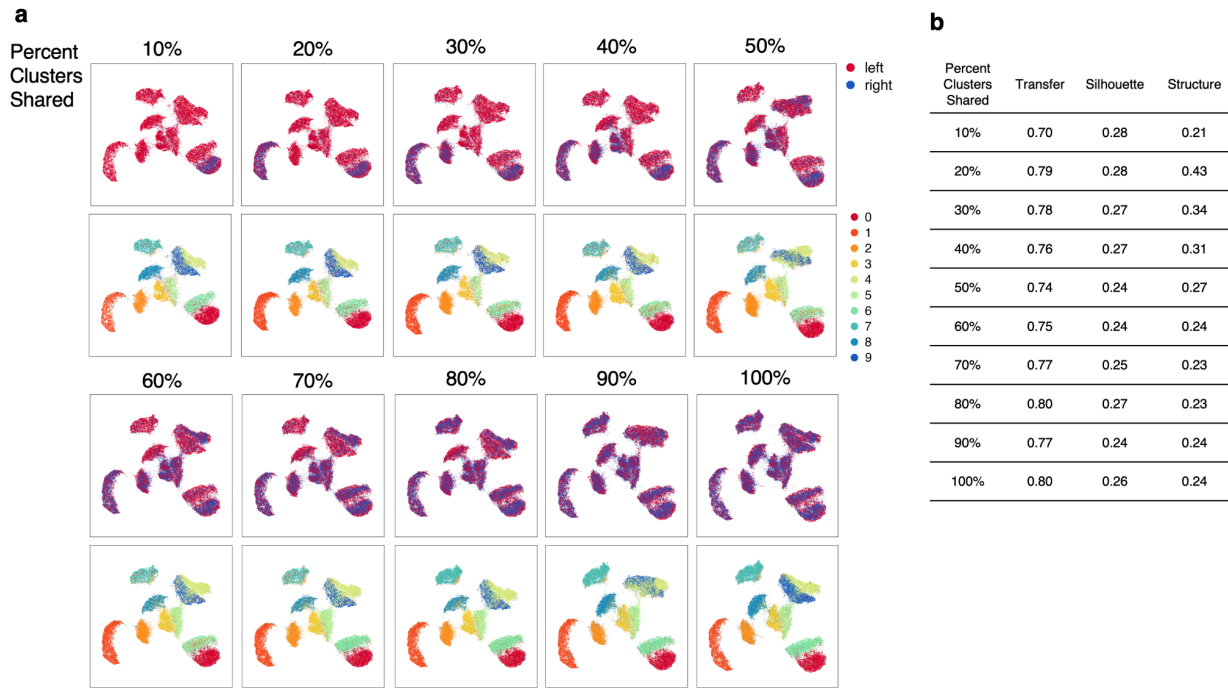
647

648

649 **Extended Data Figure 1. MultiMAP's weight parameter.** **a.** UMAP projections of the two halves
 650 of the MNIST handwritten digit images. **b.** MultiMAP embeddings as the weight parameters are
 651 varied. Each color is a different handwritten digit (0-9). When ω^1 is larger than ω^2 , the
 652 embedding more closely resembles the projection of only X^1 ; when ω^2 is larger than ω^1 , the
 653 embedding more closely resembles the projection of only X^2 . For different choices of ω^v , the
 654 datasets are well integrated in the embedding space. **c.** MultiMAP integration with varied weight
 655 parameters of published scATAC-seq¹⁵ and newly generated scRNA-seq data of the mouse
 656 spleen ($n=1$). **d.** Comparison of the MultiMAP integration of the spleen data as the weight
 657 parameter is varied -- in terms of transfer learning accuracy ("Transfer"), separation of cell type
 658 clusters as quantified by Silhouette coefficient ("Silhouette"), and preservation of high-

659 dimensional structure as measured by the Pearson correlation between distances in the high-
 660 and low-dimensional spaces (“Structure”)

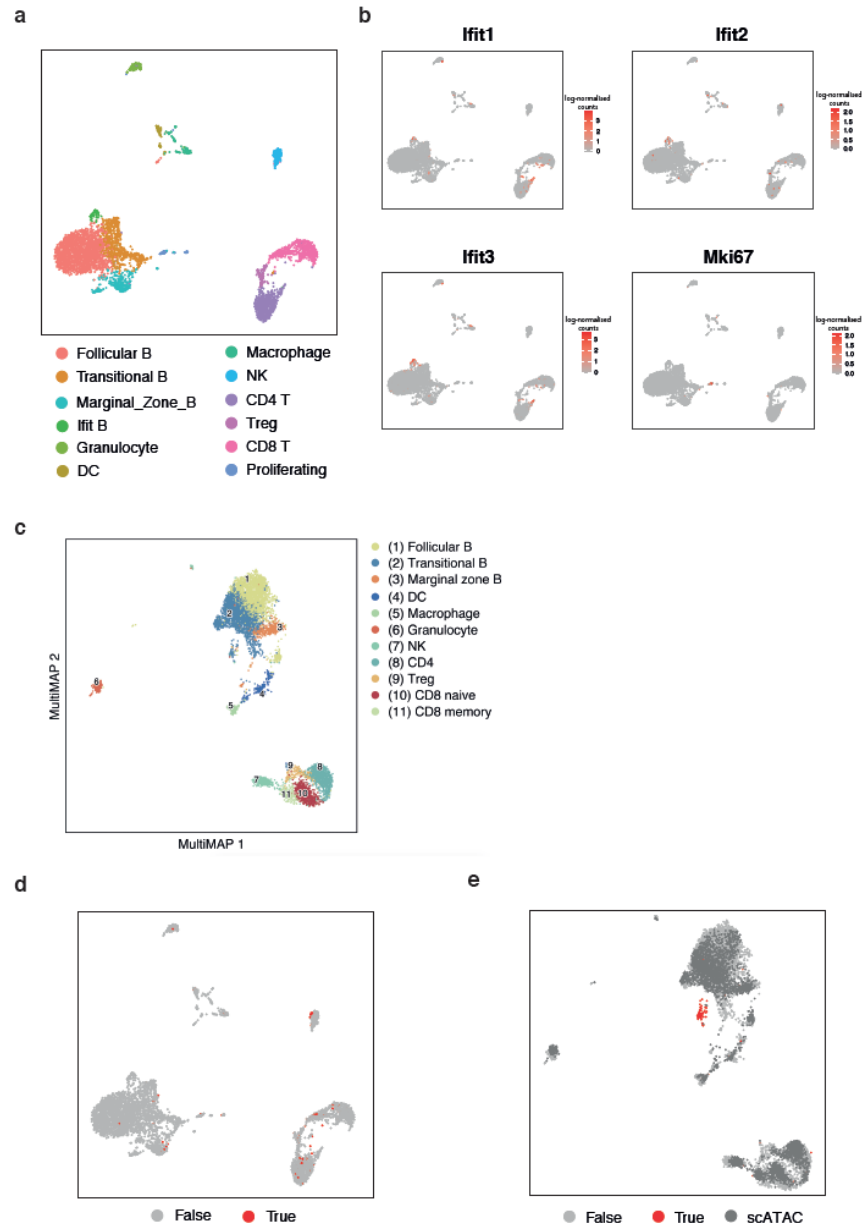
661
 662
 663



664

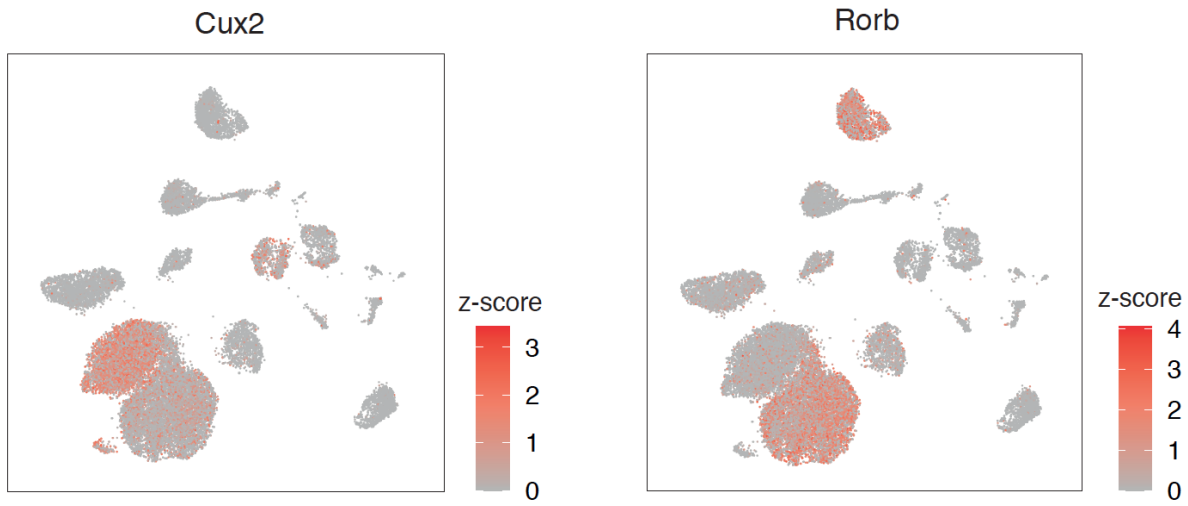
665 Extended Data Figure 2. **MultiMAP integration with non-shared clusters.** **a.** MultiMAP
 666 integration of the left and right halves of MNIST handwritten digit images with a 2 pixel wide shared
 667 region. Gaussian noise is added to the left half. MultiMAP integration is performed with a varying
 668 number of digit clusters removed from the right dataset, so that the integration ranges from one
 669 shared cluster (10%) to all clusters shared (100%). **b.** Comparison of the MultiMAP integration of
 670 the modified MNIST dataset as the percent of clusters shared is varied -- in terms of transfer
 671 learning accuracy (“Transfer”), separation of cell type clusters as quantified by Silhouette
 672 coefficient (“Silhouette”), and preservation of high-dimensional structure as measured by the
 673 Pearson correlation between distances in the high- and low-dimensional spaces (“Structure”).

674
 675



676
 677 Extended Data Figure 3. **Mouse spleen scRNA-seq and scATAC-seq data.** **a.** UMAP
 678 visualization of the mouse spleen scRNA-seq data (n=1) colored by the identified cell types. **b.**
 679 UMAP visualisation of expression levels of Ifit family genes associated with interferon response,
 680 upregulated in one specific B cell subpopulation, and the proliferation marker Mki67. **c.**
 681 MultiMAP visualization of the integrated scRNA-seq and scATAC-seq mouse spleen data (n=1)
 682 colored by the jointly identified clusters. **d, e.** UMAP (d) and MultiMAP (e) visualizations of the
 683 mouse spleen data showing cells identified as doublets (labelled “True”) using an independent
 684 pipeline (Scrublet). The MultiMAP visualisation leads to these artifactual data points being
 685 clustered in one group, highlighting the power of this method to visualise and separate data.
 686
 687

688



689

690

691

692

693

694

695

696

697

698

699

700

701

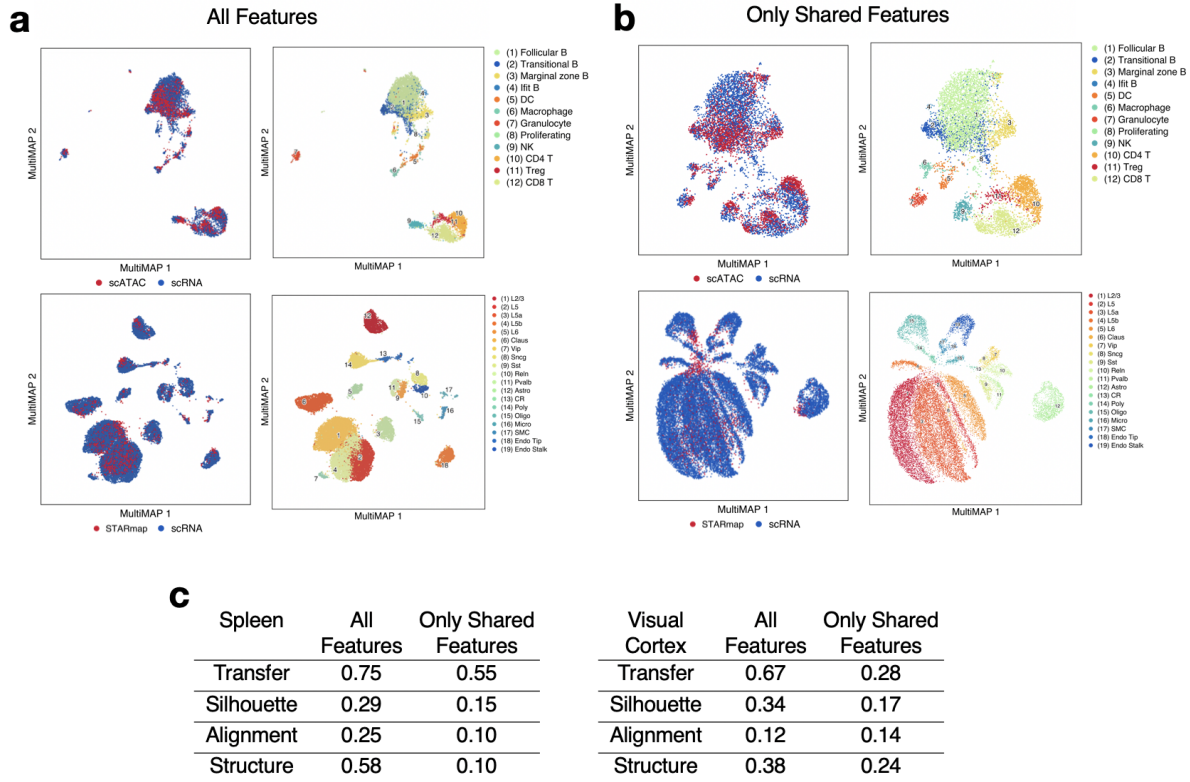
702

703

704

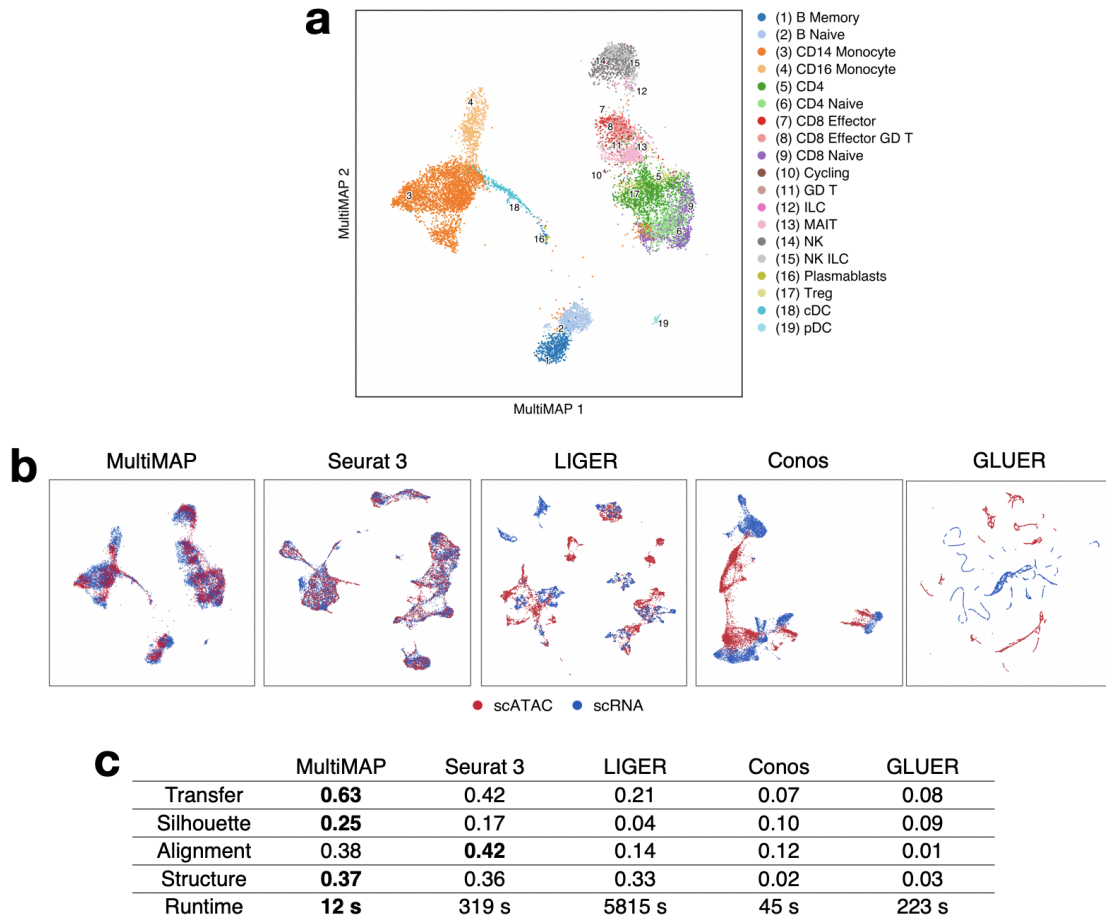
705

Extended Data Figure 4. **Marker genes of the L4 cluster identified in the scRNA-seq and STARmap integration.** MultiMAP visualisation of log-transformed gene expression of markers associated with L4 neurons. The MultiMAP integration identified L4 cells in the scRNA-seq data previously annotated as L5 neurons.



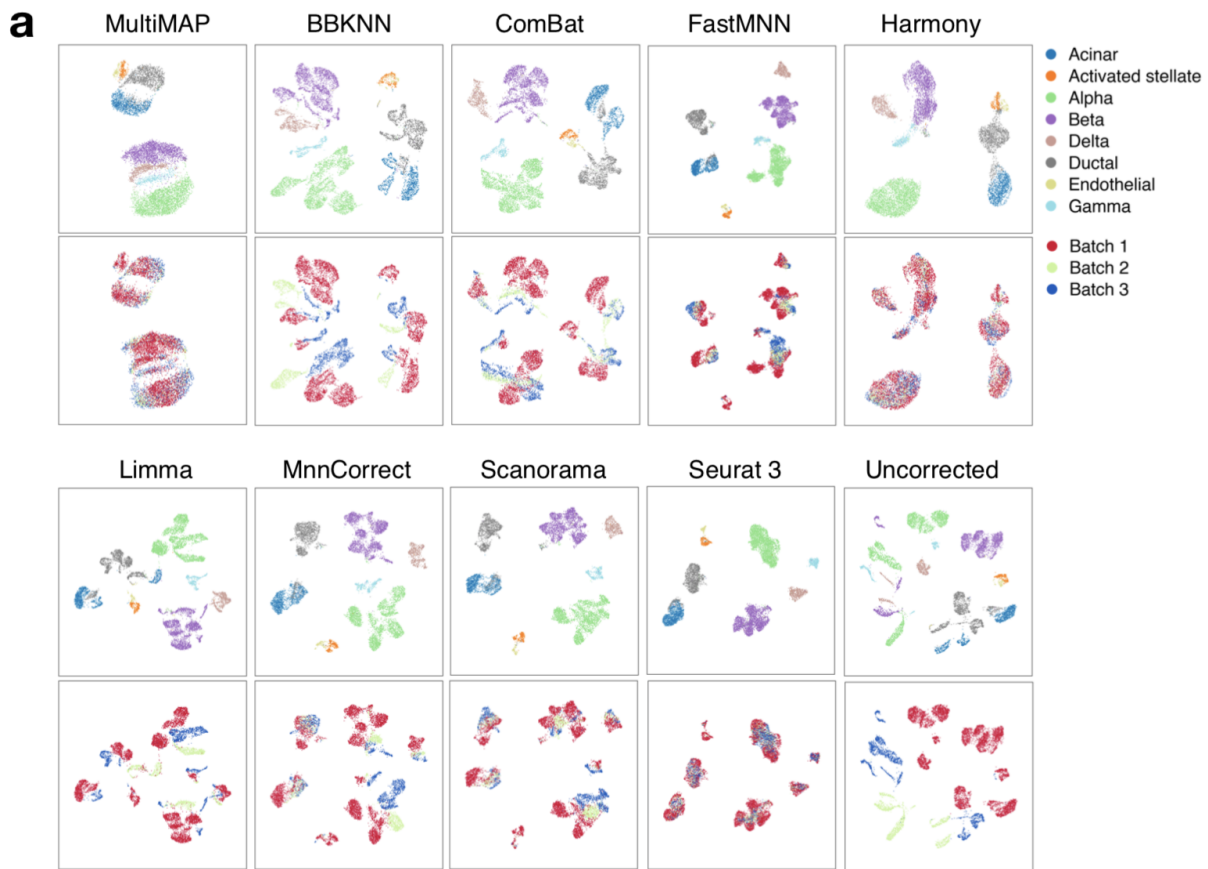
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720

Extended Data Figure 5. **MultiMAP integration with all features vs. only shared features in the spleen scRNA-seq + scATACseq, and visual cortex STARmap + scRNAseq datasets.** **a.** MultiMAP embeddings using all genes present in each dataset (intended use of MultiMAP). **b.** MultiMAP embeddings using only genes shared by all datasets in each integration. **c.** Comparison of the MultiMAP integration with all features vs. only shared features -- in terms of transfer learning accuracy ("Transfer"), separation of cell type clusters as quantified by Silhouette coefficient ("Silhouette"), mixing of different datasets as measured by fraction of nearest neighbours that belong to a different dataset ("Alignment"), and preservation of high-dimensional structure as measured by the Pearson correlation between distances in the high- and low-dimensional spaces ("Structure").



721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740

Extended Data Figure 6. **Benchmarking MultiMAP using paired PBMCs.** **a.** MultiMAP visualization of the Multiome RNA+ATAC PBMCs, colored by independently annotated cell type. **b.** Embeddings produced by alternative integration strategies, colored by omic technology. **c.** Comparison of each method in terms of transfer learning accuracy (“Transfer”), separation of cell type clusters as quantified by Silhouette coefficient (“Silhouette”), mixing of different datasets as measured by fraction of nearest neighbours that belong to a different dataset (“Alignment”), preservation of high-dimensional structure as measured by the Pearson correlation between distances in the high- and low-dimensional spaces (“Structure”), and runtime.



b

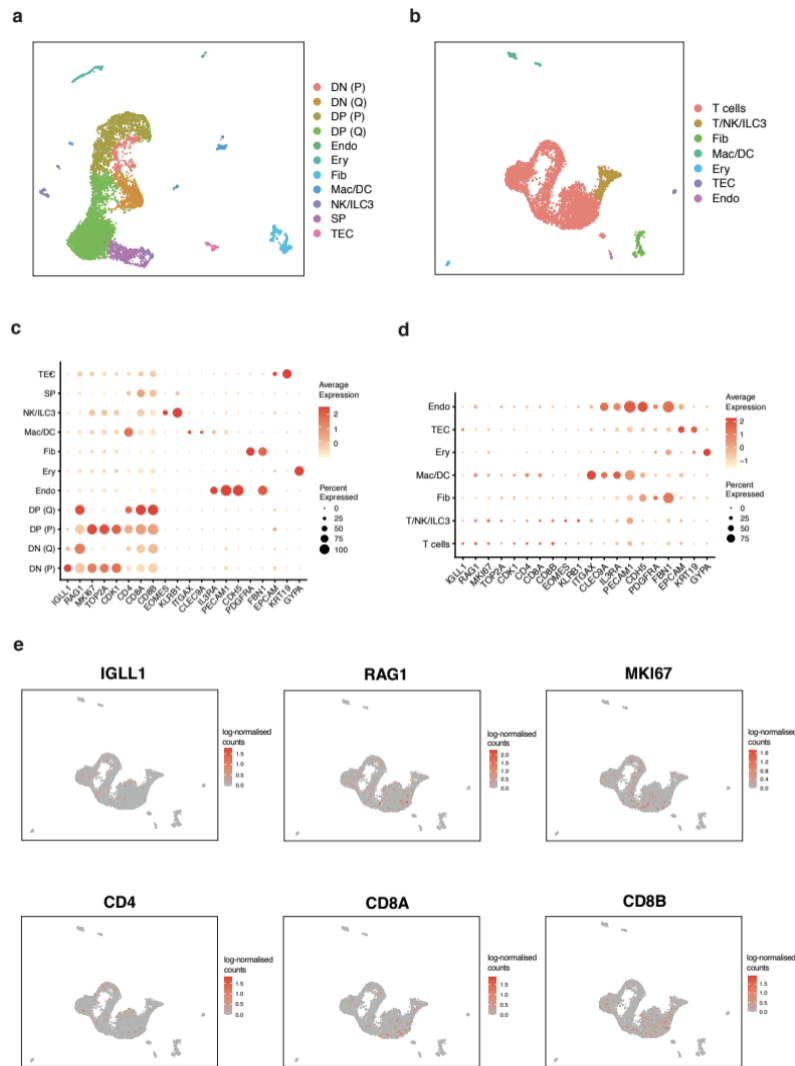
	MultiMAP	BBKNN	ComBat	FastMNN	Harmony
Silhouette	0.26	-0.27	-0.24	-0.24	-0.16
Alignment	0.29	0.01	0.09	0.17	0.28
Structure	0.57	0.59	0.52	0.53	0.51

	Limma	MnnCorrect	Scanorama	Seurat 3	Uncorrected
Silhouette	-0.24	-0.21	-0.23	-0.23	-0.31
Alignment	0.02	0.08	0.12	0.30	0.00
Structure	0.46	0.49	0.50	0.51	0.53

741
 742 Extended Data Figure 7. **Benchmarking MultiMAP against batch correction methods.** **a.**
 743 Embeddings returned by MultiMAP and batch correction methods on three scRNA-seq
 744 pancreas datasets. **b.** Comparison of separation of cell type clusters as quantified by Silhouette
 745 coefficient (“Silhouette”), mixing of different datasets as measured by fraction of nearest
 746 neighbours that belong to a different dataset (“Alignment”), preservation of high-dimensional
 747 structure as measured by the Pearson correlation between distances in the high- and low-
 748 dimensional spaces (“Structure”), and runtime.

749
 750
 751
 752

753



754

755

756

757 Extended Data Figure 8. **Fetal thymus scRNA-seq and scATAC-seq data.** **a.** UMAP

758 visualisation of the fetal thymus scRNA-seq data (n=1) colored by identified cell types shows the

759 same cell types as previously published³¹. **b.** UMAP visualisation of the fetal thymus scATAC-

760 seq data (n=1) colored by the identified cell types. **c.** Dot plot showing the z-score of the mean

761 log-transformed expression level of marker genes. **d.** Dot plot showing the z-score of the mean

762 log-transformed gene activity scores of marker genes, showing not very clear separation of T

763 cells clusters in the scATAC-seq data. **e.** UMAP visualisation of log-transformed gene activity

764 scores of markers for specific T cell subpopulations, showing that the scATAC-seq dataset does

765 not separate well the T cell clusters.

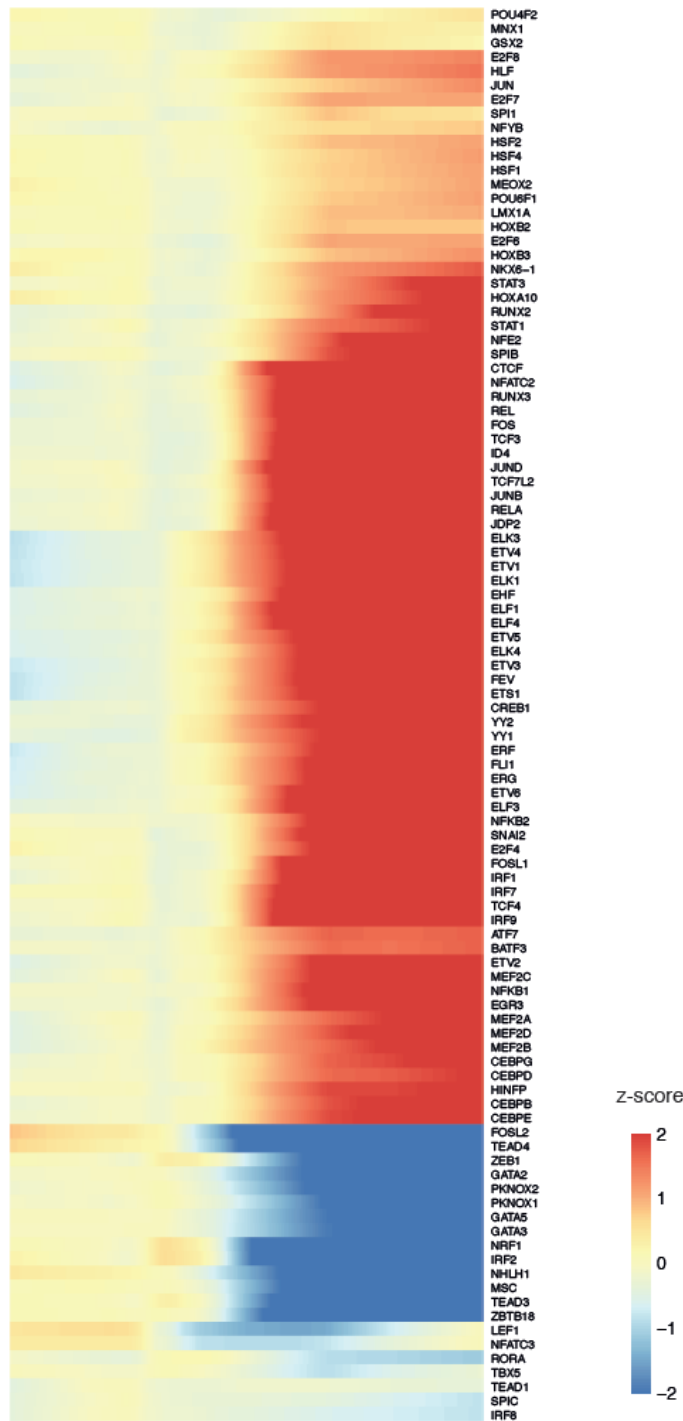
766

767

768

769

770



771

772 Extended Data Figure 9. **Chromatin accessibility of transcription factor binding sites.**

773 Smoothed heatmaps of the z-score of motif accessibility of the top 100 most variable transcription

774 factor binding sites over pseudotime. The TF binding sites that varied most in time show changes

775 in accessibility at the transition between the late DN and early DP stage of differentiation.

776

777 Methods

778

779 MultiMAP

780

781 MultiMAP (Figure 1) is a new approach for the integration and dimensionality reduction of
782 multimodal data based on a framework of Riemannian geometry and algebraic topology.
783 MultiMAP takes as input any number of datasets of potentially differing dimensions. The datasets
784 take the form \mathbf{X}^i , $i = 1, 2, \dots$, with $\mathbf{x}_j^i \in \mathbb{R}^{D_i}$ being the j 'th point in dataset \mathbf{X}^i . MultiMAP recovers
785 geodesic distances on a single latent manifold M on which all of the data is uniformly distributed.
786 The geodesic distances are calculated between data points of the same dataset by normalizing
787 distances in each dataset's ambient space X^{ii} with respect to a neighborhood distance specific to
788 the dataset, and between data points of different datasets by normalizing distances between the
789 data in a shared ambient space X^{ij} with respect to a neighborhood distance specific to the shared
790 feature space. When integrating multi-omics data with MultiMAP, the ambient spaces are the PC
791 components of each dataset's full feature space and of the shared feature space(s). These
792 neighborhood distances are the radius of a constant-radius ball B on M . These distances are then
793 used to construct a neighborhood graph (MultiGraph) on the manifold. Finally the data and
794 manifold space are projected into a low-dimensional embedding space by minimizing the cross
795 entropy of the graph in the embedding space with respect to the graph in the manifold space.
796 Specifically, this optimization minimizes cross entropy of a fuzzy set–representation $(\nu, \{\mathbf{x}_j^i\})$ of
797 the graph in the embedding space with respect to a fuzzy set–representation $(\mu, \{\mathbf{x}_j^i\})$ of the graph
798 in the manifold space. MultiMAP allows the user to modify the weight ω^i of each dataset in the
799 cross entropy loss, allowing the user to modulate the contribution of each dataset to the layout.
800 Integrated analysis can be performed on the embedding or the graph, and the embedding also
801 provides an integrated visualization. An extended description of MultiMAP, including
802 mathematical background, is in the Supplementary information.

803

804 Synthetic Data

805

806 MultiMAP was applied to two synthetic examples of multimodal data, in order to study the
807 technique in a controlled setting.

808

809 The first synthetic setting is schematized in Figure 2a. This setting consists of one dataset (\mathbf{X}^1) of
810 10,000 points sampled randomly from the canonical 3D “Swiss roll” surface (generated with
811 sklearn in Python), and a second dataset (\mathbf{X}^2) of 10,000 points sampled randomly from a 2D
812 rectangle. The two datasets can be considered multimodal data because they have different
813 feature spaces but describe a similar rectangular manifold. In addition, we are given the position
814 along the manifold of 1% of the data. Distances between data in the different datasets are
815 calculated for 1% of the data as the absolute differences between these positions. These
816 distances are supplied to MultiMAP. The purpose of this setting is to determine if MultiMAP can
817 integrate data in a nonlinear fashion and operate on datasets of different dimensionality.

818

819 The second synthetic setting is schematized in Figure 2c. This setting consists of two datasets
820 based on the MNIST database³⁴ which comprises 70,000 28x28 pixel grayscale images of
821 handwritten digits 0-9. The first dataset (\mathbf{X}^1) consists of the 28x15 pixel left half of each of images
822 flattened into a 420 dimensional vector. The second dataset (\mathbf{X}^2) consists of the 28x15 pixel right
823 half of each of 70,000 digit images, also flattened into a 420 dimensional vector. Added to the first
824 dataset is Gaussian noise with a mean of zero and a standard deviation equal to the maximum
825 pixel value. The two halves overlap by a 28x2 pixel region. Distances between data in the different
826 datasets are calculated in this shared space and supplied to MultiMAP. The two datasets can be
827 considered multimodal because they have different feature spaces but describe a similar
828 population of digit images. The purpose of this setting is to determine if MultiMAP can effectively
829 leverage features unique to certain datasets. The thin overlapping region of the two halves is not
830 enough information to create a good embedding of the data. Many distinct digits are similar in this
831 thin central sliver, and hence they should cluster together in the feature space of the two pixel
832 overlap. Indeed, in a UMAP projection of the data in the shared feature space of this overlap, the
833 clusters of different digits are not as well separated as in the UMAP projections of each half
834 (Figure 2c). A multimodal integration strategy that effectively leverages all features would use the
835 features unique to each half to separate different digits, and the shared space to bring the same
836 digits from each dataset close together.

837

838 **Acquisition and processing of human fetal thymic tissue**

839

840 The tissue sample used for this study was obtained with written informed consent from the
841 participant in accordance with the guidelines in The Declaration of Helsinki 2000. The human fetal
842 tissue was obtained from the MRC/Wellcome Trust-funded Human Developmental Biology
843 Resource (HDBR, <http://www.hdbbr.org>) with appropriate maternal written consent and approval
844 from the Newcastle and North Tyneside NHS Health Authority Joint Ethics Committee
845 (08/H0906/21+5). HDBR is regulated by the UK Human Tissue Authority (HTA; www.hta.gov.uk)
846 and operates in accordance with the relevant HTA Codes of Practice.

847

848 The developmental age was estimated from measurements of foot length and heel-to-knee
849 length, and compared against a standard growth chart³⁵. A piece of skin was collected from every
850 sample for Quantitative Fluorescence-Polymerase Chain Reaction analysis using markers for the
851 sex chromosomes and the following autosomes: 13, 15, 16, 18, 21, 22. The sample was of normal
852 karyotype.

853

854 The tissue was processed immediately after isolation using enzymatic digestion. Tissue was
855 transferred to a sterile 10mm² tissue culture dish and cut into <1mm³ segments before being
856 transferred to a 50mL conical tube. Tissues were digested with 1.6mg/mL collagenase type IV
857 (Worthington) in RPMI (Sigma-Aldrich) supplemented with 10%(v/v) heat-inactivated fetal bovine
858 serum (FBS; Gibco), 100U/mL penicillin (Sigma-Aldrich), 0.1mg/mL streptomycin (Sigma-
859 Aldrich), and 2mM L-glutamine (Sigma-Aldrich) for 30 minutes at 37°C with intermittent shaking.
860 Digested tissue was passed through a 100µm filter, and cells collected by centrifugation (500g for
861 5 minutes at 4°C). Cells were treated with 1X red blood cells (RBC lysis buffer (eBioscience) for
862 5 minutes at room temperature and washed once with a flow buffer (PBS containing 5%(v/v) FBS

863 and 2mM EDTA) prior to cell counting. For scATAC-seq, cells were taken forward for nuclei
864 isolation following 10X Genomics guidelines. Briefly, cells were centrifuged (300g for 5 minutes),
865 added the lysis buffer (Tris-HCl (pH 7.4) 10mM; NaCl 10Mm; MgCl₂ 3mM; Tween-20 0.1%; NP-
866 40 0.1%; Digitonin 0.01%; BSA 1%) and incubated on ice for 3 minutes (time optimized for
867 thymus). Following the incubation, cells were washed (Tris-HCl (pH 7.4) 10mM; NaCl 10Mm;
868 MgCl₂ 3mM; BSA 1%; Tween-20 0.1%) and centrifuged (300g for 5 minutes) and nuclei were
869 resuspended in Diluted Nuclei Buffer (10X Genomics). Isolated nuclei were high-quality with well-
870 resolved edges and no evidence of blebbing. The final nuclei concentration was determined prior
871 to loading using a hemocytometer.

872

873 **Single-cell RNA and ATAC sequencing of human thymus**

874

875 scRNA-seq targeting 5,000 cells per sample was performed using the Chromium Controller (10x
876 Genomics). Single-cell cDNA synthesis, amplification, and sequencing libraries were generated
877 using the Single Cell 5' Reagent Kit following the manufacturer's instructions. The libraries from
878 up to eight loaded channels were multiplexed together and sequenced on an Illumina HiSeq 4000.

879

880 scATAC-seq targeting 5,000 cells was performed using Chromium Single Cell ATAC Library and
881 Gel Bead kit (10x Genomics). The libraries from up to eight loaded channels were multiplexed
882 together and sequenced on an Illumina HiSeq 4000.

883

884 **Computational processing and analysis of the human fetal thymus single cell genomics** 885 **data**

886

887 scRNA-seq data were aligned and quantified using the Cell Ranger Single-Cell Software Suite
888 (version 2.0, 10x Genomics) against the GRCh38 human reference genome provided by Cell
889 Ranger. The scRNA-seq data was preprocessed using Seurat 3. Cells with fewer than 500
890 detected genes and more than 10% mitochondrial gene expression content were removed.
891 Ribosomal genes, cell cycle genes³¹ and genes associated with dissociation-induced effects³⁶
892 were removed. Clusters were identified using a community identification algorithm as
893 implemented in the Seurat 'FindClusters' function, using 30 principal components (PCs) and
894 annotated using canonical cell-type markers from³¹.

895

896 The scATAC-seq data was aligned and preprocessed using CellRanger (10x Genomics).
897 SnapATAC³⁷ was used for quality control, preprocessing, and generating cell-by-bin and log-
898 normalized gene activity matrices. The binarized cell-by-bin matrix was used as input for term
899 frequency-inverse document frequency (TF-IDF) weighting, using term frequency and smoothed
900 inverse document frequency as the weighting scheme. Weighted data were reduced to 30
901 dimensions using singular-value decomposition (SVD). Clustering and UMAP visualization were
902 performed using Seurat 3. chromVar³⁸ was used to discover transcription factor dynamics and
903 variation in their motif accessibility.

904

905 The 50 dimension reduced scATAC-seq and the 50 dimension reduced scRNA-seq data were
906 supplied as input to MultiMAP. A shared feature space with both the scATAC-seq and scRNA-

907 seq was constructed by removing genes from each dataset that were not present in the other,
908 and then reducing the space to 50 dimensions using PCA. This shared space was supplied as
909 input to MultiMAP, allowing the calculation of distances between cells from different datasets. The
910 parameters of MultiMAP were all set to their default values, including the weight parameter for
911 the scRNA-seq set to 0.8 and for ATAC-seq set to 0.2, on account of the higher-quality scRNA-
912 seq.

913
914 The Leiden algorithm³⁹ was applied directly to the MultiGraph to jointly cluster all cells. The
915 clusters were then annotated using canonical cell-type markers from ³¹. Diffusion pseudotime
916 (DPT)⁴⁰ was used for trajectory inference. The MultiGraph was supplied as input to the DPT
917 function in SCANPY⁴¹. DPT was performed only on cells annotated as T cells. Cells were removed
918 if they were positioned away from T cell clusters and close to Fibroblasts and Erythrocytes on the
919 MultiMAP plot, as this likely indicated that they were incorrectly annotated. tradeSeq⁴² was used
920 to identify genes whose expression changes significantly along the trajectory.

921

922 **Acquisition and processing of human PBMCs**

923

924 PBMCs from two donors were acquired from a LeukoLab (Clinical division of AllCells). Frozen
925 PBMC samples were thawed quickly at 37 °C in a water bath. Two pools made for technical
926 duplicates with ~500,000 cells for each donor per pool (50/50). Nuclei isolation, transposition,
927 ATAC-seq and Gene Expression (GEX) sequencing libraries construction performed according
928 to the manufacturer's Demonstrated protocol (CG000365 Rev A; 10X Genomics) and Next GEM
929 Single Cell Multiome ATAC and Gene Expression user guide (CG000338 Rev A; 10X Genomics).
930 One lane per pool with a 3,000 targeted nuclei recovery was loaded on a Chromium Next GEM
931 Chip J. ATAC-seq and GEX indexed libraries were sequenced on a NovaSeq 6000 SP Flowcell
932 according to the 10X Genomics recommendations, aiming for a minimum of 50,000 PE reads per
933 cell for both types (ATAC-Seq and GEX) libraries.

934

935 **Computational processing and analysis of the human PBMCs Multiome ATAC+RNA data**

936

937 snRNA-seq and snATAC-seq data were aligned and quantified using the Cell Ranger ARC suite
938 (10x Genomics) against the GRCh38 human reference genome provided by Cell Ranger. The
939 snRNA-seq data was preprocessed using Seurat 3. Cells with fewer than 500 detected genes
940 and more than 20% mitochondrial gene expression content were removed. Clusters were
941 identified using a community identification algorithm as implemented in the Seurat 'FindClusters'
942 function, using 30 principal components (PCs) and annotated using canonical cell-type markers.

943

944 SnapATAC³⁷ was used for quality control, preprocessing, and generating cell-by-bin and log-
945 normalized gene activity matrices for the snATAC-seq data. The binarized cell-by-bin matrix was
946 used as input for term frequency-inverse document frequency (TF-IDF) weighting, using term
947 frequency and smoothed inverse document frequency as the weighting scheme. Weighted data
948 were reduced to 30 dimensions using singular-value decomposition (SVD). Clustering and UMAP
949 visualization were performed using Seurat 3. We used the

950

951 The 50 dimension reduced snATAC-seq and the 50 dimension reduced snRNA-seq data were
952 supplied as input to MultiMAP. A shared feature space with both the snATAC-seq and snRNA-
953 seq was constructed by removing genes from each dataset that were not present in the other,
954 and then reducing the space to 50 dimensions using PCA. This shared space was supplied as
955 input to MultiMAP, allowing the calculation of distances between cells from different datasets. The
956 parameters of MultiMAP were all set to their default values, including the weight parameter for
957 the snRNA-seq set to 0.8 and for snATAC-seq set to 0.2, on account of the higher-quality snRNA-
958 seq.

959
960

961

962 **Single-cell RNA sequencing of mouse spleen and data processing**

963

964 The mice were maintained under specific pathogen-free conditions at the Wellcome Trust
965 Genome Campus Research Support Facility (Cambridge, UK). These animal facilities are
966 approved by and registered with the UK Home Office. All procedures were in accordance with the
967 Animals (Scientific Procedures) Act 1986. The protocols were approved by the Animal Welfare
968 and Ethical Review Body of the Wellcome Trust Genome Campus.

969 The spleen from a 6-month-old C57BL/6Jax mouse was removed. The splenocytes were isolated
970 by passing the spleen through a 70 µm cell strainer (Fisher Scientific 10788201) into 30 ml ice-
971 cold 1X DPBS (Thermo Fisher 14190169) with 2 mM EDTA and 0.5% (w/v) BSA (Sigma A9418)
972 using the plunger of a 2-ml syringe. Cells were spun down at 500 g for 7 minutes at 4 degree.
973 Then the supernatant was removed, and the cell pellet resuspended in 5 ml 1X RBC lysis buffer
974 (Thermo Fisher 00-4300-54). The cell suspension was vigorously vortexed for 5 seconds and left
975 on the bench for 5 minutes to lyse the red blood cells. Then 45 ml ice-cold 1X DPBS was added,
976 and cells were spun down at 500 g for 7 minutes at 4 degrees. The supernatant was removed,
977 and 30 ml ice-cold 1X DPBS with 0.1% BSA was used to resuspend the cell pellet. The cell
978 suspension was passed through a Miltenyi 30 µm Pre-Separation Filter (Miltenyi 130-041-407),
979 and the cell number was determined using the C-chip counting chamber (VWR DHC-N01). The
980 cells were spun down again, and the cell pellet resuspended in ice-cold 1X DPBS with 0.1% BSA
981 to reach a concentration of 1,000,000 cells per ml. The splenocytes were then loaded on the 10x
982 Chromium Controller, aiming to recover ~ 5000 cells (Targeted Cell Recovery 5000 cells). cDNA
983 and a sequencing library were made according to 10x Single Cell 3' Reagent Kits v2 manual. The
984 library was sequenced on an Illumina HiSeq 4000 machine.

985 The resulting scRNA-seq data were preprocessed using CellRanger (10x Genomics) and
986 downstream analysis were performed using the Seurat 3 workflow. Cells with fewer than 200
987 detected genes and more than 10% mitochondrial gene expression content were filtered out.
988 Downstream analyses such as normalization, clustering and visualization were performed using
989 Seurat 3. Clusters were identified using the community identification algorithm as implemented in
990 the Seurat 'FindClusters' function, using 20 PCs. Clusters were annotated using canonical cell-
991 type markers from the original study¹⁵. Scrublet⁴³ was used for doublet detection.

992

993 **Acquisition and processing of previously published datasets**

994

995 The mouse spleen scATAC-seq data was obtained from ArrayExpress (E-MTAB-6714) and
996 preprocessed using the code provided by Chen et al. ¹⁵([https://github.com/dbrg77/plate_scATAC-](https://github.com/dbrg77/plate_scATAC-seq)
997 [seq](https://github.com/dbrg77/plate_scATAC-seq)). Briefly, reads from all cells were merged, and open chromatin regions were identified by
998 peak calling with MACS2 ⁴⁴. Latent semantic indexing analysis was used for dimensionality
999 reduction of the resulting cell-by-bin matrix. The binary cell-by-bin accessibility was used as input
1000 for TF-IDF weighting, using term frequency and smoothed inverse document frequency as the
1001 weighting scheme. Weighted data were reduced to 50 dimensions using SVD. SnapATAC³⁷ was
1002 used to generate gene activity count matrices, which were then log-normalized. The 50 dimension
1003 reduced accessibility of the scATAC-seq and the 50 dimension reduced gene expression of the
1004 scRNA-seq data were supplied as input to MultiMAP. The 50-dimension reduced accessibility of
1005 the scATAC-seq and the 50-dimension reduced gene expression of the scRNA-seq data were
1006 supplied as input to MultiMAP. A shared feature space with both the scATAC-seq and scRNA-
1007 seq was constructed by removing genes from each dataset that were not present in the other,
1008 and then reducing the space to 50 dimensions using PCA. This shared space was supplied as
1009 input to MultiMAP, allowing the calculation of distances between cells from different datasets. The
1010 parameters of MultiMAP were all set to their default values, including the weight parameter for
1011 the scRNA-seq set to 0.8 and for ATAC-seq set to 0.2 due to the higher quality scRNA-seq.
1012 The Leiden algorithm was applied directly to the MultiGraph to jointly cluster all cells. Harmonic
1013 function-based node classification was performed directly on the MultiGraph to predict cell types
1014 of the scATAC-seq cells given the cell types of the scRNA-seq cells⁴⁵.

1015

1016 Human hematopoiesis scRNA-seq and scATAC-seq data were downloaded from
1017 <https://github.com/GreenleafLab/MPAL-Single-Cell-2019>. The scRNA-seq consists of 6
1018 experimental batches, and the scATAC-seq consists of 10 experimental batches. Severe batch
1019 effects were observed, so this data was considered to consist of 16 separate datasets for the
1020 integration with MultiMAP. The scRNA-seq data was preprocessed using Seurat 3, and each
1021 batch was log-normalized and reduced to 50 dimensions with PCA. The cell-by-bin peak
1022 accessibility was used as provided by the authors. The binary cell-by-bin accessibility was used
1023 as input for TF-IDF weighting, using term frequency and smoothed inverse document frequency
1024 as the weighting scheme. Separately for each batch, the weighted data were reduced to 50
1025 dimensions using SVD. Gene activities of the ATAC data were calculated using Cicero⁴⁶ and log-
1026 normalized. To integrate all of the data at once, all 16 datasets were provided as input to MultiMAP
1027 in the form of the 50 dimension reduced accessibility data of the scATAC-seq and the 50
1028 dimension reduced gene expression of the scRNA-seq. Shared feature spaces containing two
1029 datasets were constructed by removing genes from each of the datasets that were not present in
1030 the other, and then reducing the space to 50 dimensions using PCA. These shared spaces were
1031 supplied as input to MultiMAP to calculate distances between cells from different datasets. The
1032 parameters of MultiMAP were all set to their default values, including the weight parameter for
1033 the scRNA-seq set to 0.8 and for ATAC-seq set to 0.2 due to the higher-quality scRNA-seq data.

1034

1035 scRNA-seq data of the mouse frontal cortex acquired with Drop-seq was obtained from
1036 dropviz.org. STARmap data of the mouse visual cortex was downloaded from
1037 <https://www.starmapresources.com/data/>. Each dataset was separately preprocessed with
1038 Seurat 3¹¹, log-normalized, and reduced to 50 dimensions with PCA. Both 50 dimensional
1039 reduced datasets were supplied as input to MultiMAP. A shared feature space with both the
1040 STARmap and scRNA-seq data was constructed by removing genes from each dataset that were
1041 not present in the other, and then reducing the space to 50 dimensions using PCA. This shared
1042 space was supplied as input to MultiMAP to calculate distances between cells from different
1043 datasets. The parameters of MultiMAP were all set to their default values, including the weight
1044 parameter for the scRNA-seq set to 0.8 and for Drop-seq set to 0.2, on account of higher-quality,
1045 tighter clusters generally observed in the scRNA-seq.

1046
1047 scRNA-seq, scATAC-seq, and snmC-seq data from the mouse primary cortex²⁰ was downloaded
1048 from the Neuroscience Multi-omics Archive (NeMO). The scRNA-seq was preprocessed using
1049 Seurat 3, log-normalised, and reduced to 50 dimensions with PCA. The binary cell-by-bin
1050 accessibility and gene activity count matrix of the scATAC-seq were obtained with SnapATAC³⁷.
1051 The gene activity count data was log-normalized. Latent semantic indexing analysis was used for
1052 dimensionality reduction of the scATAC-seq accessibility. The binary cell-by-bin accessibility was
1053 used as input for TF-IDF weighting, using term frequency and smoothed inverse document
1054 frequency as weighting scheme. Weighted data were reduced to 50 dimensions using SVD. The
1055 DNA methylation data was preprocessed as described in ⁴⁷, using the provided scripts. Briefly,
1056 after mapping, the methyl-cytosine counts and total cytosine counts were calculated in two sets
1057 of genome regions for each cell: the non-overlapping 100 kb bins tiling the mm10 genome, which
1058 was used for dimensionality reduction, and gene body regions ± 2 kb, which is used for the joint
1059 alignment. Posterior mCH and mCG rates were calculated based on beta-binomial distribution for
1060 the non-overlapping 100kb bins matrix. The top 3000 highly variable features were taken and the
1061 data was reduced to 50 dimensions with PCA. Because gene body mCH proportions are
1062 negatively correlated with gene expression level, the direction of the methylation data was
1063 reversed by subtracting all values from the maximum methylation value¹². The 50 dimensional
1064 reduced scRNA-seq, scATAC-seq, and snmC-seq were supplied as input to MultiMAP. Shared
1065 feature spaces containing each pair of two datasets and all three datasets together were
1066 constructed by removing genes from each of the datasets that were not present in the other, and
1067 then reducing the space to 50 dimensions using PCA. These shared spaces were supplied as
1068 input to MultiMAP, allowing the calculation of distances between cells from different datasets. The
1069 parameters of MultiMAP were all set to their default values. The weight parameter for the scRNA-
1070 seq set to 0.8 and for the other omics set to 0.2, on account of the higher-quality scRNA-seq data.

1071 **Benchmarking**

1072
1073
1074 Benchmarking of MultiMAP, Seurat 3, LIGER, Conos and GLUER was performed using a variety
1075 of multi-omic data including the scRNA-seq and scATAC-seq data of the spleen, scRNA-seq and
1076 STARmap of the visual cortex, and the scRNA-seq, scATAC-seq, and snmC-seq of the primary
1077 cortex. These datasets were chosen because they all have cell type annotations supplied in their
1078 original publications, which was used to independently validate the integration.

1079

1080 The scRNA-seq and STARmap data was log-normalised using Seurat 3 and then used as an
1081 input for all integration methods, except GLUER where the raw data was used as an input and
1082 preprocessed using the SCANPY workflow. The scATAC-seq data was preprocessed as
1083 described above and the log-normalised gene activity matrix was used as an input for all
1084 integration methods. Seurat 3, LIGER, Conos and GLUER were executed as detailed in their
1085 tutorials, with all parameters set to their default values. Latent Semantic Indexing was used as
1086 the dimensionality reduction technique for the scATAC-seq data for weighting anchors in Seurat
1087 3. CCA was used as the dimensionality reduction technique for the scRNA-seq and STARmap
1088 data for weighting anchors in Seurat 3.

1089

1090 A diversity of performance metrics was used. After integration, label transfer of the cell type
1091 annotations from the scRNA-seq to each other omic was performed by setting the cell type of a
1092 query cell to the most frequent type among its 5 nearest labeled neighbors. The balanced
1093 accuracy of the label transfer (“Transfer”) was calculated using the annotations from the original
1094 publications as the ground truth. A high accuracy indicates that the same cell types from different
1095 modalities are near each other in the integrated embedding. After integration, the average
1096 Silhouette score⁴⁸ (“Silhouette”) across all cells was calculated using the cell type annotations
1097 from the original publications as the cluster labels. We note that the Silhouette score is not
1098 affected by the number of clusters as we use the same cell type labels, and hence number of
1099 clusters, for each integration method. A higher Silhouette score indicates the embedding is better
1100 separating distinct cell types. The degree of alignment (“Alignment”) of the different datasets in
1101 the integrated embedding was calculated as the proportion of each cell's 5 nearest neighbors that
1102 originated in a different dataset, averaged over all cells. This metric was also used in ¹². A higher
1103 value of the alignment score indicates that the different datasets are more evenly mixed in the
1104 integrated embedding. The degree to which the embedding preserves the high-dimensional
1105 structure (“Structure”) of each dataset was calculated as the Pearson correlation between all
1106 pairwise distances in the high-dimensional spaces and the corresponding distances in the
1107 embedding. A higher correlation indicates that the embedding is more faithful to the high-
1108 dimensional structure. All of these performance metrics were also calculated in the shared feature
1109 space of the datasets to be integrated, to get baseline values of the metrics prior to the application
1110 of any integration strategy.

1111

1112 The wall-clock runtime of each method on each dataset was recorded. Additionally, to
1113 characterize the runtimes of the methods on a wide range of dataset sizes, the integration
1114 methods were run on datasets ranging from 1,000 to 500,000 cells. To produce these datasets
1115 we subsampled the mouse primary cortex scRNA-seq and scATAC-seq data²⁰ using geometric
1116 sketching³³. The datasets were subsampled so that there are equal number of cells in the scRNA-
1117 seq and scATAC-seq data until 100,000 cells. Since the scATAC-seq data had 81,196 cells in
1118 total, for the 500,000 cells comparison, we used an scRNA-seq of 418,804 cells. All methods
1119 were run with 3.1 GHz Intel i7 cores and 218 GB RAM.

1120

1121

1122

1123 References

- 1124 1. Becht, E. *et al.* Dimensionality reduction for visualizing single-cell data using UMAP. *Nat.*
1125 *Biotechnol.* (2018) doi:10.1038/nbt.4314.
- 1126 2. Stoeckius, M. *et al.* Simultaneous epitope and transcriptome measurement in single cells.
1127 *Nature Methods* vol. 14 865–868 (2017).
- 1128 3. Peterson, V. M. *et al.* Multiplexed quantification of proteins and transcripts in single cells.
1129 *Nature Biotechnology* vol. 35 936–939 (2017).
- 1130 4. Klemm, S. L., Shipony, Z. & Greenleaf, W. J. Chromatin accessibility and the regulatory
1131 epigenome. *Nat. Rev. Genet.* **20**, 207–220 (2019).
- 1132 5. Karemaker, I. D. & Vermeulen, M. Single-Cell DNA Methylation Profiling: Technologies and
1133 Biological Applications. *Trends Biotechnol.* **36**, 952–965 (2018).
- 1134 6. Mayr, U., Serra, D. & Liberali, P. Exploring single cells in space and time during tissue
1135 development, homeostasis and regeneration. *Development* **146**, (2019).
- 1136 7. Regev, A. *et al.* The Human Cell Atlas. *Elife* **6**, (2017).
- 1137 8. HuBMAP Consortium. The human body at cellular resolution: the NIH Human Biomolecular
1138 Atlas Program. *Nature* **574**, 187–192 (2019).
- 1139 9. Efremova, M. & Teichmann, S. A. Computational methods for single-cell omics across
1140 modalities. *Nat. Methods* **17**, 14–17 (2020).
- 1141 10. Lähnemann, D. *et al.* Eleven grand challenges in single-cell data science. *Genome Biol.* **21**,
1142 31 (2020).
- 1143 11. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* vol. 177 1888–
1144 1902.e21 (2019).
- 1145 12. Welch, J. D. *et al.* Single-Cell Multi-omic Integration Compares and Contrasts Features of
1146 Brain Cell Identity. *Cell* vol. 177 1873–1887.e17 (2019).
- 1147 13. Lopez, R. *et al.* A joint model of unpaired data from scRNA-seq and spatial transcriptomics

- 1148 for imputing missing gene expression measurements. *arXiv [cs.LG]* (2019).
- 1149 14. GradientBased Learning Applied to Document Recognition. *Intelligent Signal Processing*
1150 (2009) doi:10.1109/9780470544976.ch9.
- 1151 15. Chen, X., Miragaia, R. J., Natarajan, K. N. & Teichmann, S. A. A rapid and robust method
1152 for single cell chromatin accessibility profiling. *Nat. Commun.* **9**, 5345 (2018).
- 1153 16. Granja, J. M. *et al.* Single-cell multiomic analysis identifies regulatory programs in mixed-
1154 phenotype acute leukemia. *Nat. Biotechnol.* **37**, 1458–1465 (2019).
- 1155 17. Saunders, A. *et al.* Molecular Diversity and Specializations among the Cells of the Adult
1156 Mouse Brain. *Cell* **174**, 1015–1030.e16 (2018).
- 1157 18. Wang, X. *et al.* Three-dimensional intact-tissue sequencing of single-cell transcriptional
1158 states. *Science* **361**, (2018).
- 1159 19. Brodmann, K. *Brodmann's: Localisation in the Cerebral Cortex.* (Springer Science &
1160 Business Media, 2007).
- 1161 20. Yao, Z. *et al.* An integrated transcriptomic and epigenomic atlas of mouse primary motor
1162 cortex cell types. 2020.02.29.970558 (2020) doi:10.1101/2020.02.29.970558.
- 1163 21. Yamawaki, N., Borges, K., Suter, B. A., Harris, K. D. & Shepherd, G. M. G. A genuine layer
1164 4 in motor cortex with prototypical synaptic circuit connectivity. *Elife* **3**, e05422 (2014).
- 1165 22. Barkas, N. *et al.* Joint analysis of heterogeneous single-cell RNA-seq dataset collections.
1166 *Nat. Methods* **16**, 695–698 (2019).
- 1167 23. Peng, T., Chen, G. M. & Tan, K. GLUER: integrative analysis of single-cell omics and
1168 imaging data by deep neural network. doi:10.1101/2021.01.25.427845.
- 1169 24. Muraro, M. J. *et al.* A Single-Cell Transcriptome Atlas of the Human Pancreas. *Cell Syst* **3**,
1170 385–394.e3 (2016).
- 1171 25. Segerstolpe, Å. *et al.* Single-Cell Transcriptome Profiling of Human Pancreatic Islets in
1172 Health and Type 2 Diabetes. *Cell Metab.* **24**, 593–607 (2016).
- 1173 26. Baron, M. *et al.* A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas

- 1174 Reveals Inter- and Intra-cell Population Structure. *Cell Syst* **3**, 346–360.e4 (2016).
- 1175 27. Chazarra-Gil, R., van Dongen, S., Kiselev, V. Y. & Hemberg, M. Flexible comparison of
1176 batch correction methods for single-cell RNA-seq using BatchBench. *Nucleic Acids Res.*
1177 (2021) doi:10.1093/nar/gkab004.
- 1178 28. Roels, J. *et al.* Distinct and temporary-restricted epigenetic mechanisms regulate human $\alpha\beta$
1179 and $\gamma\delta$ T cell development. *Nat. Immunol.* **21**, 1280–1292 (2020).
- 1180 29. Jia, G. *et al.* Single cell RNA-seq and ATAC-seq analysis of cardiac progenitor cell
1181 transition states and lineage settlement. *Nat. Commun.* **9**, 4877 (2018).
- 1182 30. Chen, H. *et al.* Single-cell trajectories reconstruction, exploration and mapping of omics
1183 data with STREAM. *Nat. Commun.* **10**, 1903 (2019).
- 1184 31. Park, J.-E. *et al.* A cell atlas of human thymic development defines T cell repertoire
1185 formation. *Science* **367**, (2020).
- 1186 32. Hosokawa, H. & Rothenberg, E. V. How transcription factors drive choice of the T cell fate.
1187 *Nature Reviews Immunology* (2020) doi:10.1038/s41577-020-00426-6.
- 1188 33. Hie, B., Cho, H., DeMeo, B., Bryson, B. & Berger, B. Geometric Sketching Compactly
1189 Summarizes the Single-Cell Transcriptomic Landscape. *Cell Syst* **8**, 483–493.e7 (2019).
- 1190 34. Lecun, Y., Bottou, L., Bengio, Y. & Haffner, P. Gradient-based learning applied to document
1191 recognition. *Proceedings of the IEEE* vol. 86 2278–2324 (1998).
- 1192 35. Hern, W. M. Correlation of fetal age and measurements between 10 and 26 weeks of
1193 gestation. *Obstet. Gynecol.* **63**, 26–32 (1984).
- 1194 36. van den Brink, S. C. *et al.* Single-cell sequencing reveals dissociation-induced gene
1195 expression in tissue subpopulations. *Nat. Methods* **14**, 935 (2017).
- 1196 37. Fang, R. *et al.* Fast and Accurate Clustering of Single Cell Epigenomes Reveals Cis-
1197 Regulatory Elements in Rare Cell Types. doi:10.1101/615179.
- 1198 38. Schep, A. N., Wu, B., Buenrostro, J. D. & Greenleaf, W. J. chromVAR: inferring
1199 transcription-factor-associated accessibility from single-cell epigenomic data. *Nat. Methods*

- 1200 **14**, 975–978 (2017).
- 1201 39. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities
1202 in large networks. *Journal of Statistical Mechanics: Theory and Experiment* vol. 2008
1203 P10008 (2008).
- 1204 40. Haghverdi, L., Büttner, M., Wolf, F. A., Buettner, F. & Theis, F. J. Diffusion pseudotime
1205 robustly reconstructs lineage branching. *Nat. Methods* **13**, 845 (2016).
- 1206 41. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression
1207 data analysis. *Genome Biol.* **19**, 15 (2018).
- 1208 42. Van den Berge, K. *et al.* Trajectory-based differential expression analysis for single-cell
1209 sequencing data. *Nat. Commun.* **11**, 1201 (2020).
- 1210 43. Wolock, S. L., Lopez, R. & Klein, A. M. Scrublet: Computational Identification of Cell
1211 Doublets in Single-Cell Transcriptomic Data. *Cell Syst* **8**, 281–291.e9 (2019).
- 1212 44. Grytten, I. *et al.* Graph Peak Caller: Calling ChIP-seq peaks on graph-based reference
1213 genomes. *PLoS Comput. Biol.* **15**, e1006731 (2019).
- 1214 45. Zhu, X., Ghahramani, Z. & Lafferty, J. D. Semi-supervised learning using gaussian fields
1215 and harmonic functions. in *Proceedings of the 20th International conference on Machine*
1216 *learning (ICML-03)* 912–919 (2003).
- 1217 46. Pliner, H. A. *et al.* Cicero Predicts cis-Regulatory DNA Interactions from Single-Cell
1218 Chromatin Accessibility Data. *Mol. Cell* **71**, 858–871.e8 (2018).
- 1219 47. Kozareva, V. *et al.* A transcriptomic atlas of the mouse cerebellum reveals regional
1220 specializations and novel cell types. doi:10.1101/2020.03.04.976407.
- 1221 48. Rousseeuw, P. J. Silhouettes: A graphical aid to the interpretation and validation of cluster
1222 analysis. *Journal of Computational and Applied Mathematics* vol. 20 53–65 (1987).

1223

1224