



Data integration and inference of gene regulation using single-cell temporal multimodal data with scTIE

Yingxin Lin, Tung-Yu Wu, Xi Chen, et al.

Genome Res. published online December 15, 2023

Access the most recent version at doi:[10.1101/gr.277960.123](https://doi.org/10.1101/gr.277960.123)

P<P	Published online December 15, 2023 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

Data integration and inference of gene regulation using single-cell temporal multimodal data with scTIE

Yingxin Lin^{1,2,3§}, Tung-Yu Wu^{4§}, Xi Chen^{4§}, Sheng Wan⁵, Brian Chao⁶, Jingxue
Xin⁴, Jean Y.H. Yang^{1,2,3}, Wing H. Wong^{4,7,8*}, 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.

³Laboratory of Data Discovery for Health Limited (D24H), Science Park, Hong Kong SAR,
China

⁴Department of Statistics, Stanford University, CA, USA.

⁵Institute of Electronics, National Yang Ming Chiao Tung University, Hsinchu, Taiwan.

⁶Department of Electrical Engineering, Stanford University, CA, USA.

⁷Department of Biomedical Data Science, Stanford University, CA, USA.

⁸Bio-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.

Keywords

Single cell multiome, Temporal data integration, Context-specific gene regulatory network

Abstract

Single-cell technologies offer unprecedented opportunities to dissect gene regulatory mechanisms in context-specific ways. Although there are computational methods for extracting gene regulatory relationships from scRNA-seq and scATAC-seq data, the data integration problem, essential for accurate cell type identification, has been mostly treated as a standalone challenge.

24 Here we present scTIE, a unified method that integrates temporal multimodal data and infers
25 regulatory relationships predictive of cellular state changes. scTIE uses an autoencoder to em-
26 bed cells from all time points into a common space using iterative optimal transport, followed
27 by extracting interpretable information to predict cell trajectories. Using a variety of synthetic
28 and real temporal multimodal datasets, we demonstrate scTIE achieves effective data integration
29 while preserving more biological signals than existing methods, particularly in the presence of
30 batch effects and noise. Furthermore, on the exemplar multiome dataset we generated from dif-
31 ferentiating mouse embryonic stem cells over time, we demonstrate scTIE captures regulatory
32 elements highly predictive of cell transition probabilities, providing new potentials to understand
33 the regulatory landscape driving developmental processes.

34 **Introduction**

35 In eukaryotic cells, gene expressions are intricately regulated through complex interactions of
36 transcription factors (TFs), various regulatory elements and target genes. Deciphering the func-
37 tions of gene regulatory networks (GRNs) in shaping cell identity and cell fate is one of the
38 central quests in understanding the mapping from genomic blueprints to phenotypes. Over the
39 past decades, much effort has been devoted to developing statistical and computational meth-
40 ods for inferring GRNs from tissue-level bulk data containing genome-wide profiling of gene
41 expression, TF binding, and 3D chromatin structure. More recently, the advent of single-cell se-
42 quencing technologies has propelled the study of GRNs into a new era, in which context-specific
43 regulation mechanisms can be investigated. Unlike global GRNs, which inherently aggregate
44 gene interactions over all biological conditions present in a given dataset, context-specific GRNs
45 are tailored to a particular biological setting. These specialized networks detail the regulatory
46 interactions that occur in unique circumstances, such as within specific cell types, lineages, tis-
47 sues, or under certain environmental conditions. Alongside new opportunities, the sparse and
48 noisy nature of these single-cell data also brings new challenges to the statistical and computa-
49 tional analyses.

50
51 A growing number of methods have been developed to extract GRNs from data generated by
52 assays of single-cell RNA-sequencing (scRNA-seq) and single-cell transposase-accessible chro-
53 matin sequencing (scATAC-seq). Most of these methods infer the relationships between TFs
54 and target genes by estimating their interactions with *cis*-regulatory elements (CREs) as an inter-

55 mediate, using information including TF motif enrichment, marginal or conditional correlations
56 between genes and CRE accessibility, and physical proximity between different elements (Duren
57 et al. 2022; Y Jiang et al. 2022; Kartha et al. 2022; Tran et al. 2022; L Zhang et al. 2022). These
58 methods typically work with multimodal data that provide joint profiling of scRNA-seq and
59 scATAC-seq from the same cells, or unpaired data from a matched population of cells, possibly
60 measured over a time course. However, they do not directly address the data integration problem
61 accompanying such data, in which noise, sparsity, and batch effects can obscure identification of
62 cell types and affect the downstream inference of context-specific GRNs. Furthermore, to com-
63 pare how GRNs dynamically evolve in developmental data, features (e.g., genes, CREs) that are
64 different between time points (or pseudotime points) are identified using differential expression
65 (DE) / accessibility (DA) analyses. While this captures marginal correlations, the features found
66 are not necessarily predictive of the developmental changes.

67

68 On a separate front, an increasing number of computational methods have been proposed to
69 perform data integration for single-cell multiomics data from unpaired measurements (Cao and
70 Gao 2022; Gong et al. 2021; Y Lin et al. 2022; Z Zhang et al. 2022). As more technologies capa-
71 ble of multimodal profiling start to emerge (S Chen et al. 2019; Ma et al. 2020; Plongthongkum
72 et al. 2021), integration methods designed for paired data (Argelaguet et al. 2020; Ashuach et
73 al. 2023; Y Hao et al. 2021; Jin et al. 2020) have also attracted significant research interests.
74 However, most of these integration methods do not directly address the immediate downstream
75 problem of inferring GRNs; one exception is GLUE (Cao and Gao 2022), although the GRNs
76 inferred there remain global and not context-specific. One difficulty lies in the fact that most of
77 these methods rely on finding a low-dimensional representation of the datasets across modalities
78 and data batches, and how to extract interpretable biological signals from blackbox methods such
79 as neural networks is a challenging problem. Neural networks offer a conceptual advantage over
80 methods built on linear models, including cross correlation analysis and non-negative matrix fac-
81 torization, as their superior representation power can capture complex nonlinear interactions in
82 the feature space. However, this comes with the drawback that the relationships between the
83 measured features (e.g., genes) and cellular phenotypes in trained models become more difficult
84 to interpret. Although alternative architectures have been proposed involving linearizing part of
85 the neural network (Svensson et al. 2020), a tradeoff remains between the network's representa-
86 tion power and interpretability.

87

88 Here, we propose scTIE, an autoencoder-based method for integrating multimodal profil-
89 ing of scRNA-seq and scATAC-seq data over a time course and inferring context-specific GRNs.
90 Unlike existing GRN inference methods that study cell type-specific or condition-specific GRNs,
91 scTIE focuses on cellular state transitions and aims to infer GRNs that predict cellular changes
92 along a developmental process. To the best of our knowledge, scTIE provides the first unified
93 framework for the integration of temporal data and the inference of context-specific GRNs that
94 predict cell fates. We achieve this through three main innovations in the architecture design of
95 the autoencoder and the interpretation of a blackbox neural network method. Firstly, scTIE uses
96 iterative optimal transport (OT) fitting to align cells in similar states between different time points
97 and estimate their transition probabilities. scTIE incorporates OT into the loss function of the
98 autoencoder so that the alignment of cells is updated iteratively throughout training to achieve
99 a desirable balance between time point alignment and cell type separation. This is in contrast
100 to many widely used applications of OT in trajectory inference of scRNA-seq data (Schiebinger
101 et al. 2019; Forrow and Schiebinger 2021), where most of the methods solve OT only once
102 on suitably constructed cell distance matrices. Secondly, scTIE removes the need for selecting
103 highly variable genes (HVGs) as input through a pair of coupled batchnorm layers to account
104 for large variations in gene expression levels, making it more robust and generalizable. Thirdly,
105 scTIE provides the means to extract interpretable features from the common embedding space by
106 linking the developmental trajectories of cell representations to their measured features (genes
107 and peaks). We formulate a trajectory prediction problem using the estimated transition prob-
108 abilities from OT and use gradient-based saliency mapping (P Yang et al. 2021; Ciortan and
109 Defrance 2021) to identify genes and peaks that are potentially driving the cellular state changes.
110 Compared to most GRN inference methods, which focus on developing new ways to construct
111 network relationships among features selected through DE/DA analysis, the main innovation of
112 scTIE lies in selecting these informative features based on their ability to predict cellular changes.

113 To demonstrate the performance of scTIE on developmental data, we have chosen to focus
114 on multimodal time-course data, as this emerging form of data provides better opportunities to
115 understand the key transcriptional regulatory activities driving a developmental process. To as-
116 sess scTIE's integration performance against other existing methods, we constructed a variety of
117 synthetic datasets using a mouse early organogenesis multiome dataset. Furthermore, we gen-
118 erated an exemplar dataset comprising paired scRNA-seq and scATAC-seq measurements from
119 $\sim 11,000$ differentiating mouse embryonic stem cells (mESCs) over a time course. Using these
120 datasets, our primary aim was to assess scTIE's ability to integrate multimodal developmental

121 data for better cell type identification and uncover key regulatory elements predictive of cell fate
122 in a unified framework.

123 **Results**

124 **Overview of scTIE**

125 scTIE consists of two main steps. In the first step, scTIE uses modality-specific encoders and
126 decoders to project high dimensional input data from all time points into a lower dimensional
127 common embedding space and reconstruct them in the original space (Fig. 1A). Each encoder-
128 decoder pair is designed to preserve the original information in the input data with minimal
129 information loss, with appropriate loss functions to guide the integration process. For scATAC-
130 seq, accessibility peaks are used as input without conversion to gene activity scores. The encoder
131 and decoder for scRNA-seq use an additional pair of coupled batchnorm layers to handle hetero-
132 geneity in gene expression levels and achieve high-fidelity reconstruction of the signals without
133 the need for selecting HVGs. Between consecutive time points, scTIE models cell trajectories
134 using the principle of OT based on the current embeddings and computes an OT loss using the
135 transport cost matrix. The OT loss is incorporated into the total loss function to update the em-
136 bedded features, aligning cells by their estimated transition probabilities in the trajectories; the
137 cost matrix itself is also updated iteratively throughout training. In addition to the OT loss, a
138 modality alignment loss is used to ensure the projected feature vectors from the two modalities
139 (RNA and ATAC) are close in distance for the same cell.

140 In the second step, scTIE finetunes the learned embeddings to build a supervised model for
141 predicting cellular transition probabilities for user-selected subgroups of cells (Fig. 1B). Genes
142 and peak regions highly predictive of the cellular transitions are selected by backpropagating the
143 gradients, allowing us to construct GRNs responsible for developmental changes.

144 We demonstrate the advantages of scTIE on a number of synthetic and real datasets. On
145 synthetic datasets constructed from a mouse early organogenesis multiome dataset, our results
146 reveal that scTIE can effectively align cells from different time points and mitigate batch effects,
147 achieving an optimal balance between time alignment, modality alignment, and cell type sepa-
148 ration. Furthermore, our analysis of an exemplar multiome dataset from differentiating mESCs
149 show its superior capacity to capture biological signals from each modality and achieve better
150 day alignment when compared to other methods, resulting in identification of distinct cell sub-

151 populations. Finally, using developmental transitions from anterior primitive streak as a case
152 study, we demonstrate scTIE’s ability to construct lineage-specific GRNs consisting of regula-
153 tory elements with a high predictive power of cell fate and identify key regulatory signals that
154 would be missed by DE or DA-based analysis.

155 **scTIE outperforms existing methods in integrating temporal multimodal** 156 **data.**

157 We first evaluated the data integration performance of scTIE against recent methods designed to
158 integrate paired multimodal data, including Seurat (Y Hao et al. 2021), scAI (Jin et al. 2020),
159 multiVI (Ashuach et al. 2023) and MOFA (Argelaguet et al. 2020). We generated four syn-
160 thetic datasets by introducing batch effects and noise into a mouse early organogenesis multiome
161 dataset (Argelaguet et al. 2022) (Fig. 2A, Supplementary Fig. S1-S5). As shown in the UMAP
162 plots of the data with synthetic batch effects introduced in RNA and noise introduced in ATAC
163 (Fig. 2A), scTIE effectively removed the batch effects while also better revealing the cell type
164 signals.

165
166 Next, we compared the performance of these methods from three aspects quantitatively,
167 namely batch effect removal, time point alignment and their ability to capture cell type sig-
168 nals (Fig. 2B-C). We quantify the quality of batch removal and time point alignment using
169 purity scores, which calculate the proportion of cells from the same batch/sampling time among
170 neighbors of given cells. A lower purity score indicates a better mixing of batch/time points.
171 We measured the cell type preservation using adjusted rand index (ARI) with the cell type an-
172 notations provided in the original paper as the ground truth. An ideal embedding should mix
173 well cells from different batches and different time points, while maintaining well-separated cell
174 types. These three metrics are summarized in Fig. 2C, where scTIE encloses the largest area,
175 thus outperforming the other methods in the overall performance (Supplementary Fig. S1). Fur-
176 thermore, scTIE’s superior performance is robust against the number of neighbors used in the
177 purity score calculation (Supplementary Fig. S2). We observe similar trends across the other
178 three synthetic scenarios, where scTIE consistently exhibits better performance than the other
179 methods (Supplementary Fig. S3-S5). Together, we demonstrate the superiority of scTIE in data
180 integration, enabling better capture of biological signals through batch effect removal and time
181 point alignment.

182 **scTIE enables identification of cellular subpopulations via modality and**
183 **time point alignment with robust performance.**

184 Encouraged by scTIE's performance in data integration, we next generated a temporal single-cell
185 multimodal dataset and leveraged scTIE for the integration of cells across time points and anno-
186 tation of cell types. We performed single-cell multiome sequencing from mESCs treated with
187 Activin A/Lithium Chloride and measured on Day 2, 4 and 6, using the 10x Chromium Single
188 Cell Multiome platform. After quality control filtering (Supplementary Fig. S6), we obtained
189 high quality measurements of RNA and ATAC from a total of 11,440 cells, with a median detec-
190 tion of 4,130 genes expressed per cell and a median of 11,267 peaks detected per cell.

191
192 By clustering on the joint embeddings produced by scTIE, we identified 17 clusters with
193 either distinct transcription or chromatin accessibility profiles that include cell types from all
194 the three germ layers as well as from extra-embryonic layers of embryonic development (Fig.
195 3A-C). We annotated these clusters based on the key markers identified in the two previous stud-
196 ies (Mittnenzweig et al. 2021; Pijuan-Sala et al. 2019) (Fig. 3C), and confirmed them by label
197 transfer using a public reference (Y Lin et al. 2020; Mittnenzweig et al. 2021) (Supplementary
198 Fig. S7). Further explorations of the motif enrichment of regions with DA in specific clusters
199 highlight the cluster-specific TFs of the annotated cell types (Fig. 3D-E). Additionally, we quan-
200 titatively assessed the clustering results using evaluation metrics. Our findings demonstrate that
201 compared with the existing methods, scTIE better preserves biological signals in each modality
202 and achieves better alignment in days, further supporting our annotation of the cells using the
203 integrated data from scTIE (Supplementary Figs. S8-S9). Furthermore, we performed the same
204 training and clustering procedure on two pseudo-replicates constructed by randomly splitting the
205 data into two halves and demonstrated the consistency of the cell type annotation results. The
206 UMAP visualizations for these two subsets are mostly consistent with an overall accuracy rate
207 of 81% across cell types (Supplementary Fig. S10-S11).

208
209 Notably, scTIE identifies three distinct clusters of definitive endoderm (Cluster 3, 4 and 7)
210 (Supplementary Fig. S12A). We find that Cluster 4 uniquely expresses several Wnt pathway
211 direct targets (*Vcan*, *Nrcam* and *Ccnd2*) and Wnt TF (LEF1), and has lower expressions in Wnt
212 inhibitors *Dkk1* and some definitive endoderm markers (*Hhex* and *Sox17*) (Supplementary Fig.
213 S12B). The activation of Wnt signaling of this group of cells could be linked to primordial lung

214 specification progenitors (Ikonomidou et al. 2020). Cluster 3 and Cluster 7 have similar expression
215 profiles to each other. Compared with Cluster 3, we find Cluster 7 with majority of cells from
216 Day 6 has lower expressions in Nodal signaling genes *Nodal* and *Tdgf1*, but higher expressions in
217 genes that negatively regulate the Nodal pathway (*Cer1* and *Lefty1*) (Supplementary Fig. S12B).

218

219 An inspection of the epiblast subsets further demonstrates that scTIE enables cellular sub-
220 population identification (Supplementary Fig. S13A). We find that one of the epiblast clusters
221 (Cluster 12) has up-regulation of genes related to Hypoxia (*Adm*, *Anxa2*, *Ddit4* and *Gbe1*), which
222 could enhance the definitive endoderm differentiation, as suggested in (LF Chu et al. 2016; Pim-
223 ton et al. 2015) (Supplementary Fig. S13B). In addition, we find that Cluster 1 is enriched
224 with anterior epiblast markers (*Pou3f1*, *Enpp3*, *Pten* and *Slc7a3*), while Cluster 10 highly ex-
225 presses posterior epiblast markers (*Lhx1*, *Ifitm1*) (Supplementary Fig. S13B) (Peng et al. 2016),
226 with down-regulation of the TFs POU5F1 and SOX2 but up-regulation of the TFs FOXA1 and
227 FOXA2 (Supplementary Fig. S13C).

228

229 Finally, we examine the stability of our results in both modality alignment and cluster iden-
230 tification, with respect to key tuning parameters in scTIE, including the weight of OT in the loss
231 function, the number of nodes in hidden layer and the updating frequency of OT. We find that
232 the weight of the OT loss is an important parameter to reach a balance between the alignment of
233 modalities and time points, with a larger weight resulting in a better alignment in time points but
234 poorer performance in modality integration (Supplementary Fig. S14A, E, Supplementary Fig.
235 S15A). In this sense, the choice of this parameter can be guided by the performance in modality
236 alignment, since the pairing information for all cells is known and serves as the ground truth. The
237 two other tuning parameters have a small impact on our results (Supplementary Fig. S14B-D,
238 F-H, Supplementary Fig. S15B-D).

239

240 Together, we demonstrate that scTIE is able to capture distinct cellular subpopulations by
241 preserving information from both epigenomic and transcriptomic profiles, while also aligning
242 the cells from different time points.

243 **scTIE embeddings capture interpretable biological features.**

244 To interpret the embedding space projected by scTIE, we deconvoluted the latent representation
245 by backpropagating the gradient of each dimension in the embedding layer with respect to gene
246 and peak input, followed by ranking the features. We then computed the enrichment scores of
247 the cell type marker list for the feature rankings of each embedding dimension (see Methods).
248 We find that each dimension exhibits distinct patterns of enrichment of cell type markers, and
249 at the same time the cell types from the same lineage share similar enrichment patterns across
250 the dimensions, indicating that scTIE captures diverse and biologically meaningful information
251 from the data (Fig. 4A). We further observe that the enrichment results of RNA and ATAC share
252 similar patterns, illustrating that scTIE is able to link the transcriptomic profiles with the chro-
253 matin accessibility through the common embeddings (Fig. 4A).

254

255 The embedding gradients can be further interpreted in terms of known biological functions,
256 based on their Gene Ontology (GO) enrichment. As illustrated in Fig. 4B, we find that the
257 embedding dimensions enriched with definitive endoderm cell type markers can be associated
258 with different pathways. We observe that dimension 39 is uniquely enriched with Activin recep-
259 tor signaling, as confirmed by the top ranking genes including *Lefty1*, *Fst*, and *Nodal* from this
260 pathway (Fig. 4C). Consistently, the nearest genes of the top ranking peaks also include genes
261 associated with the Activin pathway, such as *Nodal*, *Lefty1* and *Fgf9*. Since treatment by Actin-
262 vin is a key component of our differentiation protocol (see Methods), it is comforting to see that
263 the relevance of this pathway is captured by the fitted model. Together, we demonstrate that sc-
264 TIE is able to project the two modalities into a joint embedding space that captures interpretable
265 biological signals of the data.

266

267 Lastly, we find that the above results are robust to the choice of dimension size (i.e., number
268 of nodes in the embedding layer). We trained scTIE and performed the same gradient calcula-
269 tions with the number of dimensions set to 32 and 96 (vs. the current choice of 64) and found
270 qualitatively similar enrichment patterns (Supplementary Fig. S16). Selected embeddings also
271 show enrichment of GO pathways related to the definitive endoderm development, similar to Fig.
272 4B (Supplementary Fig. S17).

273 **scTIE uncovers cell fate-specific regulatory networks.**

274 scTIE constructs lineage-defining GRNs by combining information across different dimensions
275 of the embedding layer to predict the cell transition probabilities between time points. As a case
276 study, we investigate the transitions of cells from anterior primitive streak on earlier days into en-
277 doderm, mesoderm, as well as remaining as anterior primitive streak on later days. The primitive
278 streak is a transient embryonic structure which marks bilateral symmetry, helps confer anterior-
279 posterior spatial information during gastrulation, and initiates germ layer formation (Mikawa et
280 al. 2004). A distinct group of cells located at anterior primitive streak, the node, forms the axial
281 mesodermal structures and definitive endoderm cells (Hoodless et al. 2001).

282

283 In each of the above three possible cell fates, we fine-tuned the trained embeddings using a
284 prediction layer with weight regularization and backpropagate the gradients from the prediction
285 layer to select the top 200 genes and 500 peak regions as the most predictive features of the
286 lineage. Compared with the conventional approach that uses DE / DA analysis to select the top
287 features, scTIE selects genes and peak regions with significantly better prediction performance
288 (Fig. 5A). The superior prediction performance is consistent across a range of tuning parameters,
289 including the regularization weights and the number of top features, evaluated via cross valida-
290 tion (Supplementary Fig. S18). To demonstrate the benefit of jointly modeling RNA and ATAC
291 data, we considered the alternative approach of only integrating the RNA data across the time
292 points. Then, we used the same gradient approach to select the top genes for each of the three
293 lineages and selected top peaks by physical proximity and correlation. The predictive power of
294 these features decreased compared to joint modeling (Supplementary Fig. S19). Conceptually,
295 we note that joint modeling allows us to train a separate autoencoder for the ATAC modality and
296 backpropagate the gradients from the prediction layer to select the most informative peaks for
297 predicting the transition probabilities. Thus, the framework of scTIE is capable of jointly finding
298 the most predictive features in both modalities.

299

300 In addition, we assessed the stability of the gradient analysis through subsampling. For both
301 the definitive endoderm and mesoderm lineages, we randomly subsampled 60% of the cells con-
302 sidered for each lineage analysis, finetuned the trained neural network and calculated the feature
303 gradients for RNA and ATAC the same way as we did for the full data. The correlations of gra-
304 dients between the subsampling approach (averaged over 50 repetitions) and the full set demon-

305 strate a high level of agreement across all genes and peaks used as input to scTIE (Supplementary
306 Fig. S20).

307

308 To annotate the top peaks, we overlapped the selected peaks with the published enhancer
309 database from 12 tissues of seven developmental stages from 11.5 days after conception until
310 birth (Gorkin et al. 2020), quantified by the Jaccard index. We find that the top peaks associated
311 with mesoderm transition potential are enriched with facial prominence and limb enhancers at
312 E11.5, while endoderm transition-related peaks identified by scTIE show higher enrichment and
313 distinct overlap with stomach enhancers at E14.5, E15.5 and P0 (Fig. 5B). In contrast, the peaks
314 selected by DA analysis show enrichments in tissues that are much less specific to predicted lin-
315 eages of mesoderm or endoderm (Supplementary Fig. S21). Together, these results illustrate that
316 scTIE is able to identify peaks that are specific to lineage transition.

317

318 The identification of genes and peaks that are predictive of cell transition further allows us
319 to infer GRN for each of the lineages: anterior primitive streak, endoderm and mesoderm (see
320 Methods). In the GRN of anterior primitive streak (Fig. 5C, left panel), we identified a few TFs
321 that play key roles in jointly governing anterior mesendoderm and the node development (LHX1,
322 OTX2 and SMAD4) (GC Chu et al. 2004; Costello et al. 2015), as well as a TF related to axial
323 mesendoderm morphogenesis and patterning (MIXL1) (Hart et al. 2002). When focusing on the
324 endoderm GRN (Fig. 5C, middle panel), we find that besides identifying TFs that are central
325 regulators for the formation of definitive endoderm development (SOX17, GATA4, GATA6, and
326 GSC) (F Li et al. 2011; Fisher et al. 2017; Bossard and Zaret 1998; Kanai-Azuma et al. 2002;
327 Heslop et al. 2021), scTIE also captures TFs that are associated with early mesendoderm dif-
328 ferentiation (RUNX1) (VanOudenhove et al. 2016) and morphogenetic movement (LHX1) (Tam
329 and Loebel 2007).

330

331 Lastly, we examined the mesoderm GRN (Fig. 5C, right panel) which identifies a few key
332 TFs (HHEX, SOX17, SMAD3, ZIC3, TWIST1 and NFAT5) that are associated with mesoderm
333 lineages. Notably, most of these TFs have insignificant p-values under DE analysis (Table S1),
334 illustrating that scTIE captures key regulatory signals in this lineage that would be missed oth-
335 erwise. More specifically, the mesoderm GRN highlights TFs that are associated with cardiac
336 development such as ZIC3 in early mesodermal patterning (Z Jiang et al. 2013; Sutherland et
337 al. 2013); HHEX that is involved in mediating the SOX17 for cardiac mesoderm formation in

338 mESC (Y Liu et al. 2014) and NFAT5 for cardiomyogenic during mesodermal induction through
339 regulating the canonical Wnt pathway (Adachi et al. 2012). We also identify TFs that are essen-
340 tial for mesoderm formation and patterning (SMAD3) (Dunn et al. 2004) and cranial mesoderm
341 development (TWIST1) (Bildsoe et al. 2016).

342 **Discussion**

343 While the rapidly increasing collection of single-cell multiomics data provides a wealth of infor-
344 mation for examining context-specific regulatory mechanisms, accurate characterization of cell
345 identities remains the first hurdle to be overcome in such tasks. scTIE provides a unified frame-
346 work for the integration and joint modeling of temporal multimodal data and the subsequent
347 visualization, cell type identification and inference of key regulatory modules predictive of the
348 developmental transitions of cells. Incorporating OT into the training of an autoencoder, scTIE
349 alternates between updating the alignment of cells at different time points and using the current
350 alignment for training the projections into the common embedding space, thus achieving a better
351 balance between integrating time points and maintaining cell type specific signals. As we have
352 demonstrated on the real and synthetic datasets, scTIE outperforms existing paired methods in
353 terms of integration performance.

354

355 Different from existing integration methods that also utilize the notion of a common em-
356 bedding space, scTIE directly exploits the information in this space produced by the nonlinear
357 projections of a neural network, linking it to interpretable features such as genes and peak re-
358 gions. scTIE extracts context-specific gene regulatory relationships through the identification of
359 features that are predictive of cell transition probabilities, which quantify how likely a collection
360 of cells on earlier days will transit to a certain cell state on later days, relative to other cells. These
361 sets of cells can be flexibly defined, allowing users to investigate any cell transition process of
362 interest. In addition to cell transition probabilities derived from OT, the current framework can
363 also be adapted to select features that are predictive of other types of response variables, such
364 as pseudotime and perturbation, which potentially enables the construction of differential GRN
365 under continuous cell differentiation and in perturbed conditions.

366

367 scTIE is designed for temporal multimodal data, which is ideal for studying single-cell ge-
368 nomics in developmental trajectories. Paired measurements from the same cells remove the need

369 for computational pairing, which can introduce errors into the downstream GRN analysis if cells
370 of different cell types are paired, and the issue of cell type imbalance between different modalities.
371 The integration of unpaired developmental data across multiple time points remains an open
372 problem itself. For datasets taken from a matched population, a loss function performing global
373 alignment between modalities, such as the one used in Z Zhang et al. 2022, can be potentially
374 incorporated into the training of scTIE. However, the problem is more challenging if cells are
375 sampled at different time points or develop at a different rate across the modalities, and we will
376 pursue this in future work.

377

378 Although a large number of methods exist for inferring pseudotime ordering of cells from a
379 static snapshot of a developmental process, pseudotime inference assumes that a continuum of
380 cellular states is observed at the sampled time, and thus may not capture the entire transition process
381 (Tritschler et al. 2019). An interesting extension would be combining pseudotime inference
382 and experimental time points to create a finer temporal resolution. However, we note that this
383 would also increase the computation time of scTIE, since iterative OT estimation is performed
384 between consecutive time points; efficient and accurate OT algorithms remain an active area of
385 research.

386

387 We have focused on scRNA-seq and scATAC-seq as common modalities from multimodal
388 profiling technologies. Other modalities such as methylation and protein levels (Mimitou et al.
389 2021; Swanson et al. 2021; Y Wang et al. 2021) can be easily incorporated into scTIE through appropriate
390 encoder-decoder pairs. Since transcriptional regulation involves interactions of protein
391 complexes, histone modifications and other microenvironmental factors, we expect the addition
392 of such information will allow us to build a more accurate prediction model for cellular state
393 changes. Furthermore, emerging single-cell perturbation assays (Rubin et al. 2019) can either be
394 used to validate the top candidates found in our predictive model, or built into the neural network
395 architecture as a prior knowledge graph (Cao and Gao 2022).

396

397 In summary, scTIE provides an integrative framework for analyzing temporal multimodal
398 data, which is an emerging form of data we expect will become more readily available as interests
399 in characterizing GRNs at single-cell resolution continue to rise. On real and synthetic
400 developmental datasets, scTIE is shown to provide effective integration of cells from all time
401 points and select key regulatory elements with superior performance in predicting cellular state

402 changes. We envision that advances in single-cell technologies generating new forms of tempo-
403 ral data will enable us to further expand the functionalities of scTIE, paving the way towards a
404 holistic understanding of cellular transitions and responses in development and disease.

405 **Methods**

406 **Synthetic data construction**

407 The 10x Genomics multiome data of mouse early organogenesis, along with its cell type an-
408 notation, was obtained from the Gene Expression Omnibus database under accession number
409 GSE205117 (Argelaguet et al. 2022). The dataset comprises 59,132 cells from a time course of
410 mouse embryonic development, spanning 5 time points from E7.5 to E8.75.

411 To construct synthetic data that could be processed by most of the methods within their
412 computational capacity, we subset the data to 24,188 cells by selecting only one sample at each
413 time point. We filtered out genes expressed in less than 1% of cells and peaks expressed in less
414 than 5% of cells, resulting in 15,754 genes and 81,108 peaks. To introduce noise and batch effects
415 to the data, we used the `downsampleReads()` function in the `DropletUtils` R package (R Core
416 Team 2022) to downsample the reads. We generated five synthetic scenarios: (1) subsample 10%
417 for all cells in ATAC; (2) subsample 10% for all cells in ATAC and 50% for all cells in RNA;
418 (3) subsample 50% for half of cells in RNA to create the synthetic batch effect in the data; (4)
419 subsample 50% for half of cells in both RNA and ATAC to create the synthetic batch effect in the
420 data; and (5) subsample 10% for all cells in ATAC, subsample 50% for half of the cells in RNA
421 and 25% for the other half of the cells.

422 **mESC data generation**

423 **Cell culture**

424 Mouse embryonic stem cell line R1 was obtained from ATCC. The cells were first expanded on
425 an MEF feeder layer previously irradiated. Then, subculturing was carried out on 0.1% bovine
426 gelatin-coated tissue culture plates. The cells were propagated in mESC medium consisting of
427 Knockout DMEM supplemented with 15% Knockout Serum Replacement, 100 μ M nonessen-
428 tial amino acids, 0.5 mM beta-mercaptoethanol, 2 mM GlutaMax, and 100 U/mL Penicillin-
429 Streptomycin with the addition of 1,000 U/mL of LIF (ESGRO, Millipore).

430 **Cell differentiation**

431 mESCs were differentiated using the hanging drop method (X Wang and P Yang 2008). Trypsinized
432 cells were suspended in chemically defined medium CDM (F Li et al. 2011) to a concentration

433 of 37,500 cells/mL. CDM consists of 75% Iscove's modified Dulbecco's medium (IMDM, Invit-
434 rogen), 25% Ham's F12 medium (Invitrogen), 1X N2 supplements (Invitrogen), 0.05% bovine
435 serum albumin (BSA, Invitrogen), 2 mM Glutamax-1 (Invitrogen), 0.5 mM ascorbic acid (Sigma-
436 Aldrich), and 4.5×10^4 M MTG (Sigma-Aldrich). 20 μ L drops (\sim 750 cells per drop) were then
437 placed on the lid of a bacterial plate and the lid was upside down. After 48 h incubation at 37°C
438 incubator with 5% CO₂, Embryoid bodies (EBs) formed at the bottom of the drops were collected
439 and placed in the well of a 6-well ultra-low attachment plate (Corning) with fresh CDM medium
440 containing 50 ng/mL Activin A (R&D Systems, 338-AC-050/CF) and 2 mM Lithium Chloride
441 (LiCl, Sigma-Aldrich) for up to 6 days, with the medium being changed daily.

442 **Single cell multiome library**

443 We followed 10x Genomics single cell multiome library preparation protocol. The EBs were
444 collected at Day 2, 4, and 6 after Activin A/Lithium Chloride treatment. For each time point,
445 the cells were first treated with StemPro Accutase Cell Dissociation Reagent (Thermo Fisher
446 Scientific) at 37°C for 10-15 min with pipetting. Single cell suspension was obtained by passing
447 through 37 μ M cell strainer (STEMCELL Technologies) twice. After measuring cell concentra-
448 tion, approximately 1 million of cells were centrifuged at 300 rcf for 5 min. Nuclei were isolated
449 by following the protocol provided by 10x Genomics (Nuclei isolation for single cell multi-
450 ome ATAC + Gene expression sequencing, CG00365, Rev A). The final nuclei concentration
451 was adjusted to 3000 cell/ μ L in 1X Nuclei Buffer (10x Genomics). The sample was immedi-
452 ately submitted to Stanford Genomics Service Center (SGSC) for single cell sorting using 10x
453 Chromium Controller (target cells: 5000 per replicate, total 2-3 replicates per time point). The
454 single cell multiome library was generated using Chromium Next GEM Single Cell Multiome
455 ATAC + Gene Expression Reagent Bundle Kit (10x Genomics, PN-1000283).

456 **Data preprocessing**

457 10x Genomics Cell Ranger arc v2.0.0 was used to process the raw FASTQ files for each multiome
458 single-cell dataset separately. The reference genome and transcriptome for alignment and annota-
459 tion was version arc-mm10-2020-A-2.0.0. To integrate all filtered count matrices for scRNA-seq
460 and scATAC-seq from different replicates and time points, the cellranger-arc aggr command was
461 applied with default depth normalization method.

462 Next, we performed quality control on the cell level. We removed cells based on the following

463 criteria in scRNA-seq: (1) with the total number of UMI (nUMI) less than 6000 on Day 2, 3000
 464 on Day 4 and Day 6; (2) with nUMI greater than 100,000; (3) with the number of genes less
 465 than 2000 on Day 2, 1800 on Day 4 and 1500 on Day 6 and (4) mitochondrial reads greater than
 466 25%. We further removed cells based on the following criteria in scATAC-seq: (1) with less than
 467 500 total ATAC fragments and (2) with less than 500 peaks detected. After quality control, we
 468 retained 11440 cells (Day 2: 2896 cells; Day 4: 2796 cells and Day 6: 5748 cells). We then
 469 performed the quality control on the feature level, removing the genes that are not expressed in
 470 any cells and the peaks that are expressed at least 5% of cells, resulting in 26717 genes and 61744
 471 peaks as input in scTIE.

472 **Architecture and training of scTIE**

473 scTIE uses an autoencoder structure to project high dimensional feature vectors (i.e., gene ex-
 474 pression levels and accessibility peaks) from all time points into a lower dimensional common
 475 embedding space and reconstruct the features in the original high dimensional space. Each
 476 modality has its own encoder and decoder (Table 1). For RNA, the architecture has an addi-
 477 tional pair of coupled batchnorm layers, where the final reconstructed output uses the moving
 478 average μ and standard deviation σ stored in the first batchnorm layer of the encoder to perform
 479 rescaling. This accounts for the high variability in gene expression levels without the need for
 480 selecting HVGs, and allows us to significantly improve the performance in reconstruction cor-
 481 relation, modality and day alignment, and clustering quality (Supplementary Fig. S22). The
 482 pairing between feature vectors from the same cell is enforced through a modality loss function
 483 minimizing their distance in the embedding space. An OT matrix is used to construct cell trajec-
 484 tories between each pair of consecutive time points. In contrast to existing methods using OT for
 485 trajectory inference, we integrate an OT loss into the autoencoder training process and estimate
 486 the OT matrix iteratively throughout. A larger weight on the OT loss leads to better alignment
 487 between days (Supplementary Fig. S15A).

488 Let $X^{(t,s)}$ denote the data matrix from time point t and modality s , where $t = 1, \dots, T$ and
 489 $s = 1, 2$ for RNA and ATAC respectively. Each time point t provides measurements for N_t
 490 cells; thus in this case, $X^{(t,1)} \in \mathbb{R}^{D_1 \times N_t}$ with $D_1 =$ number of genes and $X^{(t,2)} \in \mathbb{R}^{D_2 \times N_t}$ with
 491 $D_2 =$ number of peak regions. In each iteration, a mini-batch of data is sampled by taking equal-
 492 sized subsets of cells from each time point, that is, $\mathcal{B} = \{\mathcal{B}^{(t)}\}_{t=1}^T$, where each subset $\mathcal{B}^{(t)}$ has B
 493 cells. Three loss functions are applied to the mini-batch.

1. *Reconstruction loss.* (f_s, g_s) represents the encoder-decoder pair for modality s . Compared with the architecture for ATAC, the RNA part has a pair of coupled batchnorm layers, starting with a batchnorm layer in the encoder to remove scale variations in genes and prevent the gradients from being dominated by a small number of highly expressed genes (Table 1). Let $x_i^{(t,1)}$ denote the gene expression vector from cell i at time t and $\tilde{x}_i^{(t,1)}$ denote the normalized output from the first batchnorm layer, then $\tilde{x}_i^{(t,1)} = (x_i^{(t,1)} - \mu)/\sigma$, where μ and σ are the moving average and standard deviation of the genes saved in the batchnorm layer throughout training. The reconstruction loss is applied to the normalized data and the output from the decoder, defined as

$$L_{\text{recon}}^{(1)} = \frac{1}{TB} \sum_{t=1}^T \sum_{i \in \mathcal{B}^{(t)}} \|\tilde{x}_i^{(t,1)} - g_1(f_1(x_i^{(t,1)}))\|_2^2.$$

494 For ATAC, the first layer in the encoder is a fully connected layer and the reconstruction
495 loss is computed on the input $x_i^{(t,2)}$ and output $g_2(f_2(x_i^{(t,2)}))$ as usual. The overall L_{recon} is
496 the sum of $L_{\text{recon}}^{(1)}$ and $L_{\text{recon}}^{(2)}$.

497 2. *Optimal transport loss.* We leverage OT to effectively align cells from all time points in
498 the embedding space. For notational convenience, we will suppress the dependence on
499 modality s for now, with understanding that the following steps are performed for each
500 modality. For any two adjacent time points t and $t + 1$, a transport cost matrix $C^{(t,t+1)} \in$
501 $\mathbb{R}^{N_t \times N_{t+1}}$ can be computed using the current embeddings, where the (k, l) -th entry is given
502 by $C^{(t,t+1)}(k, l) = \|f(x_k^{(t)}) - f(x_l^{(t+1)})\|_2$ for the k -th cell from t and the l -th cell from
503 $t + 1$. With the cost matrix, Waddington-OT (Schiebinger et al. 2019) is then used as the
504 algorithm to estimate a transport matrix $\gamma^{(t,t+1)} \in \mathbb{R}^{N_t \times N_{t+1}}$. Each row in $\gamma^{(t,t+1)}$ sums to
505 1, representing the transition probabilities of a cell in time step t to all the other cells in
506 time step $t + 1$. Given T time steps, we need to maintain a total of $T - 1$ transport matrices
507 throughout the autoencoder training process. For a given mini-batch \mathcal{B} in each iteration,
508 a submatrix version of $C^{(t,t+1)}$ is computed using the rows and columns specified in \mathcal{B}
509 and is denoted by $\tilde{C}^{(t,t+1)}$. Similarly, a mini-batch version $\tilde{\gamma}^{(t,t+1)}$ of $\gamma^{(t,t+1)}$ is calculated
510 by taking the appropriate submatrix and rescaling the rows to unit sum. The batch-wise
511 feature alignment loss (for each modality s) is defined as

$$L_{\text{ot}} = \frac{1}{T-1} \sum_{t=1}^T \left(\sum_{k=1}^B \sum_{l=1}^B (\tilde{C}^{(t,t+1)} \odot \tilde{\gamma}^{(t,t+1)})(k, l) \right),$$

512 where \odot is the Hadamard product. The final L_{ot} is the sum over modalities s .

513 3. *Modality alignment loss*. For each mini-batch, the modality alignment loss is simply de-
514 fined as the L2 distance between feature vectors from the same cell in the embedding space,
515 which is to be minimized:

$$L_{\text{modality}} = \frac{1}{TB} \sum_{t=1}^T \sum_{i \in \mathcal{B}^{(t)}} \|f_1(x_i^{(t,1)}) - f_2(x_i^{(t,2)})\|_2^2.$$

516 The total loss in each iteration is $L = \lambda_{\text{recon}} L_{\text{recon}} + \lambda_{\text{ot}} L_{\text{ot}} + L_{\text{modality}}$ where the λ 's are tuning
517 parameters controlling the relative weighting of the losses. For every K epochs, the transport
518 matrices (for each modality s) $\gamma_s^{(t,t+1)}$, $1 \leq i \leq T - 1$ are updated by computing OT on the
519 current embedding features.

520 The functionalities of each loss function in L :

- 521 1. The **reconstruction loss** preserves the original data signals (i.e., distinct cell type signals)
522 at each time point by encouraging the autoencoders to learn a low-dimensional embedding
523 that can reproduce the data input.
- 524 2. The **OT loss** aligns embeddings between consecutive time points by calculating an align-
525 ment cost function derived from the estimated transition probabilities. To reduce the align-
526 ment cost, cell pairs with high transition probabilities should be near each other in the
527 embedding space and vice versa. The transition probabilities and embeddings are itera-
528 tively refined. Additionally, the loss aids in mitigating batch effects, as OT can cross-align
529 cells from different batches when mapping cells between consecutive time points. As we
530 demonstrated on the synthetic data with batch effects in both RNA and ATAC data (Sup-
531 plementary Fig. S4), the pre-training stage (see Training details below), which only trains
532 the RNA autoencoder using the OT loss and the reconstruction loss, already removes most
533 of the batch effects in RNA data (Supplementary Fig. S23).
- 534 3. The **modality alignment loss** makes use of the pairing information between RNA and
535 ATAC so that the final embeddings take into account signals in both modalities.

536 Training details

537 scTIE took a collection of peak matrices from scATAC-seq data and raw counts matrices from
538 scRNA-seq data from multiple time points as input. For ATAC, the peak matrices were trans-

539 formed to binary matrices, where one represents any non-zero original values. For RNA, the
 540 raw count matrices were sized-factor normalized and then log-transformed. For the overall mul-
 541 timodal training, we first pre-trained the RNA autoencoder f_1, g_1 for 500 epochs (excluding
 542 L_{modality}). Then, we fixed the weights of the pretrained RNA model to train the ATAC model
 543 for 300 epochs with the overall loss L . Finally, the two models were jointly trained for 200
 544 epochs using the full algorithm as detailed in Algorithm 1. The final joint embeddings were
 545 calculated by taking the averages of $f_1(x_i^{(t,1)})$ and $f_2(x_i^{(t,2)})$ for each cell i from time t , followed
 546 by computing the final $\gamma^{(t,t+1)}$ from the joint embeddings. Throughout training, we used Adam
 547 as the optimizer with learning rate set to 0.1, batch size $B = 256$, tuning parameters $\lambda_{\text{recon}} = 1$,
 548 $\lambda_{\text{ot}} = 0.1$, and OT was updated every 10 epochs.

549 We note here that due to the pre-training of the RNA autoencoder, the biological signals
 550 utilized by scTIE to produce the common embeddings were mostly driven by the RNA modality.
 551 However, complementary signals from scATAC-seq still play a role in generating the embeddings
 552 since the modality alignment loss is affected by both RNA and ATAC positions in the embedding
 553 space. Pre-training with RNA signals is essential for stable training of the neural network because
 554 (i) the RNA modality generally contains stronger signals for cell type identification and (ii) the
 555 dimension of ATAC input (number of peaks) is much larger than that of the RNA modality
 556 (number of genes).

Algorithm 1 Multimodal OT Autoencoder (two-modality case)

Data matrices $X^{(t,s)}$, training iterations M , batch size B , autoencoder f_1, g_1, f_2, g_2 with weights
 θ , learning rate α , loss weight tuning parameters $\lambda_{\text{recon}}, \lambda_{\text{ot}}$, OT update frequency K .

Initialize all $\gamma_s^{(t,t+1)}, 1 \leq t \leq T - 1$ matrices with zero matrices.

for $iteration = 1, 2, \dots, M$ **do**

 Sample cells $\mathcal{B} = \{\mathcal{B}^{(t)}\}_{t=1}^T$, where each subset $\mathcal{B}^{(t)}$ has B cells.

 Compute $L_{\text{recon}}, L_{\text{ot}}, L_{\text{modality}}$

 Compute $L = \lambda_{\text{recon}}L_{\text{recon}} + \lambda_{\text{ot}}L_{\text{ot}} + L_{\text{modality}}$

 Perform gradient descent step on autoencoder weights $\theta \leftarrow \theta - \alpha \nabla_{\theta} L$

if $M \% K == 0$ **then**

 Update $\gamma_s^{(t,t+1)}, 1 \leq t \leq T - 1, s = 1, 2$ using current embeddings.

end if

end for

557 **Estimation of long-range transition probabilities**

558 Long-range transition probabilities can be estimated by multiplying the transport matrices. For
559 example, $\gamma^{(t,t+2)}$ can be calculated as $\gamma^{(t,t+1)}\gamma^{(t+1,t+2)}$. An alternative approach is to compute
560 $\gamma^{(t,t+2)}$ directly from OT. However, since OT interpolates between two observed datasets by
561 finding the shortest path in the space of distributions, one has to implicitly assume that the cells
562 do not change their expression or accessibility by large amounts over the two time points. It is
563 generally recommended that long-range time couplings are estimated by multiplying the gamma
564 matrices (Schiebinger et al. 2019). On the mESC dataset, these two ways of estimation give
565 positively correlated results, with the mode of correlations lying around 0.6 (Supplementary Fig.
566 S24).

567 **Cell type annotation of mESC data**

568 **Cell clustering of scTIE**

569 To identify the clusters on the common embedding of scTIE, we first constructed a shared nearest
570 neighbor graph using `buildSNNGraph` in R package `scrAn` (Lun et al. 2016) (v 1.23.0), with
571 the number of nearest neighbor set as 15 with weighted scheme set as `jaccard`. Next we performed
572 Leiden community detection (Traag et al. 2019) on the shared nearest graph with resolution 1.8
573 and number of iterations 50, implemented in R package `leidenAlg` (v 1.0.3), resulting in 17
574 clusters in total.

575 **Motif enrichment**

576 We used `Signac` (Stuart et al. 2021) to calculate the over-represented motif of each cluster
577 based on the differential accessible peaks. The motif position frequency matrices are obtained
578 from `Cis-BP` (Weirauch et al. 2014). We used `limma-trend` (Ritchie et al. 2015) to perform
579 differential accessibility analysis between the cells in one cluster and the remaining cells, where
580 the top 500 peaks of each cluster with log fold change greater than 0.1 and adjusted p-value less
581 than 0.001 are selected. We then performed the motif enrichment analysis using `FindMotifs`
582 to find motifs over-represented in the selected set of peaks.

583 **Benchmarking and evaluation metrics**

584 **Settings used in other methods**

585 We benchmarked the performance of scTIE against four other methods designed for single-cell
586 paired multimodal data integration: Seurat, scAI, MultiVI and MOFA. We compared scTIE's
587 performance in terms of visualisation of the latent space, alignment of the days and clustering in
588 the latent space against these methods.

- 589 • **Seurat.** R package Seurat v4.1.0 (Y Hao et al. 2021) was used. We ran Seurat (WNN)
590 using `FindMultiModalNeighbors`, with the reduction list input as the first 50 com-
591 ponents of LSI reduced dimension of scATAC-seq (with the first dimension excluded) and
592 50 top PCs of scRNA-seq, with other parameters set as default.
- 593 • **scAI.** R package scAI v1.0.0 (Jin et al. 2020) was used. We ran scAI using `run_scAI` by
594 setting the rank of the inferred factor set as 64 and `nrun = 5`, with other parameters set
595 as default.
- 596 • **MultiVI.** Python package scvi v0.15.0 (Ashuach et al. 2023) was used. We ran Mul-
597 tiVI using `MULTIVI` by setting the `fully_paired = True`, `n_hidden = 256` and
598 `n_latent = 64`, with other parameters set as default. The model was then trained with
599 `max_epochs = 200`.
- 600 • **MOFA.** R package MOFA2 v1.7.0 (Argelaguet et al. 2020) was used. We ran MOFA using
601 `run_mofa` by setting the number of factors as 64, with other parameters set as default.

602 **Benchmarking of mESC data**

603 **Modality alignment:** We used two metrics to measure scTIE's performance in the alignment of
604 the two modalities, namely FOSCTTM and paired data proportion.

- 605 • **FOSCTTM.** FOSCTTM refers to Fraction of Samples Closer than True Match, which is
606 first introduced in MMD-MA (J Liu et al. 2019) to quantify the alignment of multi-omics
607 data. To evaluate the modal alignment of scTIE using FOSCTTM, we first calculated the
608 Euclidean distance between the ATAC embedding and RNA embedding. Then for each
609 modality we calculated one FOSCTTM score, which summarizes the proportion of cells
610 that are closer to the ground truth matched cells based on the distance matrix. Finally we

611 summarized the FOSCTTM scores from the two modalities into one score by taking the
612 average.

613 • *Paired data proportion.* Paired data proportion (used in Cobolt (Gong et al. 2021)) calcu-
614 lated the proportion of cells whose ground truth matched cells are included within a certain
615 number of neighbors, based on the Euclidean distance between the ATAC embedding and
616 RNA embedding. We varied the number of neighbors from 1 to the total number of cells
617 in the data.

618 **Day alignment:** We quantified the alignment of data sampled on different days using neighbor-
619 hood purity using `neighborPurity` in R package `bluster` (v1.5.1), which calculated the
620 proportion of cells from the same day among a certain number of neighbors, based on the UMAP
621 coordinates generated from the common latent embeddings.

622

623 **Comparison with single-modality clustering:** We benchmarked clustering results from scTIE
624 against other paired data integration methods by evaluating how similar the results are compared
625 to clustering dimension-reduced scRNA-seq (PCA space) or scATAC-seq (LSI space) alone. On
626 the latent space of each method or the dimension-reduced space from scRNA-seq or scATAC-
627 seq, we performed Leiden clustering on the shared nearest neighbor graphs constructed, with
628 the same parameter settings as mentioned in Section *Cell clustering*. Note that for Seurat, we
629 performed Leiden clustering directly on the weighted nearest neighbor graph it outputs. We used
630 two metrics to quantify the results, Adjusted Rand Index and silhouette coefficient.

631 • *Adjusted Rand Index (ARI).* We computed the ARI scores of clustering results from each
632 data integration method and clustering results from scRNA-seq or scATAC-seq alone.

633 • *Silhouette coefficient.* For each clustering result, we computed the silhouette coefficient
634 based on the Euclidean distance calculated from the UMAP coordinates generated from
635 the dimension-reduced scRNA-seq or scATAC-seq.

636 For both metrics, higher values indicate a method better captures the clustering information in a
637 single modality.

638 **Benchmarking of synthetic data**

639 We benchmarked the data integration performance of scTIE with the other paired data integration
640 methods in terms of three evaluation metrics: (1) ARI scores of the cell type annotation provided

641 by the original study and the Leiden clustering results from each method; (2) neighborhood purity
642 of days; and (3) neighborhood purity of batch for scenarios with synthetic batch effects.

643 **Enrichment analysis for embedding dimensions**

644 Upon completion of training, scTIE has projected the high dimensional feature vectors (genes
645 and peaks) into a 64 dimensional embedding space. Treating each dimension as a representation
646 unit, for each cell type, we backpropagate the gradient of each unit with respect to gene and peak
647 input to select features with the largest impact. More specifically, for each cell in cell type G ,
648 we pass its gene expression vector through the autoencoder to obtain its embedding vector y and
649 compute $\frac{\partial y_j}{\partial x_i}$ for each dimension j and gene input node i . The gradients are averaged over all cells
650 in G to obtain the mean gradient for each gene. We then take the variability of gene expression
651 into account by multiplying each mean gradient by its corresponding gene standard deviation,
652 so that the final gradients are equivalent to gradients after the first batchnorm layer. Finally, we
653 rank the genes by their gradient values and calculate the enrichment scores of the top 200 genes
654 from the DE analysis of cell type G , where the DE analysis is performed using `limma-trend`
655 (Ritchie et al. 2015) between the cells in one cluster and the remaining cells. Similar steps are
656 performed for the peaks and the top 500 peaks are selected for enrichment score calculation.

657 We used `fgsea` function in the R package `fgsea` (Korotkevich et al. 2021) to perform the
658 gene set enrichment analysis (GSEA) on the pathways related to mouse embryonic stem cells (as
659 listed in **Fig. 4B**). Significant pathways are defined with adjusted p-value less than 0.05.

660 **GRN inference**

661 **Selecting features with high predictive power**

662 By building a prediction framework on the obtained transition probabilities, scTIE selects genes
663 and peaks jointly with high predictive power for developmental outcomes. In the mESC data, we
664 consider how a group of cells from earlier days, denoted as G_0 , develops into two other groups
665 G_1 and G_2 on later days.

The transition probabilities are obtained from $\gamma^{(t,t+1)}$ ($t = 1, 2$ in our data) so that each cell i
in G_0 is associated with a probability vector (p_{i1}, p_{i2}) indicating its probabilities of becoming G_1
and G_2 (See Section *Cell transition probability calculation*). We finetune a one-layer classifier
on the pretrained features in the embedding space of cells in G_0 to predict their transition prob-

abilities. A simple linear classifier is sufficient to partition the cell feature space into G_1 and G_2 when the pretrained features are representative enough. Concretely, let q be the linear classifier and \mathcal{B} be a mini-batch of cells from G_0 of size B , we employ a batch-wise KL divergence loss defined

$$L_{kl} = \frac{1}{B} \sum_{j \in \mathcal{B}} D_{KL}(q(f(x_j)) || P_j),$$

666 where f is the trained encoder, $P_j = (p_{j1}, p_{j2})$. This loss enforces the classifier q to output
 667 transition probability distributions close to those in P_j 's. We also include the modality alignment
 668 loss L_{modality} , with weight default set as 0.1. The classifier is trained with Adam setting learning
 669 rate to 0.001, training epochs to 200, batch size to 256 and L1 regularization.

670 After training, gradients from the two classification nodes are backpropagated to each gene
 671 (or peak) input the same way as in computing embedding gradients. The gene gradients are then
 672 scaled by multiplying with the gene-wise standard deviations. A positive gradient for gene (or
 673 peak) j with respect to the node for G_1 means increasing the input feature value tend to increase
 674 the cells' probabilities of becoming G_1 , while a negative value indicates more contribution to G_2 .
 675 The final feature ranking is based on the average gradients by repeating this procedure 20 times
 676 with different seeds.

677 **Selection of G_0, G_1, G_2**

678 As a case study in this paper, we focus on the transition of cells from anterior primitive streak on
 679 Day 2 and Day 4 into endoderm, mesoderm, as well as remaining as anterior primitive streak on
 680 Day 4 and Day 6.

First, we considered the cells that are annotated as anterior primitive streak (Cluster 6) on Day 2 and Day 4 as G_0 . G_1 and G_2 are then selected from the cells on Day 4 and Day 6 that are more likely to be the descendants of G_0 , as quantified by the descendant scores. The descendant scores are defined similarly as in WOT (Schiebinger et al. 2019). Recall $\gamma^{(t,t+1)}$ is the N_t by N_{t+1} transition probability matrix between time points t and $t + 1$, let $s_t \in R^{N_t}$ be the vector of descendant scores for all cells at time point t , then we can calculate

$$s_{t+1} = s_t \gamma^{(t,t+1)}, \text{ where } s_t(i) = \begin{cases} \frac{1}{|G_0|}, & \text{if cell } i \text{ is in } G_0, \\ 0, & \text{otherwise.} \end{cases}.$$

681 This formula can then be pushed forward again to calculate the descendant scores for the next
 682 time point $t + 2$, and so on. For all cells in G_0 at time point t (here $t = 1$ or 2), we calculated

683 the descendant scores s_{t+k} of all cells at the later time point $t + k$, for $k = 1, \dots, T - t$. We
 684 then considered the cells with descendant scores greater than the median of all cells at a certain
 685 time point as the potential descendants, i.e., cells with $s_{t+k}(i) > \text{median}(s_{t+k})$. Among these
 686 descendant cells, we selected three pairs of G_1 and G_2 corresponding to the three cell fates we
 687 have analyzed: G_1 that are annotated as (1) anterior primitive streak or (2) definitive endoderm
 688 or (3) mesoderm; for each selection of G_1, G_2 always represents the remaining descendant cells.

689 Cell transition probability calculation

For each cell $i \in G_0$ on Day t , and G_1, G_2 on Day $k \in K$, where $K = \{k : t < k \leq T\}$, the transition probability vector $(p_{i1}^{(t)}, p_{i2}^{(t)})$ are calculated as the following,

$$\begin{aligned}
 p_{i1}^{(t,k)} &= \sum_{y \in G_1} \gamma^{(t,k)}(i, y), \\
 p_{i2}^{(t,k)} &= \sum_{y \in G_2} \gamma^{(t,k)}(i, y), \\
 p_{ij}^{(t,k)} &= \frac{p_{ij}^{(t)}}{\sum_j p_{ij}^{(t)}}, j = 1, 2, \\
 p_{ij}^{(t)} &= \frac{1}{|K|} \sum_k p_{ij}^{(t,k)}.
 \end{aligned}$$

690 (p_{i1}, p_{i2}) is then the concatenated vector of $(p_{i1}^{(t)}, p_{i2}^{(t)})$.

691 Evaluation of cell transition probability prediction

692 To evaluate the predictive power of the selected features to the transition probability, we per-
 693 formed support vector machine (SVM) with radial kernel to predict the transition probability
 694 using Day 2 and 4 anterior primitive streak gene expression of the top selected genes and peak
 695 matrix of the top selected peaks. The performance are quantified by root mean squared error
 696 (RMSE) from a 20 repeated 5 fold cross validation. We benchmarked the predictive power of the
 697 features selected by gradients with different regularization weights (0, 1, 10, 100), against the
 698 features selected by DE/DA analysis using limma-trend (Ritchie et al. 2015).

699 Gene regulatory network construction

700 To construct the gene regulatory network for each cell fate (anterior primitive streak, definitive
 701 endoderm and mesoderm), we focus on the top 500 genes based on the gradient ranking. For each

702 gene, we consider the open chromatin regions that are within 250kb upstream and downstream of
703 its transcription start site (TSS) as well as ranked top 2000 according to the gradients as the distal
704 candidate functional regions, which results in 396, 404 and 339 gene-peak pairs for the three
705 cell fates respectively. We next filtered the pairs based on the gene-peak correlation, calculated
706 from the meta-cells. The meta-cells were constructed using the following strategies: we first
707 randomly selected 100 cells from the anterior primitive streak cells on Day 2. For each cell, we
708 looked for its 5 nearest neighbors based on the euclidean distances of the common embeddings
709 and aggregated them as a meta-cell. Then, we calculated Pearson's correlation of the gene-peak
710 pairs for these 100 meta-cells. This procedure is repeated 20 times and the gene-peak pairs with
711 an absolute average correlation greater than 0.2 are retained (APS: 35, DE: 38 and MES: 17 pairs
712 remained).

713 To link the peak region with the TF, we identified the enriched TF using *matchMotifs* func-
714 tion in R package *motifmatchr* of the peaks from the selected gene-peak pairs based on Cis-BP
715 database (Weirauch et al. 2014). We only consider if the TF are the top 500 genes. Finally, by
716 linking the TF-region and peak-gene relationships, we construct the TF-gene regulatory networks
717 that are associated cell fate probabilities.

718 In the alternative approach of only integrating RNA across time, we selected the peaks that
719 are within 250kb upstream and downstream of the transcription start site (TSS) of the top ranking
720 genes, with Pearson's correlation greater than 0.2.

721 **Data access**

722 All raw and processed sequencing data generated in this study have been submitted to the NCBI
723 Gene Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) under accession number
724 GSE223041.

725 scTIE is available at <https://github.com/SydneyBioX/scTIE> and as part of the
726 Supplementary Code.

727 **Acknowledgements**

728 We would like to thank Michael Blanco and Dhananjay Wagh from Stanford Genomics Ser-
729 vice Center (SGSC) for their kind help on the preparation of 10x Genomics single cell mul-
730 tiome libraries. We also want to thank Xuhuai Ji from SGSC for providing sequencing ser-
731 vices. The Illumina HiSeq 4000 was purchased using a NIH S10 Shared Instrumentation Grant
732 (S10OD018220). The Illumina NovaSeq 6000 was also purchased using a NIH S10 Shared
733 Instrumentation Grant (1S10OD02521201). The authors gratefully acknowledge the following
734 funding sources: Research Training Program Tuition Fee Offset and Stipend Scholarship and
735 Chen Family Research Scholarship to Y.L.; AIR@innoHK programme of the Innovation and
736 Technology Commission of Hong Kong to J.Y.H.Y. and Y.L.; the UT Austin Harrington Faculty
737 Fellowship to Y.X.R.W; NIH grants R01 HG010359 and P50 HG007735 to W.H.W.

738 *Author contributions:* T.W., W.H.W. and Y.X.R.W. conceived and designed this project; X.C.
739 performed the mESC multiome experiment; Y.L., T.W., S.W., B.C., and J.X. performed data
740 preprocessing, model development, and evaluation of results; J.Y.H.Y., W.H.W. and Y.X.R.W.
741 supervised the execution; Y.L., B.C., J.X., J.Y.H.Y., W.H.W. and Y.X.R.W. wrote the manuscript.
742 All authors read and approved the manuscript.

743 **Competing interests**

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

745 **References**

- 746 Adachi A, Takahashi T, Ogata T, Imoto-Tsubakimoto H, Nakanishi N, Ueyama T, and Matsubara
747 H. 2012. NFAT5 regulates the canonical Wnt pathway and is required for cardiomyogenic
748 differentiation. *Biochemical and biophysical research communications*. **426**: 317–323.
- 749 Argelaguet R, Arnol D, Bredikhin D, Deloro Y, Velten B, Marioni JC, and Stegle O. 2020.
750 MOFA+: a statistical framework for comprehensive integration of multi-modal single-cell
751 data. *Genome biology*. **21**: 1–17.
- 752 Argelaguet R, Lohoff T, Li JG, Nakhuda A, Drage D, Krueger F, Velten L, Clark SJ, and Reik W.
753 2022. Decoding gene regulation in the mouse embryo using single-cell multi-omics. *bioRxiv*.
754 2022–06.
- 755 Ashuach T, Gabitto MI, Koodli RV, Saldi GA, Jordan MI, and Yosef N. 2023. MultiVI: deep
756 generative model for the integration of multimodal data. *Nature Methods*. **20**: 1222–1231.
- 757 Bildsoe H, Fan X, Wilkie EE, Ashoti A, Jones VJ, Power M, Qin J, Wang J, Tam PP, and Loebel
758 DA. 2016. Transcriptional targets of TWIST1 in the cranial mesoderm regulate cell-matrix
759 interactions and mesenchyme maintenance. *Developmental biology*. **418**: 189–203.
- 760 Bossard P and Zaret KS. 1998. GATA transcription factors as potentiators of gut endoderm dif-
761 ferentiation. *Development*. **125**: 4909–4917.
- 762 Cao ZJ and Gao G. 2022. Multi-omics single-cell data integration and regulatory inference with
763 graph-linked embedding. *Nature Biotechnology*. 1–9.
- 764 Chen S, Lake BB, and Zhang K. 2019. High-throughput sequencing of the transcriptome and
765 chromatin accessibility in the same cell. *Nature biotechnology*. **37**: 1452–1457.
- 766 Chu LF, Leng N, Zhang J, Hou Z, Mamott D, Vereide DT, Choi J, Kendziorski C, Stewart R, and
767 Thomson JA. 2016. Single-cell RNA-seq reveals novel regulators of human embryonic stem
768 cell differentiation to definitive endoderm. *Genome biology*. **17**: 1–20.
- 769 Chu GC, Dunn NR, Anderson DC, Oxburgh L, and Robertson EJ. 2004. Differential require-
770 ments for Smad4 in TGF β -dependent patterning of the early mouse embryo.
- 771 Ciortan M and Defrance M. 2021. Explainability methods for differential gene analysis of single
772 cell RNA-seq clustering models. *bioRxiv*.
- 773 Costello I, Nowotschin S, Sun X, Mould AW, Hadjantonakis AK, Bikoff EK, and Robertson EJ.
774 2015. Lhx1 functions together with Otx2, Foxa2, and Ldb1 to govern anterior mesendoderm,
775 node, and midline development. *Genes & development*. **29**: 2108–2122.

- 776 Dunn NR, Vincent SD, Oxburgh L, Robertson EJ, and Bikoff EK. 2004. Combinatorial activities
777 of Smad2 and Smad3 regulate mesoderm formation and patterning in the mouse embryo.
- 778 Duren Z, Chang F, Naqing F, Xin J, Liu Q, and Wong WH. 2022. Regulatory analysis of sin-
779 gular cell multiome gene expression and chromatin accessibility data with scREG. *Genome*
780 *biology*. **23**: 1–19.
- 781 Fisher J, Pulakanti K, Rao S, and Duncan S. 2017. GATA6 is essential for endoderm formation
782 from human pluripotent stem cells. *Biology Open*. **6**: 1084–1095.
- 783 Forrow A and Schiebinger G. 2021. LineageOT is a unified framework for lineage tracing and
784 trajectory inference. *Nature communications*. **12**: 1–10.
- 785 Gong B, Zhou Y, and Purdom E. 2021. Cobolt: integrative analysis of multimodal single-cell
786 sequencing data. *Genome biology*. **22**: 1–21.
- 787 Gorkin DU, Barozzi I, Zhao Y, Zhang Y, Huang H, Lee AY, Li B, Chiou J, Wildberg A, Ding B,
788 et al. 2020. An atlas of dynamic chromatin landscapes in mouse fetal development. *Nature*.
789 **583**: 744–751.
- 790 Hao Y, Hao S, Andersen-Nissen E, Mauck III WM, Zheng S, Butler A, Lee MJ, Wilk AJ, Darby
791 C, Zager M, et al. 2021. Integrated analysis of multimodal single-cell data. *Cell*. **184**: 3573–
792 3587.
- 793 Hart AH, Hartley L, Sourris K, Stadler ES, Li R, Stanley EG, Tam PP, Elefanty AG, and Robb L.
794 2002. Mixl1 is required for axial mesendoderm morphogenesis and patterning in the murine
795 embryo.
- 796 Heslop JA, Pournasr B, Liu JT, and Duncan SA. 2021. GATA6 defines endoderm fate by control-
797 ling chromatin accessibility during differentiation of human-induced pluripotent stem cells.
798 *Cell reports*. **35**: 109145.
- 799 Hoodless PA, Pye M, Chazaud C, Labbé E, Attisano L, Rossant J, and Wrana JL. 2001. FoxH1
800 (Fast) functions to specify the anterior primitive streak in the mouse. *Genes & development*.
801 **15**: 1257–1271.
- 802 Ikonomidou L, Herriges MJ, Lewandowski SL, Marsland R, Villacorta-Martin C, Caballero IS,
803 Frank DB, Sanghrajka RM, Dame K, Kańduła MM, et al. 2020. The in vivo genetic program
804 of murine primordial lung epithelial progenitors. *Nature communications*. **11**: 1–17.
- 805 Jiang Y, Harigaya Y, Zhang Z, Zhang H, Zang C, and Zhang NR. 2022. Nonparametric single-
806 cell multiomic characterization of trio relationships between transcription factors, target genes,
807 and cis-regulatory regions. *Cell Systems*. **13**: 737–751.

- 808 Jiang Z, Zhu L, Hu L, Slesnick TC, Pautler RG, Justice MJ, and Belmont JW. 2013. *Zic3* is
809 required in the extra-cardiac perinodal region of the lateral plate mesoderm for left–right
810 patterning and heart development. *Human molecular genetics*. **22**: 879–889.
- 811 Jin S, Zhang L, and Nie Q. 2020. scAI: an unsupervised approach for the integrative analysis of
812 parallel single-cell transcriptomic and epigenomic profiles. *Genome biology*. **21**: 1–19.
- 813 Kanai-Azuma M, Kanai Y, Gad JM, Tajima Y, Taya C, Kurohmaru M, Sanai Y, Yonekawa H,
814 Yazaki K, Tam PP, et al. 2002. Depletion of definitive gut endoderm in *Sox17*-null mutant
815 mice.
- 816 Kartha VK, Duarte FM, Hu Y, Ma S, Chew JG, Lareau CA, Earl A, Burkett ZD, Kohlway
817 AS, Lebofsky R, et al. 2022. Functional inference of gene regulation using single-cell multi-
818 omics. *Cell genomics*. **2**: 100166.
- 819 Korotkevich G, Sukhov V, Budin N, Shpak B, Artyomov MN, and Sergushichev A. 2021. Fast
820 gene set enrichment analysis. *BioRxiv*. 060012.
- 821 Li F, He Z, Li Y, Liu P, Chen F, Wang M, Zhu H, Ding X, Wangenstein KJ, Hu Y, et al. 2011.
822 Combined activin A/LiCl/Noggin treatment improves production of mouse embryonic stem
823 cell-derived definitive endoderm cells. *Journal of cellular biochemistry*. **112**: 1022–1034.
- 824 Lin Y, Cao Y, Kim HJ, Salim A, Speed TP, Lin DM, Yang P, and Yang JYH. 2020. scClas-
825 sify: sample size estimation and multiscale classification of cells using single and multiple
826 reference. *Molecular systems biology*. **16**: e9389.
- 827 Lin Y, Wu TY, Wan S, Yang JY, Wong WH, and Wang Y. 2022. scJoint integrates atlas-scale
828 single-cell RNA-seq and ATAC-seq data with transfer learning. *Nature Biotechnology*. **40**:
829 703–710.
- 830 Liu J, Huang Y, Singh R, Vert JP, and Noble WS 2019. Jointly embedding multiple single-
831 cell omics measurements. In: *Algorithms in bioinformatics: International Workshop, WABI,*
832 *proceedings. WABI (Workshop)*. Vol. 143. NIH Public Access.
- 833 Liu Y, Kaneda R, Leja TW, Subkhankulova T, Tolmachov O, Minchiotti G, Schwartz RJ, Bara-
834 hona M, and Schneider MD. 2014. *Hhex* and *Cer1* mediate the *Sox17* pathway for cardiac
835 mesoderm formation in embryonic stem cells. *Stem cells*. **32**: 1515–1526.
- 836 Lun AT, McCarthy DJ, and Marioni JC. 2016. A step-by-step workflow for low-level analysis of
837 single-cell RNA-seq data with Bioconductor. *F1000Research*. **5**:
- 838 Ma S, Zhang B, LaFave LM, Earl AS, Chiang Z, Hu Y, Ding J, Brack A, Kartha VK, Tay T, et al.
839 2020. Chromatin potential identified by shared single-cell profiling of RNA and chromatin.
840 *Cell*. **183**: 1103–1116.

- 841 Mikawa T, Poh AM, Kelly KA, Ishii Y, and Reese DE. 2004. Induction and patterning of the
842 primitive streak, an organizing center of gastrulation in the amniote. *Developmental dynam-*
843 *ics: an official publication of the American Association of Anatomists.* **229**: 422–432.
- 844 Mimitou EP, Lareau CA, Chen KY, Zorzetto-Fernandes AL, Hao Y, Takeshima Y, Luo W,
845 Huang TS, Yeung BZ, Papalexi E, et al. 2021. Scalable, multimodal profiling of chromatin
846 accessibility, gene expression and protein levels in single cells. *Nature biotechnology.* **39**:
847 1246–1258.
- 848 Mittnenzweig M, Mayshar Y, Cheng S, Ben-Yair R, Hadas R, Rais Y, Chomsky E, Reines N,
849 Uzonyi A, Lumerman L, et al. 2021. A single-embryo, single-cell time-resolved model for
850 mouse gastrulation. *Cell.* **184**: 2825–2842.
- 851 Peng G, Suo S, Chen J, Chen W, Liu C, Yu F, Wang R, Chen S, Sun N, Cui G, et al. 2016. Spatial
852 transcriptome for the molecular annotation of lineage fates and cell identity in mid-gastrula
853 mouse embryo. *Developmental cell.* **36**: 681–697.
- 854 Pijuan-Sala B, Griffiths JA, Guibentif C, Hiscock TW, Jawaid W, Calero-Nieto FJ, Mulas C,
855 Ibarra-Soria X, Tyser RC, Ho DLL, et al. 2019. A single-cell molecular map of mouse gas-
856 trulation and early organogenesis. *Nature.* **566**: 490–495.
- 857 Pimton P, Lecht S, Stabler CT, Johannes G, Schulman ES, and Lelkes PI. 2015. Hypoxia en-
858 hances differentiation of mouse embryonic stem cells into definitive endoderm and distal
859 lung cells. *Stem cells and development.* **24**: 663–676.
- 860 Plongthongkum N, Diep D, Chen S, Lake BB, and Zhang K. 2021. Scalable dual-omics profiling
861 with single-nucleus chromatin accessibility and mRNA expression sequencing 2 (SNARE-
862 Seq2). *Nature Protocols.* **16**: 4992–5029.
- 863 R Core Team 2022. R: A Language and Environment for Statistical Computing. R Foundation
864 for Statistical Computing. Vienna, Austria.
- 865 Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, and Smyth GK. 2015. limma powers
866 differential expression analyses for RNA-sequencing and microarray studies. *Nucleic acids*
867 *research.* **43**: e47–e47.
- 868 Rubin AJ, Parker KR, Satpathy AT, Qi Y, Wu B, Ong AJ, Mumbach MR, Ji AL, Kim DS, Cho
869 SW, et al. 2019. Coupled single-cell CRISPR screening and epigenomic profiling reveals
870 causal gene regulatory networks. *Cell.* **176**: 361–376.
- 871 Schiebinger G, Shu J, Tabaka M, Cleary B, Subramanian V, Solomon A, Gould J, Liu S, Lin
872 S, Berube P, et al. 2019. Optimal-transport analysis of single-cell gene expression identifies
873 developmental trajectories in reprogramming. *Cell.* **176**: 928–943.

- 874 Stuart T, Srivastava A, Madad S, Lareau CA, and Satija R. 2021. Single-cell chromatin state
875 analysis with Signac. *Nature methods*. **18**: 1333–1341.
- 876 Sutherland MJ, Wang S, Quinn ME, Haaning A, and Ware SM. 2013. Zic3 is required in the mi-
877 grating primitive streak for node morphogenesis and left–right patterning. *Human molecular*
878 *genetics*. **22**: 1913–1923.
- 879 Svensson V, Gayoso A, Yosef N, and Pachter L. 2020. Interpretable factor models of single-cell
880 RNA-seq via variational autoencoders. *Bioinformatics*. **36**: 3418–3421.
- 881 Swanson E, Lord C, Reading J, Heubeck AT, Genge PC, Thomson Z, Weiss MD, Li Xj, Savage
882 AK, Green RR, et al. 2021. Simultaneous trimodal single-cell measurement of transcripts,
883 epitopes, and chromatin accessibility using TEA-seq. *Elife*. **10**: e63632.
- 884 Tam PP and Loebel DA. 2007. Gene function in mouse embryogenesis: get set for gastrulation.
885 *Nature Reviews Genetics*. **8**: 368–381.
- 886 Traag VA, Waltman L, and Van Eck NJ. 2019. From Louvain to Leiden: guaranteeing well-
887 connected communities. *Scientific reports*. **9**: 1–12.
- 888 Tran A, Yang P, Yang JY, and Ormerod JT. 2022. scREMOTE: Using multimodal single cell
889 data to predict regulatory gene relationships and to build a computational cell reprogramming
890 model. *NAR genomics and bioinformatics*. **4**: lqac023.
- 891 Tritschler S, Büttner M, Fischer DS, Lange M, Bergen V, Lickert H, and Theis FJ. 2019. Con-
892 cepts and limitations for learning developmental trajectories from single cell genomics. *De-*
893 *velopment*. **146**: dev170506.
- 894 VanOudenhove JJ, Medina R, Ghule PN, Lian JB, Stein JL, Zaidi SK, and Stein GS. 2016.
895 Transient RUNX1 expression during early mesendodermal differentiation of hESCs promotes
896 epithelial to mesenchymal transition through TGFB2 signaling. *Stem Cell Reports*. **7**: 884–
897 896.
- 898 Wang X and Yang P. 2008. In vitro differentiation of mouse embryonic stem (mES) cells using
899 the hanging drop method. *JoVE (Journal of Visualized Experiments)*. e825.
- 900 Wang Y, Yuan P, Yan Z, Yang M, Huo Y, Nie Y, Zhu X, Qiao J, and Yan L. 2021. Single-cell
901 multiomics sequencing reveals the functional regulatory landscape of early embryos. *Nature*
902 *communications*. **12**: 1–14.
- 903 Weirauch MT, Yang A, Albu M, Cote AG, Montenegro-Montero A, Drewe P, Najafabadi HS,
904 Lambert SA, Mann I, Cook K, et al. 2014. Determination and inference of eukaryotic tran-
905 scription factor sequence specificity. *Cell*. **158**: 1431–1443.

- 906 Yang P, Huang H, and Liu C. 2021. Feature selection revisited in the single-cell era. *Genome*
907 *Biology*. **22**: 1–17.
- 908 Zhang L, Zhang J, and Nie Q. 2022. DIRECT-NET: An efficient method to discover cis-regulatory
909 elements and construct regulatory networks from single-cell multiomics data. *Science Ad-*
910 *vances*. **8**: eabl7393.
- 911 Zhang Z, Yang C, and Zhang X. 2022. scDART: integrating unmatched scRNA-seq and scATAC-
912 seq data and learning cross-modality relationship simultaneously. *Genome Biology*. **23**: 1–28.

913 **Figure legends**

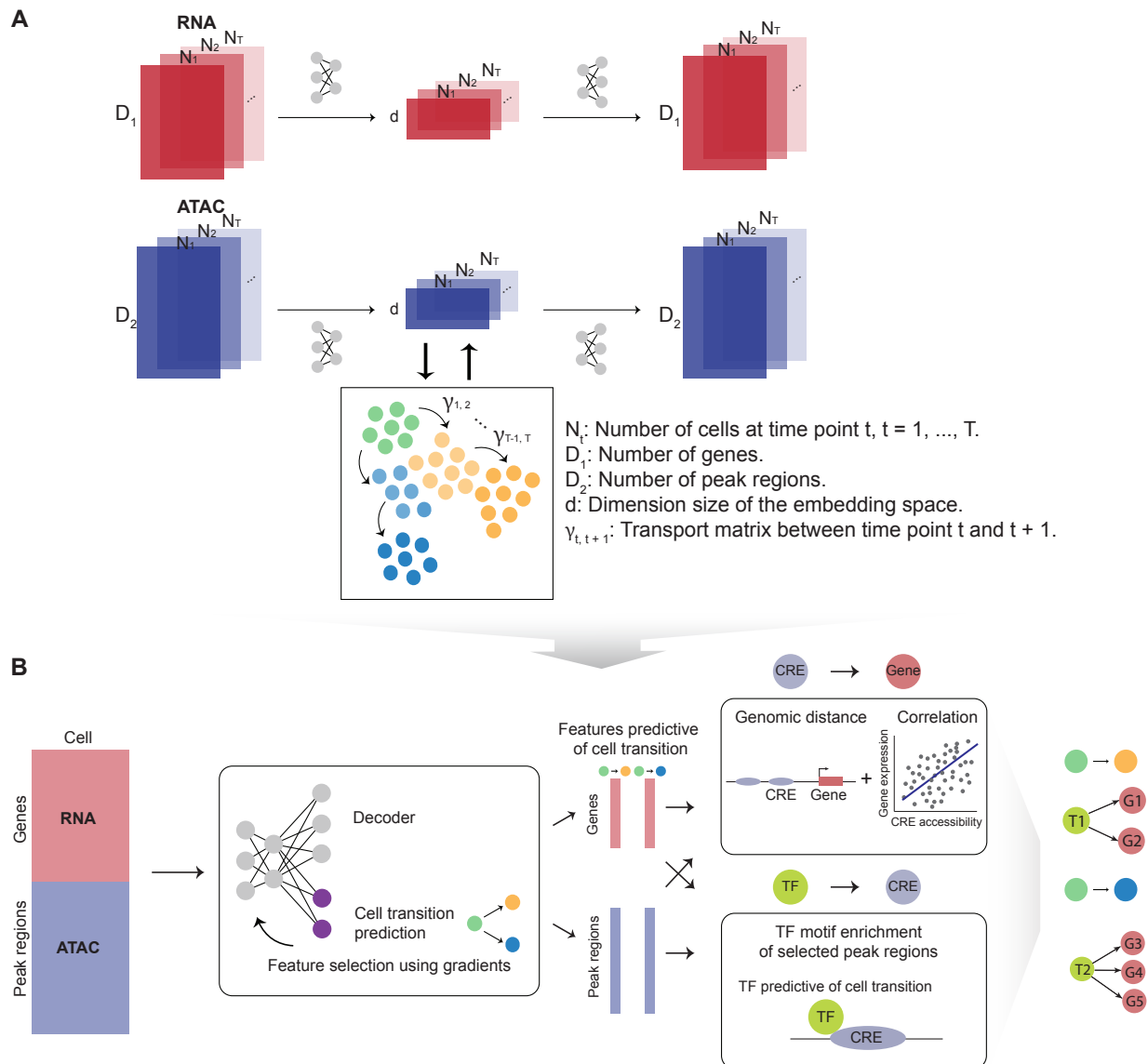


Figure 1: Overview of scTIE, a unified framework for the integration of temporal data and the inference of context-specific GRNs that predict cell fates. The input of scTIE consists of the gene expression matrix of scRNA-seq and peak matrix of scATAC-seq from single-cell multiome data over a time course. scTIE consists of two main steps: (A) In the first step, each cell, represented by a pair of gene and peak vectors, is projected into a common embedding space by separate encoders and decoders. The two modalities and time points are aligned by appropriate loss functions, while the transition probability matrix between cells from consecutive time points are iteratively estimated; (B) In the second step, users have the ability to select specific subgroups of cells whose transitions are of interest, fine-tune the previously trained neural network, identify features that are predictive of transition probabilities, and construct the corresponding GRNs.



Figure 2: Performance benchmarking for integrating temporal multimodal data. (A) Joint visualization using UMAP of the synthetic dataset with batch effect in RNA and noise in ATAC, colored by cell type annotations (first row), sampling days (second row) and synthetic batch information (third row). Each dot represents a cell in the embedding space. (B) Bar plots showing the evaluation metrics of different data integration methods, including ARI values for clustering with annotations (left); 1 - average purity scores of sampling days with the number of neighbors equal to 50 (middle) and 1 - average purity scores of the synthetic batch with the number of neighbors equal to 50 (right). Higher values indicate better agreement with annotations and mixing of batches/days. (C) Radar plot summarizing the three evaluation metrics shown in (B), where each line represents the performance of one method, and each axis represents an evaluation metric, starting from the minimum value of all methods. It is noted that scAI was not included in this benchmarking due to its long computational time (> 2 days).

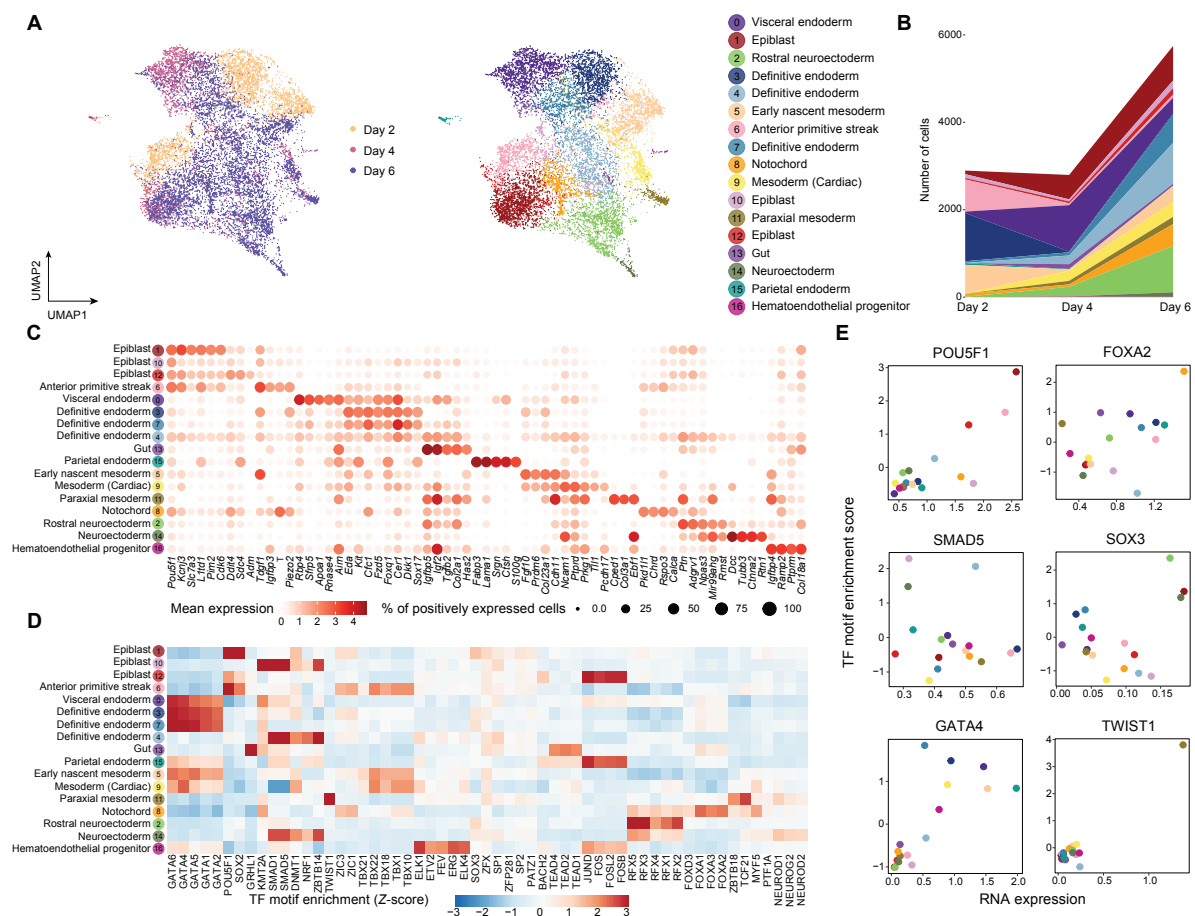


Figure 3: Integration and cell type identification of the mESC dataset by scTIE. (A) Joint visualization of the mESC dataset using UMAP, colored by sampling day and cell type annotations. Each dot represents a cell in the embedding space. (B) Cell type compositions per time point. (C) Dot plots of mean expression of RNA data. Rows represent cell types and columns indicate each gene. The color scale represents the expression level, and the size indicates proportion of positively expressed cells. The five most significantly expressed genes for each cluster are included. (D) Heatmap of the TF motif enrichment (Z-scores) of ATAC data. Rows represent cell types and columns indicate TFs. The five most significantly enriched TFs for each cluster are included. (E) Scatter plots of the mean RNA expression levels by clusters (x-axis) and the average TF motif enrichment scores of ATAC (y-axis) for the selected TFs. The dots are colored by the cell type annotations, with color legend consistent with Fig. 3A.

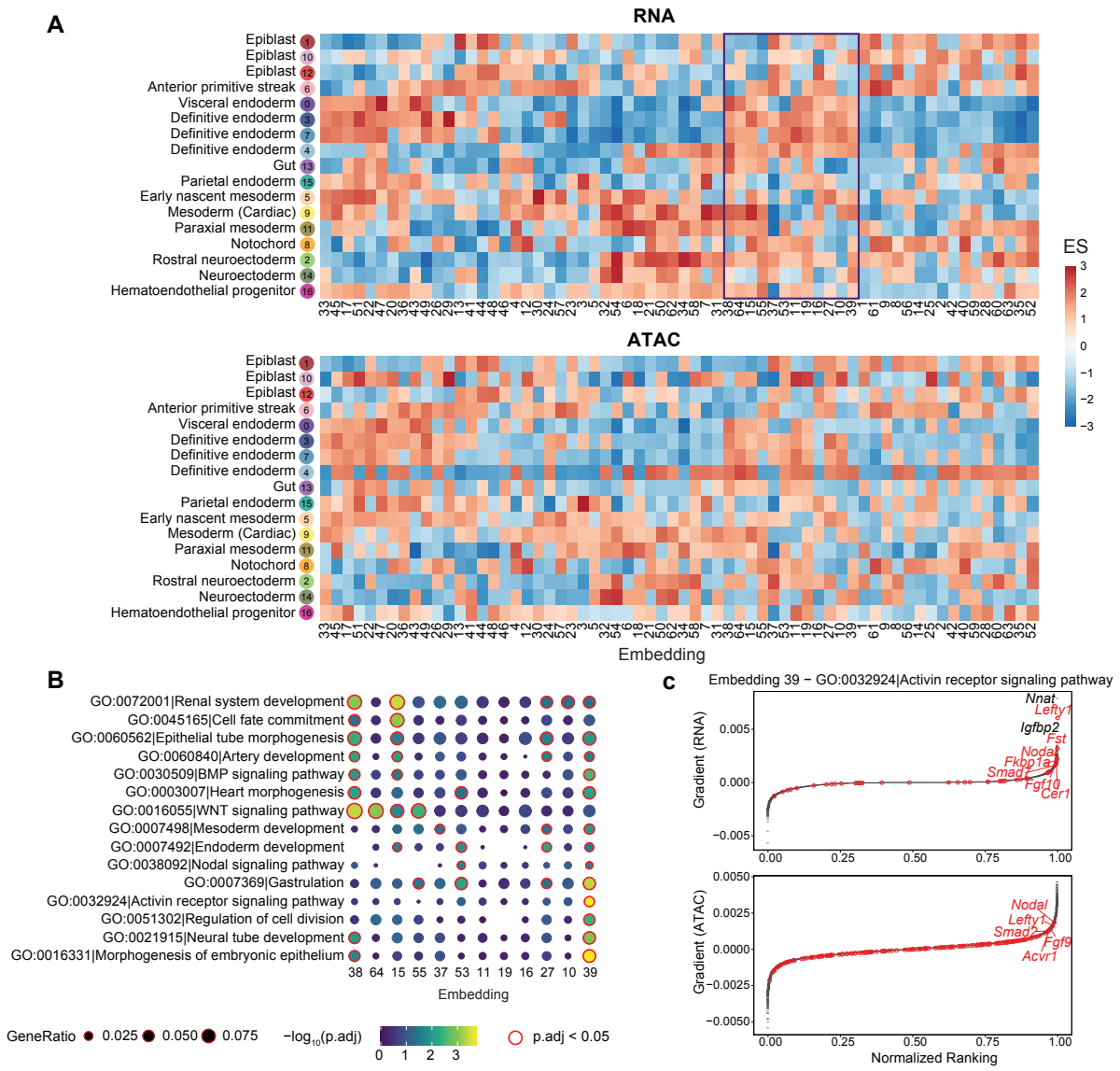


Figure 4: Biological signals in the mESC dataset captured by each embedding dimension of scTIE. (A) Enrichment scores of the gradient ranking in each embedding dimension using the RNA (top panel) and ATAC (bottom panel) marker list for each cell type. (B) Gene Ontology enrichment of selected pathways on the gradient ranking of a subset of embedding dimensions. (C) Gradient rankings for RNA (top panel) and ATAC (bottom panel) of embedding dimension 39, where genes/peaks are ranked based on the gradient values. The labeled points are genes in the selected gene set (Activin receptor signaling pathway).

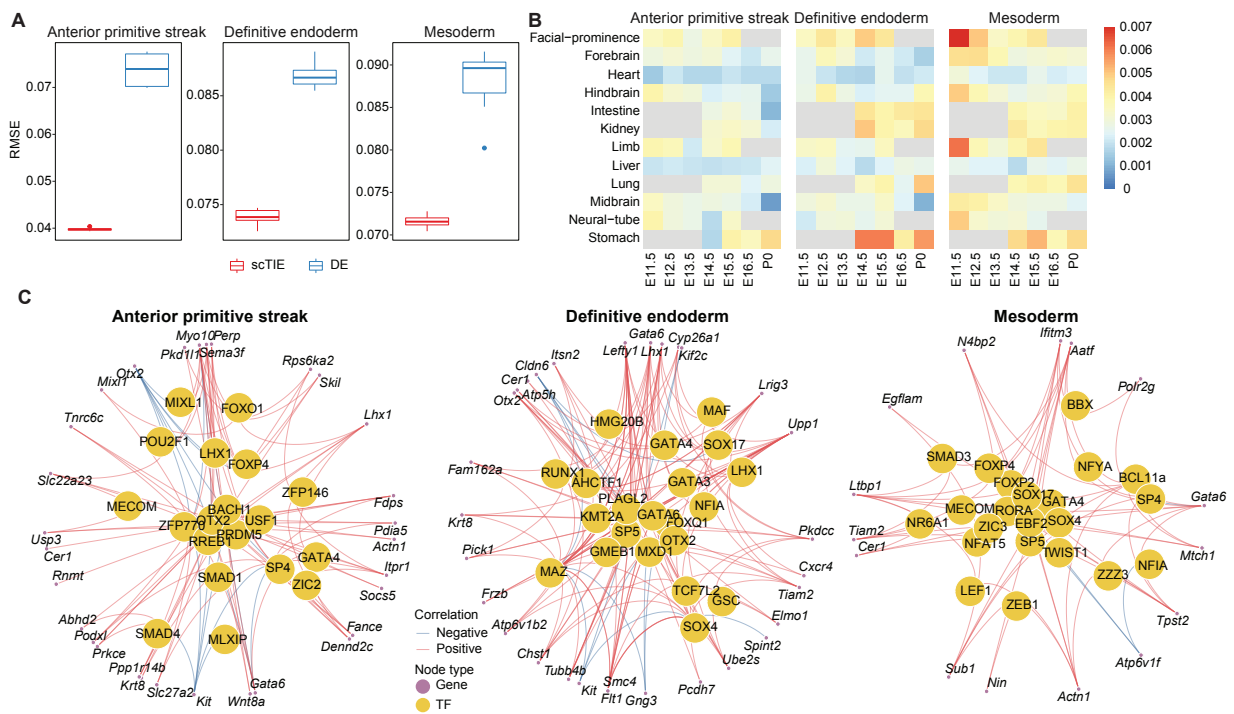
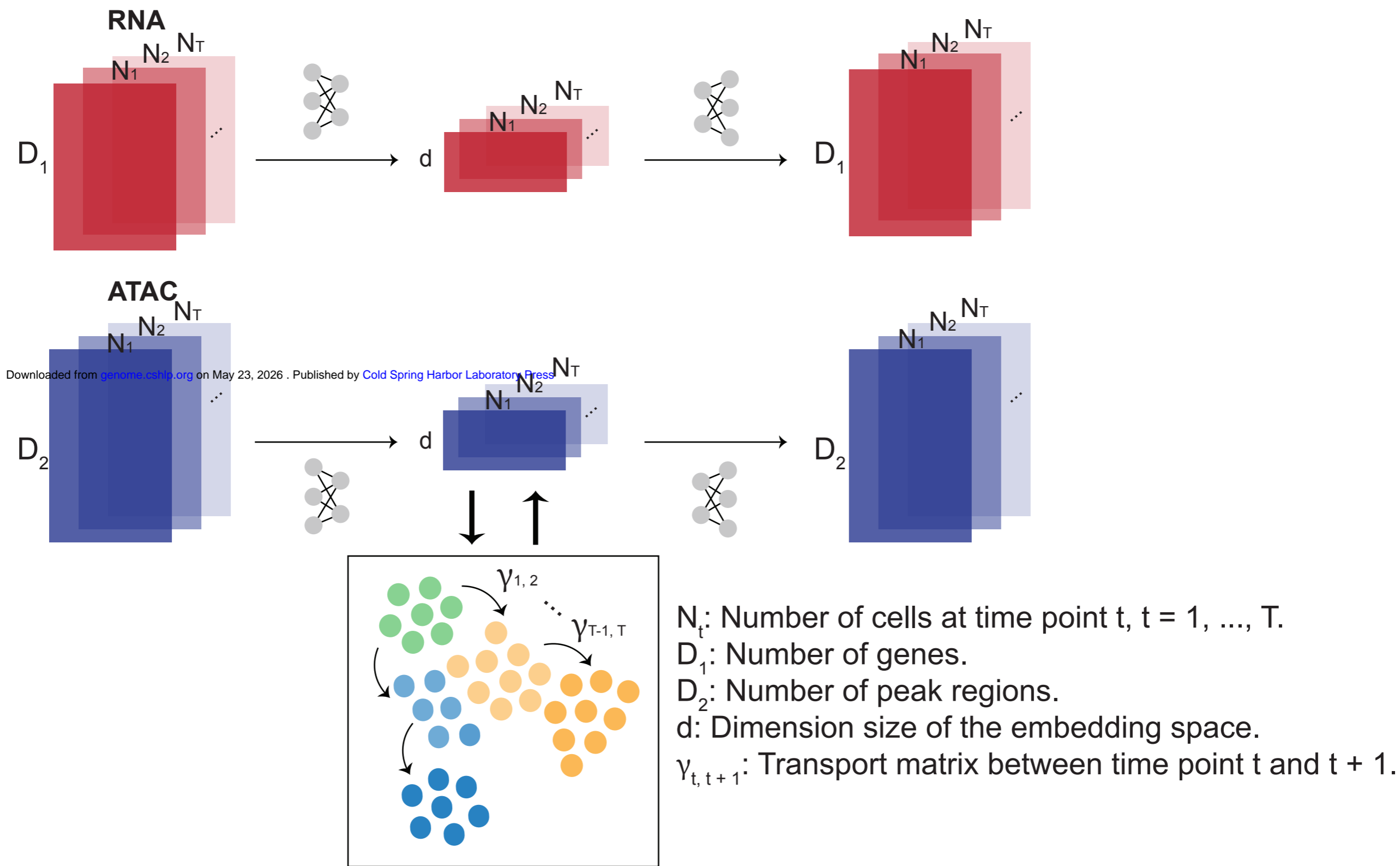
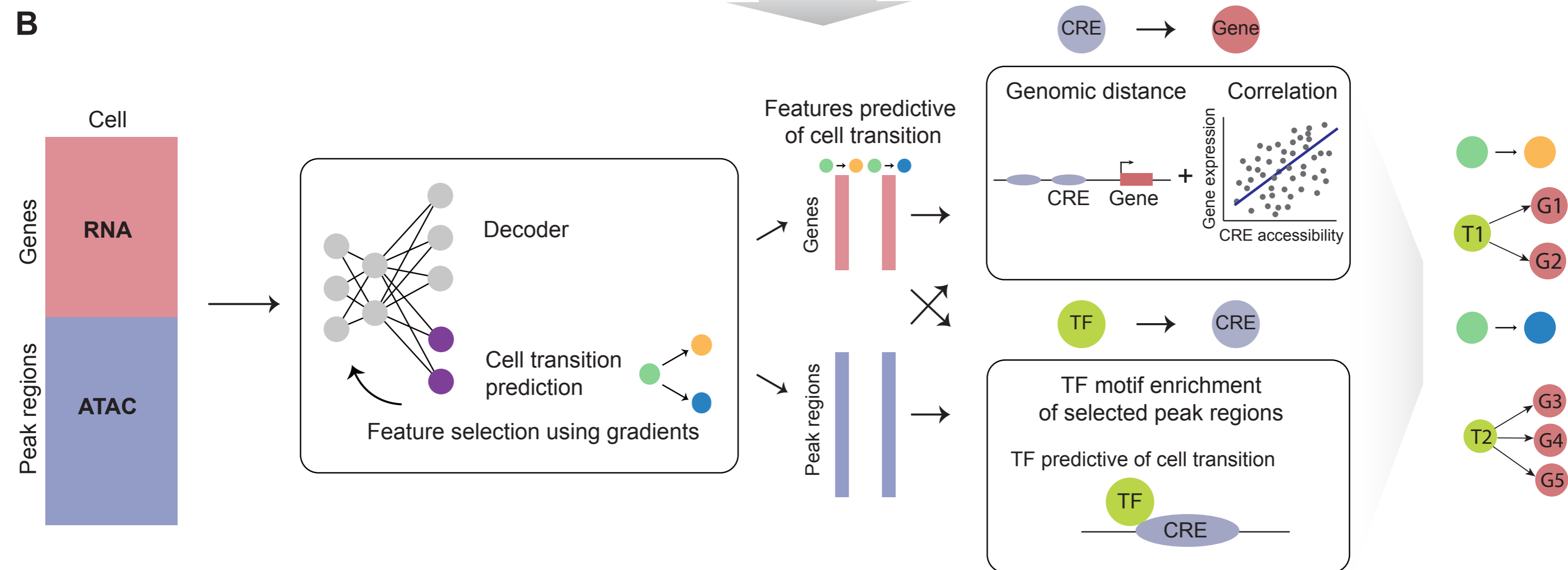
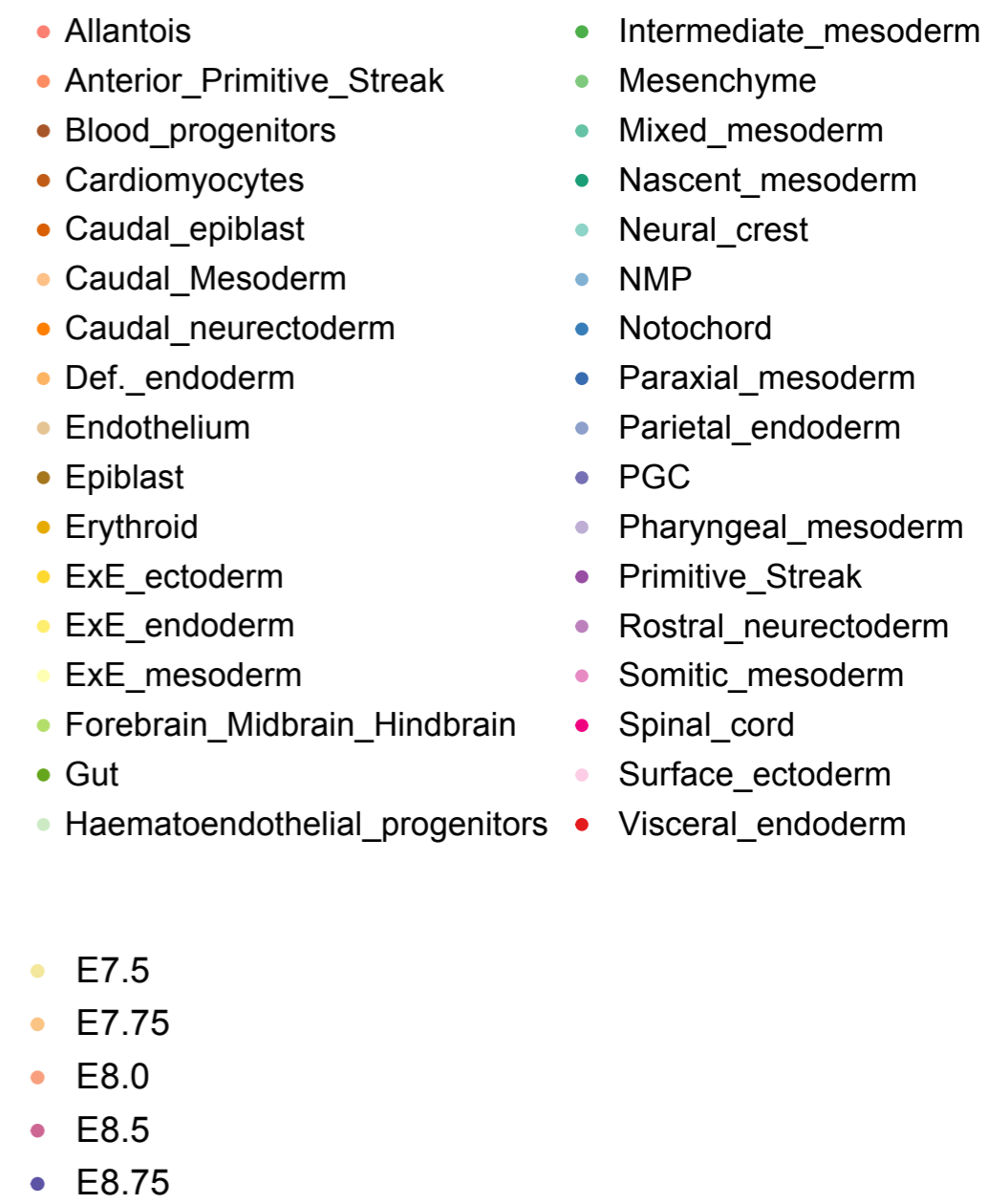
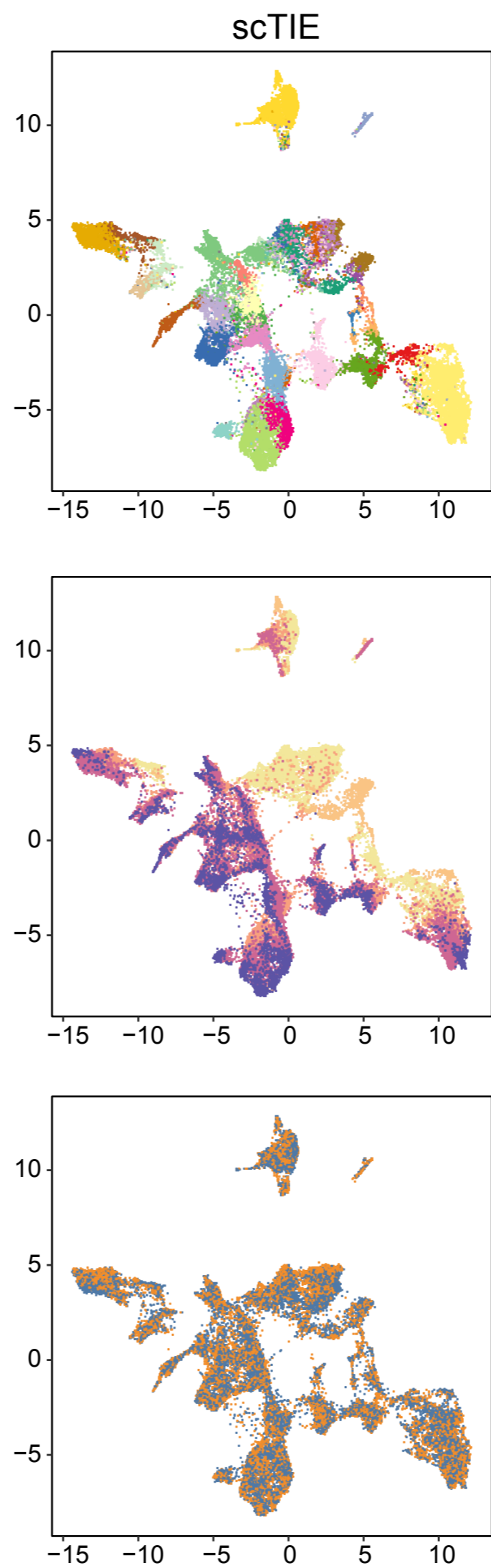
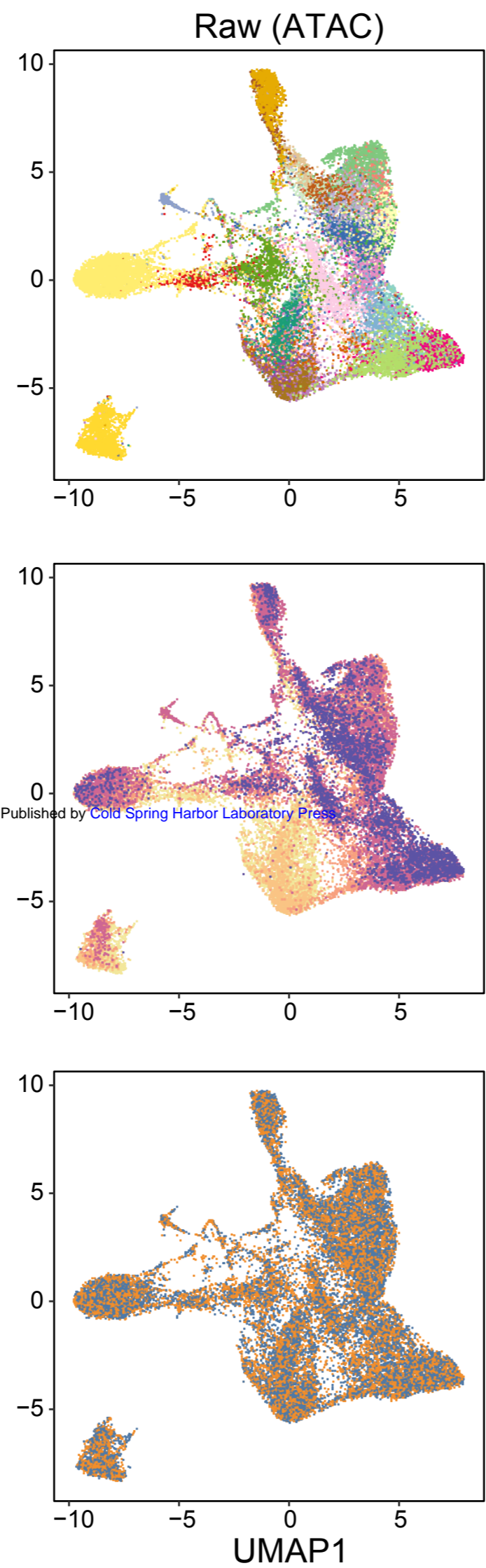
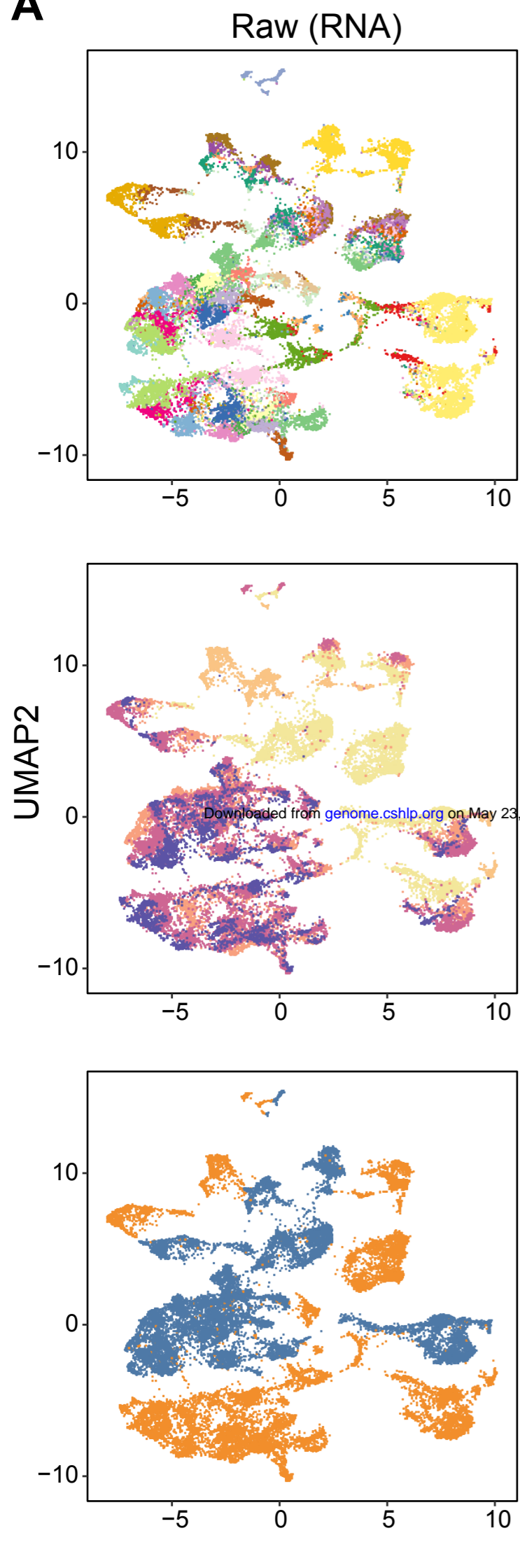


Figure 5: Lineage-specific regulatory elements selected by scTIE and the corresponding GRNs. (A) Performance of top genes and peaks selected by each method in predicting cell fate probabilities. (B) Similarity of top gradient peaks with enhancers of 12 tissues at seven developmental stages from known enhancer databases. (C) GRN of three cell fates.

Encoder	Encoder
Batchnorm (26717)	Batchnorm (61744)
Linear (26717, 1000)	Linear (61744, 1000)
Batchnorm (1000)	Batchnorm (1000)
LeakyReLU (0.2)	LeakyReLU (0.2)
Linear (1000, 1000)	Linear (1000, 1000)
Batchnorm (1000)	Batchnorm (1000)
LeakyReLU (0.2)	LeakyReLU (0.2)
Linear (1000, 64)	Linear (1000, 64)
Decoder	Decoder
Linear (64, 500)	Linear (64, 500)
Batchnorm (500)	Batchnorm (500)
LeakyReLU (0.2)	LeakyReLU (0.2)
Linear (500, 1000)	Linear (500, 1000)
Batchnorm (1000)	Batchnorm (1000)
LeakyReLU (0.2)	LeakyReLU (0.2)
Linear (1000, 26717)	Linear (1000, 61744)
Batchnorm (26717)	
Multiply by σ and add μ	

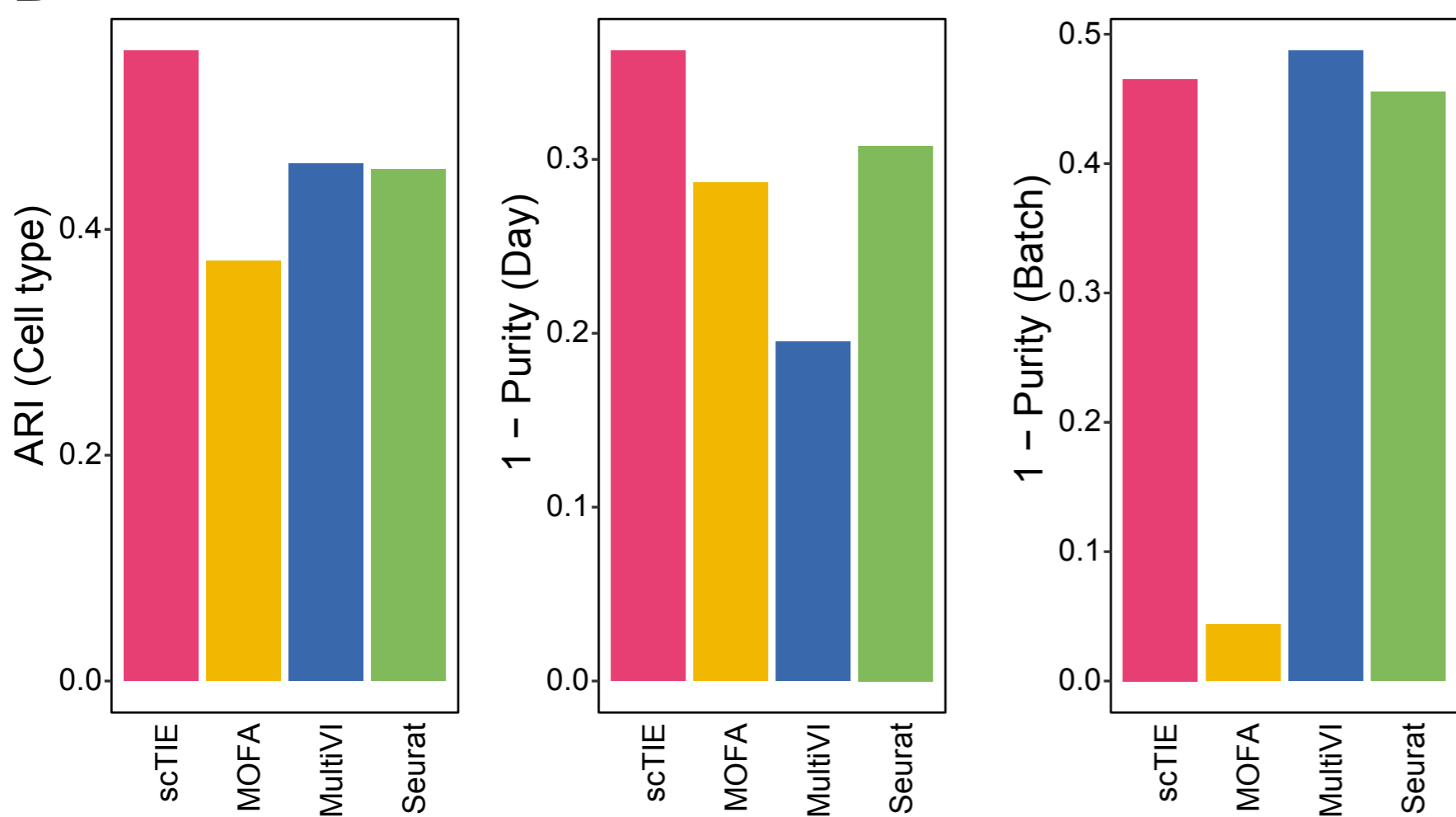
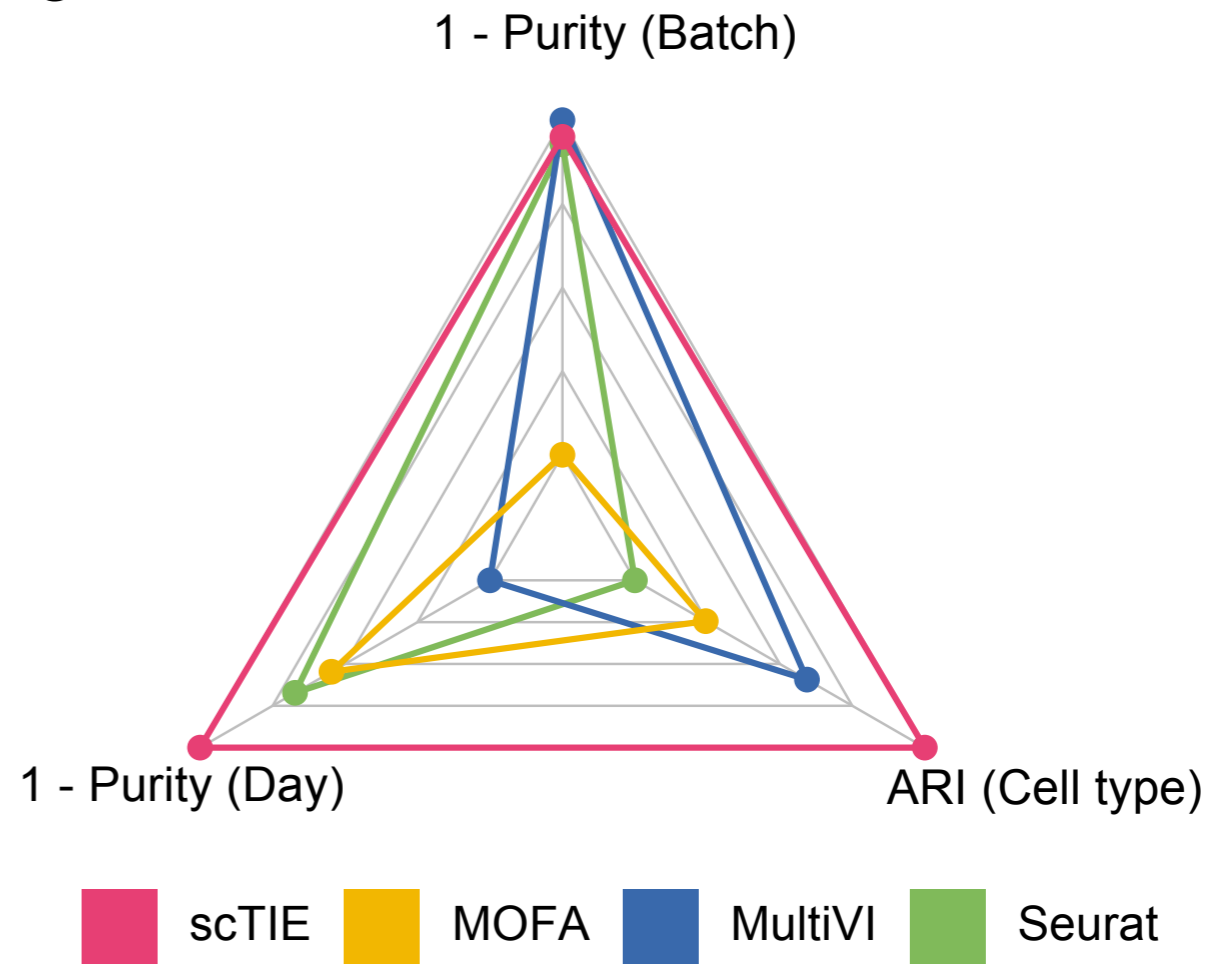
Table 1: Autoencoder architecture for RNA (left) and ATAC (right).

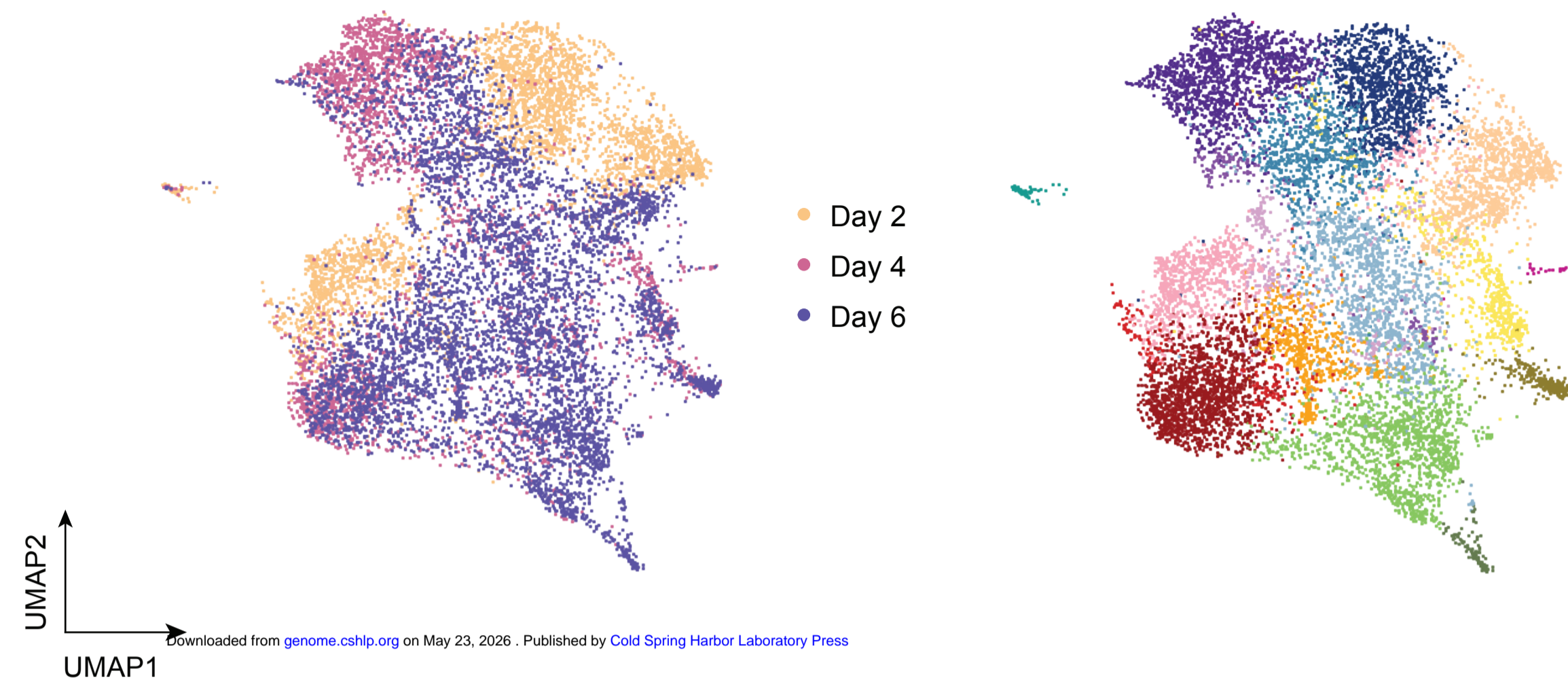
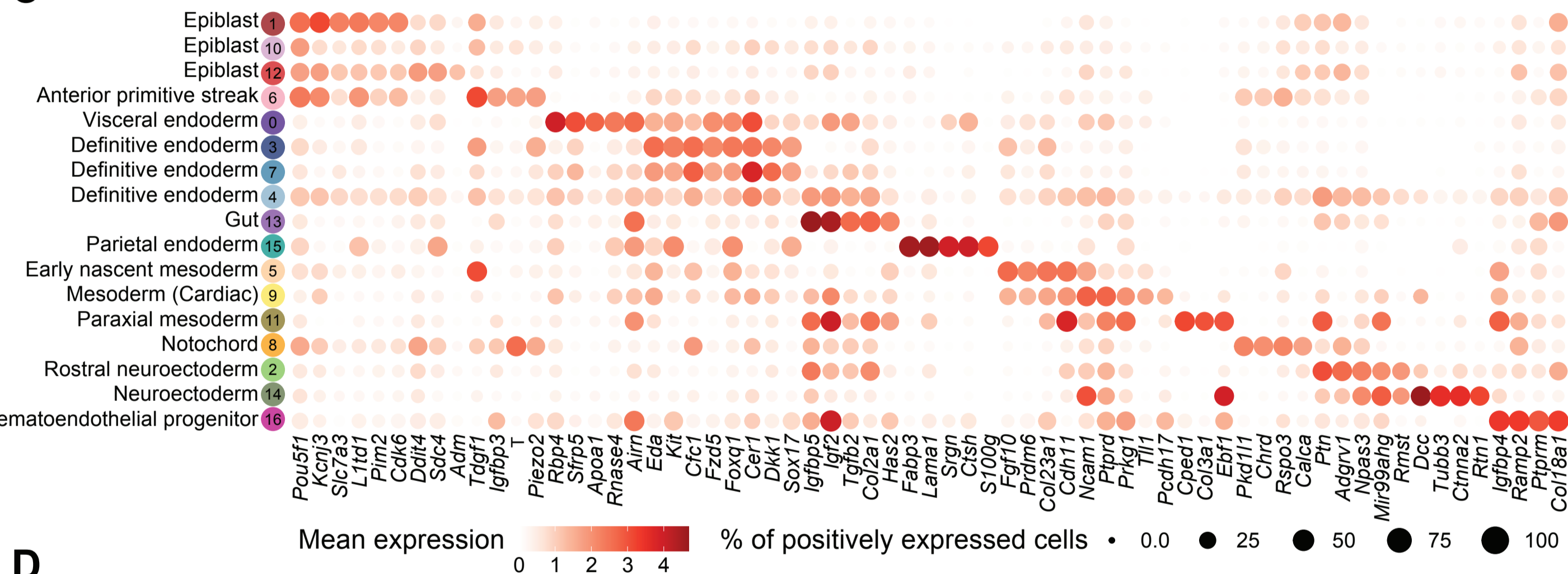
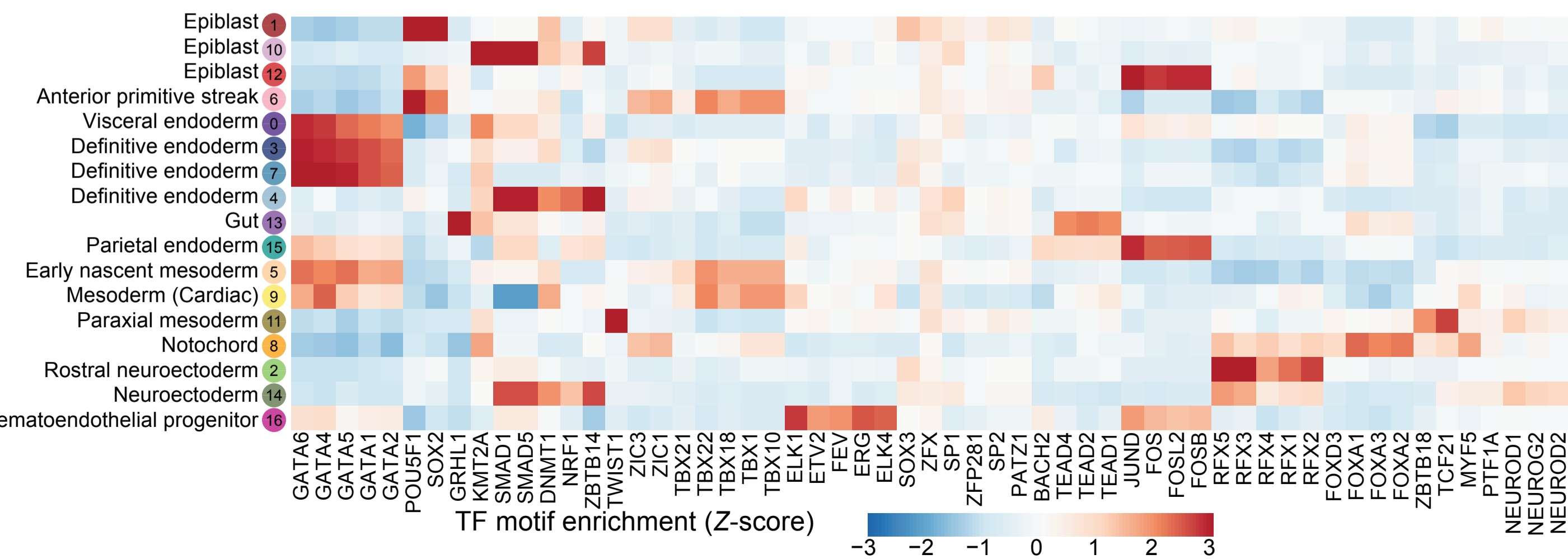
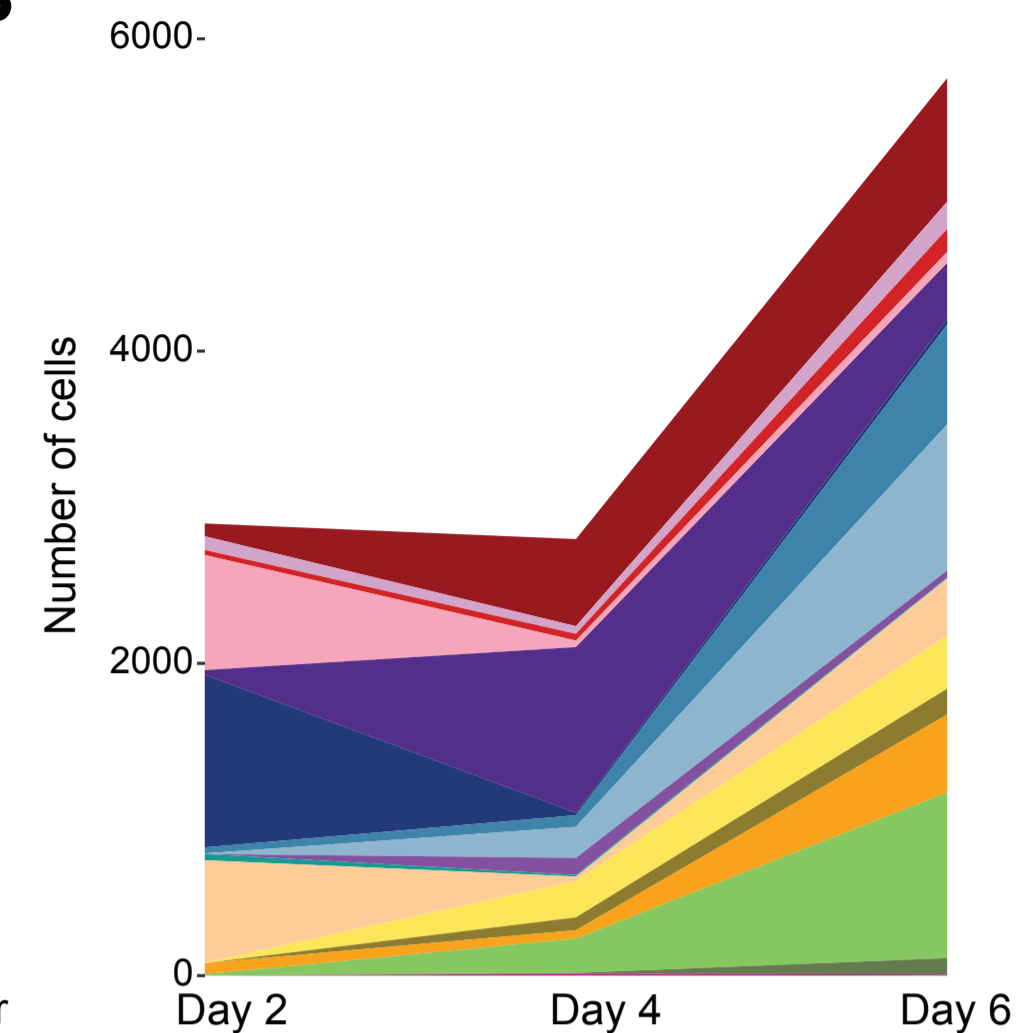
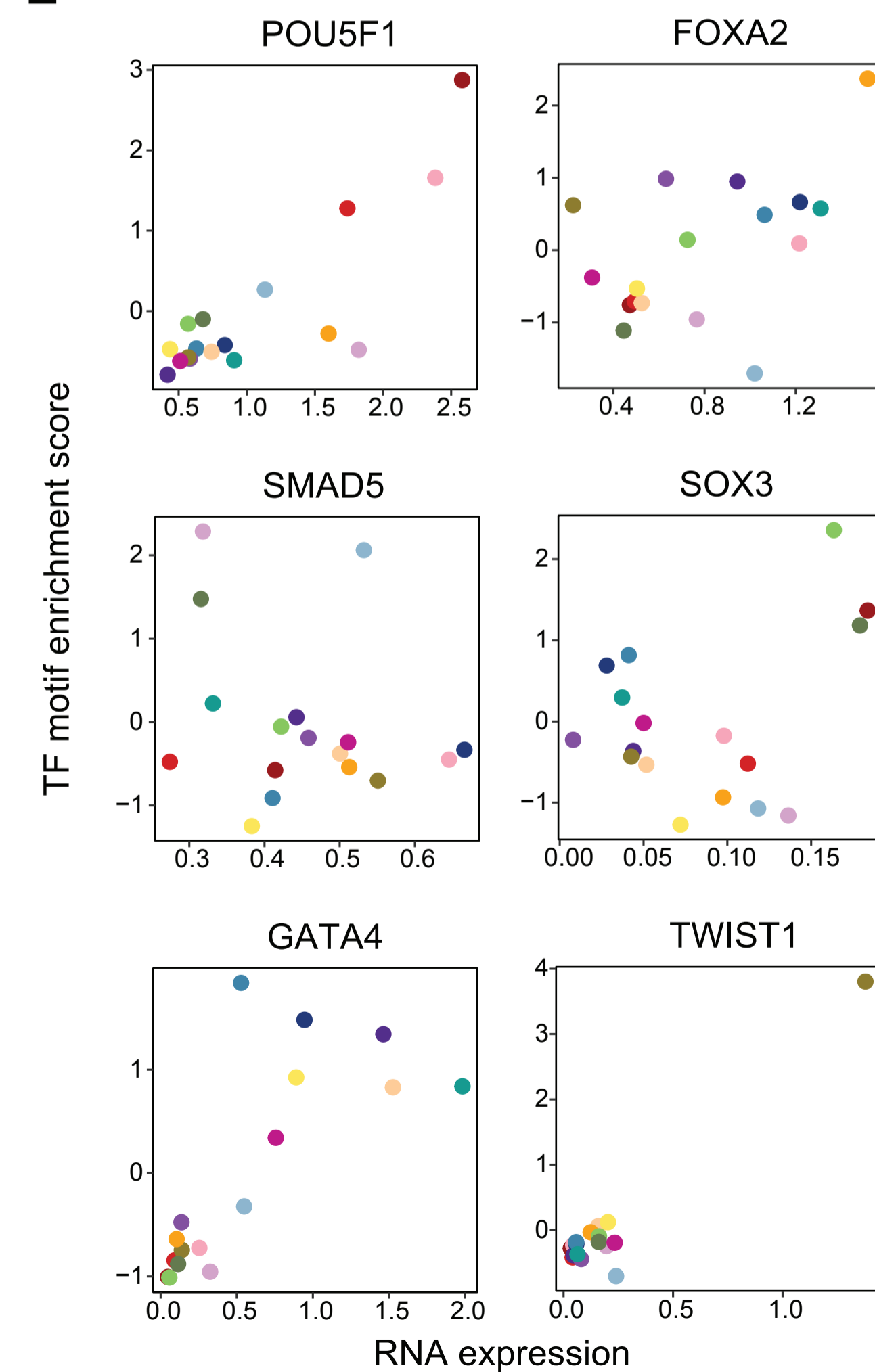
A**B**

A

UMAP2

● Batch 1
● Batch 2

B**C**

A**C****D****B****E**

- 0 Visceral endoderm
- 1 Epiblast
- 2 Rostral neuroectoderm
- 3 Definitive endoderm
- 4 Definitive endoderm
- 5 Early nascent mesoderm
- 6 Anterior primitive streak
- 7 Definitive endoderm
- 8 Notochord
- 9 Mesoderm (Cardiac)
- 10 Epiblast
- 11 Paraxial mesoderm
- 12 Epiblast
- 13 Gut
- 14 Neuroectoderm
- 15 Parietal endoderm
- 16 Hematoendothelial progenitor

