

1 Multi-omics integration and regulatory inference for unpaired single-
2 cell data with a graph-linked unified embedding framework

3 Zhi-Jie Cao¹, Ge Gao^{1,*}

4 ¹ Biomedical Pioneering Innovation Center (BIOPIC), Beijing Advanced Innovation Center for Genomics
5 (ICG), Center for Bioinformatics (CBI), and State Key Laboratory of Protein and Plant Gene Research at
6 School of Life Sciences, Peking University, Beijing, 100871, China

7 * To whom correspondence should be addressed. Tel: +86-010-62755206; Email: gaog@mail.cbi.pku.edu.cn

1 Abstract

2 With the ever-increasing amount of single-cell multi-omics data accumulated during the past years,
3 effective and efficient computational integration is becoming a serious challenge. One major obstacle
4 of unpaired multi-omics integration is the feature discrepancies among omics layers. Here, we
5 propose a computational framework called GLUE (graph-linked unified embedding), which utilizes
6 accessible prior knowledge about regulatory interactions to bridge the gaps between feature spaces.
7 Systematic benchmarks demonstrated that GLUE is accurate, robust and scalable. We further
8 employed GLUE for various challenging tasks, including triple-omics integration, model-based
9 regulatory inference and multi-omics human cell atlas construction (over millions of cells) and found
10 that GLUE achieved superior performance for each task. As a generalizable framework, GLUE
11 features a modular design that can be flexibly extended and enhanced for new analysis tasks. The full
12 package is available online at <https://github.com/gao-lab/GLUE> for the community.

1 Introduction

2 Recent technological advances in single-cell sequencing have enabled the probing of regulatory
3 maps through multiple omics layers, such as chromatin accessibility (scATAC-seq^{1,2}), DNA
4 methylation (snmC-seq³, sci-MET⁴) and the transcriptome (scRNA-seq^{5,6}), offering a unique
5 opportunity to unveil the underlying regulatory bases for the functionalities of diverse cell types⁷.
6 While simultaneous assays are emerging recently⁸⁻¹¹, different omics are usually measured
7 independently and produce unpaired data, which calls for effective and efficient *in silico* multi-omics
8 integration^{12,13}.

9
10 Computationally, one major obstacle faced when integrating unpaired multi-omics data is the distinct
11 feature spaces of different modalities (e.g., accessible chromatin regions in scATAC-seq vs. genes in
12 scRNA-seq)¹⁴. A quick fix is to convert multimodality data into one common feature space based on
13 prior information and apply single-omics data integration methods¹⁵⁻¹⁷. Such explicit “feature
14 conversion” is straightforward, but has been reported to result in significant information loss¹⁸.
15 Algorithms based on coupled matrix factorization circumvent explicit conversion but hardly handle
16 more than two omics layers^{19,20}. An alternative option is to match cells from different omics layers
17 via nonlinear manifold alignment, which removes the requirement of prior knowledge completely
18 and could reduce inter-modality information loss in theory^{21,22}; however, this technique has mostly
19 been applied to continuous, trajectory-like manifolds rather than atlases.

20
21 The ever-increasing volume of data is another serious challenge²³. Recently developed technologies
22 can routinely generate datasets at the scale of millions of cells²⁴⁻²⁶, whereas current integration
23 methods have only been applied to datasets with much smaller volumes^{15,17,19-22}. To catch up with
24 the growth in data throughput, computational integration methods should be designed with
25 scalability in mind.

26
27 Hereby, we introduce GLUE (graph-linked unified embedding), a modular framework for integrating
28 unpaired single-cell multi-omics data and inferring regulatory interactions simultaneously. By
29 modeling the regulatory interactions across omics layers explicitly, GLUE bridges the gaps between
30 various omics-specific feature spaces in a biologically intuitive manner. Systematic benchmarks and
31 case studies demonstrate that GLUE is accurate, robust and scalable for heterogeneous single-cell
32 multi-omics data. Furthermore, GLUE is designed as a generalizable framework that allows for easy

1 extension and quick adoption to particular scenarios in a modular manner. GLUE is publicly
 2 accessible at <https://github.com/gao-lab/GLUE>.

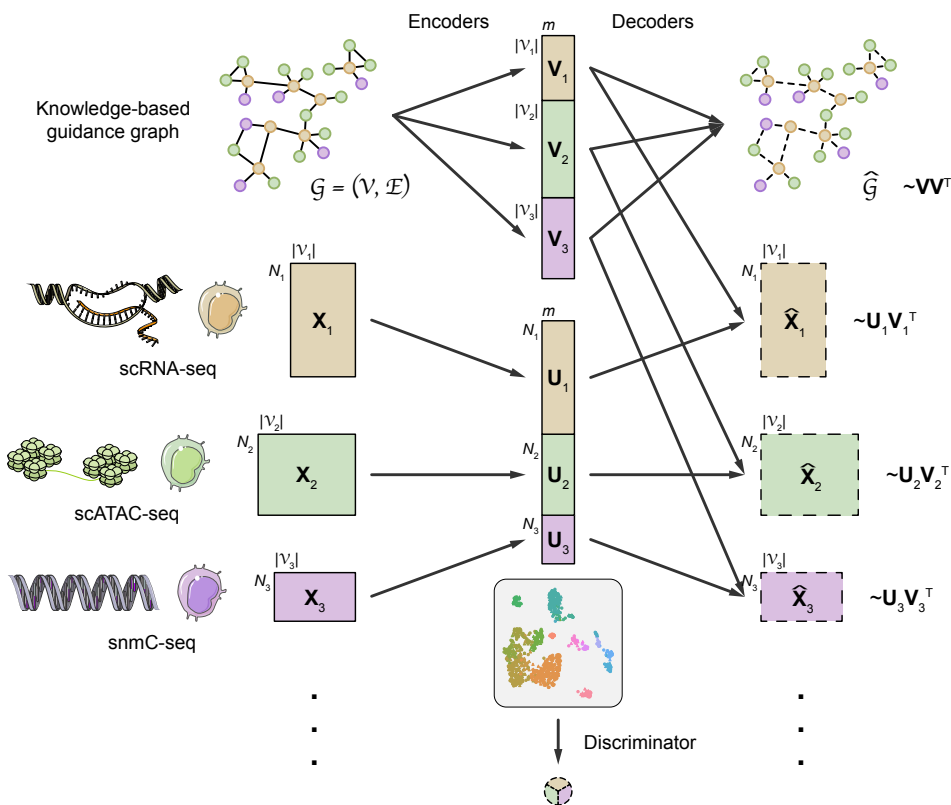
3

4 Results

5 Integrating unpaired single-cell multi-omics data via graph-guided embeddings

6 Inspired by previous works, we model cell states as low-dimensional cell embeddings learned
 7 through variational autoencoders^{27, 28}. Given their intrinsic differences in biological nature and assay
 8 technology, each omics layer is equipped with a separate autoencoder that uses a probabilistic
 9 generative model tailored to the layer-specific feature space (Fig. 1, Methods).

10



11

12 Fig. 1 Architecture of the GLUE framework.

13 GLUE employs omics-specific variational autoencoders to learn low-dimensional cell embeddings from each omics
 14 layer. The data dimensionality and generative distribution can differ across omics layers, but the cell embedding
 15 dimensions are shared. A graph variational autoencoder is used to learn feature embeddings from the prior
 16 knowledge-based guidance graph; these embeddings are then used as data decoder parameters. The feature
 17 embeddings effectively link the omics-specific autoencoders to ensure a consistent embedding orientation. Last, an
 18 omics discriminator is employed to align the cell embeddings of different omics layers via adversarial learning.

1 Taking advantage of prior biological knowledge, we propose the use of a knowledge-based graph
2 (“guidance graph”) that explicitly models cross-layer regulatory interactions for linking layer-
3 specific feature spaces; the vertices in the graph correspond to the features of different omics layers,
4 and edges represent signed regulatory interactions. For example, when integrating scRNA-seq and
5 scATAC-seq data, the vertices are genes and accessible chromatin regions (i.e., ATAC peaks), and a
6 positive edge can be connected between an accessible region and its putative downstream gene.
7 Then, adversarial multimodal alignment is performed as an iterative optimization procedure, guided
8 by feature embeddings encoded from the graph²⁹ (Fig. 1, Methods). Notably, when the iterative
9 process converges, the graph can be refined with inputs from the alignment procedure and used for
10 data-oriented regulatory inference (see below for more details).

11

12 **Systematic benchmarks demonstrate superior alignment accuracy and robustness over existing** 13 **methods**

14 We first benchmarked GLUE against multiple popular unpaired multi-omics integration methods^{15-17,}
15 ^{21, 30} using gold-standard datasets generated by recent simultaneous scRNA-seq and scATAC-seq
16 technologies^{8, 9, 31}.

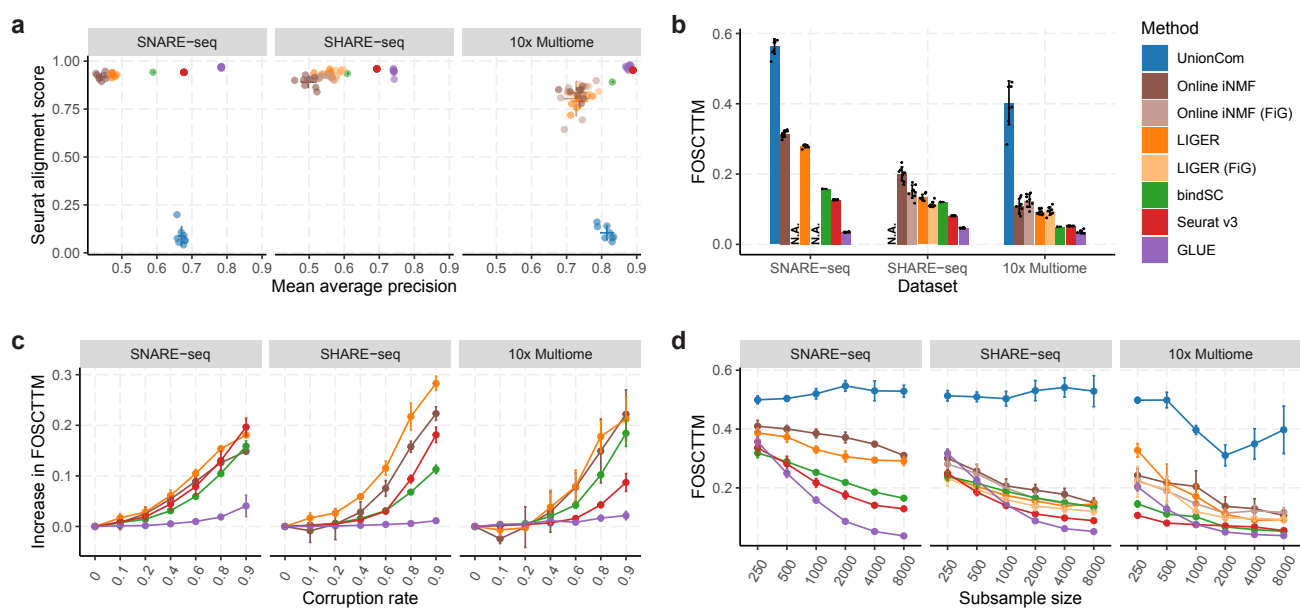
17

18 At the cell type level, an integration method should match the corresponding cell types from
19 different omics layers, producing cell embeddings where the cell types are clearly distinguishable
20 and the omics layers are well mixed. Compared to other methods, GLUE achieved overall the
21 highest cell type resolution (as quantified by mean average precision) and layer mixing (as quantified
22 by the Seurat alignment score³²) simultaneously (Fig. 2a); these results were also validated by
23 UMAP visualization of the aligned cell embeddings (Supplementary Fig. 1-3).

24

25 An optimal integration method should produce accurate alignments not only at the cell type level but
26 also at finer scales. Exploiting the ground truth cell-to-cell correspondence between scRNA-seq and
27 scATAC-seq, we further quantified single-cell level alignment error via the FOSCTTM (fraction of
28 samples closer than the true match) metric³³. On all three datasets, GLUE achieved the lowest
29 FOSCTTM, decreasing the alignment error by large margins compared to the second-best method on
30 each dataset (Fig. 2b, the decreases were 3.6-fold for SNARE-seq, 1.7-fold for SHARE-seq, and 1.4-
31 fold for 10x Multiome).

32



1

2 **Fig. 2 Systematic benchmarks on gold-standard datasets.**

3 **a**, Cell type resolution (quantified by mean average precision) vs. omics layer mixing (quantified by Seurat alignment
 4 score) for different integration methods. FiG (fragments in genes) is an alternative feature conversion method
 5 recommended by online iNMF and LIGER (Methods). Online iNMF and LIGER could not run with FiG conversion
 6 on the SNARE-seq data because the raw ATAC fragment file was not available. UnionCom failed to run on the
 7 SHARE-seq dataset due to memory overflow. **b**, Single-cell level alignment error (quantified by FOSCTTM) of
 8 different integration methods. **c**, Increases in FOSCTTM at different prior knowledge corruption rates for integration
 9 methods that rely on prior feature interactions. **d**, FOSCTTM values of different integration methods on subsampled
 10 datasets. The error bars indicate mean \pm s.d.

11 During the evaluation described above, we adopted a standard schema (ATAC peaks were linked to
 12 RNA genes if they overlapped in the gene body or proximal promoter regions) to construct the
 13 guidance graph for GLUE and to perform feature conversion for other conversion-based methods.

14 Given that our current knowledge about the regulatory interactions is still far from perfect, a useful
 15 integration method must be robust to such inaccuracies. Thus, we further assessed the methods'
 16 robustness to corruption of regulatory interactions by randomly replacing varying fractions of
 17 existing interactions with nonexistent ones. For all three datasets, GLUE exhibited the smallest
 18 performance changes even at high corruption rates (Fig. 2c), suggesting its superior robustness.

19

20 Given its neural network-based nature, GLUE may suffer from undertraining when working with
 21 small datasets. Thus, we repeated the evaluations using subsampled datasets of various sizes. GLUE
 22 remained the top-ranking method with as few as 2,000 cells, but the alignment error increased more
 23 steeply when the data volume decreased to less than 1,000 cells (Fig. 2d). Additionally, we also
 24 noted that the performance of GLUE was robust for a wide range of hyperparameter settings
 25 (Supplementary Fig. 4).

26

1 **GLUE enables effective triple-omics integration**

2 Benefitting from a modular design and scalable adversarial alignment, GLUE readily extends to
3 more than two omics layers. As a case study, we used GLUE to integrate three distinct omics layers
4 of neuronal cells in the adult mouse cortex, including gene expression³⁴, chromatin accessibility³⁵,
5 and DNA methylation³.

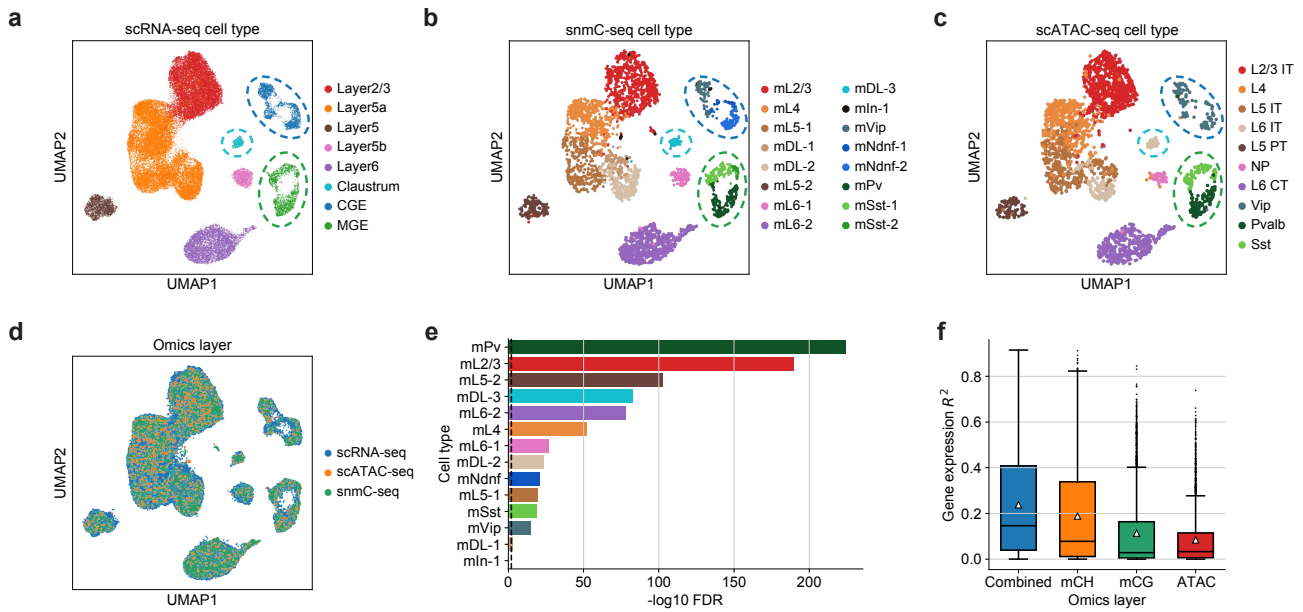
6

7 Unlike chromatin accessibility, gene body DNA methylation generally shows a negative correlation
8 with gene expression in neuronal cells³⁶. GLUE natively supports the mixture of regulatory effects
9 by modeling edge signs in the guidance graph. Such a strategy avoids data inversion, which is
10 required by previous methods^{16, 17} and can break data sparsity and the underlying distribution. For
11 the triple-omics guidance graph, we linked gene body mCH and mCG levels to genes via negative
12 edges, while the positive edges between accessible regions and genes remained the same.

13

14 The GLUE alignment successfully revealed a shared manifold of cell states across the three omics
15 layers (Fig. 3a-d). We observed highly significant marker overlap (Fig. 3e, three-way Fisher's exact
16 test³⁷, $FDR < 2 \times 10^{-15}$) for 12 out of the 14 mapped cell types (Supplementary Fig. 5, 6, Methods),
17 indicating reliable alignment. Interestingly, we found that GLUE alignment helped improve the
18 effects of cell typing in all omics layers, including the further partitioning of the scRNA-seq "MGE"
19 cluster into *Pvalb*⁺ ("mPv") and *Sst*⁺ ("mSst") subtypes (highlighted with green circles/flows in Fig.
20 3, Supplementary Fig. 5), the partitioning of the scRNA-seq "CGE" cluster and scATAC-seq "Vip"
21 cluster into *Vip*⁺ ("mVip") and *Ndnf*⁺ ("mNdnf") subtypes (highlighted with dark blue circles/flows
22 in Fig. 3, Supplementary Fig. 5), and the identification of snmC-seq "mDL-3" cells and a subset of
23 scATAC-seq "L6 IT" cells as claustrum cells (highlighted with light blue circles/flows in Fig. 3,
24 Supplementary Fig. 5).

25



1

2 **Fig. 3 Triple-omics integration of the mouse cortex.**

3 **a-c**, UMAP visualizations of the integrated cell embeddings for **a**, scRNA-seq, **b**, snmC-seq, and **c**, scATAC-seq,
4 colored by the original cell types. Cells aligning with “mPv” and “mSst” are highlighted with green circles. Cells
5 aligning with “mNdnf” and “mVip” are highlighted with dark blue circles. Cells aligning with “mDL-3” are
6 highlighted with light blue circles. **d**, UMAP visualizations of the integrated cell embeddings for all cells, colored by
7 omics layers. **e**, Significance of marker gene overlap for each cell type across all three omics layers (three-way
8 Fisher’s exact test³⁷). The dashed vertical line indicates that FDR = 0.01. We observed highly significant marker
9 overlap (FDR < 2×10^{-15}) for 12 out of the 14 cell types, indicating reliable alignment. For the remaining 2 cell types,
10 “mDL-1” had marginally significant marker overlap with FDR = 0.001, while the “mIn-1” cells in snmC-seq did not
11 properly align with the scRNA-seq or scATAC-seq cells. **f**, Coefficient of determination (R^2) for predicting gene
12 expression based on each epigenetic layer as well as the combination of all layers. The box plots indicate the medians
13 (centerlines), means (triangles), 1st and 3rd quartiles (hinges), and minima and maxima (whiskers).

14 Such triple-omics integration also sheds light on the quantitative contributions of different epigenetic
15 regulation mechanisms (Methods). Among mCH, mCG and chromatin accessibility, we found that
16 the mCH level had the highest predictive power for gene expression in cortical neurons (average R^2
17 = 0.188). When all epigenetic layers were considered, the expression predictability increased further
18 (average R^2 = 0.238), suggesting the presence of nonredundant contributions (Fig. 3f). Among the
19 neurons of different layers, DNA methylation (especially mCH) exhibited slightly higher
20 predictability for gene expression in deeper layers than in superficial layers, whereas the reverse
21 situation held for chromatin accessibility (Supplementary Fig. 7a). Across all genes, the
22 predictability of gene expression was generally correlated among the different epigenetic layers
23 (Supplementary Fig. 7b). We also observed varying associations with gene characteristics. For
24 example, mCH had higher expression predictability for longer genes, which was consistent with
25 previous studies^{17,38}, while chromatin accessibility contributed more to genes with higher expression
26 variability (Supplementary Fig. 7c).

27

1 **Model-based regulatory inference with GLUE**

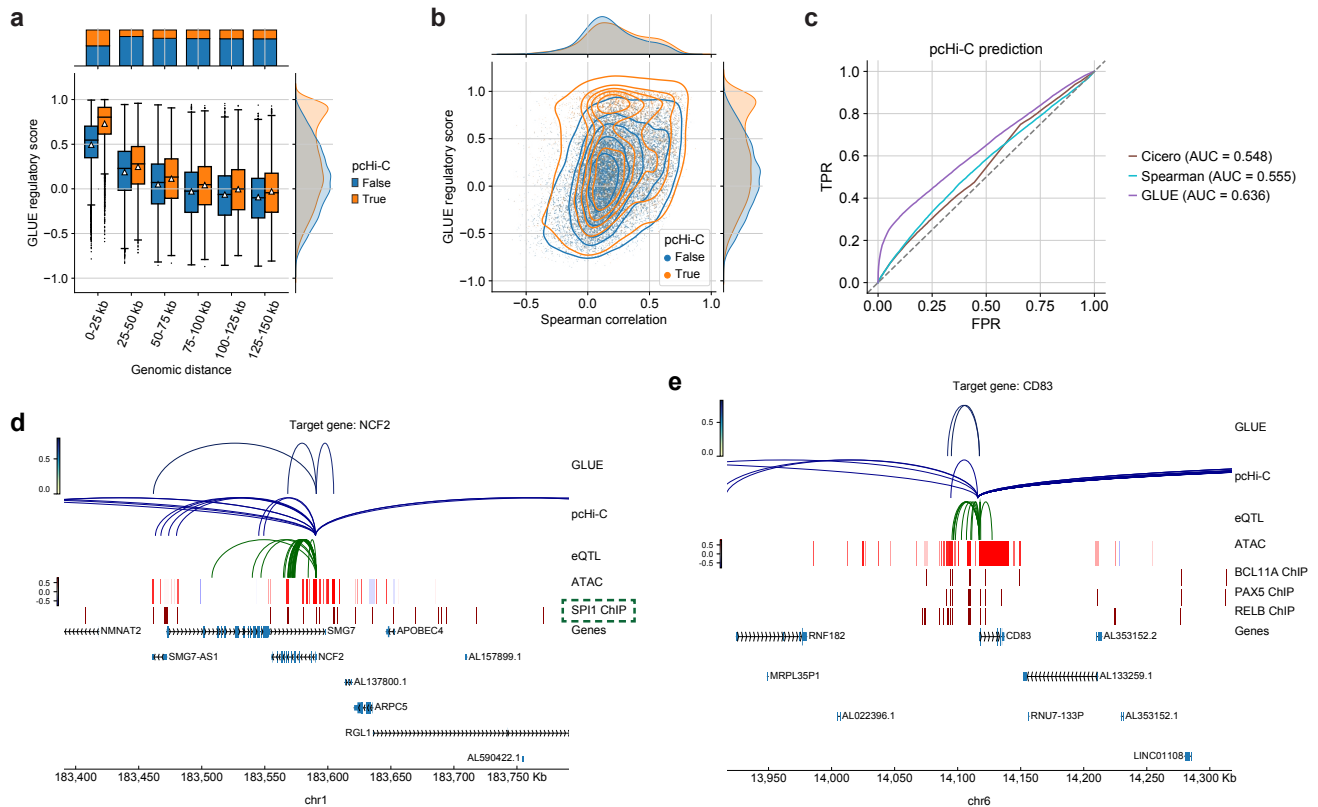
2 The incorporation of a graph explicitly modeling regulatory interactions in GLUE further enables a
3 Bayesian-like approach that combines prior knowledge and observed data for posterior regulatory
4 inference. Specifically, since the feature embeddings are designed to reconstruct the knowledge-
5 based guidance graph and single-cell multi-omics data simultaneously (Fig. 1), their cosine
6 similarities should reflect information from both aspects, which we adopt as “regulatory scores”.

7
8 As a demonstration, we employed the official PBMC (peripheral blood mononuclear cell) Multiome
9 dataset from 10x³¹ and fed it to GLUE as unpaired scRNA-seq and scATAC-seq data. To capture
10 remote cis-regulatory interactions, we employed a long-range guidance graph connecting ATAC
11 peaks and RNA genes within 150 kb windows weighted by a power-law function that models
12 chromatin contact probability^{39, 40} (Methods). Visualization of cell embeddings confirmed that the
13 GLUE alignment was correct and accurate (Supplementary Fig. 8a, b). As expected, we found that
14 the regulatory score was negatively correlated with genomic distance (Fig. 4a) and positively
15 correlated with the empirical peak-gene correlation (computed with paired cells, Fig. 4b), with
16 robustness across different random seeds (Supplementary Fig. 8c).

17
18 To further assess whether the score reflected actual cis-regulatory interactions, we compared it with
19 external evidence, including pcHi-C⁴¹ and eQTL⁴². The GLUE regulatory score was higher for pcHi-
20 C-supported peak-gene pairs in all distance ranges (Fig. 4a) and was a better predictor of pcHi-C
21 interactions than empirical peak-gene correlations (Fig. 4b, c), as well as Cicero⁴⁰, the
22 coaccessibility-based regulatory prediction method (Fig. 4c). The same held for eQTL
23 (Supplementary Fig. 8d-f).

24
25 The GLUE framework also allows additional regulatory evidence, such as pcHi-C, to be
26 incorporated intuitively via the guidance graph. Thus, we further trained new models with a
27 composite guidance graph containing distance-weighted interactions as well as pcHi-C- and eQTL-
28 supported interactions (Supplementary Fig. 9). While the multi-omics alignment was insensitive to
29 these changes, the GLUE-derived TF-target gene network (Methods) showed more significant
30 agreement with manually curated connections in the TRRUST v2 database⁴³ than individual
31 evidence-based networks (Supplementary Fig. 9e, Supplementary Fig. 10, Supplementary Table 3).

32



1

2 **Fig. 4 Model-based regulatory inference in PBMC.**

3 **a**, GLUE regulatory scores for peak-gene pairs across different genomic ranges, grouped by whether they had pcHi-C
 4 C support. The box plots indicate the medians (centerlines), means (triangles), 1st and 3rd quartiles (hinges), and
 5 minima and maxima (whiskers). **b**, Comparison between the GLUE regulatory scores and the empirical peak-gene
 6 correlations computed on paired cells. Peak-gene pairs are colored by whether they had pcHi-C support. **c**, ROC
 7 (receiver operating characteristic) curves for predicting pcHi-C interactions based on different peak-gene association
 8 scores. **d**, **e**, GLUE-identified cis-regulatory interactions for **d**, *NCF2*, and **e**, *CD83*, along with individual regulatory
 9 evidence. SPI1 (highlighted with a green box) is a known regulator of *NCF2*.

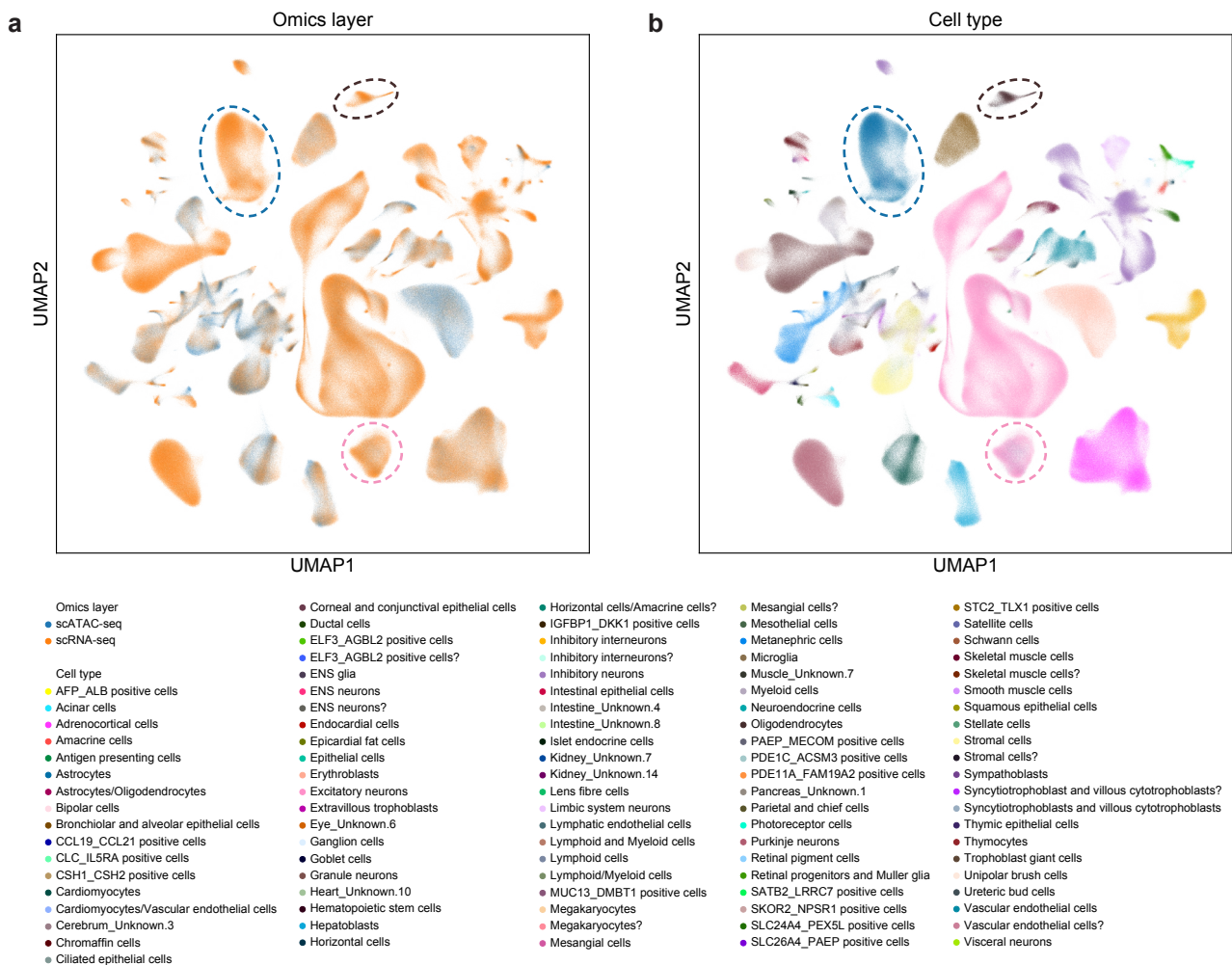
10 Interestingly, we noticed that the GLUE-inferred cis-regulatory interactions could provide new hints
 11 about the regulatory mechanisms of known TF-target pairs. For example, SPI1 is a known regulator
 12 of the *NCF2* gene, and both are highly expressed in monocytes (Supplementary Fig. 11a, b). GLUE
 13 identified three remote regulatory peaks for *NCF2* with various pieces of evidence, i.e., ~120 kb
 14 downstream, ~25 kb downstream, and ~20 kb upstream from the TSS (transcription start site) (Fig.
 15 4d), all of which were bound by SPI1. Meanwhile, most putative regulatory interactions were
 16 previously unknown. For example, *CD83* was linked with two regulatory peaks (~25 kb upstream
 17 from the TSS), which were enriched for the binding of three TFs (BCL11A, PAX5, and RELB; Fig.
 18 4e). While *CD83* was highly expressed in both monocytes and B cells, the inferred TFs showed more
 19 constrained expression patterns (Supplementary Fig. 11c-f), suggesting that its active regulators
 20 might differ per cell type. Supplementary Fig. 12 shows more examples of GLUE-inferred regulatory
 21 interactions.

1 Atlas-scale integration over millions of cells with GLUE

2 As technologies continue to evolve, the throughput of single-cell experiments is constantly
 3 increasing. Recent studies have generated human cell atlases for gene expression²⁵ and chromatin
 4 accessibility²⁶ containing millions of cells. The integration of these atlases poses a significant
 5 challenge to computational methods due to the sheer volume of data, extensive heterogeneity, low
 6 coverage per cell, and unbalanced cell type compositions, and has yet to be accomplished at the
 7 single-cell level.

8
 9 Implemented as parametric neural networks with minibatch optimization, GLUE delivers superior
 10 scalability with a sublinear time cost, promising its applicability at the atlas scale (Supplementary
 11 Fig. 13a). Using an efficient multistage training strategy for GLUE (Methods), we successfully
 12 integrated the gene expression and chromatin accessibility data into a unified multi-omics human cell
 13 atlas (Fig. 5).

14



15

1 **Fig. 5 Integration of a multi-omics human cell atlas.**

2 UMAP visualizations of the integrated cell embeddings, colored by **a**, omics layers, and **b**, cell types. The pink circles
3 highlight cells labeled as “Excitatory neurons” in scRNA-seq but “Astrocytes” in scATAC-seq. The blue circles
4 highlight cells labeled as “Astrocytes” in scRNA-seq but “Astrocytes/Oligodendrocytes” in scATAC-seq. The brown
5 circles highlight cells labeled as “Oligodendrocytes” in scRNA-seq but “Astrocytes/Oligodendrocytes” in scATAC-
6 seq.

7 While the aligned atlas was largely consistent with the original annotations²⁶ (Supplementary Fig.
8 13c-e), we also noticed several discrepancies. For example, cells originally annotated as
9 “Astrocytes” in scATAC-seq were aligned to an “Excitatory neurons” cluster in scRNA-seq
10 (highlighted with pink circles/flows in Supplementary Fig. 13). Further inspection revealed that
11 canonical radial glial (RG) markers such as *PAX6*, *HES1*, and *HOPX*^{44, 45} were actively transcribed
12 in this cluster, both in the RNA and ATAC domain (Supplementary Fig. 14), with chromatin
13 priming⁹ also detected at both neuronal and glial markers (Supplementary Fig. 15-17), suggesting
14 that the cluster consists of multipotent neural progenitors (likely RGs) rather than excitatory neurons
15 or astrocytes as originally annotated. GLUE-based integration also resolved several scATAC-seq
16 clusters that were ambiguously annotated. For example, the “Astrocytes/Oligodendrocytes” cluster
17 was split into two halves and aligned to the “Astrocytes” and “Oligodendrocytes” clusters of scRNA-
18 seq (highlighted with blue and brown circles/flows in Supplementary Fig. 13, respectively), which
19 was also supported by marker expression and accessibility (Supplementary Fig. 16, 17). These
20 results demonstrate the unique value of atlas-scale multi-omics integration.

21

22 **Discussion**

23 Combining omics-specific autoencoders with graph-based coupling and adversarial alignment, we
24 designed and implemented the GLUE framework for unpaired single-cell multi-omics data
25 integration with superior accuracy and robustness. By modeling regulatory interactions across omics
26 layers explicitly, GLUE uniquely supports model-based regulatory inference for unpaired multi-
27 omics datasets, exhibiting even higher reliability than regular correlation analysis on paired datasets
28 (notably, in a Bayesian interpretation, the GLUE regulatory inference can be seen as a posterior
29 estimate, which can be continuously refined upon the arrival of new data). Furthermore, benefitting
30 from a neural network-based design, GLUE enables notable scalability for whole-atlas alignment
31 over millions of unpaired cells, which remains a serious challenge for *in silico* integration. In fact,
32 we also attempted to perform integration using online iNMF, which was the only other method
33 capable of integrating the data at full scale, but the result was far from optimal (Supplementary Fig.

1 18a, b). Meanwhile, an attempt to integrate the data as aggregated metacells (Methods) via the
2 popular Seurat v3 method also failed (Supplementary Fig. 18c, d).

3
4 Unpaired multi-omics integration, also referred to as diagonal integration¹⁴, shares some conceptual
5 similarities with batch effect correction⁴⁶, as both call for the alignment of unpaired cells in certain
6 data representations. Nonetheless, the former is significantly more challenging because of the
7 distinct, omics-specific feature spaces. While completely unsupervised multi-omics integration has
8 been proposed^{21, 22}, such an approach is exceedingly difficult and has largely been limited to aligning
9 continuous trajectories. For general-case multi-omics integration, additional prior knowledge is
10 necessary. At the omics feature level, presumed feature interactions have been used via feature
11 conversion^{15-17, 30} or coupled matrix factorization^{19, 20}. While feature conversion may seem to be a
12 straightforward solution, the inevitable information loss¹⁸ can have a detrimental effect on
13 performance. Apart from the feature-converted data, Seurat v3¹⁵ and bindSC³⁰ also devised heuristic
14 strategies to utilize information in the original feature space, which probably explains their improved
15 performance than methods that do not^{16, 17}. At the cell level, known cell types have also been used
16 via (semi-)supervised learning^{47, 48}, but this approach incurs substantial limitations in terms of
17 applicability since such supervision is typically unavailable and in many cases serves as the purpose
18 of multi-omics integration *per se*²⁶. Notably, one of these methods was proposed with a similar
19 autoencoder architecture and adversarial alignment⁴⁸, but it relied on matched cell types or clusters to
20 orient the alignment. In fact, GLUE shares more conceptual similarity with the coupled matrix
21 factorization methods, but with superior accuracy, robustness and scalability, which mostly benefits
22 from its deep generative model-based design.

23
24 We note that the current framework also works for integrating omics layers with shared features, by
25 using either the same vertex or connected surrogate vertices for each shared feature in the guidance
26 graph. In particular, the integration between scRNA-seq and spatial transcriptomics^{49, 50} could be
27 naturally implemented in this way. After the integration, genes not detected in the spatial
28 transcriptome could be further imputed via cross-layer translation, through a combination of the
29 spatial transcriptomics encoder and the scRNA-seq decoder.

30
31 As a generalizable framework, GLUE features a modular design, where the data and graph
32 autoencoders are independently configurable.

- 33 ● The data autoencoders in GLUE are customizable with appropriate generative models that
34 conform to omics-specific data distributions. In the current work, we used the negative binomial

1 distribution for scRNA-seq and scATAC-seq, and the zero-inflated log-normal distribution for
2 snmC-seq (Methods). Nevertheless, generative distributions can be easily reconfigured to
3 accommodate other omics layers, such as protein abundance⁵¹ and histone modification⁵², and to
4 adopt new advances in data modeling techniques⁵³.

5 ● The guidance graphs used in GLUE have currently been limited to multipartite graphs,
6 containing only edges between features of different layers. Nonetheless, graphs, as intuitive and
7 flexible representations of regulatory knowledge, can embody more complex regulatory patterns,
8 including within-modality interactions, non-feature vertices, and multi-relations. Beyond
9 canonical graph convolution, more advanced graph neural network architectures⁵⁴⁻⁵⁶ may also be
10 adopted to extract richer information from the regulatory graph.

11
12 Recent advances in experimental multi-omics technologies have increased the availability of paired
13 data^{8-11, 31}. While most of the current simultaneous multi-omics protocols still suffer from lower data
14 quality or throughput than that of single-omics methods⁵⁷, paired cells can be highly informative in
15 anchoring different omics layers and should be utilized in conjunction with unpaired cells whenever
16 available. It is straightforward to extend the GLUE framework to incorporate such pairing
17 information, e.g., by adding another loss term that penalizes the embedding distances between paired
18 cells⁵⁸. Such an extension may ultimately lead to a solution for the general case of mosaic
19 integration¹⁴.

20
21 Apart from multi-omics integration, we also note that the GLUE framework could be suitable for
22 cross-species integration, especially when distal species are concerned and one-to-one orthologs are
23 limited. Specifically, we may compile all orthologs into a GLUE guidance graph and perform
24 integration without explicit ortholog conversion. Under that setting, the GLUE approach could also
25 be conceptually connected to a recent work called SAMap⁵⁹.

26
27 We believe that GLUE, as a modular and generalizable framework, creates an unprecedented
28 opportunity towards effectively delineating gene regulatory maps via large-scale multi-omics
29 integration at single-cell resolution. The whole package of GLUE, along with tutorials and demo
30 cases, is available online at <https://github.com/gao-lab/GLUE> for the community.

1 Methods

2 The GLUE framework

3 We assume that there are K different omics layers to be integrated, each with a distinct feature set
 4 $\mathcal{V}_k, k = 1, 2, \dots, K$. For example, in scRNA-seq, \mathcal{V}_k is the set of genes, while in scATAC-seq, \mathcal{V}_k is
 5 the set of chromatin regions. The data spaces of different omics layers are denoted as $\mathcal{X}_k \subseteq \mathbb{R}^{|\mathcal{V}_k|}$
 6 with varying dimensionalities. We use $\mathbf{x}_k^{(n)} \in \mathcal{X}_k, n = 1, 2, \dots, N_k$ to denote cells from the k^{th} omics
 7 layer and $\mathbf{x}_{ki}^{(n)}, i \in \mathcal{V}_k$ to denote the observed value of feature i of the k^{th} layer in the n^{th} cell. N_k is
 8 the sample size of the k^{th} layer. Notably, the cells from different omics layers are unpaired and can
 9 have different sample sizes. To avoid cluttering, we drop the superscript (n) when referring to an
 10 arbitrary cell.

11
 12 We model the observed data from different omics layers as generated by a low-dimensional latent
 13 variable (i.e., cell embedding) $\mathbf{u} \in \mathbb{R}^m$:

$$14 \quad p(\mathbf{x}_k; \theta_k) = \int p(\mathbf{x}_k | \mathbf{u}; \theta_k) p(\mathbf{u}) d\mathbf{u} \quad \text{Eq. 1}$$

15 where $p(\mathbf{u})$ is the prior distribution of the latent variable, $p(\mathbf{x}_k | \mathbf{u}; \theta_k)$ are learnable generative
 16 distributions (i.e., data decoders), and θ_k denotes learnable parameters in the decoders. The cell
 17 latent variable \mathbf{u} is shared across different omics layers. In other words, \mathbf{u} represents the common
 18 cell states underlying all omics observations, while the observed data from each layer are generated
 19 by a specific type of measurement of the underlying cell states.

20
 21 With the introduction of variational posteriors $q(\mathbf{u} | \mathbf{x}_k; \phi_k)$ (i.e., data encoders, where ϕ_k are
 22 learnable parameters in the encoders), model fitting can be efficiently performed by maximizing the
 23 following evidence lower bounds:

$$24 \quad \mathcal{L}_{\mathcal{X}_k}(\phi_k, \psi_k) = \mathbb{E}_{\mathbf{x}_k \sim p_{\text{data}}(\mathbf{x}_k)} \left[\mathbb{E}_{\mathbf{u} \sim q(\mathbf{u} | \mathbf{x}_k; \phi_k)} \log p(\mathbf{x}_k | \mathbf{u}; \theta_k) - \text{KL}(q(\mathbf{u} | \mathbf{x}_k; \phi_k) \parallel p(\mathbf{u})) \right] \quad \text{Eq. 2}$$

25
 26 Since different autoencoders are independently parameterized and trained on separate data, the cell
 27 embeddings learned for different omics layers could have inconsistent semantic meanings unless
 28 they are linked properly.

29

1 To link the autoencoders, we propose a guidance graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, which incorporates prior
 2 knowledge about the regulatory interactions among features at distinct omics layers, where $\mathcal{V} =$
 3 $\bigcup_{k=1}^K \mathcal{V}_k$ is the universal feature set and $\mathcal{E} = \{(i, j) | i, j \in \mathcal{V}\}$ is the set of edges. Each edge is also
 4 associated with signs and weights, which are denoted as s_{ij} and w_{ij} , respectively. We require that
 5 $w_{ij} \in (0, 1]$, which can be interpreted as interaction credibility, and that $s_{ij} \in \{-1, 1\}$, which
 6 specifies the sign of the regulatory interaction. For example, an ATAC peak located near the
 7 promoter of a gene is usually assumed to positively regulate its expression, so they can be connected
 8 with a positive edge ($s_{ij} = 1$). Meanwhile, DNA methylation in the gene promoter is usually
 9 assumed to suppress expression, so they can be connected with a negative edge ($s_{ij} = -1$). In
 10 addition to the connections between features, self-loops are also added for numerical stability, with
 11 $s_{ii} = 1, w_{ii} = 1, \forall i \in \mathcal{V}$.

12
 13 We treat the guidance graph as observed variable and model it as generated by low-dimensional
 14 feature latent variables (i.e., feature embeddings) $\mathbf{v}_i \in \mathbb{R}^m, i \in \mathcal{V}$. Furthermore, differing from the
 15 previous model, we now model \mathbf{x}_k as generated by the combination of feature latent variables $\mathbf{v}_i \in$
 16 $\mathbb{R}^m, i \in \mathcal{V}_k$ and the cell latent variable $\mathbf{u} \in \mathbb{R}^m$. For convenience, we introduce the notation $\mathbf{V} \in$
 17 $\mathbb{R}^{m \times |\mathcal{V}|}$, which combines all feature embeddings into a single matrix. The model likelihood can thus
 18 be written as:

$$19 \quad p(\mathbf{x}_k, \mathcal{G}; \theta_k, \theta_G) = \int p(\mathbf{x}_k | \mathbf{u}, \mathbf{V}; \theta_k) p(\mathcal{G} | \mathbf{V}; \theta_G) p(\mathbf{u}) p(\mathbf{V}) d\mathbf{u} d\mathbf{V} \quad \text{Eq. 3}$$

20 where $p(\mathbf{x}_k | \mathbf{u}, \mathbf{V}; \theta_k)$ and $p(\mathcal{G} | \mathbf{V}; \theta_G)$ are learnable generative distributions for the omics data (i.e.,
 21 data decoders) and knowledge graph (i.e., graph decoder), respectively. θ_k and θ_G are learnable
 22 parameters in the decoders. $p(\mathbf{u})$ and $p(\mathbf{V})$ are the prior distributions of the cell latent variable and
 23 feature latent variables, respectively, which are fixed as standard normal distributions for simplicity:

$$24 \quad p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; \mathbf{0}, \mathbf{I}_m) \quad \text{Eq. 4}$$

$$25 \quad p(\mathbf{v}_i) = \mathcal{N}(\mathbf{v}_i; \mathbf{0}, \mathbf{I}_m), \quad p(\mathbf{V}) = \prod_{i \in \mathcal{V}} p(\mathbf{v}_i) \quad \text{Eq. 5}$$

26 though alternatives may also be used⁶⁰. For convenience, we also introduce the notation $\mathbf{V}_k \in$
 27 $\mathbb{R}^{m \times |\mathcal{V}_k|}$, which contains only feature embeddings in the k^{th} omics layer, and \mathbf{u}_k , which emphasizes
 28 that the cell embedding is from a cell in the k^{th} omics layer.

29
 30 The graph likelihood $p(\mathcal{G} | \mathbf{V}; \theta_G)$ (i.e., graph decoder) is defined as:

$$1 \quad p(\mathcal{G}|\mathbf{V}; \theta_{\mathcal{G}}) = \mathbb{E}_{i,j \sim p(i,j;w_{ij})} \left[\sigma(s_{ij} \cdot \mathbf{v}_i^T \mathbf{v}_j) \cdot \mathbb{E}_{j' \sim p_{\text{ns}}(j'|i)} \left(1 - \sigma(s_{ij} \cdot \mathbf{v}_i^T \mathbf{v}_{j'}) \right) \right] \quad \text{Eq. 6}$$

2 where σ is the sigmoid function and p_{ns} is a negative sampling distribution⁶¹. In other words, we
3 first sample the edges (i, j) with probabilities proportional to the edge weights and then sample
4 vertices j' that are not connected to i and treat them as if $s_{ij'} = s_{ij}$. When maximizing the graph
5 likelihood, the inner products between features are maximized or minimized (per edge sign) based on
6 the Bernoulli distribution. For example, ATAC peaks located near the promoter of a gene would be
7 encouraged to have similar embeddings to that of the gene, while DNA methylation in the gene
8 promoter would be encouraged to have a dissimilar embedding to that of the gene.

9

10 The data likelihoods $p(\mathbf{x}_k|\mathbf{u}, \mathbf{V}; \theta_k)$ (i.e., data decoders) in Eq. 3 are built upon the inner product
11 between the cell embedding \mathbf{u} and feature embeddings \mathbf{V}_k . Thus, analogous to the loading matrix in
12 principal component analysis (PCA), the feature embeddings \mathbf{V}_k confer semantic meanings for the
13 cell embedding space. As \mathbf{V}_k are modulated by interactions among omics features in the guidance
14 graph, the semantic meanings become linked. The exact formulation of data likelihood depends on
15 the omics data distribution. For example, for count-based scRNA-seq and scATAC-seq data, we
16 used the negative binomial (NB) distribution:

$$17 \quad p(\mathbf{x}_k|\mathbf{u}, \mathbf{V}; \theta_k) = \prod_{i \in \mathcal{V}_k} \text{NB}(\mathbf{x}_{k_i}; \boldsymbol{\mu}_i, \boldsymbol{\theta}_i) \quad \text{Eq. 7}$$

$$18 \quad \text{NB}(\mathbf{x}_{k_i}; \boldsymbol{\mu}_i, \boldsymbol{\theta}_i) = \frac{\Gamma(\mathbf{x}_{k_i} + \boldsymbol{\theta}_i)}{\Gamma(\boldsymbol{\theta}_i)\Gamma(\mathbf{x}_{k_i} + 1)} \left(\frac{\boldsymbol{\mu}_i}{\boldsymbol{\theta}_i + \boldsymbol{\mu}_i} \right)^{\mathbf{x}_{k_i}} \left(\frac{\boldsymbol{\theta}_i}{\boldsymbol{\theta}_i + \boldsymbol{\mu}_i} \right)^{\boldsymbol{\theta}_i} \quad \text{Eq. 8}$$

$$19 \quad \boldsymbol{\mu}_i = \text{Softmax}_i(\boldsymbol{\alpha} \odot \mathbf{V}_k^T \mathbf{u} + \boldsymbol{\beta}) \cdot \sum_{j \in \mathcal{V}_k} \mathbf{x}_{k_j} \quad \text{Eq. 9}$$

20 where $\boldsymbol{\mu}, \boldsymbol{\theta} \in \mathbb{R}_+^{|\mathcal{V}_k|}$ are the mean and dispersion of the NB distribution, respectively, and $\boldsymbol{\alpha} \in$
21 $\mathbb{R}_+^{|\mathcal{V}_k|}, \boldsymbol{\beta} \in \mathbb{R}^{|\mathcal{V}_k|}$ are scaling and bias factors. \odot is the Hadamard product. Softmax_i represents the i^{th}
22 dimension of the softmax output. The set of learnable parameters is $\theta_k = \{\boldsymbol{\theta}, \boldsymbol{\alpha}, \boldsymbol{\beta}\}$. Analogously,
23 many other distributions can also be supported, as long as we can parameterize the means of the
24 distributions by feature-cell inner products.

25

26 For efficient inference and optimization, we introduce the following factorized variational posterior:

$$27 \quad q(\mathbf{u}, \mathbf{V}|\mathbf{x}_k, \mathcal{G}; \phi_k, \phi_{\mathcal{G}}) = q(\mathbf{u}|\mathbf{x}_k; \phi_k) \cdot q(\mathbf{V}|\mathcal{G}; \phi_{\mathcal{G}}) \quad \text{Eq. 10}$$

28 The graph variational posterior $q(\mathbf{V}|\mathcal{G}; \phi_{\mathcal{G}})$ (i.e., graph encoder) is modeled as diagonal-covariance
29 normal distributions parameterized by a graph convolutional network⁶²:

$$1 \quad q(\mathbf{V}|\mathcal{G}; \phi_{\mathcal{G}}) = \prod_{i \in \mathcal{V}} q(\mathbf{v}_i|\mathcal{G}; \phi_{\mathcal{G}}) \quad \text{Eq. 11}$$

$$2 \quad q(\mathbf{v}_i|\mathcal{G}; \phi_{\mathcal{G}}) = \mathcal{N}\left(\mathbf{v}_i; \text{GCN}_{\mu_i}(\mathcal{G}; \phi_{\mathcal{G}}), \text{GCN}_{\sigma_i^2}(\mathcal{G}; \phi_{\mathcal{G}})\right) \quad \text{Eq. 12}$$

3 where $\phi_{\mathcal{G}}$ represents the learnable parameters in the GCN encoder.

4

5 The variational data posteriors $q(\mathbf{u}|\mathbf{x}_k; \phi_k)$ (i.e., data encoders) are modeled as diagonal-covariance
6 normal distributions parameterized by multilayer perceptron (MLP) neural networks:

$$7 \quad q(\mathbf{u}|\mathbf{x}_k, \mathbf{V}_k; \phi_k) = \mathcal{N}\left(\mathbf{u}; \text{MLP}_{k,\mu}(\mathbf{x}_k; \phi_k), \text{MLP}_{k,\sigma^2}(\mathbf{x}_k; \phi_k)\right) \quad \text{Eq. 13}$$

8 where ϕ_k is the set of learnable parameters in the MLP encoder of the k^{th} omics layer.

9

10 Model fitting can then be performed by maximizing the following evidence lower bound:

$$11 \quad \sum_{k=1}^K \mathbb{E}_{\mathbf{x}_k \sim p_{\text{data}}(\mathbf{x}_k)} \left[\mathbb{E}_{\mathbf{u} \sim q(\mathbf{u}|\mathbf{x}_k; \phi_k), \mathbf{V} \sim q(\mathbf{V}|\mathcal{G}; \phi_{\mathcal{G}})} \log p(\mathbf{x}_k|\mathbf{u}, \mathbf{V}; \theta_k) p(\mathcal{G}|\mathbf{V}; \theta_{\mathcal{G}}) \right. \\ \left. - \text{KL}\left(q(\mathbf{u}|\mathbf{x}_k; \phi_k) q(\mathbf{V}|\mathcal{G}; \phi_{\mathcal{G}}) \parallel p(\mathbf{u}) p(\mathbf{V})\right) \right] \quad \text{Eq. 14}$$

12 which can be further rearranged into the following form:

$$13 \quad K \cdot \mathcal{L}_{\mathcal{G}}(\theta_{\mathcal{G}}, \phi_{\mathcal{G}}) + \sum_{k=1}^K \mathcal{L}_{x_k}(\theta_k, \phi_k, \phi_{\mathcal{G}}) \quad \text{Eq. 15}$$

14 where we have

$$15 \quad \mathcal{L}_{x_k}(\theta_k, \phi_k, \phi_{\mathcal{G}}) = \mathbb{E}_{\mathbf{x}_k \sim p_{\text{data}}(\mathbf{x}_k)} \left[\mathbb{E}_{\mathbf{u} \sim q(\mathbf{u}|\mathbf{x}_k; \phi_k), \mathbf{V} \sim q(\mathbf{V}|\mathcal{G}; \phi_{\mathcal{G}})} \log p(\mathbf{x}_k|\mathbf{u}, \mathbf{V}; \theta_k) - \text{KL}\left(q(\mathbf{u}|\mathbf{x}_k; \phi_k) \parallel p(\mathbf{u})\right) \right] \quad \text{Eq. 16}$$

$$16 \quad \mathcal{L}_{\mathcal{G}}(\theta_{\mathcal{G}}, \phi_{\mathcal{G}}) = \mathbb{E}_{\mathbf{V} \sim q(\mathbf{V}|\mathcal{G}; \phi_{\mathcal{G}})} \log p(\mathcal{G}|\mathbf{V}; \theta_{\mathcal{G}}) - \text{KL}\left(q(\mathbf{V}|\mathcal{G}; \phi_{\mathcal{G}}) \parallel p(\mathbf{V})\right) \quad \text{Eq. 17}$$

17 Below, for convenience, we denote the union of all encoder parameters as $\phi = (\cup_{k=1}^K \phi_k) \cup \phi_{\mathcal{G}}$ and

18 the union of all decoder parameters as $\theta = (\cup_{k=1}^K \theta_k) \cup \theta_{\mathcal{G}}$.

19

20 To ensure the proper alignment of different omics layers, we use the adversarial alignment strategy^{28,}

21 ⁶³. A discriminator D with a K -dimensional softmax output is introduced, which predicts the omics

22 layers of cells based on their embeddings \mathbf{u} . The discriminator D is trained by minimizing the

23 multiclass classification cross entropy:

$$24 \quad \mathcal{L}_D(\phi, \psi) = -\frac{1}{K} \sum_{k=1}^K \mathbb{E}_{\mathbf{x}_k \sim p_{\text{data}}(\mathbf{x}_k)} \mathbb{E}_{\mathbf{u} \sim q(\mathbf{u}|\mathbf{x}_k; \phi_k)} \log D_k(\mathbf{u}; \psi) \quad \text{Eq. 18}$$

25 where D_k represents the k^{th} dimension of the discriminator output and ψ is the set of learnable

26 parameters in the discriminator. The data encoders can then be trained in the opposite direction to

1 fool the discriminator, ultimately leading to the alignment of cell embeddings from different omics
2 layers⁶⁴.

3

4 The overall training objective of GLUE thus consists of:

5
$$\min_{\psi} \lambda_D \cdot \mathcal{L}_D(\phi, \psi) \quad \text{Eq. 19}$$

6
$$\max_{\theta, \phi} \lambda_D \cdot \mathcal{L}_D(\phi, \psi) + \lambda_G K \cdot \mathcal{L}_G(\theta_G, \phi_G) + \sum_{k=1}^K \mathcal{L}_{x_k}(\theta_k, \phi_k, \phi_G) \quad \text{Eq. 20}$$

7 The two hyperparameters λ_D and λ_G control the contributions of adversarial alignment and graph-
8 based feature embedding, respectively. We use stochastic gradient descent (SGD) to train the GLUE
9 model. Each SGD iteration is divided into two steps. In the first step, the discriminator is updated
10 according to objective Eq. 19. In the second step, the data and graph autoencoders are updated
11 according to Eq. 20. The RMSprop optimizer with no momentum term is employed to ensure the
12 stability of adversarial training.

13

14 **Implementation details**

15 We applied linear dimensionality reduction using canonical methods such as PCA (for scRNA-seq)
16 or LSI (latent semantic indexing, for scATAC-seq) as the first transformation layers of the data
17 encoders (note that the decoders were still fitted in the original feature spaces). This effectively
18 reduced model size and enabled a modular input, so advanced dimensionality reduction or batch
19 effect correction methods can also be used instead as preprocessing steps for GLUE integration.

20

21 To ensure stable alignment, we used batch normalization in the data encoder layers and employed
22 additive noise annealing. Specifically, noise $\epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \tau \cdot \mathbf{I}_m)$ was added to the cell embeddings \mathbf{u}
23 before passing to the discriminator. The parameter τ controls the noise level, which starts at $\tau = 1$
24 and decreases linearly per epoch until reaching 0 (i.e., noise annealing). The number of annealing
25 epochs was set automatically based on the data size and learning rate to match a learning progress
26 equivalent to 4,000 iterations at a learning rate of 0.002.

27

28 During model training, 10% of the cells were used as the validation set. In the final stage of training,
29 the learning rate would be reduced by factors of 10 if the validation loss did not improve for
30 consecutive epochs. Training would be terminated if the validation loss still did not improve for
31 consecutive epochs. The patience for learning rate reduction, training termination, and the maximal

1 number of training epochs were automatically set based on the data size and learning rate to match a
2 learning progress equivalent to 1,000, 2,000, and 16,000 iterations at a learning rate of 0.002,
3 respectively.

4
5 For all benchmarks and case studies with GLUE, we used the default hyperparameters unless
6 explicitly stated. The set of default hyperparameters is presented in Supplementary Fig. 4.

8 **Systematic benchmarks**

9 UnionCom²¹ and GLUE were executed using the Python packages “unioncom” (v0.3.0) and “scglue”
10 (v0.1.1), respectively. Online iNMF¹⁶, LIGER¹⁷, bindSC³⁰, and Seurat v3¹⁵ were executed using the
11 R packages “rliger” (v1.0.0), “rliger” (v1.0.0), “bindSC” (v1.0.0), and “Seurat” (v4.0.2),
12 respectively. For each method, we used the default hyperparameter settings and data preprocessing
13 steps as recommended. For the scRNA-seq data, 2,000 highly variable genes were selected using the
14 Seurat “vst” method. To construct the guidance graph, we connected ATAC peaks with RNA genes
15 via positive edges if they overlapped in either the gene body or proximal promoter regions (defined
16 as 2 kb upstream from the TSS). For the methods that require feature conversion (online iNMF,
17 LIGER, bindSC, and Seurat v3), we converted the scATAC-seq data to gene-level activity scores by
18 summing up counts in the ATAC peaks connected to specific genes in the guidance graph. Notably,
19 online iNMF and LIGER also recommend an alternative way of ATAC feature conversion, i.e.,
20 directly counting ATAC fragments falling in gene body and promoter regions without resorting to
21 ATAC peaks ([https://htmlpreview.github.io/?https://github.com/welch-](https://htmlpreview.github.io/?https://github.com/welch-lab/liger/blob/master/vignettes/Integrating_scRNA_and_scATAC_data.html)
22 [lab/liger/blob/master/vignettes/Integrating_scRNA_and_scATAC_data.html](https://htmlpreview.github.io/?https://github.com/welch-lab/liger/blob/master/vignettes/Integrating_scRNA_and_scATAC_data.html)), which we abbreviate
23 as FiG (fragments in genes). We also tested the FiG feature conversion method with online iNMF
24 and LIGER.

25
26 MAP (mean average precision) was used to evaluate the cell type resolution. Supposing that the cell
27 type of the i^{th} cell is $y^{(i)}$ and that the cell types of its K ordered nearest neighbors are
28 $y_1^{(i)}, y_2^{(i)}, \dots, y_K^{(i)}$, the MAP is then defined as follows:

$$29 \quad \text{MAP} = \frac{1}{N} \sum_{i=1}^N \text{AP}^{(i)} \quad \text{Eq. 21}$$

$$AP^{(i)} = \begin{cases} \frac{\sum_{k=1}^K \mathbb{1}_{y^{(i)}=y_k} \cdot \frac{\sum_{j=1}^k \mathbb{1}_{y^{(i)}=y_j}}{k}}{\sum_{k=1}^K \mathbb{1}_{y^{(i)}=y_k}}, & \text{if } \sum_{k=1}^K \mathbb{1}_{y^{(i)}=y_k} > 0 \\ 0, & \text{otherwise} \end{cases} \quad \text{Eq. 22}$$

where $\mathbb{1}_{y^{(i)}=y_k}$ is an indicator function that equals 1 if $y^{(i)} = y_k$ and 0 otherwise. For each cell, AP (average precision) computes the average cell type precision up to each cell type-matched neighbor, and MAP is the average AP across all cells. We set K to 1% of the total number of cells in each dataset. MAP has a range of 0 – 1, and higher values indicate better cell type resolution.

SAS (Seurat alignment score) was used to evaluate the extent of mixing among distinct omics layers and was computed as described in the original paper³²:

$$SAS = 1 - \frac{\bar{x} - \frac{k}{N}}{k - \frac{k}{N}} \quad \text{Eq. 23}$$

where \bar{x} is the average number of cells from the same omics layer among the k nearest neighbors (different layers were first subsampled to the same number of cells as the smallest layer), and N is the number of omics layers. We set k to 1% of the subsampled cell number. SAS has a range of 0 – 1, and higher values indicate better mixing.

FOSCTTM (fraction of samples closer than the true match)³³ was used to evaluate the single-cell level alignment accuracy. It was computed on two datasets with known cell-to-cell pairings. Suppose that each dataset contains N cells, and that the cells are sorted in the same order, i.e., the i^{th} cell in the first dataset is paired with the i^{th} cell in the second dataset. Denote \mathbf{x} and \mathbf{y} as the cell embeddings of the first and second dataset, respectively. The FOSCTTM is then defined as:

$$FOSCTTM = \frac{1}{2N} \left(\sum_{i=1}^N \frac{n_1^{(i)}}{N} + \sum_{i=1}^N \frac{n_2^{(i)}}{N} \right) \quad \text{Eq. 24}$$

$$n_1^{(i)} = |\{j | d(\mathbf{x}_j, \mathbf{y}_i) < d(\mathbf{x}_i, \mathbf{y}_i)\}| \quad \text{Eq. 25}$$

$$n_2^{(i)} = |\{j | d(\mathbf{x}_i, \mathbf{y}_j) < d(\mathbf{x}_i, \mathbf{y}_i)\}| \quad \text{Eq. 26}$$

where $n_1^{(i)}$ and $n_2^{(i)}$ are the number of cells in the first and second dataset, respectively, that are closer to the i^{th} cell than their true matches in the opposite dataset. d is the Euclidean distance.

For the baseline benchmark, each method was run 8 times with different random seeds, except for bindSC, which has a deterministic implementation and was run only once. For the guidance

1 corruption benchmark, we removed the specified proportions of existing peak-gene interactions and
 2 added equal numbers of nonexistent interactions, so the total number of interactions remained
 3 unchanged. Of note, feature conversion was also repeated using the corrupted guidance graphs. The
 4 corruption procedure was repeated 8 times with different random seeds. For the subsampling
 5 benchmark, the scRNA-seq and scATAC-seq cells were subsampled in pairs (so FOSCTTM could
 6 still be computed). The subsampling process was also repeated 8 times with different random seeds.

7

8 For the systematic scalability test (Supplementary Fig. 13a), all methods were run on a Linux
 9 workstation with 40 CPU cores (two Intel Xeon Silver 4210 chips), 250 GB of RAM, and NVIDIA
 10 GeForce RTX 2080 Ti GPUs. Only a single GPU card was used when training GLUE.

11

12 Triple-omics integration

13 The scRNA-seq and scATAC-seq data were handled as previously described (see the section
 14 “Systematic benchmarks”). Due to low coverage per single-C site, the snmC-seq data were converted
 15 to average methylation levels in gene bodies. The mCH and mCG levels were quantified separately,
 16 resulting in 2 features per gene. The gene methylation levels were normalized by the global
 17 methylation level per cell. An initial dimensionality reduction was performed using PCA (see the
 18 section “Implementation details”). For the triple-omics guidance graph, the mCH and mCG levels
 19 were connected to the corresponding genes with negative edges.

20

21 The normalized methylation levels were positive, with dropouts corresponding to the genes that were
 22 not covered in single cells. As such, we used the zero-inflated log-normal (ZILN) distribution for the
 23 data decoder:

$$24 \quad p(\mathbf{x}_k | \mathbf{u}, \mathbf{V}; \theta_k) = \prod_{i \in \mathcal{V}_k} \text{ZILN}(\mathbf{x}_{k_i}; \boldsymbol{\mu}_i, \boldsymbol{\sigma}_i, \boldsymbol{\delta}_i) \quad \text{Eq. 27}$$

$$25 \quad \text{ZILN}(\mathbf{x}_{k_i}; \boldsymbol{\mu}_i, \boldsymbol{\sigma}_i, \boldsymbol{\delta}_i) = \begin{cases} \frac{1 - \boldsymbol{\delta}_i}{\mathbf{x}_{k_i} \boldsymbol{\sigma}_i \sqrt{2\pi}} \exp\left(-\frac{(\log \mathbf{x}_{k_i} - \boldsymbol{\mu}_i)^2}{2\boldsymbol{\sigma}_i^2}\right), & \mathbf{x}_{k_i} > 0 \\ \boldsymbol{\delta}_i, & \mathbf{x}_{k_i} = 0 \end{cases} \quad \text{Eq. 28}$$

$$26 \quad \boldsymbol{\mu}_i = \boldsymbol{\alpha} \odot \mathbf{V}_k^T \mathbf{u} + \boldsymbol{\beta} \quad \text{Eq. 29}$$

27 where $\boldsymbol{\mu} \in \mathbb{R}^{|\mathcal{V}_k|}$, $\boldsymbol{\sigma} \in \mathbb{R}_+^{|\mathcal{V}_k|}$, $\boldsymbol{\delta} \in (0, 1)^{|\mathcal{V}_k|}$ are the log-scale mean, log-scale standard deviation and
 28 zero-inflation parameters of the ZILN distribution, respectively, and $\boldsymbol{\alpha} \in \mathbb{R}_+^{|\mathcal{V}_k|}$, $\boldsymbol{\beta} \in \mathbb{R}^{|\mathcal{V}_k|}$ are scaling
 29 and bias factors.

1

2 To unify the cell type labels, we performed a nearest neighbor-based label transfer with the snmC-
3 seq dataset as a reference. The 5 nearest neighbors in snmC-seq were identified for each scRNA-seq
4 and scATAC-seq cell in the aligned embedding space, and majority voting was used to determine the
5 transferred label. To verify whether the alignment was correct, we tested for significant overlap in
6 cell type marker genes. The features of all omics layers were first converted to genes. Then, for each
7 omics layer, the cell type markers were identified using the one-vs.-rest Wilcoxon rank-sum test with
8 the following criteria: FDR < 0.05 and log-fold change > 0 for scRNA-seq/scATAC-seq; FDR <
9 0.05 and log-fold change < 0 for snmC-seq. The significance of marker overlap was determined by
10 the three-way Fisher's exact test³⁷.

11

12 To perform correlation and regression analysis after the integration, we clustered all cells from the
13 three omics layers using fine-scale k-means (k = 200). Then, for each omics layer, the cells in each
14 cluster were aggregated into a metacell by summing their expression/accessibility counts or
15 averaging their DNA methylation levels. Therefore, the metacells were established as paired
16 samples, based on which feature correlation and regression analyses could be conducted.

17

18 **Model-based cis-regulatory inference**

19 To ensure consistency of cell types, we first selected the overlapping cell types between the 10x
20 Multiome and pcHi-C data. The remaining cell types included T cells, B cells and monocytes. The
21 eQTL data were used as is, because they were not cell type-specific. For scRNA-seq, we selected
22 6,000 highly variable genes. For the initial regulatory inference, the guidance graph was constructed
23 by connecting RNA genes with ATAC peaks within 150 kb of the gene promoters (defined as 2 kb
24 upstream from the TSS); the graph was weighted by a power-law function $w = (d + 1)^{-0.75}$ (d is
25 the genomic distance in kb), which has been proposed to model the probability of chromatin
26 contact^{39, 40}.

27

28 To incorporate the regulatory evidence of pcHi-C and eQTL, we anchored all evidence to that
29 between the ATAC peaks and RNA genes. A peak-gene pair was considered supported by pcHi-C if
30 (1) the gene promoter was within 1 kb of a bait fragment, (2) the peak was within 1 kb of an other-
31 end fragment, and (3) significant contact was identified between the bait and the other-end fragment
32 in pcHi-C. The pcHi-C-supported peak-gene interactions were weighted by multiplying the
33 promoter-to-bait and the peak-to-other-end power-law weights (see above). If a peak-gene pair was

1 supported by multiple pcHi-C contacts, the weights were summed and clipped at a maximum of 1. A
2 peak-gene pair was considered supported by eQTL if (1) the peak overlapped an eQTL locus and (2)
3 the locus was associated with the expression of the gene. The eQTL-supported peak-gene
4 interactions were assigned weights of 1. The composite guidance graph was constructed by adding
5 the pcHi-C- and eQTL-supported interactions to the previous distance-based interactions, allowing
6 for multi-edges.

7
8 For regulatory inference, only peak-gene pairs within 150 kb in distance were considered. The
9 GLUE training process was repeated 4 times with different random seeds. For each repeat, the peak-
10 gene regulatory score was computed as the cosine similarity between the feature embeddings. The
11 final regulatory inference was obtained by averaging the regulatory scores across the 4 repeats.

12 13 **TF-target gene regulatory inference**

14 We employed the SCENIC workflow⁶⁵ to construct a TF-gene regulatory network from the inferred
15 peak-gene regulatory interactions. Briefly, the SCENIC workflow first constructs a gene
16 coexpression network based on the scRNA-seq data, and then uses external cis-regulatory evidence
17 to filter out false positives. SCENIC accepts cis-regulatory evidence in the form of gene rankings per
18 TF, i.e., genes with higher TF enrichment levels in their regulatory regions are ranked higher. To
19 construct the rankings based on our inferred peak-gene interactions, we first overlapped the
20 ENCODE TF ChIP peaks⁶⁶ with the ATAC peaks and counted the number of ChIP peaks for each
21 TF in each ATAC peak. Since different genes can have different numbers of connected ATAC
22 peaks, and the ATAC peaks vary in length (longer peaks can contain more ChIP peaks by chance),
23 we devised a sampling-based approach to evaluate TF enrichment. Specifically, for each gene, we
24 randomly sampled 1,000 sets of ATAC peaks that matched the connected ATAC peaks in both
25 number and length distribution. We counted the numbers of TF ChIP peaks in these random ATAC
26 peaks as null distributions. For each TF in each gene, an empirical P value could then be computed
27 by comparing the observed number of ChIP peaks to the null distribution. Finally, we ranked the
28 genes by the empirical P values for each TF, producing the cis-regulatory rankings used by SCENIC.
29 Since peak-gene-based inference is mainly focused on remote regulatory regions, proximal
30 promoters could be missed. As such, we provided SCENIC with both the above peak-based and
31 proximal promoter-based cis-regulatory rankings.

32

1 **Integration for the human multi-omics atlas**

2 The scRNA-seq and scATAC-seq atlases have highly unbalanced cell type compositions, which is
3 primarily caused by differences in organ sampling sizes (Supplementary Fig. 13b). Although cell
4 types are unknown during real-world analyses, organ sources are typically available and can be
5 utilized to help balance the integration process. To perform organ-balanced data preprocessing, we
6 first subsampled each omics layer to match the organ compositions. For the scRNA-seq data, 4,000
7 highly variable genes were selected using the organ-balanced subsample. Then, for the initial
8 dimensionality reduction, we fitted PCA (scRNA-seq) and LSI (scATAC-seq) on the organ-balanced
9 subsample and applied the projection to the full data. The PCA/LSI coordinates were used as the first
10 transformation layer in the GLUE data encoders (see the section “Implementation details”), as well
11 as for metacell aggregation (see below). The guidance graph was constructed as described previously
12 (see the section “Systematic benchmarks”).

13
14 The two atlases consist of large numbers of cells but with low coverage per cell. To alleviate dropout
15 and increase the training speed simultaneously, we designed a multistage training strategy, where the
16 GLUE model was pretrained on aggregated metacells and then fine-tuned on the original single cells.
17 Specifically, in the first stage, we clustered the cells in each omics layer using fine-scaled k-means (k
18 = 100,000 for scRNA-seq and $k = 40,000$ for scATAC-seq). To balance the organ compositions at
19 the same time, k-means centroids were fitted on the previous organ-balanced subsample and then
20 applied to the full data. The cells in each k-means cluster were aggregated into a metacell by
21 summing their expression/accessibility counts and averaging their PCA/LSI coordinates. GLUE was
22 then pretrained on the aggregated metacells without adversarial alignment, which roughly oriented
23 the cell embeddings but did not actually align them. To better utilize the large data size, the hidden
24 layer dimensionality was doubled to 512 from the default 256.

25
26 In the second stage, GLUE was fine-tuned on the full single-cell data with weighted adversarial
27 alignment (see below). As shown in previous work²⁸, pure adversarial alignment amounts to
28 minimizing a generalized form of Jensen–Shannon divergence among the cell embedding
29 distributions of different omics layers:

$$30 \quad \frac{1}{K} \sum_{k=1}^K \text{KL} \left(q_k(\mathbf{u}) \parallel \frac{1}{K} \sum_{k=1}^K q_k(\mathbf{u}) \right) \quad \text{Eq. 30}$$

31 where $q_k(\mathbf{u}) = \mathbb{E}_{\mathbf{x}_k \sim p_{\text{data}}(\mathbf{x}_k)} q(\mathbf{u} | \mathbf{x}_k; \phi_k)$ represents the marginal cell embedding distribution of the
32 k^{th} layer. Without other loss terms, Eq. 30 converges at perfect alignment, i.e., when $q_i(\mathbf{u}) =$

1 $q_j(\mathbf{u}), \forall i \neq j$. This can be problematic when cell type compositions differ dramatically across
 2 different layers, as was the case here. To address this issue, we added cell-specific weights $w^{(n)}$ to
 3 the discriminator loss in Eq. 18:

$$4 \quad \mathcal{L}_D(\phi, \psi) = -\frac{1}{K} \sum_{k=1}^K \frac{1}{W_k} \sum_{n=1}^{N_k} w^{(n)} \cdot \mathbb{E}_{\mathbf{u} \sim q(\mathbf{u}|\mathbf{x}_k^{(n)}; \phi_k)} \log D_k(\mathbf{u}; \psi) \quad \text{Eq. 31}$$

5 where the normalizer $W_k = \sum_{n=1}^{N_k} w^{(n)}$. The adversarial alignment still amounts to minimizing

6 Eq. 30 but with weighted marginal cell embedding distributions $q_k(\mathbf{u}) =$

7 $\frac{1}{W_k} \sum_{n=1}^{N_k} w^{(n)} q(\mathbf{u}|\mathbf{x}_k^{(n)}; \phi_k)$. By assigning appropriate weights to balance the cell distributions across

8 different layers, the optimum of $q_i(\mathbf{u}) = q_j(\mathbf{u}), \forall i \neq j$ could be much closer to the desired

9 alignment. For example, a viable choice for the cell weights is organ-balancing weights. Suppose

10 that the organ proportions in scRNA-seq and scATAC-seq are f_1, f_2, \dots, f_O and g_1, g_2, \dots, g_O (O is

11 the number of organs, $\sum_{i=1}^O f_i = 1, \sum_{i=1}^O g_i = 1$), respectively. We can weight the RNA cells in the

12 i^{th} organ by $\sqrt{g_i/f_i}$ and weight the ATAC cells in the i^{th} organ by $\sqrt{f_i/g_i}$, so that the i^{th} organ has a

13 balanced accumulative contribution of $\sqrt{f_i g_i}$, regardless of omics layers. However, in case there are

14 major differences among the cell type compositions within the same organ, organ-level balancing

15 can be insufficient. As such, we designed a method to compute cell type-level balancing weights in

16 an unsupervised manner. For each omics layer, we clustered the cell embeddings using Louvain

17 clustering and matched the clusters in different layers via cosine similarity. The population size of

18 each cluster was distributed to its matched counterparts by cosine similarity. Subsequently, for each

19 omics layer, we could obtain the proportion of each cluster f_i and its effective proportion g_i in the

20 opposite layer (by normalizing the effective population size received from the opposite layer).

21 Balancing weights could then be computed as above. The weight-balanced alignment proved

22 effective in aligning the highly skewed data (Fig. 5).

23

24 For a comparison with other integration methods, we also tried online iNMF and Seurat v3. Online

25 iNMF was the only other method that could scale to millions of cells, so we applied it to the full

26 dataset. On the other hand, Seurat v3 showed the second-best accuracy in our previous benchmark.

27 We also managed to apply it to the aggregated data as used in the first stage of GLUE training, due

28 to the fact that Seurat v3 could not scale to the full dataset (Supplementary Fig. 13a).

1 Data availability

2 All datasets used in this study are already published and were obtained from public data repositories.
3 See Supplementary Table 1 for detailed information, including access codes and URLs. All
4 benchmark data are available in Supplementary Table 2.

6 Code availability

7 The GLUE framework was implemented in the “scglue” Python package, which is available at
8 <https://github.com/gao-lab/GLUE>. For reproducibility, the scripts for all benchmarks and case
9 studies were assembled using Snakemake, which is also available in the above repository.

11 Acknowledgments

12 The authors would like to thank Drs. Zemin Zhang, Xiaoliang Sunney Xie, Letian Tao, Cheng Li,
13 Jian Lu (at Peking University) and Yang Ding (at the Beijing Institute of Radiation Medicine) for
14 their helpful discussions and comments during the study, as well as authors of the datasets used in
15 this work for their kindly help.

16 This work was supported by funds from the National Key Research and Development Program
17 (2016YFC0901603), the China 863 Program (2015AA020108), and the State Key Laboratory of
18 Protein and Plant Gene Research and the Beijing Advanced Innovation Center for Genomics (ICG)
19 at Peking University. The research of G.G. was supported in part by the National Program for
20 Support of Top-notch Young Professionals.

21 Part of the analysis was performed on the Computing Platform of the Center for Life Sciences of
22 Peking University and supported by the High-performance Computing Platform of Peking
23 University.

25 Author contributions

26 G.G. conceived the study and supervised the research; Z.J.C. designed and implemented the
27 computational framework and conducted benchmarks and case studies, with guidance from G.G.;
28 Z.J.C. and G.G. wrote the manuscript.

1

2 **Competing interests**

3 The authors declare no competing interests.

1 References

- 2 1. Cusanovich, D.A. et al. Multiplex single cell profiling of chromatin accessibility by
3 combinatorial cellular indexing. *Science* **348**, 910–914 (2015).
- 4 2. Chen, X., Miragaia, R.J., Natarajan, K.N. & Teichmann, S.A. A rapid and robust method for
5 single cell chromatin accessibility profiling. *Nat. Commun.* **9**, 5345 (2018).
- 6 3. Luo, C. et al. Single-cell methylomes identify neuronal subtypes and regulatory elements in
7 mammalian cortex. *Science* **357**, 600–604 (2017).
- 8 4. Mulqueen, R.M. et al. Highly scalable generation of DNA methylation profiles in single
9 cells. *Nat. Biotechnol.* **36**, 428–431 (2018).
- 10 5. Picelli, S. et al. Smart-seq2 for sensitive full-length transcriptome profiling in single cells.
11 *Nat. Methods* **10**, 1096–1098 (2013).
- 12 6. Zheng, G.X. et al. Massively parallel digital transcriptional profiling of single cells. *Nat.*
13 *Commun.* **8**, 14049 (2017).
- 14 7. Packer, J. & Trapnell, C. Single-cell multi-omics: An engine for new quantitative models of
15 gene regulation. *Trends Genet.* **34**, 653–665 (2018).
- 16 8. Chen, S., Lake, B.B. & Zhang, K. High-throughput sequencing of the transcriptome and
17 chromatin accessibility in the same cell. *Nat. Biotechnol.* **37**, 1452–1457 (2019).
- 18 9. Ma, S. et al. Chromatin potential identified by shared single-cell profiling of RNA and
19 chromatin. *Cell* **183**, 1103–1116 (2020).
- 20 10. Clark, S.J. et al. scNMT-seq enables joint profiling of chromatin accessibility DNA
21 methylation and transcription in single cells. *Nat. Commun.* **9**, 781 (2018).
- 22 11. Wang, Y. et al. Single-cell multiomics sequencing reveals the functional regulatory landscape
23 of early embryos. *Nat. Commun.* **12**, 1247 (2021).
- 24 12. Lake, B.B. et al. Integrative single-cell analysis of transcriptional and epigenetic states in the
25 human adult brain. *Nat. Biotechnol.* **36**, 70–80 (2018).
- 26 13. Bravo Gonzalez-Blas, C. et al. Identification of genomic enhancers through spatial
27 integration of single-cell transcriptomics and epigenomics. *Mol. Syst. Biol.* **16**, e9438 (2020).
- 28 14. Argelaguet, R., Cuomo, A.S.E., Stegle, O. & Marioni, J.C. Computational principles and
29 challenges in single-cell data integration. *Nat. Biotechnol.* (2021).
- 30 15. Stuart, T. et al. Comprehensive integration of single-cell data. *Cell* **177**, 1888–1902 (2019).
- 31 16. Gao, C. et al. Iterative single-cell multi-omic integration using online learning. *Nat.*
32 *Biotechnol.* **39**, 1000–1007 (2021).
- 33 17. Welch, J.D. et al. Single-cell multi-omic integration compares and contrasts features of brain
34 cell identity. *Cell* **177**, 1873–1887 (2019).
- 35 18. Chen, H. et al. Assessment of computational methods for the analysis of single-cell ATAC-
36 seq data. *Genome Biol.* **20**, 241 (2019).
- 37 19. Duren, Z. et al. Integrative analysis of single-cell genomics data by coupled nonnegative
38 matrix factorizations. *Proc. Natl. Acad. Sci. U. S. A.* **115**, 7723–7728 (2018).
- 39 20. Zeng, W. et al. DC3 is a method for deconvolution and coupled clustering from bulk and
40 single-cell genomics data. *Nat. Commun.* **10**, 4613 (2019).
- 41 21. Cao, K., Bai, X., Hong, Y. & Wan, L. Unsupervised topological alignment for single-cell
42 multi-omics integration. *Bioinformatics* **36**, i48–i56 (2020).
- 43 22. Demetci, P., Santorella, R., Sandstede, B., Noble, W.S. & Singh, R. Gromov-Wasserstein
44 optimal transport to align single-cell multi-omics data. Preprint at
45 <https://www.biorxiv.org/content/10.1101/2020.04.28.066787> (2020).
- 46 23. Svensson, V., Vento-Tormo, R. & Teichmann, S.A. Exponential scaling of single-cell RNA-
47 seq in the past decade. *Nat. Protoc.* **13**, 599–604 (2018).

- 1 24. Kozareva, V. et al. A transcriptomic atlas of the mouse cerebellum reveals regional
2 specializations and novel cell types. Preprint at
3 <https://www.biorxiv.org/content/10.1101/2020.03.04.976407> (2020).
- 4 25. Cao, J. et al. A human cell atlas of fetal gene expression. *Science* **370**, eaba7721 (2020).
- 5 26. Domcke, S. et al. A human cell atlas of fetal chromatin accessibility. *Science* **370**, eaba7612
6 (2020).
- 7 27. Lopez, R., Regier, J., Cole, M.B., Jordan, M.I. & Yosef, N. Deep generative modeling for
8 single-cell transcriptomics. *Nat. Methods* **15**, 1053–1058 (2018).
- 9 28. Cao, Z.J., Wei, L., Lu, S., Yang, D.C. & Gao, G. Searching large-scale scRNA-seq databases
10 via unbiased cell embedding with Cell BLAST. *Nat. Commun.* **11**, 3458 (2020).
- 11 29. Kipf, T.N. & Welling, M. Variational graph auto-encoders. Preprint at
12 <https://arxiv.org/abs/1611.07308> (2016).
- 13 30. Dou, J. et al. Unbiased integration of single cell multi-omics data. Preprint at
14 <https://www.biorxiv.org/content/10.1101/2020.12.11.422014> (2020).
- 15 31. 10x Genomics. PBMC from a healthy donor, single cell multiome ATAC gene expression
16 demonstration data by Cell Ranger ARC 1.0.0. [https://support.10xgenomics.com/single-cell-](https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/1.0.0/pbmc_granulocyte_sorted_10k)
17 [multiome-atac-gex/datasets/1.0.0/pbmc_granulocyte_sorted_10k](https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/1.0.0/pbmc_granulocyte_sorted_10k) (2020).
- 18 32. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell
19 transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.*
20 **36**, 411–420 (2018).
- 21 33. Singh, R. et al. Unsupervised manifold alignment for single-cell multi-omics data. In
22 *Proceedings of the 11th ACM International Conference on Bioinformatics, Computational*
23 *Biology and Health Informatics*. (ACM, Virtual Event, USA, 2020).
- 24 34. Saunders, A. et al. Molecular diversity and specializations among the cells of the adult mouse
25 brain. *Cell* **174**, 1015–1030 (2018).
- 26 35. 10x Genomics. Fresh cortex from adult mouse brain (v1), single cell ATAC demonstration
27 data by Cell Ranger 1.1.0. [https://support.10xgenomics.com/single-cell-](https://support.10xgenomics.com/single-cell-atac/datasets/1.1.0/atac_v1_adult_brain_fresh_5k)
28 [atac/datasets/1.1.0/atac_v1_adult_brain_fresh_5k](https://support.10xgenomics.com/single-cell-atac/datasets/1.1.0/atac_v1_adult_brain_fresh_5k) (2019).
- 29 36. Mo, A. et al. Epigenomic signatures of neuronal diversity in the mammalian brain. *Neuron*
30 **86**, 1369–1384 (2015).
- 31 37. Wang, M., Zhao, Y. & Zhang, B. Efficient test and visualization of multi-set intersections.
32 *Sci. Rep.* **5**, 16923 (2015).
- 33 38. Gabel, H.W. et al. Disruption of DNA-methylation-dependent long gene repression in Rett
34 syndrome. *Nature* **522**, 89–93 (2015).
- 35 39. Dekker, J., Marti-Renom, M.A. & Mirny, L.A. Exploring the three-dimensional organization
36 of genomes: Interpreting chromatin interaction data. *Nat. Rev. Genet.* **14**, 390–403 (2013).
- 37 40. Pliner, H.A. et al. Cicero predicts cis-regulatory DNA interactions from single-cell chromatin
38 accessibility data. *Mol. Cell* **71**, 858–871 (2018).
- 39 41. Javierre, B.M. et al. Lineage-specific genome architecture links enhancers and non-coding
40 disease variants to target gene promoters. *Cell* **167**, 1369–1384 (2016).
- 41 42. Aguet, F. et al. Genetic effects on gene expression across human tissues. *Nature* **550**, 204–
42 213 (2017).
- 43 43. Han, H. et al. TRRUST v2: An expanded reference database of human and mouse
44 transcriptional regulatory interactions. *Nucleic Acids Res.* **46**, D380–D386 (2018).
- 45 44. Thomsen, E.R. et al. Fixed single-cell transcriptomic characterization of human radial glial
46 diversity. *Nat. Methods* **13**, 87–93 (2016).
- 47 45. Pollen, Alex A. et al. Molecular identity of human outer radial glia during cortical
48 development. *Cell* **163**, 55–67 (2015).
- 49 46. Tran, H.T.N. et al. A benchmark of batch-effect correction methods for single-cell RNA
50 sequencing data. *Genome Biol.* **21**, 12 (2020).

- 1 47. Stark, S.G. et al. SCIM: Universal single-cell matching with unpaired feature sets.
2 *Bioinformatics* **36**, i919–i927 (2020).
- 3 48. Yang, K.D. et al. Multi-domain translation between single-cell imaging and sequencing data
4 using autoencoders. *Nat. Commun.* **12**, 31 (2021).
- 5 49. Eng, C.-H.L. et al. Transcriptome-scale super-resolved imaging in tissues by RNA seqfish+.
6 *Nature* **568**, 235–239 (2019).
- 7 50. Rodriques, S.G. et al. Slide-seq: A scalable technology for measuring genome-wide
8 expression at high spatial resolution. *Science* **363**, 1463–1467 (2019).
- 9 51. Bandura, D.R. et al. Mass cytometry: Technique for real time single cell multitarget
10 immunoassay based on inductively coupled plasma time-of-flight mass spectrometry. *Anal.*
11 *Chem.* **81**, 6813–6822 (2009).
- 12 52. Bartosovic, M., Kabbe, M. & Castelo-Branco, G. Single-cell CUT&Tag profiles histone
13 modifications and transcription factors in complex tissues. *Nat. Biotechnol.* **39**, 825–835
14 (2021).
- 15 53. Ashuach, T., Reidenbach, D.A., Gayoso, A. & Yosef, N. PeakVI: A deep generative model
16 for single cell chromatin accessibility analysis. Preprint at
17 <https://www.biorxiv.org/content/10.1101/2021.04.29.442020> (2021).
- 18 54. Hamilton, W., Ying, Z. & Leskovec, J. Inductive representation learning on large graphs. In
19 Advances in Neural Information Processing Systems. (eds. I. Guyon et al.) 1024–1034
20 (Curran Associates, Inc., Long Beach, CA, USA, 2017).
- 21 55. Veličković, P. et al. Graph attention networks. Preprint at <https://arxiv.org/abs/1710.10903>
22 (2017).
- 23 56. Vashishth, S., Sanyal, S., Nitin, V. & Talukdar, P. Composition-based multi-relational graph
24 convolutional networks. In *Proceedings of the 8th International Conference on Learning*
25 *Representations*. (Addis Ababa, Ethiopia, 2020).
- 26 57. Stuart, T. & Satija, R. Integrative single-cell analysis. *Nat. Rev. Genet.* **20**, 257–272 (2019).
- 27 58. Amodio, M. & Krishnaswamy, S. MAGAN: Aligning biological manifolds. In *Proceedings*
28 *of the 35th International Conference on Machine Learning*. (eds. J.G. Dy & A. Krause) 215–
29 223 (PMLR, Stockholm, Sweden, 2018).
- 30 59. Tarashansky, A.J. et al. Mapping single-cell atlases throughout metazoa unravels cell type
31 evolution. *eLife* **10**, e66747 (2021).
- 32 60. Ding, J. & Regev, A. Deep generative model embedding of single-cell RNA-seq profiles on
33 hyperspheres and hyperbolic spaces. *Nat. Commun.* **12**, 2554 (2021).
- 34 61. Mikolov, T., Sutskever, I., Chen, K., Corrado, G. & Dean, J. Distributed representations of
35 words and phrases and their compositionality. In Advances in Neural Information Processing
36 Systems. (eds. C.J.C. Burges, L. Bottou, M. Welling, Z. Ghahramani & K.Q. Weinberger)
37 3111–3119 (Curran Associates, Inc., Lake Tahoe, NV, USA, 2013).
- 38 62. Kipf, T.N. & Welling, M. Semi-supervised classification with graph convolutional networks.
39 In *Proceedings of the 5th International Conference on Learning Representations*. (Toulon,
40 France, 2017).
- 41 63. Dincer, A.B., Janizek, J.D. & Lee, S.-I. Adversarial deconfounding autoencoder for learning
42 robust gene expression embeddings. *Bioinformatics* **36**, i573–i582 (2020).
- 43 64. Goodfellow, I. et al. Generative adversarial nets. In Advances in Neural Information
44 Processing Systems. (eds. Z. Ghahramani, M. Welling, C. Cortes, N.D. Lawrence & K.Q.
45 Weinberger) 2672–2680 (Curran Associates, Inc., Montreal, Quebec, Canada, 2014).
- 46 65. Aibar, S. et al. SCENIC: Single-cell regulatory network inference and clustering. *Nat.*
47 *Methods* **14**, 1083–1086 (2017).
- 48 66. Davis, C.A. et al. The encyclopedia of DNA elements (ENCODE): Data portal update.
49 *Nucleic Acids Res.* **46**, D794–D801 (2018).