

# Comprehensive View Embedding Learning for Single-cell Multimodal Integration

Zhenchao Tang<sup>1</sup>, Jiehui Huang<sup>1</sup>, Guanxing Chen<sup>1</sup>, Calvin Yu-Chian Chen<sup>1,2,3,4,5\*</sup>

<sup>1</sup>School of Intelligent Systems Engineering, Shenzhen Campus of Sun Yat-sen University

<sup>2</sup>AI for Science (AI4S) - Preferred Program, Peking University Shenzhen Graduate School

<sup>3</sup>School of Electronic and Computer Engineering, Peking University Shenzhen Graduate School

<sup>4</sup>Department of Medical Research, China Medical University Hospital

<sup>5</sup>Department of Bioinformatics and Medical Engineering, Asia University

## Abstract

**Motivation:** Advances in single-cell measurement techniques provide rich multimodal data, which helps us to explore the life state of cells more deeply. However, multimodal integration, or, learning joint embeddings from multimodal data remains a current challenge. The difficulty in integrating unpaired single-cell multimodal data is that different modalities have different feature spaces, which easily leads to information loss in joint embedding. And few existing methods have fully exploited and fused the information in single-cell multimodal data. **Result:** In this study, we propose CoVEL, a deep learning method for unsupervised integration of single-cell multimodal data. CoVEL learns single-cell representations from a comprehensive view, including regulatory relationships between modalities, fine-grained representations of cells, and relationships between different cells. The comprehensive view embedding enables CoVEL to remove the gap between modalities while protecting biological heterogeneity. Experimental results on multiple public datasets show that CoVEL is accurate and robust to single-cell multimodal integration. **Data availability:** <https://github.com/shapsider/scintegration>.

## Introduction

Single-cell measurement technology can measure the abundance of molecules at the cellular level, which allows us to understand the life state of organisms more clearly and screen drugs to target diseases (Lartigue et al. 2007; Lv et al. 2023; Tang et al. 2023). Single-cell measurement technology has developed rapidly in the past 10 years. For example, Eberwine et al. measured RNA content at the single-cell level for the first time (Eberwine et al. 1992). Tang et al. further introduced single-cell transcriptome sequencing (Tang et al. 2009). Klein et al. and Macosko et al. used droplet technology to expand the measurement scale (Klein et al. 2015; Macosko et al. 2015), CITE-seq simultaneously measures RNA expression and cell surface protein abundance (Stoeckius et al. 2017), sci-CAR and SNARE-seq can jointly measure RNA expression and chromatin accessibility (Cao et al. 2018; Chen, Lake, and Zhang 2019). Joint measurement methods can provide single-cell multimodal data and bring valuable insights into the overall understanding of biological systems. However,

single-cell multimodal (multi-omics) data often consists of unpaired cells with mismatched features across modalities (see Figure 1), making single-cell multimodal data analysis a huge challenge.

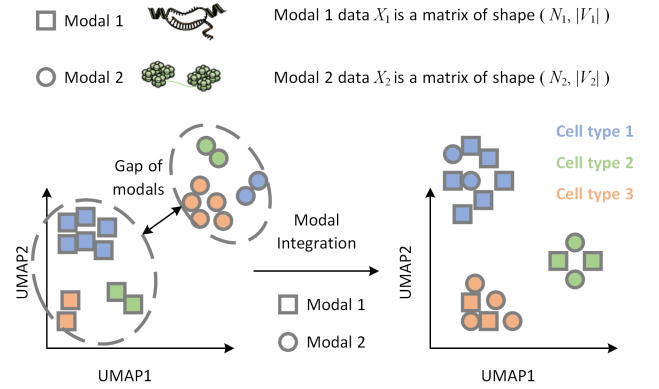


Figure 1: The goal of multimodal integration is to eliminate gaps (enhance batch remove metric) and preserve biological heterogeneity (cluster cells of the same type together, enhance biology conservation metric).

In Figure 1, the problem description of single-cell multimodal integration is: single-cell data is stored in a matrix, rows represent cells (i.e. samples), columns represent features. Different modalities correspond to different feature spaces, and there is no one-to-one correspondence between samples in different modalities (the same cell does not exist in modal 1 and modal 2, but cells of the same type may exist in modal 1 and modal 2). For the same type of cells with different modalities, since the feature spaces of two modalities are completely different, there is a significant gap in the embeddings after direct dimensionality reduction. The goal of multimodal integration is to eliminate gaps (enhance batch remove metric) and preserve biological heterogeneity (cluster cells of the same type together, enhance biology conservation metric).

A major obstacle faced when integrating unpaired single-cell multimodal data is that different modalities have different feature spaces, such as accessible chromatin regions in the ATAC modality and genes in the RNA modality. A common

\*Corresponding author, [chenyuchian@mail.sysu.edu.cn](mailto:chenyuchian@mail.sysu.edu.cn).

approach is to project multimodal data into a common feature space and then correct for differences between modalities, but early projection operations can lead to information loss (Stuart et al. 2019; Welch et al. 2019; Korsunsky et al. 2019). Algorithms based on matrix factorization can avoid projection operations, but cannot handle more than two omics data (Gao et al. 2021). There are also methods to match data from different modalities via nonlinear manifold alignment, which can reduce information loss between modalities, however, they cannot align uniformly distributed data (Cao et al. 2020; Cao, Hong, and Wan 2022). Recently, multimodal deep learning methods have been applied to bioinformatics field (Baltrušaitis, Ahuja, and Morency 2018; Steyaert et al. 2023). However, few existing methods have fully exploited and fused the information in single-cell multimodal data. CLUE designs multiple auxiliary tasks to supervise VAE to capture the relationship between single-cell multi-omics data (Tu et al. 2022). GLUE achieves multimodal integration by building a guidance graph, that is, to obtain joint embedding (Cao and Gao 2022). The guidance graph can learn the regulatory relationship between modalities. However, they lack the exploration of cell representation learning, because from the field of single cell representation learning, scBERT uses Transformer to learn cell fine-grained features (Yang et al. 2022a), and Concerto uses contrastive learning and builds relationships between cells (Yang et al. 2022b). Note that some single-cell representation learning methods target RNA modality, and these methods cannot handle multimodal data. But their ideas are worth considering, such as fine-grained representation learning and contrastive learning have achieved excellent performance in RNA modality. From the embedding point of view, all the above methods learn representations of three views respectively: regulatory between modalities, cell fine-grained representation, and contrastive learning between cells. However, neither current single-cell representation learning methods nor single-cell multimodal integration methods consider fusing these three view representations.

In this study, we propose a **Comprehensive View Embedding Learning** method CoVEL. CoVEL learns embedding from three views on single-cell multimodal data, such as: the regulatory relationship between different modalities, named graph-linked embedding; the relationship between single-cell fine-grained features in each modality, named single-cell fine-grained embedding; the representation of a single cell sample in each modality, named contrastive cell embedding. By learning graph-linked embedding, CoVEL can model the regulatory relationship across modalities, and bridge the gap between feature spaces under different modalities with biological knowledge. CoVEL combines generative methods and contrastive learning to unsupervisedly learn single-cell fine-grained embedding and contrastive cell embedding on multimodal data, ensuring that the gap between modalities is removed and biological heterogeneity is preserved. From the perspective of representation learning, contrastive learning tries to find information between data, while generative methods focus on learning information within data (Liu et al. 2021). Therefore, CoVEL is a comprehensive self-supervised learning method. Experimental results on public

datasets demonstrate that CoVEL is accurate and robust to single-cell multimodal integration.

## Related work

### Single-cell representation learning

Single-cell representation learning can be divided into three views: regulatory between modalities, cell fine-grained representation, and contrastive learning between cells. scGCN and HGT learn high-order relationships between cells based on graphs and infer biological networks (Song, Su, and Zhang 2021; Ma et al. 2023). scBERT leverages large-scale language models to learn fine-grained representations of single-cell transcriptomes (RNA modality) (Yang et al. 2022a). TOSICA combines pathway information with Transformer to learn fine-grained representations and assists in the discovery of markers (Chen et al. 2023). scMDC optimized VAE with ZIBN, also obtained good fine-grained cell representations (Lin et al. 2022). Concerto and scDML utilize contrastive learning to eliminate batch effects in the single-cell transcriptome (RNA modality) embedding space and avoid the loss of original rare cell type information (Yang et al. 2022b; Yu et al. 2023). Some single-cell representation learning methods target RNA modality, and these methods cannot handle multimodal data. But their ideas are worth considering, such as fine-grained representation learning and contrastive learning have achieved excellent performance in RNA modality. Therefore, CoVEL combines the three types of representation learning to integrate multimodal data according to multi-view context.

### Single-cell multimodal integration

Single-cell multimodal integration methods can be classified into semi-supervised and unsupervised. Semi-supervised methods utilize paired multimodal observations. Seurat v4 uses paired multimodal observations to interconnect other single-modal data (Hao et al. 2022). MultiVI uses multiple VAEs to learn multimodal data and align joint embedding (Ashuach et al. 2021). CLUE introduces auxiliary tasks to establish associations in multimodal spaces (Tu et al. 2022). scMoGNN extracts cell-modal connections with GNNs and uses known single-modal embedding to learn multimodal joint embedding (Wen et al. 2022). Unsupervised methods are used to integrate multimodal measurements without any paired information. Seurat v3 (Stuart et al. 2019), LIGER (Welch et al. 2019), and Harmony (Korsunsky et al. 2019) perform feature transformation first, and then correct the differences between modalities in the common space, but feature transformation will cause significant information loss. The matrix factorization-based method iNMF (Gao et al. 2021) can avoid the projection operation, but cannot handle more than two omics data. bindSC is limited by RNA and ATAC modalities and lacks scalability (Dou et al. 2022). UnionCom (Cao et al. 2020) and Pamona (Cao, Hong, and Wan 2022) match data from different modalities via nonlinear manifold alignment, which cannot align uniformly distributed data and are limited by the number of cell types. Scanorama (Hie, Bryson, and Berger 2019) and GLUE (Cao and Gao 2022)

are graph-based methods, Scanorama is suitable for batch integration, and GLUE requires additional expert knowledge as well as complex training strategies. CoVEL is not limited to the number of cell types and complex training strategies, and supports unsupervised integration of three and more modalities.

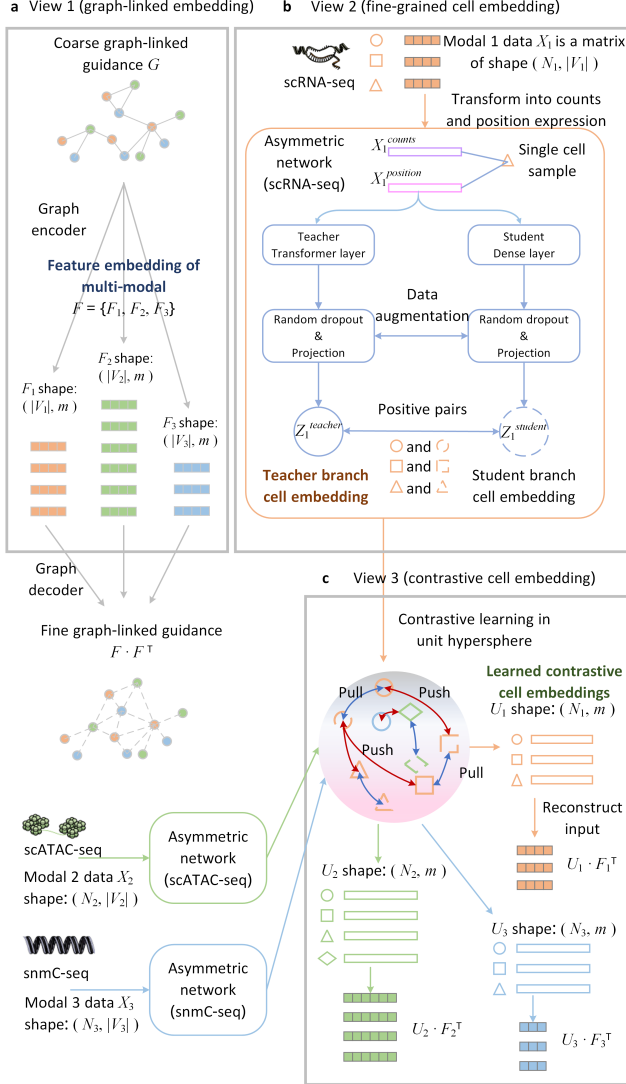


Figure 2: The overall architecture of CoVEL, and view 3 embedding is the embedding of cells in the multimodal joint space. **a**: graph-linked embedding learning (view 1 embedding, blue bold). **b**: fine-grained cell embedding learning (view 2 embedding, orange bold). **c**: contrastive cell embedding learning (view 3 embedding, green bold).

## Method

For single-cell multimodal data, CoVEL learning includes embeddings from three views, graph-linked embedding, which describes the relationship between different modalities; fine-grained cell embedding, which describes the fine-grained representations of cell in a specified modality; contrastive

cell embedding, which is an optimized cell representation based on contrastive learning. We will introduce the learning details of three embeddings respectively.

### Graph-linked embedding learning

Assume there are  $K$  modalities, each modality has a different feature set  $V_k$ ,  $k = 1, \dots, K$ . For RNA modality, the feature set is a collection of genes, and for ATAC modality, the feature set is a collection of chromatin regions. The data for mode  $k$  is  $X_k \subseteq \mathbb{R}^{N_k \times |V_k|}$ ,  $N_k$  is the number of cells. Let  $x_k^{(n)} \in X_k$ ,  $n = 1, 2, \dots, N_k$  represents a cell in the modal  $k$ . We follow GLUE (Cao and Gao 2022) and use the features of the modality as nodes to build a Coarse graph-linked guidance  $G = (\{V_1, \dots, V_K\}, E)$ . The guidance graph has  $\sum_{k=1}^K |V_k|$  nodes,  $E$  is the edge set, which is the interaction between different modal features. Edge  $(i, j)$ ,  $i, j \in \{V_1, \dots, V_K\}$  has weight  $w_{ij} \in (0, 1]$  represents the interaction confidence, and edge  $(i, j)$  has sign  $s_{ij} \in \{-1, 1\}$ . For example ATAC often positively regulate gene ( $s_{ij} = 1$ ,  $i, j$  represents gene and the corresponding chromatin region respectively). The construction of guidance graph is easily available, see section Experimental configuration.

We model the guidance graph with graph VAE (Kipf and Welling 2016) to get the modal feature embedding (view 1 embedding in Figure 2a)  $f_i \in \mathbb{R}^m$  for node  $i$ , where  $f_i \subseteq F$ ,  $i \in V = \{V_1, \dots, V_K\}$ . Graph encoder is 2-layer GAT (Veličković et al. 2017):

$$q(F|G) = \prod_{i=1}^{|V|} \mathcal{N}(f_i | \mu_i, \sigma_i^2) \\ = \prod_{i=1}^{|V|} \mathcal{N}(f_i | \text{GAT}_{\mu_i}(G), \text{GAT}_{\sigma_i^2}(G)), \quad (1)$$

where  $u_i^\top, \sigma_i^\top \in \mathbb{R}^m$  represent mean and variance of node  $i$  embedding. Node  $i$  embedding (view 1 embedding) is  $f_i = \text{GAT}_{\mu_i}(G) + \epsilon \times \text{GAT}_{\sigma_i^2}(G)$ ,  $\epsilon \sim \mathcal{N}(0, I)$ . Graph decoder calculates the probability of an edge between any two nodes in the guidance graph, and obtains Fine graph-linked guidance:

$$p(\hat{G}|F) = \prod_{i=1}^{|V|} \prod_{j=1}^{|V|} p(\hat{G}|f_i, f_j) \\ = \prod_{i=1}^{|V|} \prod_{j=1}^{|V|} \text{sigmoid}(f_i \cdot f_j^\top). \quad (2)$$

Therefore, Fine graph-linked guidance is  $\hat{G} = F \cdot F^\top$ . The loss function includes the distance between Fine graph-linked guidance and Coarse graph-linked guidance, and the KL divergence of feature embedding distribution and normal distribution:

$$L_{\text{view1}} = -\mathbb{E}_{F \sim q(F|G)} [\log p(\hat{G}|F)] \\ + KL[q(F|G) || \mathcal{N}(0, I)]. \quad (3)$$

By learning graph-linked embedding, we can integrate the feature information of different modalities, and the obtained feature embedding  $F_1 \in \mathbb{R}^{|V_1| \times m}, \dots, F_K \in \mathbb{R}^{|V_K| \times m}$  contains interaction information of different modalities.

### Fine-grained cell embedding learning

Fine-grained learning can enhance the model’s recognition of fine-grained discriminative features in each cell. We use Teacher transformer layer to learn single cell fine-grained embedding. The Asymmetric network of each modality contains a Teacher Transformer layer respectively, see Figure 2b. The Asymmetric network corresponding to each modality is parameter-unshared, see Figure 2 lower left. Taking RNA modality as example, a cell  $x_1^{(n)} \in X_1$  ( $n = 1, \dots, N_1$ ) is transformed into  $x_1^{(n),position} \in X_1^{position}$  and  $x_1^{(n),counts} \in X_1^{counts}$ .  $x_1^{(n),counts}$  is gene expression of cell  $n$ .  $x_1^{(n),position}$  is gene position embedding of cell  $n$ , we can use gene2vec (Du et al. 2019) to encode each gene, or use view 1 embedding  $F_1$  as position embedding. Similarly, view 1 embedding  $F_k$  can be flexibly used as position embedding for modality  $k$ . Position embedding provide context information for features.

We use Performer (Choromanski et al. 2021) to learn fine-grained cell embedding. For RNA modality, Performer considers the relationship between any pair of genes. Input is  $x_1^{(n),in} = \text{concat}(x_1^{(n),position}, x_1^{(n),counts})$ , and  $x_1^{(n),in} \in \mathbb{R}^{|V_1| \times c}$ ,  $c$  is the dimension number after concatenating. The attention module  $\text{Attn}(\cdot)$  is calculated as:

$$q = x_1^{(n),in} W_q, k = x_1^{(n),in} W_k, v = x_1^{(n),in} W_v, \quad (4)$$

$$z = \text{Performer}(q, k, v) W_o, \quad (5)$$

where,  $W_q, W_k, W_v, W_o \in \mathbb{R}^{c \times c}$  are projection matrix,  $z \in \mathbb{R}^{|V_1| \times c}$  is output of attention module. When constructing Teacher Transformer layer, we adopt module  $\text{MLP}(\cdot)$  with two linear transformations and GELU activation function to provide nonlinear transformation. Then we add layer normalization  $\text{LN}(\cdot)$  and residual connection. Therefore, the  $l$ -th layer of Teacher Transformer layer is:

$$z_l = \text{MLP}(\text{LN}(\text{Attn}(\text{LN}(z_{l-1})) + z_{l-1})) + \text{Attn}(\text{LN}(z_{l-1})) + z_{l-1}. \quad (6)$$

We set 2 layers Teacher Transformer. One cell output is  $z_{teacher} \in \mathbb{R}^{|V_1| \times c}$ . In addition, for Student Dense layer, we set a simple fully connected network  $\text{Dense}(\cdot)$ , input is  $x_1^{(n),in}$ , output is  $z_{student} = \text{Dense}(x_1^{(n),in}) \in \mathbb{R}^{|V_1| \times c}$ .

For contrastive learning, we use  $z_{teacher}$  and  $z_{student}$  to construct positive sample pairs. With random Dropout and Projection, we transform  $z_{teacher}$  and  $z_{student}$  into a pair of vectors  $z_{teacher} \in \mathbb{R}^m$  and  $z_{student} \in \mathbb{R}^m$ , respectively. For  $N_1$  samples, we obtain two corresponding positive sample matrices  $Z_1^{teacher} \in \mathbb{R}^{N_1 \times m}$  and  $Z_1^{student} \in \mathbb{R}^{N_1 \times m}$ , where, the cells with the same row in two matrices constitute a positive sample pair. Similarly, for modality  $k$ , we have  $Z_k^{teacher} \in \mathbb{R}^{N_k \times m}$  and  $Z_k^{student} \in \mathbb{R}^{N_k \times m}$ . And  $Z_k^{teacher}$  is single cell fine-grained embedding (view 2 embedding) in modality  $k$ .

### Contrastive cell embedding learning

Contrastive learning on the unit hypersphere space. For a cell, we treat all other cells and their corresponding positive samples as negative samples for this cell. We separate cells from negative samples and bring cells closer to positive samples. Let cell  $j$  in modality  $k$  has a positive pair  $z_j^{teacher} \in Z_k^{teacher}$  and  $z_j^{student} \in Z_k^{student}$ , we define the distance between them is:

$$s_{j,j^+} = \frac{(z_j^{teacher})^\top \cdot z_j^{student}}{\tau \|z_j^{teacher}\| \|z_j^{student}\|}, \quad (7)$$

where,  $j, j^+$  is a positive pair,  $\tau$  is an adjustable coefficient.

For distance between cell  $j$  and negative sample  $r$ . We have:

$$s_{j,r} = \frac{(z_j^{teacher})^\top \cdot z_r^{student}}{\tau \|z_j^{teacher}\| \|z_r^{student}\|} + \frac{(z_j^{teacher})^\top \cdot z_r^{teacher}}{\tau \|z_j^{teacher}\| \|z_r^{teacher}\|}. \quad (8)$$

We randomly sample  $N_{batch}$  cells from a dataset containing all modalities  $\{N_1, \dots, N_K\}$ . Based on Teacher and Student branches,  $2N_{batch}$  data points can be obtained. Therefore, given a cell from  $N_{batch}$ , we have  $2N_{batch} - 2$  data points as negative samples. The contrastive loss is:

$$L_{view3} = -\frac{1}{2N_{batch}} \sum_{j=1}^{N_{batch}} \log \frac{e^{s_{j,j^+}}}{\sum_{r=1, r \neq j}^{N_{batch}} [e^{s_{j^+,r}} + e^{s_{j,r}}]}. \quad (9)$$

With contrastive learning, we can obtain representations that contain more discriminative information, and these representations are all in the same feature space. We use the optimized Teacher branch embedding of contrastive learning as joint embedding, that is, view 3 embedding  $U_k \in \mathbb{R}^{N_k \times m}$ ,  $k = 1, \dots, K$ , see Figure 2c.

We combine view 1 embedding with view 3 embedding to reconstruct the input. This is a simple decoding process:

$$\hat{X}_k = U_k \cdot F_k^\top. \quad (10)$$

The reconstruction input is used for learning view 2 embedding (single cell fine-grained embedding):

$$L_{view2} = \sum_{k=1}^K \|X_k - \hat{X}_k\|. \quad (11)$$

The overall learning goal of CoVEL is to minimize the sum of loss functions  $L$ :

$$L = L_{view1} + L_{view2} + L_{view3}. \quad (12)$$

## Experiment

### Dataset

**Mouse skin dataset** (Ma et al. 2020): It is from SHARE-seq technology, which can simultaneously measure chromatin accessibility and gene expression in the same single-cell, the dataset contains 32,231 jointly measured cells of mouse



Table 1: Comparison of different methods on test set. Abbreviations: Batch remove; Bio.c, Biology conservation.

Methods	Chen-2019		10x-Multiome		Muto-2021		Yao-2021	
	Batch.r	Bio.c	Batch.r	Bio.c	Batch.r	Bio.c	Batch.r	Bio.c
iNMF(Gao et al. 2021)	0.695	0.491	0.911	0.571	0.957	0.621	0.815	0.609
LIGER(Welch et al. 2019)	0.718	0.495	0.927	0.559	0.962	0.636	0.814	0.634
bindSC(Dou et al. 2022)	0.717	0.511	0.983	0.545	0.977	0.568	-	-
Harmony(Korsunsky et al. 2019)	0.734	0.507	0.984	0.559	0.971	0.625	0.972	0.556
Seurat v3(Stuart et al. 2019)	0.778	0.521	0.981	0.613	0.982	0.681	0.968	0.562
GLUE(Cao and Gao 2022)	0.856	0.574	0.989	0.602	0.969	0.638	0.989	0.604
CoVEL	0.881	0.593	0.981	0.639	0.998	0.721	0.994	0.681

skin, including RNA modality and ATAC modality. Each cell in the dataset is matched one-to-one across modalities. **Mouse cortex dataset** (Saunders et al. 2018; Luo et al. 2017): Three distinct omics layers from neuronal cells in the adult mouse cortex, including gene expression (RNA modality), chromatin accessibility (ATAC modality), and DNA methylation (snmC modality). And the data come from different technologies: Drop-seq, 10x ATAC and snmC-seq. Since there is no joint measurement, the cells in the dataset do not have a matching relationship between modalities, and each modality has different cell classification criteria, which leads challenges for the interpretability of multimodal integration. **Other datasets:** Chen-2019 used SNARE-seq technology to jointly measure 9,190 cells in mouse cortex (Chen, Lake, and Zhang 2019). 10x-Multiome used 10x-Multiome technology to jointly measure 9,631 cells of human PBMC (Cao and Gao 2022). Muto-2021 measured 44,190 cells from human kidney using snRNA-seq and snATAC-seq technologies, respectively (Muto et al. 2021). Yao-2021 measured 124,571 cells from mouse MOP using 10x RNA v3 and snATAC-seq technology (Yao et al. 2021). All datasets are preprocessed according to the standard of scanpy (Wolf, Angerer, and Theis 2018).

## Experimental configuration

The construction of guidance graph is easily available. For example, for the RNA modality and ATAC modality, we reserve the "genome coordinates" field for the features of these two modalities, and then use the "genome coordinates" field to build the relationships between cross-modal features. The relationship between them, that is, the edges of the graph. For RNA modality and snmC modality, we use the 'gene name' field to establish relationships between cross-modal features.

According to the regulatory relationship across modalities, we can obtain the edge sign, see section Graph-linked embedding learning. For example, ATAC peaks usually positively regulate gene expression ( $s_{ij}=1$ , where  $i, j$  represent gene and the corresponding chromatin region, respectively), DNA methylation usually inhibits gene expression ( $s_{ij} = -1$ , where  $i, j$  represent gene and the corresponding methylated fragment, respectively).

CoVEL trained using NVIDIA GeForce RTX A6000 with 48 GB memory. Adam optimizer with 0.001 learning rate was

used to update model parameters. The batch size was set to 16. We use six representative methods for multimodal integration as baselines: Seurat v3 (Stuart et al. 2019), GLUE (Cao and Gao 2022), Harmony (Korsunsky et al. 2019), LIGER (Welch et al. 2019), bindSC (Dou et al. 2022), iNMF (Gao et al. 2021). Evaluation metrics include Batch remove (Batch.r: graph connectivity) and Biology conservation (Bio.c: cell type average silhouette width).

## Results and discussion

### Multimodal integration

We evaluate CoVEL and baselines on the mouse skin dataset (Ma et al. 2020), where each cell is matched one-to-one across modalities, and good multimodal integration methods should align the same cells in different modalities.

We integrate the RNA modality and the ATAC modality using different methods, and visualize the joint embedding across two modalities using UMAP (McInnes, Healy, and Melville 2018), see Figure 3a. For the representative methods GLUE and Seurat v3, it is obvious that they have not aligned the embedding on some cells, and CoVEL obtains better results (blue cells in ATAC modality achieve greater alignment with orange cells in RNA modality). We quantitatively evaluated the multimodal integration performance of different methods using Batch removal and Biology conservation, each method was tested under 8 different random seeds (the random seed is sequentially set to an integer from 0 to 7). Compared with baselines, CoVEL removes the gap between modalities to the greatest extent, protects biological heterogeneity to the greatest extent, and the result is robust, see Figure 3b. For multimodal integration task, all methods learn joint embedding unsupervisedly on the entire dataset.

To evaluate the method's robustness to dataset size, we apply the methods on subsampled datasets of different sizes. We compare the performance of all methods on each subsampled dataset, see Figure 3c and Figure 3d. The results show that some methods also perform well on highly downsampled datasets. But there is a certain degree of loss in integrated performance compared to the full dataset. CoVEL can achieve excellent performance with small-scale dataset. And as the scale of the dataset increases, CoVEL will have more advantages. The integration results on more datasets are compared

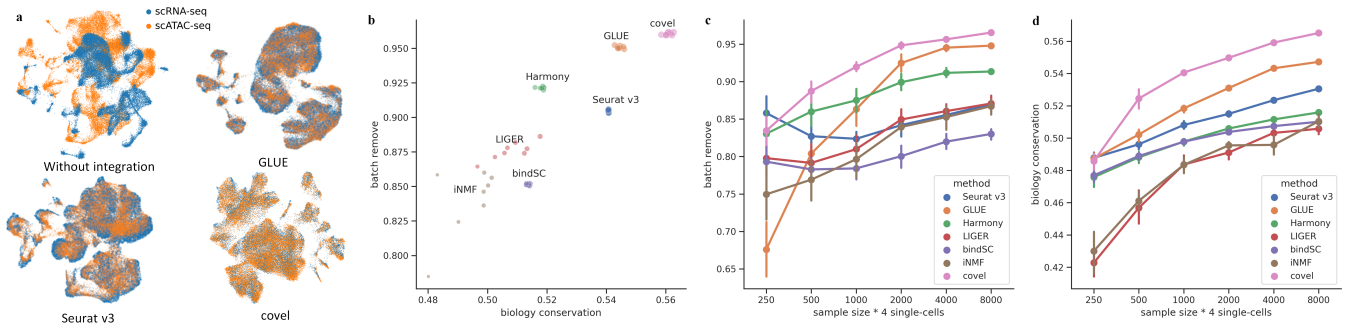


Figure 3: Multimodal integration results. **a**: UMAP visualization of different integration methods. **b**: multimodal integration performance comparison. **c** and **d**: robustness evaluation with subsampled datasets of various sizes.

in Table 1. Yao-2021 contains too many cells which cannot be processed by bindSC.

### Ablation study

We evaluate ablation study on mouse skin dataset (Ma et al. 2020). We randomly remove edges in the guidance graph to obtain graph with varying degrees of damage. For fine-grained learning, we replace Teacher Transformer layer with a Dense layer. For contrastive learning, we use view 2 embedding as the joint embedding for multimodal integration. CoVEL's ablation study results see Table 2. CoVEL is robust to damaged guidance graph, and full guidance graph can help CoVEL preserve biological heterogeneity in multimodal data. In addition, fine-grained learning and contrastive learning allow CoVEL to achieve good batch removal performance.

Table 2: Ablation study of CoVEL. Let light-gray row is reference.  $\Delta$  is the difference between the result and reference. Abbreviations: e.d, edge dropout; f.l, fine-grained learning; c.l, contrastive learning; Batch.r, Batch remove; Bio.c, Biology conservation.

e.d	f.l	c.l	Batch.r	Bio.c	$\Delta$ Batch.r	$\Delta$ Bio.c
0.0	Yes	Yes	0.968	0.562	0.0	0.0
0.3	Yes	Yes	0.966	0.558	-0.002	-0.004
0.6	Yes	Yes	0.961	0.541	-0.007	-0.021
0.0	No	Yes	0.919	0.560	-0.049	-0.002
0.0	Yes	No	0.904	0.559	-0.064	-0.003
0.0	No	No	0.881	0.537	-0.087	-0.025
0.6	No	No	0.874	0.482	-0.094	-0.080

### Interpretability of joint embedding

CoVEL supports the integration of three modalities and more. We utilize CoVEL to integrate mouse cortex dataset (Saunders et al. 2018; Luo et al. 2017). This dataset contains RNA modality, ATAC modality and snmC modality, and each modality has different cell classification criteria (the RNA modality corresponds to 8 types of cells, the ATAC modality corresponds to 10 types of cells, and the snmC

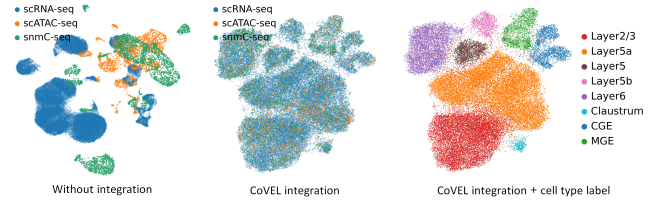


Figure 4: UMAP visualization of 3-modal integration. Left: direct dimensionality reduction to the same embedding space for data from three modalities. Middle: CoVEL eliminates the gap between modalities. Right: Adding cell type annotations to the integrated embedding space, same types of cells are clustered, demonstrating that the integration preserves biological heterogeneity.

modality corresponds to 16 types of cells). CoVEL fully integrates the 3 modalities and obtains a uniformly distributed joint embedding space, see Figure 4 middle part. When we add real cell type annotations to the joint embedding, we can find that the distribution of cell types is uniform and biological heterogeneity is preserved, see the right part of Figure 4. Note that the three modalities are essentially different descriptions of mouse cortex. Therefore, for joint embedding, we can just choose a modality (such as RNA) to represent biological heterogeneity. In Figure 4 right part, we added cell type annotations of RNA modality classification criteria. This result shows that CoVEL is correct for 3-modal integration. In the right part of Figure 4, some clusters are not dense enough (such as CGE cluster), because these cell clusters can be further divided into new subtype clusters. The joint embedding interpretability discussed below can verify this phenomenon.

There are known facts (Saunders et al. 2018; Luo et al. 2017): MGE and CGE in the RNA modality can be divided into two subtypes, and L6-IT and Vip in the ATAC modality can be divided into two subtypes. At the same time we focused on the rare subtype mDL-3. Using the snmC modality as reference dataset (snmC modality cell classification criteria: snmC modality corresponds to 16 types of cells), we use KNN to classify the joint embedding of the RNA modality and the ATAC modality into cell types under the snmC modal-

ity. At the same time, mNdnf-1 and mNdnf-2 in the snmC modality were merged into mNdnf, and mSst-1 and mSst-2 were merged into mSst. The results of cell classification after reference mapping are shown Figure 5b.

We count the flow of categories from Figure 5a to Figure 5b to obtain Figure 5c. We set the cells related to mPv and mSst to be highlighted with lightblue flow, the cells related to mNdnf and mVip to be highlighted with lightpink flow, and the cells related to mDL-3 to be highlighted with lightgreen flow. The result of Figure 5c shows that the joint embedding learned by CoVEL is consistent with known facts, which verifies the interpretability of CoVEL.

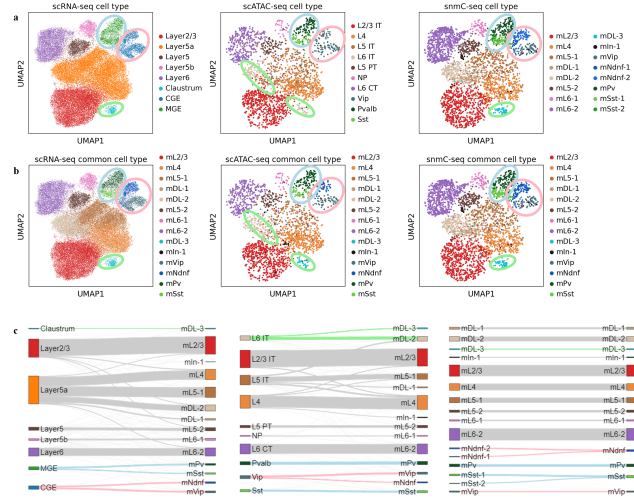


Figure 5: Interpretability analysis. **a**: use CoVEL to learn the joint embedding of three modalities, and add corresponding real annotations (cell type categories) to the embedding UMAP visualization of each modality. **b**: KNN is used to classify the joint embedding in RNA modality and ATAC modality to the cell type corresponding to snmC modality. **c**: statistics of category flow from **a** to **b** (reference mapping, RNA modality and ATAC modality as query, snmC modality as reference).

## Application of downstream task

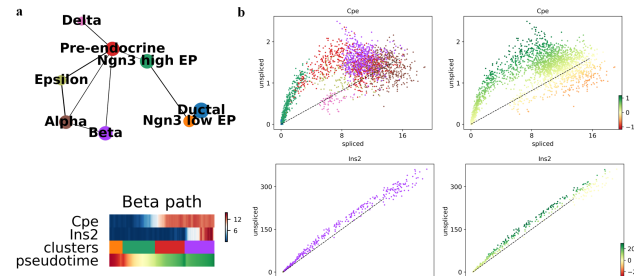


Figure 6: Verify trajectory inference. **a**: CoVEL-based trajectory inference results. **b**: The reliability of trajectory inference is verified by RNA velocity, the left corresponds to clusters, and the right corresponds to pseudotime.

In Figure 3a and Figure 4, there is a slight gap between different clusters (although the gap is small, it is enough to discriminate cell types), which indicates that the embedding

space integrated using CoVEL presents continuity. This property enables the subtle differences of cells to be reflected in trajectory inference task, thus more accurately showing the cell differentiation process. We use CoVEL integrate RNA data (Bastidas-Ponce et al. 2019) and ATAC data (Duvall et al. 2022) of mouse pancreas, and use the integrated embedding for PAGA (Wolf et al. 2019) trajectory inference. We can obtain 4 differentiation trajectories, as shown in Figure 6a: from Ngn3 low EP to Alpha, Beta, Epsilon and Delta respectively. The changes of the statistical genes Cpe and Ins2 on Beta trajectory are shown in Figure 6a lower. The reliability of trajectory inference was verified using RNA velocity (spliced and unspliced), see Figure 6b. Cpe explains the differentiation direction: Ngn3 high EP (green scatter) → Pre-endocrine (red scatter) → Beta (purple scatter). Ins2 explains separate expression in Beta cells. The consistency of Figure 6a and Figure 6b shows that CoVEL-based trajectory inference results are consistent with real cell differentiation.

## Conclusion

In this study, based on the current challenges of single-cell multimodal integration, we propose CoVEL, a deep learning method for unsupervised single-cell multimodal integration. In order to fully mine and fuse information in multimodal data, CoVEL learns single-cell representations from comprehensive views, including regulatory relationships between modalities, fine-grained representations of cells, and relationships between different cells. The comprehensive view embedding enables CoVEL to remove the gap between modalities while protecting biological heterogeneity. Experimental results on multiple public datasets show that CoVEL is accurate and robust to the single-cell multimodal integration. For the challenging unpaired 3-modal integration task, CoVEL still has good interpretability. Finally, ablation study shows that the regulatory relationship between modalities can help CoVEL preserve biological heterogeneity in multimodal data. Fine-grained representation learning and contrastive learning of cells enable CoVEL to achieve excellent batch removal performance.

## Acknowledgments

This work was supported by the National Natural Science Foundation of China (Grant No. 62176272), Research and Development Program of Guangzhou Science and Technology Bureau (No. 2023B01J1016), and Key-Area Research and Development Program of Guangdong Province (No. 2020B1111100001).

## References

- Ashuach, T.; Gabitto, M. I.; Jordan, M. I.; and Yosef, N. 2021. Multivi: deep generative model for the integration of multi-modal data. *bioRxiv*, 2021–08.
- Baltrušaitis, T.; Ahuja, C.; and Morency, L.-P. 2018. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2): 423–443.
- Bastidas-Ponce, A.; Tritschler, S.; Dony, L.; Scheibner, K.; Tarquis-Medina, M.; Salinno, C.; Schirge, S.; Burtscher, I.;

- Böttcher, A.; Theis, F. J.; et al. 2019. Comprehensive single cell mRNA profiling reveals a detailed roadmap for pancreatic endocrinogenesis. *Development*, 146(12): dev173849.
- Cao, J.; Cusanovich, D. A.; Ramani, V.; Aghamirzaie, D.; Pliner, H. A.; Hill, A. J.; Daza, R. M.; McFaline-Figueroa, J. L.; Packer, J. S.; Christiansen, L.; et al. 2018. Joint profiling of chromatin accessibility and gene expression in thousands of single cells. *Science*, 361(6409): 1380–1385.
- Cao, K.; Bai, X.; Hong, Y.; and Wan, L. 2020. Unsupervised topological alignment for single-cell multi-omics integration. *Bioinformatics*, 36(Supplement\_1): i48–i56.
- Cao, K.; Hong, Y.; and Wan, L. 2022. Manifold alignment for heterogeneous single-cell multi-omics data integration using Pamona. *Bioinformatics*, 38(1): 211–219.
- Cao, Z.-J.; and Gao, G. 2022. Multi-omics single-cell data integration and regulatory inference with graph-linked embedding. *Nature Biotechnology*, 40(10): 1458–1466.
- Chen, J.; Xu, H.; Tao, W.; Chen, Z.; Zhao, Y.; and Han, J.-D. J. 2023. Transformer for one stop interpretable cell type annotation. *Nature Communications*, 14(1): 223.
- Chen, S.; Lake, B. B.; and Zhang, K. 2019. High-throughput sequencing of the transcriptome and chromatin accessibility in the same cell. *Nature biotechnology*, 37(12): 1452–1457.
- Choromanski, K. M.; Likhoshervstov, V.; Dohan, D.; Song, X.; Gane, A.; Sarlos, T.; Hawkins, P.; Davis, J. Q.; Mohiuddin, A.; Kaiser, L.; Belanger, D. B.; Colwell, L. J.; and Weller, A. 2021. Rethinking Attention with Performers. In *International Conference on Learning Representations*.
- Dou, J.; Liang, S.; Mohanty, V.; Miao, Q.; Huang, Y.; Liang, Q.; Cheng, X.; Kim, S.; Choi, J.; Li, Y.; et al. 2022. Bi-order multimodal integration of single-cell data. *Genome biology*, 23(1): 1–25.
- Du, J.; Jia, P.; Dai, Y.; Tao, C.; Zhao, Z.; and Zhi, D. 2019. Gene2vec: distributed representation of genes based on co-expression. *BMC genomics*, 20: 7–15.
- Duvall, E.; Benitez, C. M.; Tellez, K.; Enge, M.; Pauerstein, P. T.; Li, L.; Baek, S.; Quake, S. R.; Smith, J. P.; Sheffield, N. C.; et al. 2022. Single-cell transcriptome and accessible chromatin dynamics during endocrine pancreas development. *Proceedings of the National Academy of Sciences*, 119(26): e2201267119.
- Eberwine, J.; Yeh, H.; Miyashiro, K.; Cao, Y.; Nair, S.; Finnell, R.; Zettel, M.; and Coleman, P. 1992. Analysis of gene expression in single live neurons. *Proceedings of the National Academy of Sciences*, 89(7): 3010–3014.
- Gao, C.; Liu, J.; Kriebel, A. R.; Preissl, S.; Luo, C.; Castanon, R.; Sandoval, J.; Rivkin, A.; Nery, J. R.; Behrens, M. M.; et al. 2021. Iterative single-cell multi-omic integration using online learning. *Nature biotechnology*, 39(8): 1000–1007.
- Hao, Y.; Stuart, T.; Kowalski, M.; Choudhary, S.; Hoffman, P.; Hartman, A.; Srivastava, A.; Molla, G.; Madad, S.; Fernandez-Granda, C.; et al. 2022. Dictionary learning for integrative, multimodal, and scalable single-cell analysis. *bioRxiv*, 2022–02.
- Hie, B.; Bryson, B.; and Berger, B. 2019. Efficient integration of heterogeneous single-cell transcriptomes using Scanorama. *Nature biotechnology*, 37(6): 685–691.
- Kipf, T. N.; and Welling, M. 2016. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308*.
- Klein, A. M.; Mazutis, L.; Akartuna, I.; Tallapragada, N.; Veres, A.; Li, V.; Peshkin, L.; Weitz, D. A.; and Kirschner, M. W. 2015. Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell*, 161(5): 1187–1201.
- Korsunsky, I.; Millard, N.; Fan, J.; Slowikowski, K.; Zhang, F.; Wei, K.; Baglaenko, Y.; Brenner, M.; Loh, P.-r.; and Raychaudhuri, S. 2019. Fast, sensitive and accurate integration of single-cell data with Harmony. *Nature methods*, 16(12): 1289–1296.
- Lartigue, C.; Glass, J. I.; Alperovich, N.; Pieper, R.; Parmar, P. P.; Hutchison III, C. A.; Smith, H. O.; and Venter, J. C. 2007. Genome transplantation in bacteria: changing one species to another. *science*, 317(5838): 632–638.
- Lin, X.; Tian, T.; Wei, Z.; and Hakonarson, H. 2022. Clustering of single-cell multi-omics data with a multimodal deep learning method. *Nature Communications*, 13(1): 7705.
- Liu, X.; Zhang, F.; Hou, Z.; Mian, L.; Wang, Z.; Zhang, J.; and Tang, J. 2021. Self-supervised learning: Generative or contrastive. *IEEE Transactions on Knowledge and Data Engineering*, 35(1): 857–876.
- Luo, C.; Keown, C. L.; Kurihara, L.; Zhou, J.; He, Y.; Li, J.; Castanon, R.; Lucero, J.; Nery, J. R.; Sandoval, J. P.; et al. 2017. Single-cell methylomes identify neuronal subtypes and regulatory elements in mammalian cortex. *Science*, 357(6351): 600–604.
- Lv, Q.; Chen, G.; Yang, Z.; Zhong, W.; and Chen, C. Y.-C. 2023. Meta Learning With Graph Attention Networks for Low-Data Drug Discovery. *IEEE Transactions on Neural Networks and Learning Systems*.
- Ma, A.; Wang, X.; Li, J.; Wang, C.; Xiao, T.; Liu, Y.; Cheng, H.; Wang, J.; Li, Y.; Chang, Y.; et al. 2023. Single-cell biological network inference using a heterogeneous graph transformer. *Nature Communications*, 14(1): 964.
- Ma, S.; Zhang, B.; LaFave, L. M.; Earl, A. S.; Chiang, Z.; Hu, Y.; Ding, J.; Brack, A.; Kartha, V. K.; Tay, T.; et al. 2020. Chromatin potential identified by shared single-cell profiling of RNA and chromatin. *Cell*, 183(4): 1103–1116.
- Macosko, E. Z.; Basu, A.; Satija, R.; Nemesh, J.; Shekhar, K.; Goldman, M.; Tirosh, I.; Bialas, A. R.; Kamitaki, N.; Martersteck, E. M.; et al. 2015. Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. *Cell*, 161(5): 1202–1214.
- McInnes, L.; Healy, J.; and Melville, J. 2018. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*.
- Muto, Y.; Wilson, P. C.; Ledru, N.; Wu, H.; Dimke, H.; Waikar, S. S.; and Humphreys, B. D. 2021. Single cell transcriptional and chromatin accessibility profiling redefine cellular heterogeneity in the adult human kidney. *Nature communications*, 12(1): 2190.
- Saunders, A.; Macosko, E. Z.; Wysoker, A.; Goldman, M.; Krienen, F. M.; de Rivera, H.; Bien, E.; Baum, M.; Bortolin,

- L.; Wang, S.; et al. 2018. Molecular diversity and specializations among the cells of the adult mouse brain. *Cell*, 174(4): 1015–1030.
- Song, Q.; Su, J.; and Zhang, W. 2021. scGCN is a graph convolutional networks algorithm for knowledge transfer in single cell omics. *Nature communications*, 12(1): 3826.
- Steyaert, S.; Pizurica, M.; Nagaraj, D.; Khandelwal, P.; Hernandez-Boussard, T.; Gentles, A. J.; and Gevaert, O. 2023. Multimodal data fusion for cancer biomarker discovery with deep learning. *Nature Machine Intelligence*, 1–12.
- Stoeckius, M.; Hafemeister, C.; Stephenson, W.; Houck-Loomis, B.; Chattopadhyay, P. K.; Swerdlow, H.; Satija, R.; and Smibert, P. 2017. Simultaneous epitope and transcriptome measurement in single cells. *Nature methods*, 14(9): 865–868.
- Stuart, T.; Butler, A.; Hoffman, P.; Hafemeister, C.; Papalexi, E.; Mauck III, W. M.; Hao, Y.; Stoeckius, M.; Smibert, P.; and Satija, R. 2019. Comprehensive integration of single-cell data. *Cell*, 177(7): 1888–1902.
- Tang, F.; Barbacioru, C.; Wang, Y.; Nordman, E.; Lee, C.; Xu, N.; Wang, X.; Bodeau, J.; Tuch, B. B.; Siddiqui, A.; et al. 2009. mRNA-Seq whole-transcriptome analysis of a single cell. *Nature methods*, 6(5): 377–382.
- Tang, Z.; Chen, G.; Yang, H.; Zhong, W.; and Chen, C. Y.-C. 2023. DSIL-DDI: A Domain-Invariant Substructure Interaction Learning for Generalizable Drug–Drug Interaction Prediction. *IEEE Transactions on Neural Networks and Learning Systems*.
- Tu, X.; Cao, Z.-J.; Mostafavi, S.; Gao, G.; et al. 2022. Cross-Linked Unified Embedding for cross-modality representation learning. *Advances in Neural Information Processing Systems*, 35: 15942–15955.
- Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Lio, P.; and Bengio, Y. 2017. Graph attention networks. *arXiv preprint arXiv:1710.10903*.
- Welch, J. D.; Kozareva, V.; Ferreira, A.; Vanderburg, C.; Martin, C.; and Macosko, E. Z. 2019. Single-cell multi-omic integration compares and contrasts features of brain cell identity. *Cell*, 177(7): 1873–1887.
- Wen, H.; Ding, J.; Jin, W.; Wang, Y.; Xie, Y.; and Tang, J. 2022. Graph neural networks for multimodal single-cell data integration. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 4153–4163.
- Wolf, F. A.; Angerer, P.; and Theis, F. J. 2018. SCANPY: large-scale single-cell gene expression data analysis. *Genome biology*, 19: 1–5.
- Wolf, F. A.; Hamey, F. K.; Plass, M.; Solana, J.; Dahlin, J. S.; Göttgens, B.; Rajewsky, N.; Simon, L.; and Theis, F. J. 2019. PAGA: graph abstraction reconciles clustering with trajectory inference through a topology preserving map of single cells. *Genome biology*, 20: 1–9.
- Yang, F.; Wang, W.; Wang, F.; Fang, Y.; Tang, D.; Huang, J.; Lu, H.; and Yao, J. 2022a. scBERT as a large-scale pretrained deep language model for cell type annotation of single-cell RNA-seq data. *Nature Machine Intelligence*, 4(10): 852–866.
- Yang, M.; Yang, Y.; Xie, C.; Ni, M.; Liu, J.; Yang, H.; Mu, F.; and Wang, J. 2022b. Contrastive learning enables rapid mapping to multimodal single-cell atlas of multimillion scale. *Nature Machine Intelligence*, 4(8): 696–709.
- Yao, Z.; Liu, H.; Xie, F.; Fischer, S.; Adkins, R. S.; Aldridge, A. I.; Ament, S. A.; Bartlett, A.; Behrens, M. M.; Van den Berge, K.; et al. 2021. A transcriptomic and epigenomic cell atlas of the mouse primary motor cortex. *Nature*, 598(7879): 103–110.
- Yu, X.; Xu, X.; Zhang, J.; and Li, X. 2023. Batch alignment of single-cell transcriptomics data using deep metric learning. *Nature Communications*, 14(1): 960.