

# Supplementary Materials for Integrated Analysis of Multimodal Single-Cell Data with Structural Similarity

Yingxin Cao<sup>1,5,6\*</sup>, Laiyi Fu<sup>1,2 \*</sup>, Jie Wu<sup>3</sup>, Qinke Peng<sup>2</sup>, Qing Nie<sup>4,5,6</sup>, Jing  
Zhang<sup>1 \*\*</sup>, and Xiaohui Xie<sup>1 \*\*</sup>

<sup>1</sup> Department of Computer Science, University of California, Irvine, CA, 92697, USA

<sup>2</sup> Systems Engineering Institute, School of Electronic and Information Engineering,  
Xi'an Jiaotong University, Xi'an, Shanxi, 710049, China

<sup>3</sup> Department of Biological Chemistry, University of California, Irvine, CA, 92697,  
USA

<sup>4</sup> Department of Mathematics, University of California, Irvine, CA, 92697, USA

<sup>5</sup> Center for Complex Biological Systems, University of California, Irvine, CA, 92697,  
USA

<sup>6</sup> NSF-Simons Center for Multiscale Cell Fate Research, University of California,  
Irvine, CA, 92697, USA

---

\* Equal contributions

\*\* To whom correspondence should be addressed

## Table of Contents

1	Results on Hyperparameter Stability .....	3
2	Cell Type Specific Marker Genes .....	4
3	Extra Results on Clustering .....	6
4	Results on Share-seq Dataset .....	7
5	Results on Batch Correction .....	8
6	Results on Motif Enrichment Analysis .....	9

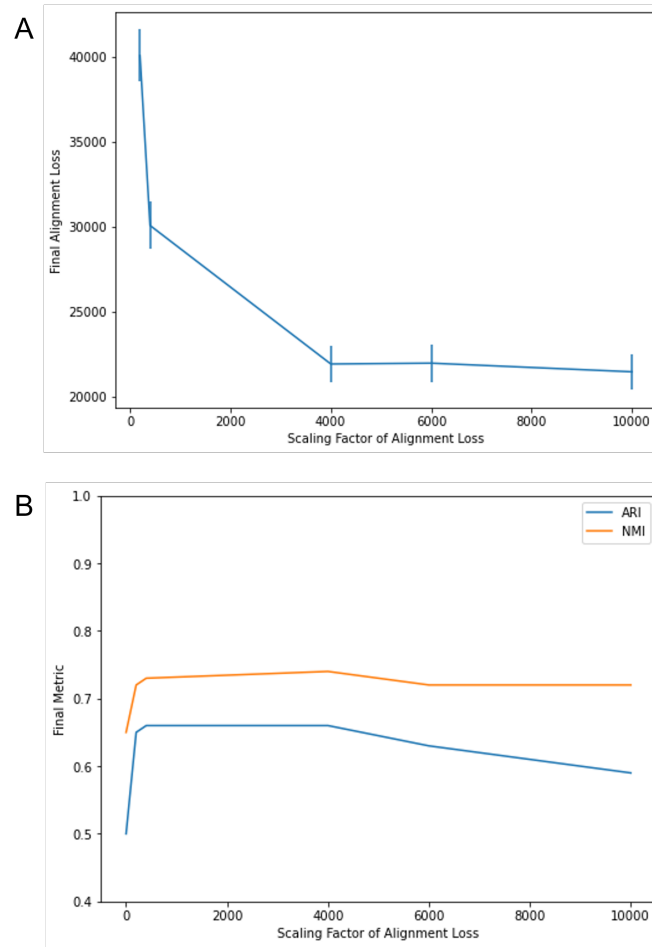
### List of Tables

S1	List of cell-type specific marker genes used to visualize expressions. . . .	4
S2	Mean expressions of markers on cells clustered by different methods. . .	5

### List of Figures

S1	Hyperparameter Stability .....	3
S2	Marker Gene Activities in the PBMC Dataset .....	4
S3	Marker Gene Expressions in the PBMC Dataset .....	5
S4	Clustering scores .....	6
S5	Comparisons of embeddings .....	7
S6	Results on Share-seq dataset .....	7
S7	Batch Effect Correction .....	8
S8	Clustering result on batches with unique cell types. ....	8
S9	Motif enrichment scores on imputed data .....	9

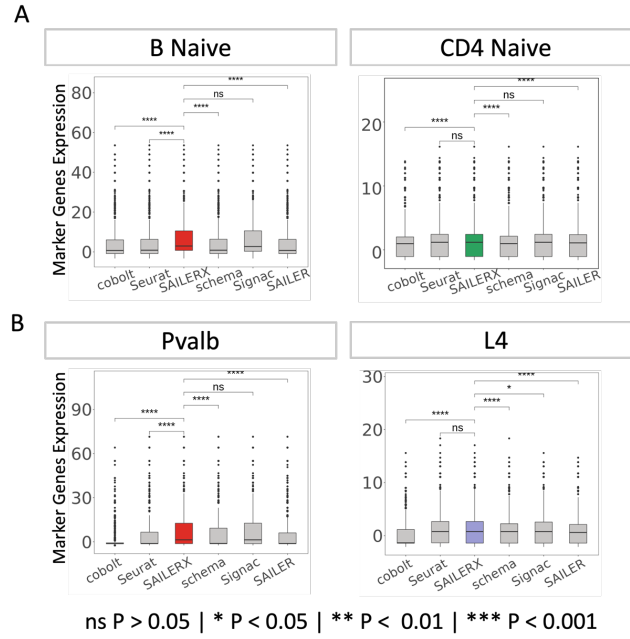
## 1 Results on Hyperparameter Stability



**Fig. S1:** Results on hyperparameter stability. (A) The alignment loss  $L_{Local}$  decreases as the scaling factor increases. (B) Clustering metrics ARI and NMI as a function of scaling factor.



sion of marker genes of different methods. Pairwise t-tests between SAILERX and other methods indicate whether the marker genes from SAILERX show significantly higher expression than those from other methods. T-test p-values are indicated by ns (p-value > 0.05, i.e., not significant), \* (p-value < 0.05), \*\* (p-value < 0.01) and \*\*\* (p-value < 0.001). The plots show that the marker genes show overall higher expression in corresponding cell clusters discovered by SAILERX than those by other methods.

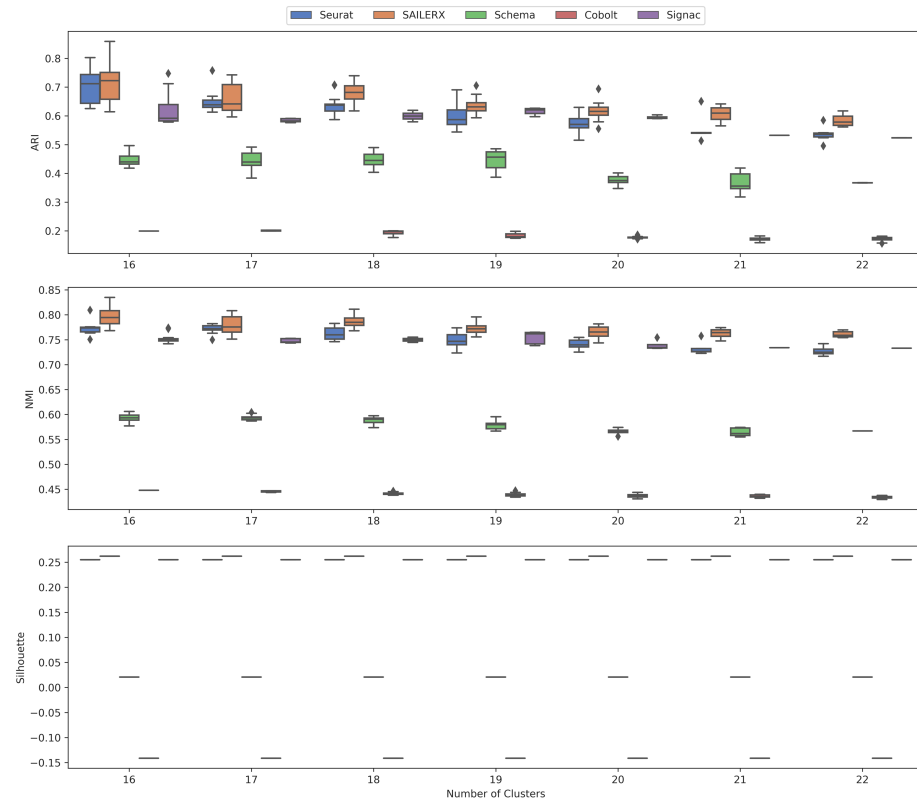


**Fig. S3:** Comparing the expression of marker genes in clusters derived by different methods. (A) Mean expression of marker genes of B Naive cells and CD4 Naive cells from PBMC 10k dataset. (B) Mean expression of marker genes of Pvalb cells and L4 cells from SNARE-seq dataset.

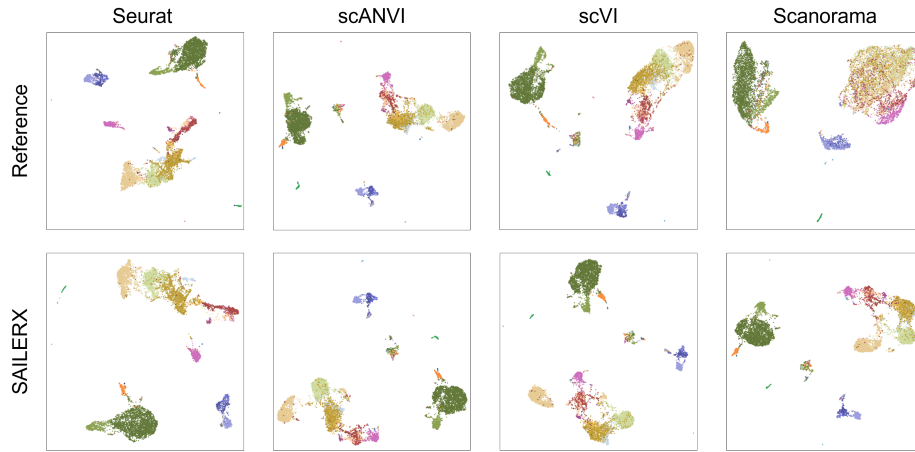
**Table S2:** Mean expressions of markers on cells clustered by different methods.

Cell Type	SAILERX	Seurat	Signac	Cobolt	Schema	SAILER
Pvalb	<b>7.10</b>	4.25	7.07	0.72	5.01	3.71
L4	<b>1.05</b>	<b>1.11</b>	0.97	-0.01	0.92	0.77
B naive	<b>6.29</b>	3.60	6.09	3.55	3.68	3.60
CD4	<b>1.36</b>	1.34	1.34	1.05	1.07	1.23

### 3 Extra Results on Clustering

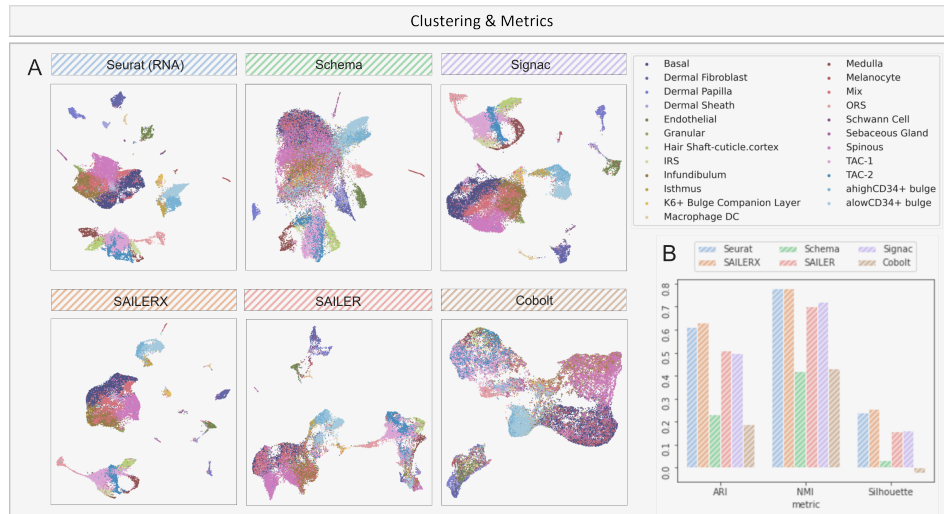


**Fig. S4:** Clustering scores of PBMC 10k dataset by different number of identified clusters.



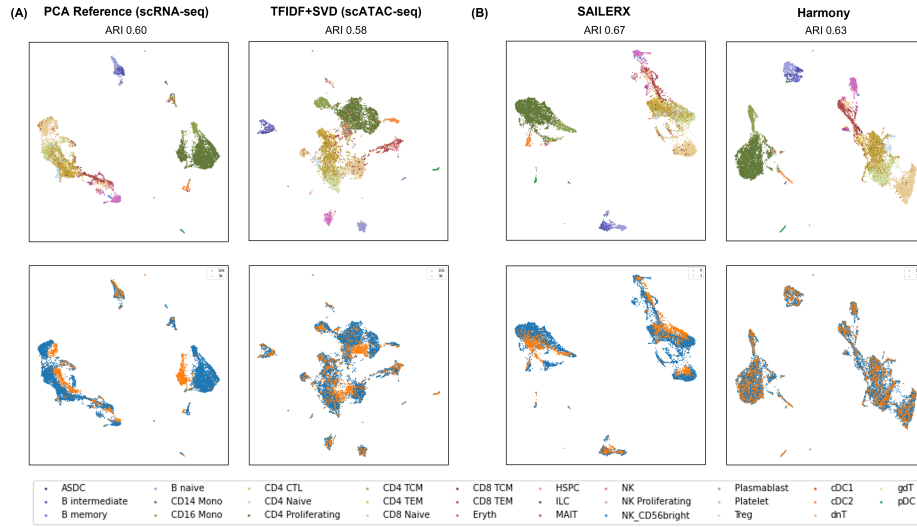
**Fig. S5:** UMAP Visualizations of reference embeddings vs SAILERX embeddings. Top row: UMAP visualizations of reference gene expression embeddings generated by different methods. Bottom row: joint embeddings generated by SAILERX after training.

#### 4 Results on Share-seq Dataset



**Fig. S6:** Results on Share-seq dataset. Cells colored by ground truth label. (A) UMAP visualizations of embeddings on mouse skin Share-seq dataset generated by different methods. (B) Quantitative metrics of ARI, NMI and Silhouette Score on clustering.

## 5 Results on Batch Correction



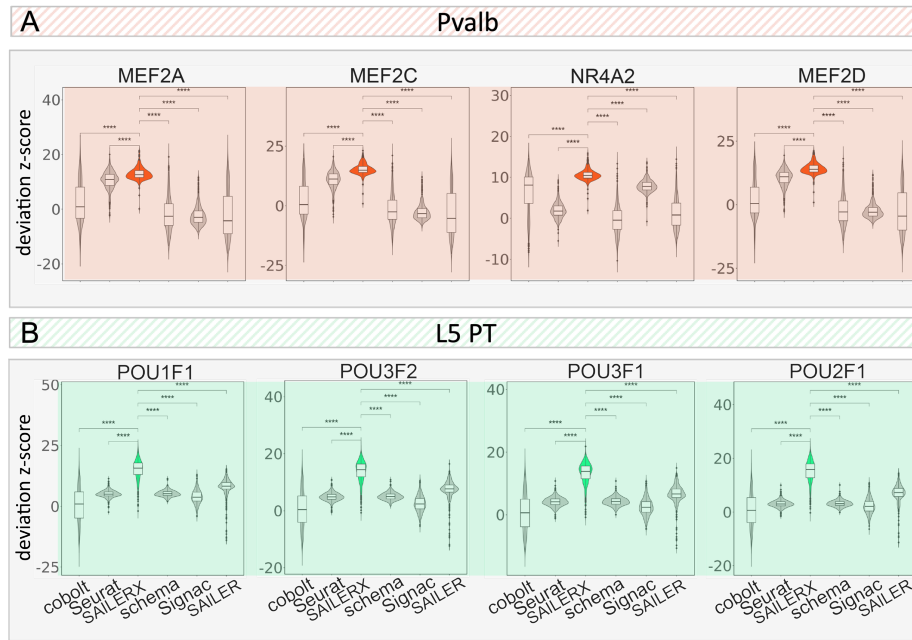
**Fig. S7:** Results on batch effect correction on PBMC 10k and 3k datasets. (A) UMAP Visualizations of PCA (left) embedding on gene expression modality and TF-IDF + SVD (right) embedding on chromatin accessibility modality before batch effect corrections. (B) UMAP visualization of embeddings after batch effect correction. Top row: colored by cell types; Bottom row: colored by batches.



**Fig. S8:** UMAP visualizations of the embedding generated by SAILERX. Left: colored by cell types; Right: colored by batches.



## 6 Results on Motif Enrichment Analysis



**Fig. S9:** Motif deviation z-scores on cells identified as (A) Pvalb and (B) L5 PT by different methods from SNARE-seq imputed data. The data is imputed through SAILERX. For each cell type, four enriched motifs are selected. Pairwise t-tests are performed between SAILERX and all other methods. Three-stars refers to differential significance between two methods (p-value less than 0.05).