

# 1 **scGREAT: Graph-based regulatory element analysis tool for single-cell**

## 2 **multi-omics data**

3 Chaozhong Liu<sup>1</sup>, Linhua Wang<sup>1</sup>, Zhandong Liu<sup>2,3,\*</sup>

4

5 1. Graduate Program in Quantitative and Computational Biosciences, Baylor College of Medicine,

6 Houston, USA

7 2. Jan and Dan Duncan Neurological Research Institute at Texas Children's Hospital, Houston, USA

8 3. Department of Pediatrics, Baylor College of Medicine, Houston, USA

9 \* Corresponding author: zhandonl@bcm.edu

10

## 11 **Abstract**

12 **Motivation:** With the development in single-cell multi-omics sequencing technology and data  
13 integration algorithms, we have entered the single-cell multi-omics era. Current multi-omics  
14 analysis algorithms failed to systematically dissect the heterogeneity within the datasets when  
15 inferring cis-regulatory events. Thus, there is a need for cis-regulatory element inferring  
16 algorithms that considers the cellular heterogeneity.

17 **Results:** Here, we propose scGREAT, a single-cell multi-omics regulatory state analysis Python  
18 package with a rapid graph-based correlation measurement  $L$ . The graph-based correlation  
19 method assigns each cell a local  $L$  index, pinpointing specific cell groups of certain regulatory  
20 states. Such single-cell resolved regulatory state information enables the heterogeneity analysis  
21 equipped in the package. Applying scGREAT to the 10X Multiome PBMC dataset, we  
22 demonstrated how it could help subcluster cell types, infer regulation-based pseudo-time  
23 trajectory, discover feature modules, and find cluster-specific regulatory gene-peak pairs.

1 Besides, we showed that global L index, which is the average of all local L values, is a better  
2 replacement for Pearson's r in ruling out confounding regulatory relationships that are not of  
3 research interests.

4 **Availability:** <https://github.com/ChaozhongLiu/scGREAT>

5

6

## 7 **1. Introduction**

8 Chromatin accessibility refers to the level of physical compaction of chromatin (Minnoye *et al.*,  
9 2021). While most of the chromatin regions are densely arranged and hard to be accessed, ~2-3%  
10 of total DNA sequences are accessible to transcriptional factors (TFs) that determine the cellular  
11 state (Thurman *et al.*, 2012, Klemm *et al.*, 2019). Such accessible DNA regions that regulate the  
12 transcription of genes are called cis-regulatory elements (CREs) (Li *et al.*, 2015). And it is a well-  
13 known marker that reflects the cell state, such as cell development, differentiation, and abnormal  
14 disease state. However, it is still one of the major challenges to determine the target genes of  
15 CRE activity (Preissl *et al.*, 2022) and how CREs' variations are associated with cellular states.

16

17 Recent advances in single-cell technologies and algorithms have enabled the direct link between  
18 accessible genome loci and nearby genes. Sequencing techniques like SNARE-seq (Ma *et al.*,  
19 2020), SHARE-seq (Chen *et al.*, 2019), and 10X Multiome profile the transcriptome and  
20 chromatin accessibility simultaneously on the same batch of cells, providing researchers with  
21 both modalities' information; Algorithms like GLUE (Cao and Gao, 2022) and Seurat (Stuart *et*  
22 *al.*, 2019), can integrate single-cell RNA-seq and ATAC-seq data, generating data with the two  
23 modalities aligned. However, multi-omics analysis tools available now such as Seurat v4 (Hao *et*

1 *al.*, 2021), Liger (Liu *et al.*, 2020), MUON (Bredikhin *et al.*, 2022), etc., all focus on generating  
2 an optimal dimension reduction space utilizing both modalities' information. However, these  
3 tools do not provide downstream analyses, such as CRE inference, after integrating the two  
4 modalities.

5  
6 One common method to infer CREs is the co-accessibility measurement (Preissl *et al.*, 2022). For  
7 example, Cicero utilizes scATAC-seq data to connect distal regulatory elements with target  
8 genes by computing a covariance matrix of accessible sites (Pliner *et al.*, 2018). ArchR (Granja  
9 *et al.*, 2021), another popular scATAC-seq analysis software, applies Pearson correlation  
10 (Benesty *et al.*, 2009) to infer the links between scATAC-seq and scRNA-seq data. Such  
11 correlation-based methods could fairly detect CREs in homogenous population of cells. However,  
12 since cellular heterogeneity is not modeled in their approaches, these methods could miss  
13 regulatory relations that are specific to subpopulations of cells. These cell subpopulations could  
14 represent different cell types (Grün *et al.*, 2015), development stages (Velten *et al.*, 2017), and  
15 more. Although diverse bioinformatics methods are available to either annotate cell types (Clarke  
16 *et al.*, 2021) or infer cell development trajectory (Saelens *et al.*, 2019), the origin of  
17 heterogeneity is more complicated and any attempt to divide cells into subpopulations is a  
18 simplification and information loss. Thus, a systematic detection of regulatory heterogeneity  
19 without an assumption on heterogeneity origin is needed.

20  
21 In geographical studies, researchers utilize measurement specially designed for spatial data to  
22 dissect the correlations between features. This kind of metrics considers the spatial pattern in  
23 correlation, and provides a correlation measurement for each location within its neighborhood,

1 which is called the local correlation value. Such local correlation value represents the  
2 heterogeneity of feature correlation among all the locations, and the sum of all local values is the  
3 general correlation trend, or the global correlation. We see the resemblance of spatial correlation  
4 between the geographical and single-cell multi-omics studies since multi-omics data could be  
5 represented in a low-dimensional space and constructed as a k-nearest neighboring (KNN) graph.  
6 Thus, it is feasible to develop a similar global and local correlation measurement to dissect the  
7 general correlation and cellular regulatory heterogeneity.

8 Here, we adapted the geographical spatial correlation measure (Lee, 2001) to single-cell studies  
9 and developed the Python package – scGREAT (single-cell Graph-based Regulatory Element  
10 Analysis Toolkit) to infer single-cell resolved regulatory states using local  $L$ -index values. First,  
11 we will show how the local  $L$  matrix has enabled the comprehensive study of heterogeneity in  
12 the 10X Multiome PBMC dataset, including sub-clustering cell types, inferring regulation-based  
13 pseudo-time trajectory, discovering feature modules, and finding cluster-unique regulatory  
14 relationships. Then we will demonstrate why our global  $L$  index, which averages local  $L$ -index  
15 values across all cells, is a better replacement for Pearson's  $r$  in studying the general gene-peak  
16 regulatory trend by ruling out the confounding gene-peak pairs that are not of research interests.

17

## 18 **2. Methods**

### 19 **2.1 Single-cell multi-omics data preprocessing**

20 Three datasets from common multiome techniques, including 10X Multiome, SHARE-seq, and  
21 SNARE-seq, have been employed to test the feasibility of our software. 10X Multiome  
22 Peripheral Blood Mononuclear Cells (PBMC) data were processed using Seurat v4 by following  
23 the Seurat tutorial at

1 [https://satijalab.org/seurat/articles/weighted\\_nearest\\_neighbor\\_analysis.html](https://satijalab.org/seurat/articles/weighted_nearest_neighbor_analysis.html) (Hao *et al.*, 2021).  
2 Then the scRNA-seq and scATAC-seq count matrices, nearest neighbor graph and metadata  
3 were exported and loaded in Python to create the final AnnData object (Virshup *et al.*, 2021) as  
4 the input of our package.  
5 SHARE-seq mouse skin and SNARE-seq mouse brain datasets were downloaded at GEO by  
6 GSE140203 and GSE126074. The raw count matrices were loaded and preprocessed using  
7 Scanpy (Wolf *et al.*, 2018) following the standard preprocessing pipeline ([https://scanpy-](https://scanpy-tutorials.readthedocs.io/en/latest/pbmc3k.html)  
8 [tutorials.readthedocs.io/en/latest/pbmc3k.html](https://scanpy-tutorials.readthedocs.io/en/latest/pbmc3k.html)). The final AnnData objects were saved as the  
9 input for our package. All codes are available in the GitHub repository.

10

## 11 **2.2 Graph-based correlation index $L$**

12 The graph-based correlation measurement is the core of our package. The correlation index  $L$   
13 between gene vector  $\mathbf{x} \in \mathbb{R}^n$  and peak vector  $\mathbf{y} \in \mathbb{R}^n$  combines the normal correlation score  
14 and graph dependence values between the two features. It is defined as:

$$L_{x,y} = \frac{\sum_i (\tilde{x}_i - \bar{x})(\tilde{y}_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_i (y_i - \bar{y})^2}}$$

15 where  $i \in \{1, \dots, n\}$  represents cell index,  $\bar{x}$  and  $\bar{y}$  are the numeric mean values of  $\mathbf{x}$  and  $\mathbf{y}$ ,  $\tilde{\mathbf{x}}$  and  
16  $\tilde{\mathbf{y}}$  are the spatial lag values which are composed of weighted averages of cell neighbors. Spatial  
17 lag of  $x_i$  is defined as:

$$\tilde{x}_i = \sum_j w_{ij} \cdot x_j$$

18 Where  $j$  is index of connected cell with  $i$  in the graph, and  $w_{ij}$  is their connectivity weight. Here,  
19 we take the  $K$ -nearest neighbor graph derived from the dimensional reduction results as the  
20 connectivity matrix  $\mathbf{W} \in \mathbb{R}^{n \times n}$ , where  $n$  indicates the cell number.

1 One major challenge of implementing such approach for single-cell data is the enormous number  
2 of gene-peak pairs, which is computational infeasible. To speed up the computation and allow  
3 scalability, we used a vectorized implementation. Suppose  $\mathbf{W}$  is the row-standardized (row sums  
4 are all 1) connectivity matrix,  $\mathbf{Z}$  is an z-scored ( $z = \frac{x-\mu}{\sigma}$ ) form of  $[\mathbf{x} \ \mathbf{y}] \in \mathbb{R}^{n \times 2}$ ,  $\mathbf{L}$  between  $\mathbf{x}$   
5 and  $\mathbf{y}$  can be calculated with:

$$\mathbf{L} = \begin{bmatrix} L_{x,x} & L_{x,y} \\ L_{x,y} & L_{y,y} \end{bmatrix} = \frac{\mathbf{Z}^T (\mathbf{W}^T \mathbf{W}) \mathbf{Z}}{n}$$

6 When implementing the calculation in Python, we vectorized the formula again to include all  $p$   
7 gene-peak pairs at once:

$$\mathbf{L} = \left[ L_{x_1, y_1}, \dots, L_{x_p, y_p} \right] = \frac{\text{RowSum}(\mathbf{Z}_Y^T (\mathbf{W}^T \mathbf{W}) \circ \mathbf{Z}_X^T)}{n}$$

8 where  $\mathbf{Z}_X$  and  $\mathbf{Z}_Y$  are z-scored matrices of the RNA-seq log-transformed data  $\mathbf{X} \in \mathbb{R}^{n \times p}$ , and  
9 scATAC-seq log-transformed data  $\mathbf{Y} \in \mathbb{R}^{n \times p}$ .

10 The local value of  $L$ , associated with each cell, indicates each cell's relative contribution to the  
11 global  $L$ . Local  $L$  index of cell  $i$  between gene  $\mathbf{x}$  and peak  $\mathbf{y}$  is defined as:

$$L_{x,y}^{(i)} = \frac{n \cdot (\tilde{x}_i - \bar{x})(\tilde{y}_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_i (y_i - \bar{y})^2}},$$

12 allowing it to capture a cell's association within its neighborhood.

13

14 The vectorized implementation of the local  $L$  matrix  $\in \mathbb{R}^{n \times p}$  can be calculated with:

$$\text{local } \mathbf{L} = (\mathbf{Z}_Y^T \mathbf{W}^T)^T \circ (\mathbf{W} \mathbf{Z}_X)$$

15 To speed up all the graph (spatial) dependence measurement, we also vectorized the  
16 implementation of Moran's  $I$  in the package. Details are discussed in the supplementary  
17 materials (Package Details 2.4). We evaluated the computational time under the constrain of 16

1 Gb memory and eight threads. In summary, the graph-based measurements in scGREAT are  
2 much faster than the original implementation in geographical package *esda* (Rey and Anselin,  
3 2007) (Figure 1B).

4  
5 To quantify the statistical significance of  $L$  index, we implemented a permutation test.  
6 Specifically, the paired vectors were shuffled together  $n$  times to get a reference distribution of  
7 the  $L$  index. Then the z-score of the  $L$  index under the reference distribution is calculated and  
8 taken as the simulated p-value. But in single-cell studies, data are usually sparse (Lähnemann *et*  
9 *al.*, 2020) and of large sample size (cell number), which results in all significant p-values during  
10 initial experiments. Here, we introduce a conditional permutation process. In each permutation,  
11 we only shuffle a certain percentage (default 10%) of values and permute among non-zero and  
12 zeros values (see pseudo-code in supplementary material). After permutation, the same process  
13 mentioned above will be done to get the final p-value as the significance level.

### 14 15 **2.3 Implementation of all package functions**

16 **Figure 1A** shows the analysis pipeline and all functions in scGREAT. Modules will be briefly  
17 described here, and all details can be found in the supplementary materials.

18 *Data preparation.* The input of scGREAT is scRNA-seq and scATAC-seq data that have already  
19 been aligned. It can be naturally paired from multiome sequencing techniques or aligned by  
20 integration algorithms. Users can use single-cell data for the analyses or generate pseudo-bulk  
21 data using scGREAT functions. After annotating genes with their nearby peaks, the final  
22 multiome AnnData object that combines transcriptome and chromatin accessibility data will be  
23 generated for the other modules.

1 *Graph-based measurements.* The graph-based measurements are the core of scGREAT,  
2 including Moran's  $I$ , global  $L$  index, and local  $L$  matrix. Users just need to run a one-line code to  
3 perform each measurement.

4 *Unlabeled analysis.* When users are interested in one or few clusters of cells or when no label  
5 information is available, unlabeled analysis can be done to perform sub-clustering and trajectory  
6 analysis based on the regulatory state measured in core functions. With pseudo-time inferred,  
7 self-organized map (SOM) algorithm was implemented to discover feature modules.

8 *Labeled analysis.* When users have already annotated the data with biological groups, a  
9 statistical test can be done to discover all the differentially regulated pairs among all the groups  
10 using the local  $L$  matrix calculated.

11 *Motif Enrichment Analysis.* With a list of peaks derived from any analysis, we implemented  
12 functions to prepare input and call *Homer* (Heinz *et al.*, 2010) for motif enrichment analysis.  
13 There is also a function helping extract peaks and related information from Homer results with  
14 users' motifs of interest.

15 *Visualization.* To help users understand the results and generate publishable plots, we developed  
16 visualization functions for each step and module, including pseudo-bulk quality summary, local  
17  $L$  index heatmap, feature plot in UMAP, volcano plot after statistical test, and feature module  
18 plot after SOM implementation.

19

## 20 **3. Results**

### 21 **3.1 Pseudo-bulk data shows higher correlations between gene and promoter regions**

22 The first step for the single-cell multi-omics regulatory element analysis is data preparation. In  
23 this step, generating pseudo-bulk data is recommended. Here, we measured the correlation



1 between gene expression and promoter accessibility by the graph-based  $L$  index and Pearson's  $r$   
2 in three datasets generated from common multiome techniques, including 10X Multiome,  
3 SHARE-seq, and SNARE-seq. Both methods showed a higher level of correlation in pseudo-  
4 bulk data compared with the single-cell data (Figure 2). We believe this is due to the single-cell  
5 data sparsity (Lähnemann *et al.*, 2020) that makes the data noisier than bulk ones. When  
6 generating pseudo-bulks averaging the values of several cells within a local neighborhood, the  
7 data is less noisy, and thus the results have a better correlation signal. With this observation, in  
8 the package, we provide the function to generate pseudo-bulk data (see supplementary method  
9 for details) for cis-regulatory element discovering. The following results are all based on the  
10 pseudo-bulk data generated by our package.

11

### 12 **3.2 Local $L$ index as an indicator for cell regulatory state**

13 The local  $L$  index indicates the cell's gene-peak correlation strength (magnitude) and direction  
14 (sign) within its neighborhood. For each gene-peak pair, local  $L$  is calculated for each cell to  
15 measure single-cell resolved regulatory states. By doing this for all cells and all gene-peak pairs,  
16 we generated a new layer of information representing the regulatory state besides the  
17 transcriptome and chromatin accessibility matrix. Here, we will demonstrate how local  $L$  index  
18 and the cell regulatory state matrix can help study the regulatory heterogeneity in the well-  
19 annotated 10X Multiome Peripheral Blood Mononuclear Cell (PBMC) dataset.

20

21 Supplementary Figure 1 shows examples of such regulatory states between genes and peaks.  
22 RORA and chr15:60649315-60650653 in CD4+ T cells, for example, are not highly correlated  
23 (Pearson's  $r = 0.404$ ). However, the regulatory state defined by the local  $L$  shows that the

1 correlation is high in CD4+ T naïve cells and CD4+ TEM cells, but not in the majority of CD4+  
2 TCM cells. And this heterogeneity was overlooked by Pearson correlation. RORA is known to  
3 function in T cell differentiation (Solt *et al.*, 2011). The correlated peak lies upstream of isoform  
4 b and c of RORA, and the regulatory state changes could indicate the alternative splicing and  
5 different functions in cell differentiation. Thus, within-cluster heterogeneity could also be  
6 important in understanding the cellular process besides the global correlation trend measured  
7 with Pearson's  $r$ .

8 With this cell-specific regulatory information encoded in the cells-by-regulatory pairs matrix,  
9 scGREAT enables cluster-specific CRE analysis. Under the assumption that regulatory states  
10 indicate cell state, such as cell differentiation, the local  $L$  matrix is utilized for sub-clustering and  
11 trajectory analysis in scGREAT similar to scRNA-seq analysis. Here, we used the same 10X  
12 Multiome CD4+ T cells as an example. After sub-clustering in the pseudo-bulk data, scGREAT  
13 mapped the labels back to the single-cell data, shown in Figure 3A left. The regulatory state-  
14 derived clusters split the three manually annotated cell types into more subtypes. Using the same  
15 KNN graph constructed in the sub-clustering process, trajectory analysis was performed with  
16 functions in scGREAT utilizing diffusion pseudo-time (Haghverdi *et al.*, 2016) implemented by  
17 Scanpy (Wolf *et al.*, 2018), and the pseudo-time labels were transferred back to single-cell data  
18 (Figure 3A middle). While the trajectory agrees with the pseudo-time inferred using scRNA-seq  
19 data only (Figure 3A right), the regulatory state-based results were smoother and better separated  
20 cell types CD4+ TCM and CD4+ TEM as shown in the density plots (Figure 3B).

21

22 To dissect the trajectory pattern, we implemented the self-organizing map (SOM) (Vettigli, 2018;  
23 Kohonen, 1990) in our package for feature module discovery such that features with similar

1 temporal patterns will be clustered together. The SOM algorithm is an unsupervised artificial  
2 neural network that groups all observations in a low-dimensional representation while preserving  
3 the topological structure of the data. Here, we take genes or peaks as observations and the  
4 averaged expression levels within each pseudo-time bins as features, and clustered the genes or  
5 peaks by SOM. Users can run the function multiple times to have the optimized results by  
6 changing the SOM shape (how many modules are needed and what is the similarity among all  
7 modules), learning rate, number of iterations, and regulation power sigma. After optimizing the  
8 SOM, scGREAT can visualize the results to give users an intuitive understanding how genes or  
9 peaks are changing along the pseudo-timeline. In our example of CD4+ T cells analysis, the self-  
10 organizing map grouped all the genes into nine modules with different patterns (Figure 3C).  
11 Module 5 and module 8 have similar trends in which high expression in CD4+ Naïve T cells  
12 dropped in CD4+ TCM and raised slightly in CD4+ TEM cells. We performed Gene Ontology  
13 enrichment analysis to determine the function of the two gene modules, shown in Figure 3D. The  
14 functions discovered by GO analysis are consistent with the cellular process of the inferred  
15 trajectory and CD4+ T cells, including lymphocyte differentiation and T cell activation. In our  
16 regulatory state analysis, genes and peaks are always linked. Following the detection of gene-  
17 linked peaks in module 5 and 8, we performed motif enrichment analysis with *Homer* (Heinz *et*  
18 *al.*, 2010). All top enriched motifs are from the ETS transcription factors family binding domain  
19 and RUNX1 binding domain. Such results are consistent with the previous findings that ETS  
20 family of transcription factors are related to lymphoid differentiation (Russell and Garrett-Sinha,  
21 2010). Moreover, RUNX1 functions in the development of normal hematopoiesis (Fujimoto *et*  
22 *al.*, 2007), and controls the energy and suppressive function of regulatory T-cells (Ono *et al.*,  
23 2007).

1

### 2 **3.3 Regulatory marker discovery based on the local $L$**

3 With the regulatory state matrix from local  $L$  and cluster labels, we implemented the t-test for  
4 differentially correlated gene-peak pairs in the same way as DEGs (Differentially Expressed  
5 Genes) and DARs (Differentially Accessible Regions) analysis. If a certain gene-peak pair is  
6 differentially correlated in one cluster compared with all other clusters, it is considered a  
7 potential regulatory marker worth further studying. In contrast, Pearson correlation which  
8 doesn't have such a local representation, cannot achieve the same aim with only a single value.

9 Here, we used functions (*FindAllMarkers* and *FindMarkers*) in scGREAT to find all regulatory  
10 markers in peripheral blood mononuclear cell types (Supplementary Figure 2A). With scGREAT,  
11 regulatory state differences can be visualized in a heatmap with cell labels annotated  
12 (Supplementary Figure 2B). Next, scGREAT functions can help select significant markers of  
13 each cluster by the mean difference, adjusted p-value, and feature sparsity, then visualize the  
14 differential regulatory pairs with the volcano plot (Supplementary Figure 2C). If a pair is of great  
15 interest, the local  $L$  index in pseudo-bulk data can also be mapped back to the original single-cell  
16 data (Supplementary Figure 2D, Package Details 1.5).

17

18 With this regulatory marker discovery method, we studied the regulatory changes in B cell  
19 development and CD4<sup>+</sup> T cell differentiation (Figure 4A, B). After the markers have been found,  
20 we extracted the peaks and performed motif enrichment analysis. Besides RUNX1 binding  
21 domain mentioned previously in T cell differentiation, the PU.1 binding motif was enriched in B  
22 cell development. This transcription factor is encoded by the *Spi1* gene, binds with the PU-box,  
23 and activates gene expression during myeloid and B-lymphoid cell development (le Coz *et al.*,

1 2021). Genes linked with these peaks are potentially the targets of the transcription factors. So,  
2 we collected the genes potentially regulated by PU.1 and RUNX1, and performed Gene  
3 Ontology enrichment analysis (Figure 4C, D). During B cell development, regulation of antigen  
4 receptor-mediated signaling pathway has changed, as well as the immunoglobulin subunit Fc  
5 receptor-mediated stimulatory signaling pathway. During naïve T cell differentiation, antigen  
6 processing and presentation process has changed; lymphocyte anergy functions have also  
7 changed, indicating the activation of CD4+ T cell from Central memory cells to effector memory  
8 cells. All these analyses have validated the markers discovered by our method.

9

### 10 **3.4 The global $L$ index is a better replacement for Pearson's $r$**

11 The global  $L$  index, which averages local  $L$ -index values across all cells, serves the same role as  
12 Pearson's correlation to show the general correlation trend in the data. Both Pearson's and our  
13 graph-based ( $L$ ) correlation are variants of Mantel's general cross-product association measure  
14 (Mantel, 1967). But the graph-based correlation takes the single-cell dataset as a K-nearest  
15 neighbor (KNN) graph and cares about the neighborhood patterns, which varies from Pearson  
16 correlation (see Method 2.2 for details).

17

18 When measuring the general correlation trend (no consideration in heterogeneity) within the  
19 dataset,  $L$  index and Pearson's  $r$  are consistent. This pattern was observed in all three datasets  
20 (Figure 5A), with statistically significant high correlations. This consistency was still there when  
21 analyzing a single cluster rather than the whole dataset (Supplementary Figure 3A).

22

1 Beyond the consistency, only the global  $L$  index can rule out potentially confounding regulatory  
2 pairs that are not of research interests, which is an advantage of our method over Pearson's  $r$ . To  
3 demonstrate this observation, we took the CD4+ T cell cluster from 10X Multiome data as an  
4 example. In Figure 5B, we visualized the data by UMAP and labeled the cells with their cell  
5 cycle phases and differentiation stages. It is clearly shown that cell cycle phases are mixing in  
6 the UMAP while the differentiation stages are separated off. Thus, we believe the CD4+ T cell  
7 data is about T cell differentiation and activation. And the regulatory pairs between cell cycle  
8 genes and their nearby cis-regulatory elements, as one of many examples, are confounding pairs  
9 to the research. To infer the regulatory relationships in T cell differentiation, we would like the  
10 correlation method to exclude such confounding pairs.

11 Here, by taking cell cycle genes and their nearest peaks (most likely to be promoter regions) as  
12 confounding examples, we compared the behavior of  $L$  index and Pearson's  $r$  in CD4+ T cells  
13 from 10X Multiome data, cluster 7 from SHARE-seq data, and cluster 2 from SNARE-seq data.  
14 Figure 5C shows the correlation (top) and significance test (bottom) between cell cycle genes  
15 and their nearest upstream peaks. The  $L$  index tends to have smaller values than Pearson's  $r$  and  
16 distinguished the confounding pairs through significant test. We think it is because such cell  
17 cycle-related gene expression and peak accessibility are uniformly distributed in the data  
18 (examples in Supplementary Figure 4A). Though they correlate well, there are low and high  
19 values mixing within each cell neighborhood. Such discordance in the neighborhoods will  
20 influence the graph-based  $L$  index measurement and significance test. On the contrary, the  
21 immune-function related gene TARBP1 has a higher  $L$  index with its nearby peak compared to  
22 cell cycle gene YWHAZ and its nearby peak, though the two pairs have the same Pearson's  
23 correlation  $r$  (Figure 5D, E). By taking the neighborhood consistency into the correlation

1 measurement, the global  $L$  index ruled out such confounding regulatory pair. However, Pearson  
2 correlation doesn't know how cells are distributed in the data space and cannot achieve the same  
3 effect.

4

5 In practice, this behavior of our correlation measurement could enhance the specificity for  
6 downstream analysis. For example, when gene-peak pairs are ranked by correlation and taken for  
7 Gene Set Enrichment Analysis (GSEA) (Subramanian et al., 2005), ruling out such confounding  
8 pairs will make the molecular process and pathways of research interest stand out from the pool.  
9 Here, we compared the GSEA results from  $L$  index-based and Pearson's  $r$ -based gene ranks.  
10 Because of the difference in ranks, only 96 out of 8995 gene sets were significantly enriched  
11 from Pearson's  $r$ -based ranks, while 355 were discovered in  $L$  index-based ranks (adjusted  $p$ -  
12 value  $< 0.01$ ). Then, we collected all T cells related sets and compared the results. In general, T  
13 cell-related gene sets rank higher from  $L$ -based correlation (Supplementary Figure 4C). Among  
14 the 15 significantly enriched T cell sets in  $L$  index-based results, only 4 can be discovered with  
15 Pearson's  $r$ -based ranks (Supplementary Figure 4D).

16

17

#### 18 **4. Conclusion and Discussion**

19 scGREAT is a novel single-cell multi-omics analysis toolkit with a rapid graph-based correlation  
20 measurement and related analysis functions. Compared with other correlation methods like  
21 Pearson correlation, the graph-based correlation method cares about spatial dependence and  
22 poses a specially designed significance test for single-cell data. It thus enabled the filtering of  
23 confounding regulatory pairs out of research interests and kept only regulatory pairs with

1 patterns consistent with the main biological variance. Besides, with the local  $L$  index from the  
2 correlation method, scGREAT can generate the regulatory state matrix, which is a new layer of  
3 information. With the graph-based correlation scores, scGREAT filled the gap in multi-omics  
4 regulatory analysis by enabling labeled and unlabeled analysis, functional annotation, and  
5 visualization.

6 Pearson correlation, our graph-based correlation, and any other variants of Mantel's general  
7 cross-product association measure (Mantel, 1967), are affected by dropouts in single-cell data.  
8 After scaling the data with zero means and unit variances, dropout values will turn to negatives  
9 from zero. These negative but meaningless values will somewhat bias the final correlation results.  
10 To deal with this problem, implementations in single-cell scenarios need to pay special attention  
11 to feature sparsity and dropouts in all processes. In scGREAT, first, we offer users pseudo-bulk  
12 analysis options to decrease the dropout rate. Second, we keep the global  $L$  index with no  
13 modification, but offer users feature sparsity information to measure the influence of dropouts on  
14 final values quantitatively. Then for the local  $L$  index or the regulatory state matrix, scGREAT  
15 will zero the values when it is derived from any of the two features with dropout values, making  
16 the final matrix similar to any single-cell data with dropouts. Implementing this process is  
17 beneficial for both labeled and unlabeled analysis to avoid misleading results. Nevertheless, the  
18 future direction in measuring multi-omics regulatory relationships could be to avoid the effect of  
19 dropouts.

20

21 The groundbreaking Spatial-omics techniques have arisen in recent years (Asp *et al.*, 2020; Marx,  
22 2021; Deng *et al.*, 2022). These techniques profile omics data and map them back to the two- and  
23 three-dimensional geography of tissues, allowing researchers to study cell-cell communication



1 and tissue organization (Ståhl *et al.*, 2016; Moffitt *et al.*, 2016; Deng *et al.*, 2022). Our  $L$  index is  
2 a good fit for the gene co-expression analysis in spatial omics data. The  $L$  index considers spatial  
3 patterns in correlation measurements, and the local  $L$  index provides the correlation pattern in the  
4 geography space. To fit scGREAT into spatial data analysis, we need to replace the connectivity  
5 matrix with the spatial connections-based matrix, which is feasible. Thus, our future direction  
6 will be how to apply scGREAT to spatial transcriptomics spatial co-expression analysis and  
7 multi-omics spatial correlation analysis.

8

### 9 **Data and Code Availability**

10 The data underlying this article are publicly available (10X Multiome PBMC:  
11 <https://support.10xgenomics.com/single-cell-multiome-atac->  
12 [geq/datasets/1.0.0/pbmc\\_granulocyte\\_sorted\\_10k](https://support.10xgenomics.com/single-cell-multiome-atac-geq/datasets/1.0.0/pbmc_granulocyte_sorted_10k), SHARE-seq mouse skin: GSE140203, and  
13 SNARE-seq mouse: GSE126074).

14 All codes to replicate the results are available at the GitHub repository.

15

### 16 **Funding**

17 Research reported in this publication was supported by the Eunice Kennedy Shriver National  
18 Institute of Child Health & Human Development of the National Institutes of Health under  
19 Award Number P50HD103555 for use of the Bioinformatics Core facilities. The content is solely  
20 the responsibility of the authors and does not necessarily represent the official views of the  
21 National Institutes of Health. ZL, CL and LW are also partially supported by the Chao  
22 Endowment and the Huffington Foundation. The funders had no role in study design, data  
23 collection and analysis, decision to publish, or preparation of the manuscript.

1

## 2 **Acknowledgements**

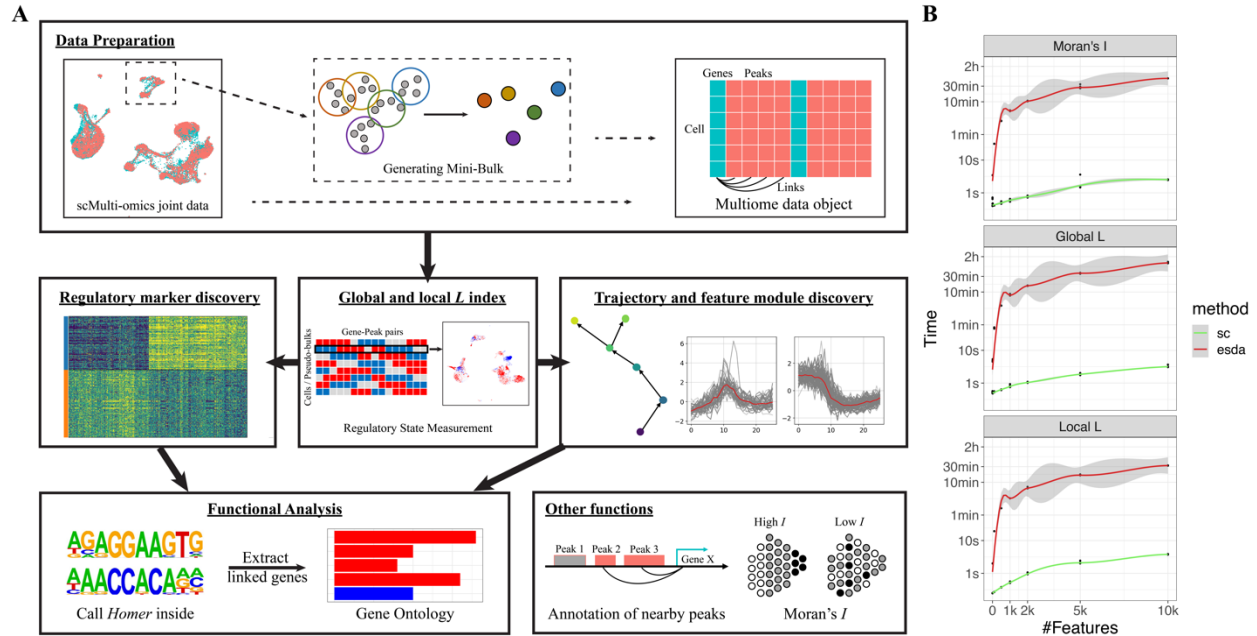
3 The authors would like to thank members of the Liu lab for their suggestions and discussion.

4

## 5 **References**

- 6 Asp, M. *et al.* (2020) Spatially Resolved Transcriptomes—Next Generation Tools for Tissue  
7 Exploration. *BioEssays*, **42**, 1900221.
- 8 Benesty, J. *et al.* (2009) Pearson Correlation Coefficient. In, Cohen, I. *et al.* (eds), *Noise Reduction*  
9 *in Speech Processing*. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 1–4.
- 10 Bredikhin, D. *et al.* (2022) MUON: multimodal omics analysis framework. *Genome Biol*, **23**, 42.
- 11 Cao, Z.-J. and Gao, G. (2022) Multi-omics single-cell data integration and regulatory inference  
12 with graph-linked embedding. *Nat Biotechnol*, **40**, 1458–1466.
- 13 Chen, S. *et al.* (2019) High-throughput sequencing of the transcriptome and chromatin  
14 accessibility in the same cell. *Nat Biotechnol*, **37**, 1452–1457.
- 15 Clarke, Z.A. *et al.* (2021) Tutorial: guidelines for annotating single-cell transcriptomic maps using  
16 automated and manual methods. *Nat Protoc*, **16**, 2749–2764.
- 17 le Coz, C. *et al.* (2021) Constrained chromatin accessibility in PU.1-mutated  
18 agammaglobulinemia patients. *Journal of Experimental Medicine*, **218**, e20201750.
- 19 Deng, Y. *et al.* (2022) Spatial profiling of chromatin accessibility in mouse and human tissues.  
20 *Nature*, **609**, 375–383.
- 21 Fujimoto, T. *et al.* (2007) Cdk6 blocks myeloid differentiation by interfering with Runx1 DNA  
22 binding and Runx1-C/EBP $\alpha$  interaction. *EMBO J*, **26**, 2361–2370.
- 23 Granja, J.M. *et al.* (2021) ArchR is a scalable software package for integrative single-cell  
24 chromatin accessibility analysis. *Nat Genet*, **53**, 403–411.
- 25 Grün, D. *et al.* (2015) Single-cell messenger RNA sequencing reveals rare intestinal cell types.  
26 *Nature*, **525**, 251–255.
- 27 Haghverdi, L. *et al.* (2016) Diffusion pseudotime robustly reconstructs lineage branching. *Nat*  
28 *Methods*, **13**, 845–848.
- 29 Hao, Y. *et al.* (2021) Integrated analysis of multimodal single-cell data. *Cell*, **184**, 3573–3587.e29.
- 30 Heinz, S. *et al.* (2010) Simple Combinations of Lineage-Determining Transcription Factors Prime  
31 cis-Regulatory Elements Required for Macrophage and B Cell Identities. *Mol Cell*, **38**, 576–  
32 589.
- 33 Klemm, S.L. *et al.* (2019) Chromatin accessibility and the regulatory epigenome. *Nat Rev Genet*,  
34 **20**, 207–220.
- 35 Kohonen, T. (1990) The self-organizing map. *Proceedings of the IEEE*, **78**, 1464–1480.
- 36 Lähnemann, D. *et al.* (2020) Eleven grand challenges in single-cell data science. *Genome Biol*, **21**,  
37 31.
- 38 Lee, S.-I. (2001) Developing a bivariate spatial association measure: An integration of Pearson's  
39  $r$  and Moran's  $I$ . *J Geogr Syst*, **3**, 369–385.

- 1 Li, Y. *et al.* (2015) The identification of cis-regulatory elements: A review from a machine  
2 learning perspective. *Biosystems*, **138**, 6–17.
- 3 Liu, J. *et al.* (2020) Jointly defining cell types from multiple single-cell datasets using LIGER. *Nat*  
4 *Protoc*, **15**, 3632–3662.
- 5 Ma, S. *et al.* (2020) Chromatin Potential Identified by Shared Single-Cell Profiling of RNA and  
6 Chromatin. *Cell*, **183**, 1103-1116.e20.
- 7 Mantel, N. (1967) The Detection of Disease Clustering and a Generalized Regression Approach.  
8 *Cancer Res*, **27**, 209–220.
- 9 Marx, V. (2021) Method of the Year: spatially resolved transcriptomics. *Nat Methods*, **18**, 9–14.
- 10 Minnoye, L. *et al.* (2021) Chromatin accessibility profiling methods. *Nature Reviews Methods*  
11 *Primers*, **1**, 10.
- 12 Moffitt, J. R. *et al.* (2016) High-throughput single-cell gene-expression profiling with multiplexed  
13 error-robust fluorescence in situ hybridization. *Proceedings of the National Academy of*  
14 *Sciences*, **113**, 11046–11051.
- 15 Ono, M. *et al.* (2007) Foxp3 controls regulatory T-cell function by interacting with AML1/Runx1.  
16 *Nature*, **446**, 685–689.
- 17 Pliner, H. A. *et al.* (2018) Cicero Predicts cis-Regulatory DNA Interactions from Single-Cell  
18 Chromatin Accessibility Data. *Mol Cell*, **71**, 858-871.e8.
- 19 Preissl, S. *et al.* (2022) Characterizing cis-regulatory elements using single-cell epigenomics. *Nat*  
20 *Rev Genet*.
- 21 Rey, S. J. and Anselin, L. (2007) PySAL: A Python Library of Spatial Analytical Methods. *Rev Reg*  
22 *Stud*, **37**, 5–27.
- 23 Russell, L. and Garrett-Sinha, L. A. (2010) Transcription factor Ets-1 in cytokine and chemokine  
24 gene regulation. *Cytokine*, **51**, 217–226.
- 25 Saelens, W. *et al.* (2019) A comparison of single-cell trajectory inference methods. *Nat*  
26 *Biotechnol*, **37**, 547–554.
- 27 Solt, L. A. *et al.* (2011) Suppression of TH17 differentiation and autoimmunity by a synthetic ROR  
28 ligand. *Nature*, **472**, 491–494.
- 29 Ståhl, P. L. *et al.* (2016) Visualization and analysis of gene expression in tissue sections by spatial  
30 transcriptomics. *Science (1979)*, **353**, 78–82.
- 31 Stuart, T. *et al.* (2019) Comprehensive Integration of Single-Cell Data. *Cell*, **177**, 1888-1902.e21.
- 32 Velten, L. *et al.* (2017) Human haematopoietic stem cell lineage commitment is a continuous  
33 process. *Nat Cell Biol*, **19**, 271–281.
- 34 Vettigli, G. (2018) MiniSom: minimalistic and NumPy-based implementation of the Self  
35 Organizing Map.
- 36 Virshup, I. *et al.* (2021) anndata: Annotated data. *bioRxiv*, 2021.12.16.473007.
- 37 Wolf, F. A. *et al.* (2018) SCANPY: large-scale single-cell gene expression data analysis. *Genome*  
38 *Biol*, **19**, 15.
- 39



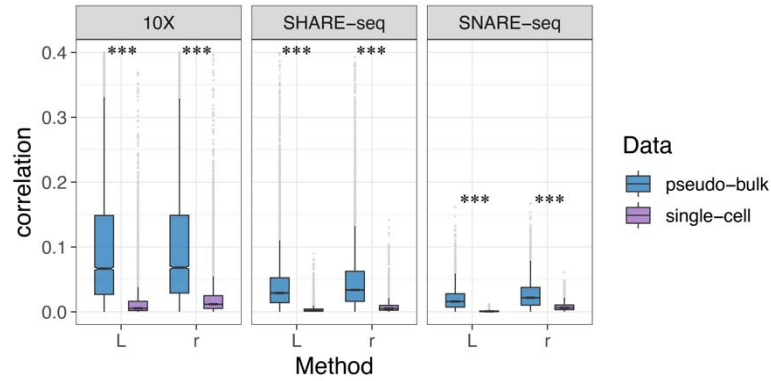
1 **Figure 1.** scGREAT overview. (A) Analysis pipeline of scGREAT. scGREAT receives as input  
2 the integrated scRNA-seq and scATAC-seq data. Users can choose to generate pseudo-bulk data  
3 or skip this step and use the original single-cell data. After linking genes with nearby peaks,  
4 scGREAT will preprocess the raw data and generate the final multiome data object for further  
5 analysis. Next, before all other analyses, the  $L$  index will be computed with user-chosen  
6 parameters. With the global  $L$  value and local  $L$  matrix, users can perform labeled analysis  
7 (regulatory marker discovery) and unlabeled analysis (sub-clustering, trajectory, feature module  
8 discovery), followed by functional analysis to help explain *Homer* output results. (B)  
9 Computation time comparison between scGREAT and *esda*. Under the limitation of 16 Gb  
10 memory and eight threads, the time consumed to calculate Moran's  $I$ , global  $L$ , and local  $L$ , was  
11 estimated with different feature numbers.

12

13

14

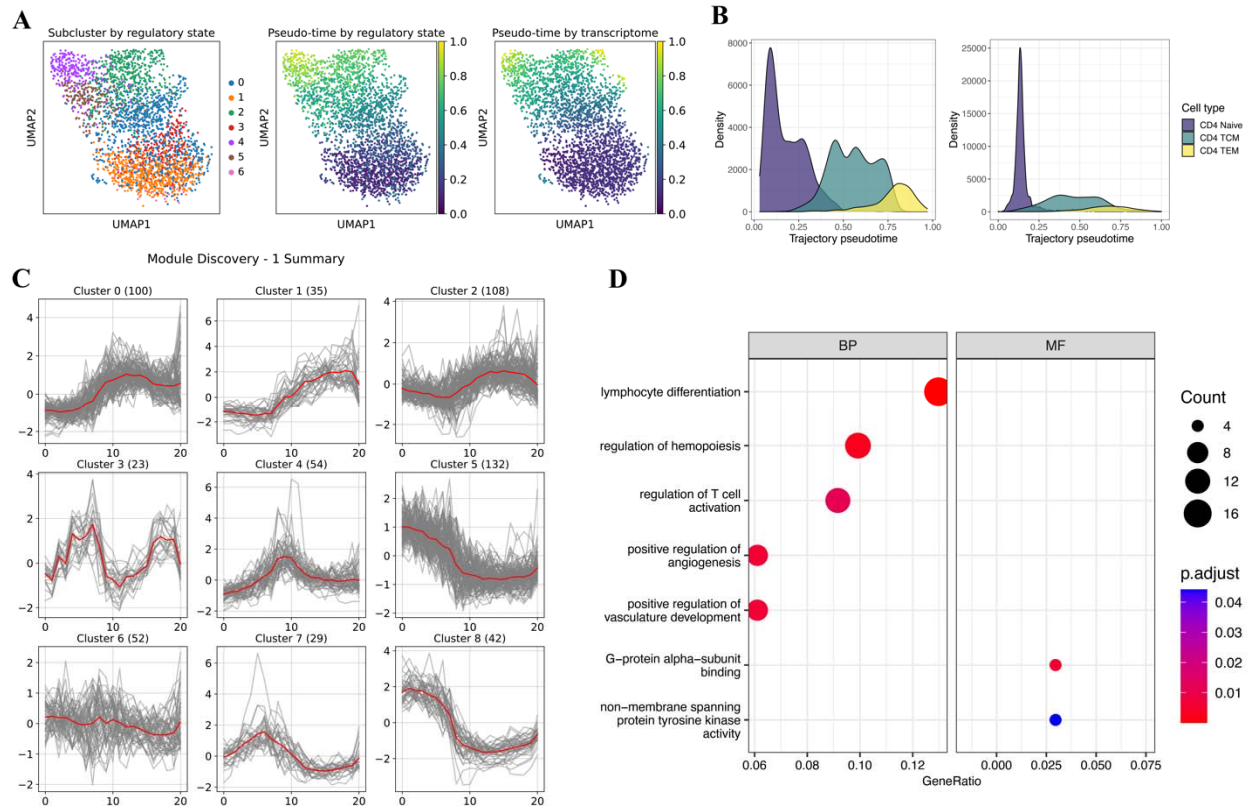
15



1 **Figure 2.** The correlation between genes and their nearest upstream peaks by  $r$  and  $L$  in pseudo-  
2 bulk and single-cell data. Three datasets were applied, including 10X Multiome PBMC,  
3 SHARE-seq mouse skin, and SNARE-seq mouse brain, to compare pseudo-bulk and single-cell  
4 data. Paired t-test was done for significance.

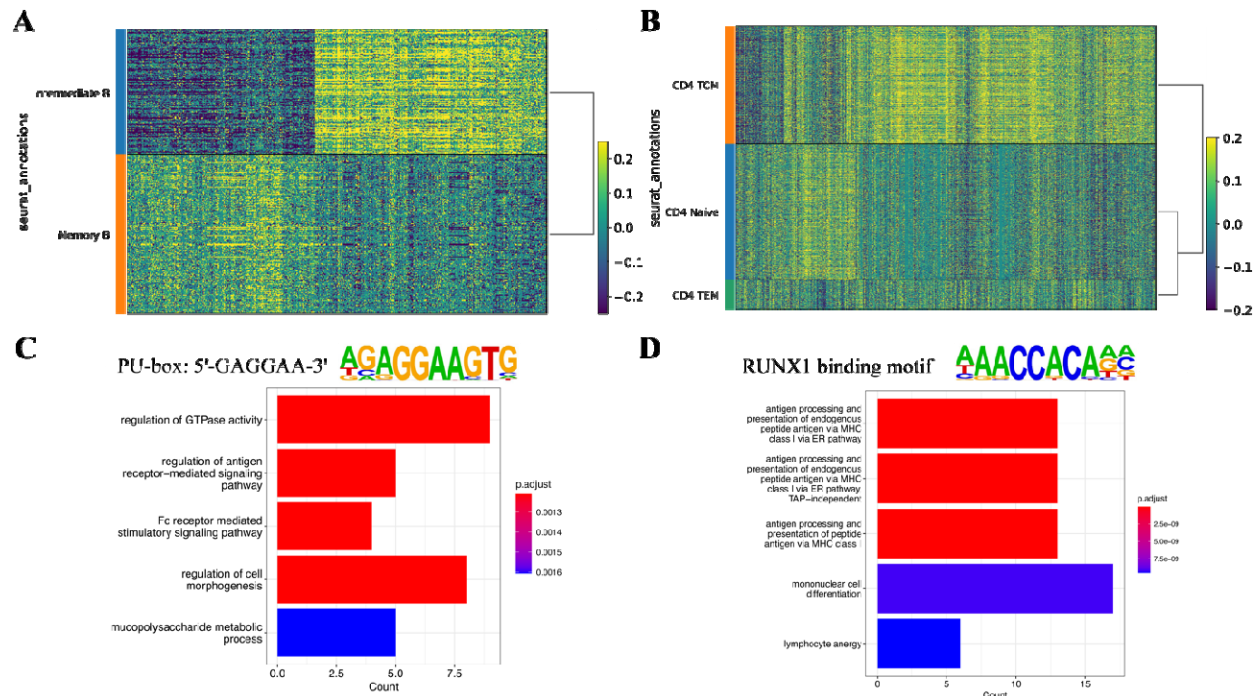
5

6



1 **Figure 3.** Heterogeneity analysis on CD4+ T cells. (A) UMAP labeled by subclusters (left),  
 2 trajectory pseudo-time (middle) derived from regulatory state, trajectory pseudo-time (right)  
 3 derived from the transcriptome. (B) Density plots showing the pseudo-time distribution in  
 4 regulatory state and transcriptome-derived trajectory, grouped by cell sub-types. (C)  
 5 Visualization of the self-organizing map (SOM) optimization results. It demonstrates the features  
 6 modules found by SOM. (D) Gene Ontology enrichment results of genes in clusters 5 and 8. Size  
 7 represents the number of genes; color represents the adjusted p-value.

8  
 9  
 10  
 11  
 12



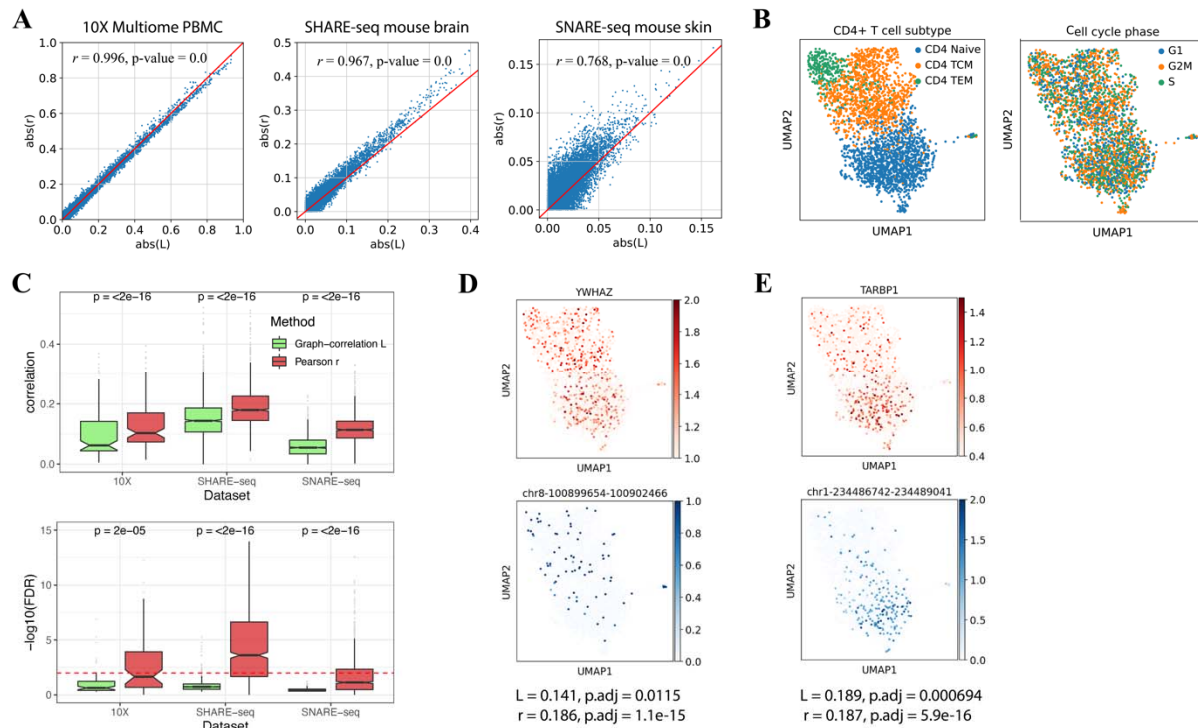
1 **Figure 4.** Markers discovery on 10X Multiome PBMC dataset. (A) Heatmap showing the gene-  
2 peak correlations in pseudo-bulk data of B cell clusters. Each row is a pseudo-bulk, each column  
3 is a gene-peak pair, and color indicates the correlation level. (B) Heatmap showing the gene-peak  
4 correlations in pseudo-bulk data of CD4+ T cell clusters. (C) Differentially correlated peaks in B  
5 cell clusters were enriched in PU-box binding motif. Genes linked with the peaks were extracted  
6 for Gene Ontology enrichment analysis to discover the functional changes. (D) Differentially  
7 correlated peaks in CD4+ T cell clusters were enriched in the RUNX1 binding motif. Genes  
8 linked with the peaks were extracted for Gene Ontology enrichment analysis to discover the  
9 functional changes.

10

11

12

13



1 **Figure 5.** Comparison between Pearson's coefficient  $r$  and global  $L$  index. (A) The consistency  
 2 between  $L$  and  $r$  in studying the general trend of regulatory relationships. The scatter plots  
 3 showed all gene-peak pairs correlation levels measured with  $r$  (y-axis) and  $L$  (x-axis). The  
 4 consistency between  $r$  and  $L$  was tested significantly by Pearson correlation (results shown in the  
 5 plots). (B) UMAP visualization of CD4+ T cells labeled by T cell subtype (left) and cell cycle  
 6 phase (right). (C) Box plots showing the correlation level (top) and significance test FDR  
 7 (bottom) measured by  $L$  and  $r$ , between cell cycle genes and their nearby peaks. (D) UMAP  
 8 visualization of cell cycle gene YWHAZ (top), and its cis-regulatory element chr8-100899654-  
 9 100902466 (bottom), (E) immune-related gene TARBP1 (top), and its cis-regulatory element  
 10 chr1-234486742-234489041 (bottom). The two pairs have similar Pearson correlations but  
 11 because of the pattern in data structure, their  $L$  index and significant test differ.

12

13