# Supplementary Information for "Unbiased integration of single cell multi-omics data"

Jinzhuang Dou[1], jdou1@mdanderson.org

Shaoheng Liang[1], sliang3@mdanderson.org

Vakul Mohanty[1], vmohanty@mdanderson.org

Xuesen Cheng[2], xuesenc@bcm.edu

Sangbae Kim[2], Sangbae.Kim@bcm.edu

Jongsu Choi[2], Jongsu.Choi@bcm.edu

Yumei Li[2], yumeil@bcm.edu

Katayoun Rezvani[4], krezvani@mdanderson.org

Rui Chen[2,3], ruichen@bcm.edu

Ken Chen[1,*]. kchen3@mdanderson.org


[1]Department of Bioinformatics and Computational Biology, The University of Texas MD Anderson Cancer Center

[2]HGSC, Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX, 77030, USA

[3]Verna and Marrs McLean Department of Biochemistry and Molecular Biology, Baylor College of Medicine,

Houston, TX, 77030, USA

[4]Department of Stem Cell Transplantation and Cellular Therapy, The University of Texas MD Anderson Cancer

Center, Houston, Texas

* Correspondence: kchen3@mdanderson.org

# Supplementary Notes

**Supplementary Note 1 Previous studies on multi-omics integration**

Recent study [1] evaluated 14 single-cell batch-effect correction/integration methods, showing that Harmony [2], LIGER [3], and Seurat3.0 [4] are the recommended methods for batch integration in general. Thus, we compare bindSC with these three available state-of-the-art methods.

*Harmony*

Harmony [2] uses an iterative clustering approach to align cells from different batches. The algorithm first combines the batches and projects the data into a dimensionally reduced space using PCA. It then uses an iterative procedure to remove the multi-dataset specific effects. In our analysis, we ran Harmony within the Seurat3.0 based on the guide (http://htmlpreview.github.io/?https://github.com/immunogenomics/harmony/blob/master/docs/SeuratV3.html).

*Seurat3.0*

Seurat [4] uses CCA to first compute the linear combinations of genes with the maximum correlation between batches then adopts mutual nearest neighbor (MNN) to align the cells between datasets based on anchor cells identified. In our analysis, we used the Seurat package version 3.0 in the *R* language environment to perform multi-omics integration. Adhering to the suggested integration workflow (https://satijalab.org/seurat/v3.2/atacseq_integration_vignette.html).

*LIGER*

LIGER [3] uses integrative non-negative matrix factorization (iNMF) to first learn a low-dimensional space where each gene is characterized by two sets of factors. The first set contains dataset-specific factors, and the second contains shared factors. The shared factor space is then used to identify similar cell types across datasets by first constructing a shared factor neighborhood

44 graph to connect cells with similar factor loading patterns. Joint clusters are then identified using

45 the Louvain community detection. Thereafter, the factor loading quantiles are normalized using

46 the largest data batch as a reference to achieve batch-correction. In our work, we followed the

47 LIGER documentation

48 (http://htmlpreview.github.io/?https://github.com/MacoskoLab/liger/blob/master/vignettes/Integr

49 ating_scRNA_and_scATAC_data.html). For preprocessing, we used the LIGER preprocessing

50 functions, where we first selected genes with high variance. We then performed iNMF-based

51 factorization using an alternating least squares algorithm, followed by data alignment using joint

52 clustering and quantile alignment.

53

54 **Supplementary Note 2 Evaluation of peak-gene correlations based on pseudo-cell profiles**

55 On the DEX-treated lung adenocarcinoma (A549) dataset, we ran bindSC to derived co-

56 embeddings. The shared nearest neighbor (SNN) graph was constructed by calculating the l-

57 nearest neighbors ($l$ = 20 by default) based on the Euclidean distance of L2-normalized co-

58 embedding coordinates. The modularity optimization technique Leiden algorithm was used to

59 group cells into interconnected clusters (termed meta-cluster) based on constructed SNN graph

60 (*resolution* = 0.5). The Leiden algorithm was performed again on each cluster with a higher

61 resolution (= 2) to further generate pseudo-cells. Finally, we got 206 pseudo-cells which included

62 a median of 27 cells from scRNA-seq and 16 cells scATAC-seq dataset (**Supplementary Fig. 6e**).

63 We observed only one cell that was modality specific (scRNA-seq) and removed it for downstream

64 analysis. The RNA-seq and ATAC-seq profiles of each pseudo-cell were aggregated respectively.

65 In this way, each pseudo-cell had paired gene expression and chromatin accessibility profiles. The

66 same strategy was used to construct pseudo-cell profiles for Seurat, LIGER, and Harmony. For

67    Seurat, LIGER, and Harmony, 41/198, 1/89, and 15/142 modality specific pseudo-cells were

68    removed, respectively. A high proportion of modality-specific pseudo-cells indicates that two

69    modalities were not well aligned in co-embeddings.

70

71    We then explored peak-gene correlation based on pseudo-cell profiles from each method. For each

72    peak-gene pair, Spearman rank correlation coefficients (SRCC) between a normalized ATAC peak

73    level and a gene expression levels of all the pseudo-cells were calculated. There are 4,759 genes

74    and 24,953 peaks in the peak-gene correlation matrix. The SRCC of each peak-gene pair calculated

75    based on 1,429 co-assayed cell profiles was used as the gold standard including 1,836,974 cis

76    peak-gene pairs and 118.7 million trans peak-gene pairs. The overall concordance between each

77    method and the gold-standard was further quantified using a single SRCC across all peak-gene

78    pairs (**Fig. 3c**). In most cases, the correlation of peak-gene may include many false positive and

79    indirect targets. We therefore focused on peak-gene pairs that were supported by Hi-C data from

80    an independent study [5]. There were 585 trans peak-gene pairs associated with the top 200 *NR3C1*

81    target binding genes identified. Among these trans peak-gene pairs, bindSC has the best agreement

82    with that from co-assayed cell profiles among all methods (**Supplementary Fig. 5**).

83

84    To explore TF-gene correlation at the pseudo-cell level, we obtained motif-based TF activity

85    matrix calculated based on peak profiles using ChromVAR [6]. The final TF activity matrix included

86    profiles for 386 TFs. Pseudo-cell level TF activity was obtained by aggregating cell profiles in

87    each pseudo-cell. The SRCC was calculated for each TF-gene pair on pseudo-cell level. Overall,

88    SRCC was 0.67 for bindSC and less than 0.59 for other methods (**Fig. 3c**).  The SRCC of TF-gene

89    pairs was higher than that from peak-gene pairs partly due to the fact that motif-based TF activity

90    was derived from genome-wide motif regions and it was less noisy than single peak region.

91

92    **Supplementary Note 3 Joint profiling of chromatin accessibility and transcription in DEX-**

93    **treated A549 cells**

94    Besides using the DEX-treated A549 cell dataset as the gold standard for method performance

95    evaluation, we performed downstream analysis to show how bindSC improved previous studies

96    by integrating transcriptomic and epigenomic datasets. Joint clustering module in bindSC defined

97    5 clusters (**Supplementary Fig. 6a**). Cells from the two technologies were well mixed together

98    within each cluster. This classification result was in good concordance with the treatment time:

99    cluster 1 consists of cells from mostly 0-hour (92%), and clusters 3-5 include cells from 1 and 3

100   hours (> 99%). Clusters 2 included cells from multiple time points and may represent transitional

101   states (**Supplementary Fig. 7b**). The list of transcription factors (TFs) that are associated with the

102   joint chromatin accessibility and gene expression changes and their activity levels across states

103   can be derived at pseudo-cell resolution, and so can the genes differentially expressed in each

104   cluster (**Supplementary Fig. 7d**). Such co-embedding yielded higher granularity in delineating

105   cell states and associated TFs than did embeddings derived from only one modality or based on

106   the treatment times.

107

108   **Supplementary Note 4 Integrating single cell epigenomic data with single cell transcriptomic**

109   **data on the mouse skin cell dataset**

110   We examined the performance of bindSC in integrating the scRNA-seq and scATAC-seq data

111   derived from mouse skin tissue. This dataset was generated using SHARE-seq [7] which included

112    34,774 cells that have joint profiles of RNA and ATAC profiles. The final ATAC-seq matrix (i.e.,

113    **Y**) includes 25,594 cells on 74,161 peaks after quality control (including removing cells with less

114    than 350 genes expressed; peaks that exist in less than 500 cells). In addition, 4,894 genes were

115    identified that were highly variable in both gene expression and gene activity profiles (i.e., both **X**

116    and **Z** includes 25,594 cells on 4,894 genes). For this evaluation, we only focused on the third

117    metric (i.e., anchoring accuracy) that represents the chance for the two instances of a co-assayed

118    cell to appear from the co-embeddings. The dimensionality $E$ was set to 15 for bindSC. BindSC

119    achieved substantially shorter anchoring distance than the other methods (**Supplementary Fig. 7**).

## Reference

1  Tran, H. T. N. *et al.* A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome biology* **21**, 1-32 (2020).

2  Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony. *Nature methods*, 1-8 (2019).

3  Rosenberg, A. B. *et al.* Single-cell profiling of the developing mouse brain and spinal cord with split-pool barcoding. *Science* **360**, 176-182 (2018).

4  Stuart, T. *et al.* Comprehensive integration of single-cell data. *Cell* **177**, 1888-1902. e1821 (2019).

5  Consortium, E. P. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57-74 (2012).

6  Schep, A. N., Wu, B., Buenrostro, J. D. & Greenleaf, W. J. chromVAR: inferring transcription-factor-associated accessibility from single-cell epigenomic data. *Nature methods* **14**, 975-978 (2017).

7  Ma, S. *et al.* Chromatin Potential Identified by Shared Single-Cell Profiling of RNA and Chromatin. *Cell* (2020).

## Supplementary Figure/Table Legends

**Fig. S1 Implementation of bindSC.** Bi-CCA iteration procedure (**a**). Implementation of divide-and-conquer SVD in bi-CCA for large matrix SVD decomposition (**b**).

**Fig. S2 Simulation settings.** Simulation of gene score matrix $\mathbf{Z}$ (**a**). Each row in $\mathbf{X}$ denotes a gene (feature) and each column a cell. MR: misalignment rate; SNR: Signal-Noise-Ratio. Previous methods including CCA, Seurat, LIGER, and Harmony take $\mathbf{X}$ and $\mathbf{Z}$ as input assuming that cell alignment is unknown (**b**). bindSC takes two parts as input: 1) $\mathbf{X}$ and $\mathbf{Z}$ with cell alignment unknown; 2) $\mathbf{X}$ and $\mathbf{Z}$ with feature alignment unknown (**c**).

**Fig. S3 Benchmarking bindSC performance on simulation datasets.** Comparison of bindSC to CCA, Seurat, LIGER, and Harmony based on Silhouette score and alignment mixing score (**a**). The dataset contains 1,000 genes and 1,000 cells equally distributed in 3 cell types. Signal-to-noise ratio (SNR) was set at 0. X-axes denote the misalignment rates (MR) between features in the two datasets, which ranges from 0 to 1. The features between two datasets have perfect match if MR = 0 and are unrelated if MR = 1. UMAP views of the co-embeddings generated by bindSC, CCA, Seurat, LIGER, and Harmony (**b**). From top to bottom are results with MR = 0.1, 0.5, and 0.9, respectively. Each point denotes one cell that is colored based on its true cell type label (red, green, or cyan).

**Fig. S4 Benchmarking bindSC performance on simulation datasets.** Comparison of bindSC to CCA, Seurat, LIGER, and Harmony based on Silhouette score and alignment mixing score (**a**). The dataset contains 1,000 genes and 1,000 cells equally distributed in 3 cell types. Signal-to-noise ratio (SNR) was set at 0.5. X-axes denote the misalignment rates (MR) between features in the two datasets, which ranges from 0 to 1. The features between two datasets have perfect match if MR = 0 and are unrelated if MR = 1. UMAP

162  views of the co-embeddings generated by bindSC, CCA, Seurat, LIGER, and Harmony **(b)**. From top to

163  bottom are results with MR = 0.1, 0.5, and 0.9, respectively. Each point denotes one cell that is colored

164  based on its true cell type label (red, green, or cyan).

165

166

167  **Fig. S5 Estimation of trans peak-gene regulatory elements supported by the Hi-C data.** We selected

168  the top 200 *NR3C1* target genes based on co-assayed cell profiles and identified 585 trans peak-gene

169  regulatory elements that were supported by the published Hi-C data [5]. X-axes are the SRCCs of peak-gene

170  pairs estimated from the co-assayed cells, which serve as the gold standard, while Y-axes are the SRCCs

171  estimated from pseudo-cells generated by each method. The overall concordance between X and Y are

172  further quantified using a single SRCC shown on the up-left corner of each subfigure. Also, the peak-gene

173  pair CFLAR@chr2:217,704,437-201,770,992 is highlighted in each subfigure. Pearson's correlation was

174  performed to produce the coefficients (R) and the P values.

175

176  **Fig. S6 Joint profiling of gene expression and chromatin accessibility data at the pseudo-cell**

177  **resolution on the A549 lung cancer cell-line.** UMAPs of cells coloring by cluster IDs obtained from

178  unsupervised clustering (meta-cluster) in the bindSC co-embedding **(a)**. Proportion of cells from the 3

179  treatment times in each of the meta-cluster **(b)**. Histogram showing the number of cells in each pseudo-cell

180  **(c)**. Heatmap showing known genes and TFs associated with glucocorticoid receptor (GR) activation

181  process **(d)**. Each row is one gene/TF and each column is one pseudo cell, grouped/colored by cluster ID.

182  Scatter plot showing the number of cells derived from the scRNA-seq and the scATAC-seq data for each

183  pseudo-cell **(e)**. Each dot denotes one pseudo-cell and the dot size denotes number of cells included in it.

184

185  **Fig. S7 Integrating single-cell RNA-seq and ATAC-seq data of a mouse skin cell atlas.** UMAP of the

186  mouse skin cells before performing integration, colored by clusters deriving from unsupervised clustering

187    of the RNA data and the ATAC data, respectively **(a)**.  Anchoring distances achieved by bindSC, Seurat,

188    LIGER and Harmony **(b)**. UMAP of cells in the multiomics co-embeddings generated by bindSC **(c)**, Seurat

189    **(d)**, LIGER **(e)**, and Harmony **(f)**, respectively. For each method, the left panel shows cells from the RNA-

190    seq data and the right panel shows cells from the ATAC-seq data.

191

192    **Fig. S8 Cell type annotation for cells in the mouse retina cell atlas**. In the heatmap, X-axes denote cluster

193    IDs in the RNA clusters, while Y-axes denotes known retinal cell-type-specific marker genes for the AC,

194    BC, cone, HC, RGC, rod, and RPC cells, respectively. The color gradient in each dot denotes the expression

195    level and the dot size denotes percentage of cells that express the gene.

196

197    **Fig. S9 UMAP visualization of mouse retina cells in the *in silico* co-embeddings generated by Seurat,**

198    **LIGER, and Harmony.** From top to bottom are the results for Seurat (**a**), LIGER (**b**), and Harmony (**c**)

199    respectively. The left panel shows cells from the RNA-seq data. The right panel shows cells from the

200    ATAC-seq data. Cells were colored based on cell types identified in **Supplementary Fig. 8**. The oval

201    regions were zoomed in **Fig. 4 g-j**.

202

203    **Fig. S10 Integrating 10x Visium spatial transcriptomics data with SMART-Seq2 scRNA-seq data**

204    **from mouse frontal cortex cells.**  Schematic representation of data used for integration **(a)**. Histology

205    image of mouse frontal cortex overlaying with cells from the 10x Visium technology **(b)**. UMAP view of

206    the RNA expression of the 1,072 spots in the 10x Visium data **(c)**. UMAP view of the transcriptomes of

207    14,249 frontal cortex cells produced by SMART-Seq2 technology **(d)**. Cell-type labels in **(d)** are derived

208    from the published SMART-Seq2 dataset.

209

210    **Fig. S11 Performance of four methods on integrating spatially resolved transcriptomic (ST) data with**

211    **dissociated scRNA-seq data from mouse frontal cortex cells.** Related to **Fig. 5a**. UMAP of cells from

212    mouse frontal cortex datasets, separated by sequencing technology with ST on the top panel and scRNA-

213    seq data on the bottom panel **(a)**. Cell-type labels are consistent with those from **Supplementary Fig. 10c-**

214    **d**. Comparison of cell-type classification based on Silhouette scores **(b)**. Comparison of dataset alignment

215    based on alignment mixing scores **(c)**. Gene expression profiles of three Lamp5-related marker genes *Lsp1*,

216    *Npy2r*, and *Dock5* from the scRNA-seq data **(d)** and the ST data **(e)**.

217

218    **Fig. S12 Cell types mapped by Seurat onto mouse brain histology images.**

219

220    **Fig. S13 Cell types mapped by LIGER onto mouse brain histology images.**

221

222    **Fig. S14 Cell types mapped by Harmony onto mouse brain histology images.**

223

224    **Fig. S15 Performance of three methods on integration of transcriptomic and proteomic data.** The

225    cluster colors for each modality are consistent with those in **Fig. 6**.

226

227    **Table S1 Summary of datasets evaluated in this study.** Also listed are the key parameters for running

228    bindSC, Seurat, LIGER, and Harmony on each dataset.

229

230    **Table S2 Simulation results with 5,000 cells**.

231

232    **Table S3 Simulation results with 10,000 cells**.