

1 **SMNN: Batch Effect Correction for Single-cell RNA-seq data via**  
2 **Supervised Mutual Nearest Neighbor Detection**

3 Yuchen Yang<sup>1\*</sup>, Gang Li<sup>2\*</sup>, Huijun Qian<sup>2</sup>, Kirk C. Wilhelmsen<sup>1,3</sup>, Yin Shen<sup>4,5</sup>, and Yun  
4 Li<sup>1,6,7,†</sup>

5 <sup>1</sup>Department of Genetics, <sup>2</sup>Statistics and Operations Research, <sup>6</sup>Biostatistics, <sup>7</sup>Computer Science, <sup>3</sup>Renaissance  
6 Computing Institute, University of North Carolina, Chapel Hill, NC 27599, USA. <sup>4</sup>Institute for Human Genetics,  
7 <sup>5</sup>Department of Neurology, University of California, San Francisco, San Francisco, CA 94143, USA

8 \* Equal contribution

## 9 **Abstract**

10 An ever-increasing deluge of single-cell RNA-sequencing (scRNA-seq) data has been generated, often involving  
11 different time points, laboratories or sequencing protocols. Batch effect correction has been recognized to be  
12 indispensable when integrating scRNA-seq data from multiple batches. A recent study proposed an effective  
13 correction method based on mutual nearest neighbors (MNN) across batches. However, the proposed MNN method  
14 is unsupervised in that it ignores cluster label information of single cells. Such cluster or cell type label information  
15 can further improve effectiveness of batch effect correction, particularly under realistic scenarios where true  
16 biological differences are not orthogonal to batch effect. Under this motivation, we propose SMNN which performs  
17 supervised mutual nearest neighbor detection for batch effect correction of scRNA-seq data. Our SMNN either takes  
18 cluster/cell-type label information as input, or, in the absence of such information, infers cell types by performing  
19 clustering of scRNA-seq data. It then detects mutual nearest neighbors within matched cell types and corrects batch  
20 effect accordingly. Our extensive evaluations in simulated and real datasets show that SMNN provides improved  
21 merging within the corresponding cell types across batches, leading to reduced differentiation across batches over  
22 MNN. Furthermore, SMNN retains more cell type-specific features after correction. Differentially expressed genes  
23 (DEGs) identified between cell types after SMNN correction are biologically more relevant, and the DEG true  
24 positive rates improve by up to 841%. SMNN is implemented in R, and freely available at  
25 <https://yunliweb.its.unc.edu/SMNN/> and <https://github.com/yycunc/SMNNcorrect>.

26

## 27 **Author summary**

28 The presence of batch effects poses grand challenges to integrative analysis of scRNA-seq data from multiple  
29 resources. One powerful tool MNN corrects batch effect of scRNA-seq data based on mutual nearest neighbors  
30 across batches. However, this method makes a critical assumption that batch effect is orthogonal to true biological  
31 differences. This assumption in practice can easily be violated. When that happens, MNN suffers from biases  
32 introduced by wrongly matched pairs of cells. To overcome this shortcoming, here we present a new method,  
33 SMNN, which performs supervised mutual nearest neighbor detection for batch effect correction. We benchmark the  
34 performance of SMNN using both simulations and real data, and demonstrate that, compared to MNN, our SMNN  
35 can better mix cells of the same type/state across batches. More importantly, SMNN can more effectively retain

36 biologically relevant features, and thereof provide improved cell type clustering and enhanced power for detecting  
37 differentially expressed genes (DEGs) between different cell types.

38

## 39 **Introduction**

40 An ever-increasing amount of single cell RNA-sequencing (scRNA-seq) data has been generated as scRNA-seq  
41 technologies mature and sequencing costs continue dropping. However, large scale scRNA-seq data, for example,  
42 those profiling tens of thousands to millions of cells (such as the Human Cell Atlas Project [1], almost inevitably  
43 involve multiple batches across time points, laboratories, or experimental protocols. The presence of batch effect  
44 renders joint analysis across batches challenging [2,3]. Batch effect, or systematic differences in gene expression  
45 profiles across batches, not only can obscure the true underlying biology, but also may lead to spurious findings.  
46 Thus, batch effect correction, which aims to mitigate the discrepancies across batches, is crucial and deemed  
47 indispensable for the analysis of scRNA-seq data across batches [4].

48 Because of its importance, a number of batch effects correction methods has been recently proposed and  
49 implemented. Most of these methods, including limma [5], ComBat [6], and svaseq [7], are regression-based.  
50 Among them, limma and ComBat explicitly model known batch effect as a blocking term. Because of the regression  
51 framework adopted, standard statistical approaches to estimate the regression coefficients corresponding to the  
52 blocking term can be conveniently employed. In contrast, svaseq is often used to detect underlying unknown factors  
53 of variation, for instance, unrecorded differences in the experimental protocols. svaseq first identifies these unknown  
54 factors as surrogate variables and subsequently corrects them. For these regression-based methods, once the  
55 regression coefficients are estimated or the unknown factors are identified, one can then regress out these batch  
56 effects accordingly, obtaining residuals that will serve as the batch-effect corrected expression matrix for further  
57 analyses. These methods have become standard practice in the analysis of bulk RNA-seq data. However, when it  
58 comes to scRNA-seq data, one key underlying assumption behind these methods, that the cell composition within  
59 each batch is identical, might not hold. Consequently, estimates of the coefficients might be inaccurate. As a matter  
60 of fact, when applied to scRNA-seq data, the corrected results derived from these methods widely adopted for bulk  
61 RNA-seq data might be even inferior to raw data without no correction, in some extreme cases [8].

62 To address the heterogeneity and high dimensionality of complex data, several dimension-reduction approaches  
63 have been adopted. An incomplete list of these strategies includes principal component analysis (PCA), autoencoder,  
64 or force-based methods such as t-distributed stochastic neighbor embedding (t-SNE) [9]. Through those dimension  
65 reduction techniques, one can project new data onto the reference dataset using a set of landmarks from the  
66 reference [8,10,11,12] to remove batch effects between any new dataset and the reference dataset. Such projection  
67 methods require the reference batch contains all the cell types across batches. As one example, Spitzer *et al.* [11]  
68 employed force-based dimension reduction and showed that leveraging a few landmark cell types from bone marrow  
69 (the most appropriate tissue in that it provides the most complete coverage of immune cell types) allowed mapping  
70 and comparing immune cells across different tissues and species. When applied to scRNA-seq data, however, these  
71 methods suffer when cells from a new batch fall out of the space inferred from the reference. Furthermore,  
72 determining the dimensionality of the low dimensional manifolds is still an open and challenging problem. To  
73 address the limitations of existing methods, a recent study proposed a new method, MNN correction, which  
74 leverages information of mutual nearest neighbors across batches. MNN has demonstrated superior performance  
75 over alternative methods [8]. MNN makes a critical assumption that true biological differences are orthogonal to  
76 batch effect. For this assumption to hold, variation from batch effect is required, at the minimum, to be much smaller  
77 than that from biological effect. Under this assumption, MNN finds, across batches, mutual nearest neighboring cells  
78 for each cell to be corrected, and then computes batch-effect correction vectors under a Gaussian kernel. However,  
79 this orthogonality assumption might not hold in real data, particularly given that different batches may easily differ  
80 in many aspects, including samples used, single cell capture method, or library preparation approach. Under non-  
81 orthogonal scenarios, MNN will not be optimal using its global (ignoring cell type information) nearest neighbor  
82 search strategy, leading to undesired correction results. For example, under the scenario depicted in **Fig 1b**, MNN  
83 leads to cluster 1 (C1) and cluster 2 (C2) mis-corrected due to mismatching single cells in the two clusters/cell-types  
84 across batches.

85 To address the above issue, here we present SMNN, a supervised version of MNN that incorporates cell type  
86 information. SMNN performs nearest neighbor searching within the same cell type, instead of global searching  
87 ignoring cell type labels (**Fig 1a**). Cell type information, when unknown *a priori*, can be inferred via clustering  
88 methods [13,14,15,16] .

89

## 90 **Results**

### 91 **SMNN Framework**

92 The motivation of our SMNN lies in the potential of single-cell cluster information to improve the accuracy of  
93 nearest neighbor (NN) identification. A preliminary clustering within each batch before any correction can provide  
94 knowledge regarding cell composition within each scRNA-seq dataset, which serves to encode the cellular  
95 correspondence across batches (**Fig 1a**). With this clustering information, we can refine the nearest neighbor  
96 searching space within a certain population of cells that are of the same or similar cell type(s) or state(s) in all the  
97 batches.

98 SMNN takes a natural two-step approach to leverage cell type label information for enhanced batch effect  
99 correction (**Fig 1c and S1 text**). First, it takes the expression matrices across multiple batches as input, and performs  
100 clustering separately for each batch. Specifically, in this first step, SMNN uses Seurat v. 3.0 [17] where dimension  
101 reduction is conducted via principal component analysis (PCA) to the default of 20 PCs, and then graph-based  
102 clustering follows on the dimension-reduced data with *resolution* parameter of 0.9 [18]. Obtaining an accurate  
103 matching of the cluster labels across batches is of paramount importance for subsequent nearest neighbor detection.  
104 SMNN requires users to specify a list of marker genes and their corresponding cell type labels to match clusters/cell  
105 types across batches. We later refer to this cell type or cluster matching as cluster harmonization across batches.  
106 Because not all cell types are necessarily shared across batches, and no prior knowledge exists regarding the exact  
107 composition of cell types in each batch, SMNN allows users to take discretion in terms of the marker genes to  
108 include, representing the cell types that, they believe, are shared across batches. Based on the marker gene  
109 information, a harmonized label is assigned to *every* cluster identified across all the batches according to two  
110 criteria: the percentage of cells in a cluster expressing a certain marker gene and the average gene expression levels  
111 across all the cells in the cluster. After harmonization, cluster labels are comparable across batches. This completes  
112 step 1 of SMNN. Note that this entire clustering step can be bypassed by feeding SMNN cluster labels that are  
113 consistent or comparable across batches.

114 With the harmonized cluster or cell type label information obtained in the first step, SMNN, in the second step,  
115 searches mutual nearest neighbors only within each matched cell type between the first batch (which serves as the  
116 reference batch) and any of the other batches (the current batch), and performs batch effect correction accordingly.  
117 Compared to MNN, where the mutual nearest neighbors are searched globally, SMNN identifies each pair of

118 neighbors from the same cell population or state, which can provide more accurate information for downstream  
119 correction. Then, following MNN, SMNN computes batch effect correction vector for each identified pair of cells  
120 and calculates the cell-specific correction vectors by exploiting a Gaussian kernel to obtain a weighted average  
121 across all the pair-specific vectors of their mutual nearest neighbors. Each cell's correction vector is further scaled  
122 according to the cell's location in the space defined by the correction vector, and standardized according to quantiles  
123 across batches, in order to eliminate the “kissing effects” phenomenon. “Kissing effects” refer to the phenomenon  
124 that naïve batch effect correction brings only the surfaces of two point-clouds in contact, rather than fully merging  
125 them [8]. At the end of the second step, SMNN returns the batch-effect corrected expression matrix for each batch,  
126 as well as the information regarding nearest neighbors between the reference batch and the current batch under  
127 correction. This step is carried out for every batch except the reference batch so that all batches are corrected to the  
128 same reference batch in the end.

129

## 130 **Simulation results**

131 Since MNN has been shown to excel alternative methods [4,8], we here focus on comparing our SMNN with MNN.  
132 We first compared the performance of SMNN to MNN in simulated data. In our simulations, SMNN demonstrates  
133 superior performance over MNN under both orthogonal and non-orthogonal scenarios (**Fig 2 and 3 and S1-3 Fig**).  
134 We show t-SNE plot for each cell type before and after MNN and SMNN correction under both the orthogonal and  
135 non-orthogonal scenarios. Under orthogonality, the two batches partially overlapped in the t-SNE plot before  
136 correction, suggesting that the variation due to batch effect was indeed much smaller than that due to biological  
137 effect. Both MNN and SMNN successfully mixed single cells from two batches (**S2 Fig**). However, for cell types 1  
138 and 3, there were still some cells from the second batch left unmixed with those from the first batch after MNN  
139 correction (**S2a and c Fig**). Under the non-orthogonal scenario, the differences between two batches were more  
140 pronounced before correction, and SMNN apparently outperformed MNN (**S3 Fig**), especially in cell type 1 (**S3a**  
141 **Fig**). Moreover, we also computed Frobenius norm distance [19] for each cell between its simulated true profile  
142 before introducing batch effects and after SMNN and MNN correction. The results showed an apparently reduced  
143 deviation from the truth after SMNN correction than MNN (**Fig 3**). These results suggest that SMNN provides  
144 improved batch effect correction over MNN under both orthogonal and non-orthogonal scenarios.

145

## 146 **Real data results**

147 For performance evaluation in real data, we first carried out batch effect correction on two hematopoietic datasets  
148 using SMNN and MNN, respectively. **Fig 4a-c** shows t-SNE plot before and after correction. Notably, both SMNN  
149 and MNN substantially mitigated discrepancy between the two datasets. Comparatively, SMNN better mixed cells  
150 of the same cell type across batches (**S4 Fig**), especially for CMP and MEP cells, which were wrongly corrected by  
151 MNN due to sub-optimal nearest neighbor search ignoring cell type information (**S5 Fig**). We also compared the  
152 distance for the cells between batch 1 and 2, and found that, compared to data before correction, both MNN and  
153 SMNN reduced the Euclidean distance between the two batches (**S6 Fig**). Moreover, SMNN further decreased the  
154 distance by up to 8.2% than MNN (2.8%, 4.3% and 8.2% for cells of type CMP, MEP and GMP, respectively).  
155 Regarding the overall variance in the two batches, our results show that, SMNN reduced the overall variance by up  
156 to 4.8% on top of MNN corrected results (**Fig 4d-f**). These results suggest improved batch effect correction by  
157 SMNN.

158

## 159 **SMNN identifies differentially expressed genes that are biologically relevant**

160 We then compared the DEGs among different cell types identified by SMNN and MNN. After correction, in the  
161 merged hematopoietic dataset, 1012 and 1145 up-regulated DEGs were identified in CMP cells by SMNN and  
162 MNN, respectively, when compared to GMP cells, while 926 and 1108 down-regulated DEGs were identified by the  
163 two methods, respectively (**Fig 5a and S7a Fig**). Of them, 736 up-regulated and 842 down-regulated DEGs were  
164 shared between SMNN and MNN corrected data. GO enrichment analysis showed that, the DEGs detected only by  
165 SMNN were overrepresented in GO terms related to blood coagulation and hemostasis, such as platelet activation  
166 and aggregation, hemostasis, coagulation and regulation of wound healing (**Fig 5b**). Similar DEG detection was  
167 carried out to detect genes differentially expressed between CMP and MEP cells. 181 SMNN-specific DEGs were  
168 identified out of the 594 up-regulated DEGs in CMP cells when compared to MEP cells (**Fig 5c**), and they were  
169 found to be enriched for GO terms involved in immune cell proliferation and differentiation, including regulation of  
170 leukocyte proliferation, differentiation and migration, myeloid cell differentiation and mononuclear cell proliferation  
171 (**Fig 5d**). Lastly, genes identified by SMNN to be up-regulated in GMP when compared to MEP cells, were found to  
172 be involved in immune processes; whereas up-regulated genes in MEP over GMP were enriched in blood

173 coagulation (**S7e-h Fig**). These cell-function-relevant SMNN-specific DEGs indicate SMNN can maintain some cell  
174 features that are missed by MNN after correction.

175 In addition, we considered two sets of “working truth”: first, DEGs identified in uncorrected batch 1;  
176 second DEGs identified in batch 2, and we compared SMNN and MNN results to both sets of working truth. The  
177 results showed that, in both comparisons (one comparison for each set of working truth), fewer DEGs were observed  
178 in SMNN-corrected batch 2, but higher TPR in each of the three cell types than those in MNN results. When  
179 compared to the uncorrected batch 1, 3.6% - 841% improvements were observed in SMNN results than MNN (**Fig 6**  
180 **and S8 and S9 Fig**). Similarly, SMNN increased the TPR by 6.2% - 54.0% on top of MNN when compared to  
181 uncorrected batch 2 (**S10-12 Fig**). Such an improvement in the accuracy of DEG identification indicates that higher  
182 amount of information regarding cell structure was retained after SMNN correction than MNN.

183 We also identified DEGs between T cells and B cells in the merged PBMC and T cell datasets after SMNN  
184 and MNN correction, respectively. Compared to B cells, 3213 and 4180 up-regulated DEGs were identified in T  
185 cells by SMNN and MNN, respectively, 2203 of which were shared between the two methods (**S13a Fig**). GO  
186 enrichment analysis showed that, the SMNN-specific DEGs were significantly enriched for GO terms relevant to the  
187 processes of immune signal recognition and T cell activation, such as T cell receptor signaling pathway, innate  
188 immune response–activating signal transduction, cytoplasmic pattern recognition receptor signaling pathway and  
189 regulation of autophagy (**S13b Fig**). In B cells, 5422 and 3462 were found to be up-regulated after SMNN and MNN  
190 correction, where 2765 were SMNN-specific (**S13c Fig**). These genes were overrepresented in GO terms involved in  
191 protein synthesis and transport, including translational elongation and termination, ER to Golgi vesicle–mediated  
192 transport, vesicle organization and Golgi vesicle budding (**S13d Fig**). These results again suggest that SMNN more  
193 accurately retains or rescues cell features after correction.

194

### 195 **SMNN more accurately identifies cell clusters**

196 Finally, we examined the ability to differentiate cell types after SMNN and MNN correction in three datasets (**S1**  
197 **Table**). In all three real datasets, ARI after SMNN correction showed 7.6 - 42.3% improvements over that of MNN  
198 (**Fig 7**), suggesting that SMNN correction more effectively recovers cell-type specific features.

199

## 200 **Discussion**



201 In this study, we present SMNN, a batch effect correction method for scRNA-seq data via supervised mutual nearest  
202 neighbor detection. Our work is built on the recently developed method MNN, which has showed advantages in  
203 batch effect correction than existing alternative methods. On top of MNN, our SMNN relaxes a strong assumption  
204 that underlies MNN: that the biological differentiations are orthogonal to batch effects [8]. When this fundamental  
205 assumption is violated, especially under the realistic scenario that the two batches are rather different, MNN tends to  
206 err when searching nearest neighbors for cells belonging to the same biological cell type across batches. Our SMNN,  
207 in contrast, explicitly considers cell type label information to perform supervised mutual nearest neighbor matching,  
208 thus empowered to extract only desired neighbors from the same cell type.

209 A notable feature of our SMNN is that it can detect and match the corresponding cell populations across  
210 batches with the help of feature markers provided by users. SMNN performs clustering within each batch before  
211 merging across batches, which can reveal basic data structure, i.e. cell composition and proportions of contributing  
212 cell types, without any adverse impact due to batch effects. Cells of each cluster are labeled by leveraging their  
213 average expression levels of certain marker(s), thus enabling us to limit the mutual nearest neighbor detection within  
214 a smaller search space (i.e., only among cells of the same or similar cell type or status). This supervised approach  
215 eliminates the correction biases incurred by pairs of cells wrongly matched across cell types. We benchmarked  
216 SMNN together with MNN on both simulated and three published scRNA datasets. Our results clearly show the  
217 advantages of SMNN in terms removing batch effects. For example, our results for the hematopoietic datasets show  
218 that SMNN better mixed cells of all the three cell types across the two batches (**Fig 4a-c**), and reduced the  
219 differentiation between the two batches by up to 8.2% on top of MNN corrected results (**Fig 4d-f and S6 Fig**),  
220 demonstrating that our SMNN method can more effectively mitigate batch effect.

221 More importantly, the wrongly matched cell pairs may wipe out the distinguishing features of cell types.  
222 This is mainly because, for a pair of cells from two different cell types, the true biological differentiations between  
223 them would be considered as technical biases and subsequently removed in the correction process. Compared to  
224 MNN, SMNN also appears to more accurately recover cell-type specific features: clustering accuracy using SMNN-  
225 corrected data increases substantially in all the three real datasets (by 7.6 to 42.3% when measured by ARI) (**Fig 7**).  
226 Furthermore, we observe power enhancement in detecting DEGs between different cell types in the data after  
227 SMNN correction than MNN (**Fig 5 and 6 and S7-12 Fig**). Specifically, the true positive rates of the DEGs  
228 identified by SMNN were improved by up to 841% and 54.0% than those by MNN when compared to the two set of

229 working truth, respectively (**Fig 6c and d and S8-12 Fig**). Moreover, GO term enrichment results show that, the up-  
230 regulated DEGs identified only in SMNN-corrected GMP and MEP cells were involved in immune process and  
231 blood coagulation, respectively (**S7f and h Fig**), which accurately reflect the major features of these two cell types  
232 [20]. Similarly, DEGs identified between T and B cells after SMNN correction are also biologically more relevant  
233 than those identified after MNN correction (**S13b and c Fig**). These results suggest that SMNN can eliminate the  
234 overcorrection between different cell types and thus maintains more biological features in corrected data than MNN.  
235 Efficient removal of batch effects at reduced cost of biological information loss, manifested by SMNN in our  
236 extensive simulated and real data evaluations, empowers valid and more powerful downstream analysis.

237 In summary, extensive simulation and real data benchmarking suggest that our SMNN can not only better  
238 rescue biological features and thereof provide improved cluster results, but also facilitate the identification of  
239 biologically relevant DEGs. Therefore, we anticipate that our SMNN is valuable for integrated analysis of multiple  
240 scRNA-seq datasets, accelerating genetic studies involving single-cell dynamics.

241

## 242 **Materials and methods**

### 243 **Simulation Framework**

244 We simulated two scenarios, orthogonal and non-orthogonal, to compare the performance of MNN and SMNN. The  
245 difference between the two scenarios lies in the directions of the true underlying batch effect vectors with respect to  
246 those of the biological effects.

247

### 248 **Baseline simulation**

249 Our baseline simulation framework, similar to that adopted in the MNN paper, contains two steps:

250 Firstly, different batches of data are independently generated from a Gaussian mixture model to represent a  
251 low dimensional biological space, with each component in the mixture corresponding to one cell type. Specifically,  
252 we consider two batches with gene expression matrix  $X_k$  and  $Y_l$ , each follows a three-component Gaussian mixture  
253 model in a three-dimensional space, representing the low (here three) dimensional biological space.

$$X_k \sim \sum_{i=1}^3 w_{1i} N(\mu_{1i}, I_3), \text{ with } \sum_{i=1}^3 w_{1i} = 1, \text{ and } w_{11}, w_{12}, w_{13} \geq 0, \text{ for } k = 1, 2, \dots, n_1 \quad (1)$$

$$Y_l \sim \sum_{j=1}^3 w_{2j} N(\mu_{2j}, I_3), \text{ with } \sum_{j=1}^3 w_{2j} = 1, \text{ and } w_{21}, w_{22}, w_{23} \geq 0, \text{ for } l = 1, 2, \dots, n_2 \quad (2)$$

254 where  $\mu_{1i}$  is the vector specifying cell-type specific means for the three cell types in the first batch, reflecting the  
 255 biological effect; similarly for  $\mu_{2j}$ ;  $n_1$  and  $n_2$  is the total number of cells in the first and second batch, respectively;  
 256  $w_{1i}$  and  $w_{2j}$  are the different mixing coefficients for the three cell types in the two batches; and  $I_3$  is the three  
 257 dimensional identity matrix with diagonal entries as ones and the rest entries as zeros. In our simulations, we set  
 258  $n_1 = 1000, n_2 = 1100$  and

$$(w_{11}, w_{12}, w_{13}) = (0.3, 0.5, 0.2) \quad (3)$$

$$(w_{21}, w_{22}, w_{23}) = (0.25, 0.5, 0.25) \quad (4)$$

259 Secondly, we project the low dimensional data with batch effect to the high dimensional gene expression  
 260 space. We map both datasets to  $G = 50$  dimensions by linear transformation using the same random Gaussian matrix  
 261  $\mathbf{P}$ , to simulate high-dimensional gene expression profiles.

$$\widetilde{X}_k = \mathbf{P}X_k, \text{ for } k = 1, 2, \dots, n_1 \quad (5)$$

$$\widetilde{Y}_l = \mathbf{P}Y_l, \text{ for } l = 1, 2, \dots, n_2 \quad (6)$$

262 Here  $\mathbf{P}$  is a  $G \times 3$  Gaussian random matrix with each entry simulated from the standard normal distribution.

263

## 264 Introduction of batch effects

265 In the MNN paper, batch effects are directly introduced in the high dimensional gene expression space. Specifically,  
 266 a Gaussian random vector  $b = (b_1, b_2, \dots, b_G)^T$  is simulated and added to the second dataset via the following:

$$X_{observed,k} = \widetilde{X}_k + \varepsilon_{1,k}, \text{ for } k = 1, 2, \dots, n_1 \quad (7)$$

$$Y_{observed,l} = \widetilde{Y}_l + b + \varepsilon_{2,l}, \text{ for } l = 1, 2, \dots, n_2 \quad (8)$$

267 where  $\widetilde{X}_k$  and  $\widetilde{Y}_l$  are projected high-dimensional gene expression profiles;  $\varepsilon_{1,k}$  and  $\varepsilon_{2,l}$  are independent random  
 268 noises added to the expression of each “gene” for each cell in the two batches.

269 In our simulations, we adopt a different approach: we introduce batch effects in the low dimensional  
 270 biological space. Specifically, we simulate a bias vector  $c = (c_1, c_2, c_3)^T$  in the biological space:

$$X_{observed,k} = \widetilde{X}_k + \varepsilon_{1,k} = \mathbf{P}X_k + \varepsilon_{1,k}, \text{ for } k = 1, 2, \dots, n_1 \quad (9)$$

$$Y_{Observed,l} = Y_{SMNN,l} + \varepsilon_{2,l} = \mathbf{P}(Y_l + c) + \varepsilon_{2,l} = \mathbf{P}Y_l + Pc + \varepsilon_{2,l}, \text{ for } l = 1, 2, \dots, n_2 \quad (10)$$

271 Comparing our simulation framework with that employed in the MNN paper, we would like to note the following:

- 272 1) For any vector  $c$ , there is a corresponding vector  $b = \mathbf{P}c$  given a fixed projection matrix  $\mathbf{P}$ . This implies that our  
 273 approach generates data that are special cases of those from MNN. In particular, since  $(b)_l = (\mathbf{P}c)_l =$   
 274  $\sum_{i=1}^G P_{li} c_i \sim N(0, \sum_{i=1}^G c_i^2)$ , if  $\sum_{i=1}^G c_i^2 = 1$ , our approach becomes equivalent to generating a standard Gaussian  
 275 random vector.
- 276 2) Our formulation allows flexible modeling of the biological effects and batch effects in the same low  
 277 dimensional biological space. Specifically,

$$\mu_{2i} = \mu_{1i} + c, \text{ for } i = 1, 2, 3. \quad (11)$$

278 Note that  $(\mu_{1j} - \mu_{1i})c = 0, \text{ for } i \neq j \in \{1, 2, 3\}$  in the orthogonal case and  $(\mu_{1j} - \mu_{1i})c \neq 0, \text{ for } i \neq j \in \{1, 2, 3\}$  in  
 279 the non-orthogonal case.

280 In summary, our simulation framework, allowing flexible manipulation of biological and batch effects in  
 281 the same low dimensional space, is effectively a special case of that adopted in the MNN paper.

## 282 **The two scenarios**

283 As aforementioned, we consider two scenarios, orthogonal case and non-orthogonal case. Orthogonality is defined  
 284 in the sense that biological differences (that is, mean difference between any two clusters/cell-types), are orthogonal  
 285 to those from batch effects.

286 Leveraging the simulation framework described before, we simulate two scenarios via the following:

- 287 1) In the orthogonal case, we set  $c = (0, 0, 2)^T$
- 288 a.  $\mu_{11} = (5, 0, 0)^T, \mu_{12} = (0, 0, 0)^T, \mu_{13} = (0, 5, 0)^T$
- 289 b.  $\mu_{21} = (5, 0, 2)^T, \mu_{22} = (0, 0, 2)^T, \mu_{23} = (0, 5, 2)^T$
- 290 2) In the non-orthogonal case, we set  $c = (0, 5, 2)^T$
- 291 a.  $\mu_{11} = (5, 0, 0)^T, \mu_{12} = (0, 0, 0)^T, \mu_{13} = (0, 5, 0)^T$
- 292 b.  $\mu_{21} = (5, 5, 2)^T, \mu_{22} = (0, 5, 2)^T, \mu_{23} = (0, 10, 2)^T$

293

## 294 **Performance evaluation**

295 MNN and SMNN share the goal to correct batch effects. Mathematically, using the notations introduced in baseline  
 296 simulation, the goal translates into de-biasing vector  $c$  (which would be effectively reduced to  $b$  in the orthogonal

297 case). Without loss of generality and following MNN, we treat the first batch as the reference and correct the second  
298 batch  $\{Y_{observed,l}: l = 1, \dots, n_2\}$  to the first batch  $\{X_{observed,k}: k = 1, \dots, n_1\}$ . Denote the corrected values from  
299 MNN and SMNN as  $\{\widehat{Y}_{MNN,l}: l = 1, \dots, n_2\}$  and  $\{\widehat{Y}_{SMNN,l}: l = 1, \dots, n_2\}$ , respectively.

300 To measure the performance of the two correction methods, we utilize the Frobenius norm [19] to define  
301 the loss function:

$$L(\widetilde{Y}, \widehat{Y}) = \|\widetilde{Y} - \widehat{Y}\|_F = \sqrt{\sum_{l=1}^{n_2} \|\widetilde{Y}_l - \widehat{Y}_l\|^2} = \sqrt{\sum_{l=1}^{n_2} \sum_{g=1}^G |\widetilde{Y}_{l,g} - \widehat{Y}_{l,g}|^2} \quad (12)$$

302 where  $\widetilde{Y} = [\widetilde{Y}_1, \dots, \widetilde{Y}_k, \dots, \widetilde{Y}_{n_2}]$ ,  $\widehat{Y} = [\widehat{Y}_1, \dots, \widehat{Y}_k, \dots, \widehat{Y}_{n_2}]$ . Note that  $\widetilde{Y}$  is the simulated true profiles introduced in  
303 equations (Error! Reference source not found.) and (Error! Reference source not found.Error! Reference  
304 source not found.Error! Reference source not found.Error! Reference source not found.) before batch effects,  
305 and noises are introduced in equations (Error! Reference source not found.) and (Error! Reference source not  
306 found.). Since MNN conducts cosine normalization to the input and the output, we use cosine-normalized  $\widetilde{Y}$  when  
307 calculating the above loss function.

308

### 309 Real data benchmarking

310 To assess the performance of SMNN in real data, we applied both SMNN and MNN to two hematopoietic scRNA-  
311 seq datasets, generated using different sequencing platforms, MARs-seq and SMART-seq2 (S1 Table) [10,21]. The  
312 first batch produced by MARs-seq consists of 1920 cells of six major cell types, and the second batch generated by  
313 SMART-seq2 contains 2730 of three cell types, where three cell types, common myeloid progenitor (CMP) cells,  
314 granulocyte-monocyte progenitors (GMP) cells and megakaryocyte-erythrocyte progenitor (MEP) cells, are shared  
315 between these two batches (here the two datasets). Batch effect correction was carried out using both MNN and  
316 SMNN, following the default instructions. Cell type labels were fed to SMNN directly according to the annotation  
317 from the original papers. To better compare the performance between MNN and SMNN, only the three cell types  
318 shared between the two batches were extracted for our downstream analyses. The corrected results of all the three  
319 cell types together, as well as for each of them separately, were visualized by t-SNE using *Rtsne* function from *Rtsne*  
320 package [9,22]. In order to qualify the mixture of single cells using both batch correction methods, we calculated: 1)

321 the distance for the cells within each cell type in batch 2 to the centroid of the corresponding cell group in batch 1;  
322 and 2) the overall variance in the two batches.

323 To measure the separation of cell types after correction, we additionally attempted to detect differentially  
324 expressed genes (DEGs) between different cell types in both SMNN and MNN corrected datasets. The corrected  
325 expression matrices of the two batches were merged and DEGs were detected by Seurat using Wilcoxon rank sum  
326 test. Genes with an adjusted p-value < 0.01 were considered as differentially expressed. Gene ontology (GO)  
327 enrichment analysis was performed for the DEGs exclusively identified by SMNN using *clusterProfiler* [23].  
328 Because there is no ground truth for DEGs, we further identified DEGs between different cell types within corrected  
329 batch 2 and then compared them to those identified in uncorrected batch 1 and uncorrected batch 2, which  
330 supposedly are not affected by the choice of batch effect correction method. True positive rate (TPR) was computed  
331 for each comparison.

332 Additionally, we also performed batch effect correction on another two tissues/cell lines, pancreas [24,25]  
333 human peripheral blood mononuclear cells (PBMCs) [26], again using both SMNN and MNN. Single cell clustering  
334 was applied to batch-effects corrected gene expression matrices following the pipeline described in MNN paper.  
335 Cell type labels before correction were considered as ground truth and Adjusted Rand Index (ARI) [27] was  
336 employed to measure the clustering similarity before and after correction:

$$ARI(L_q, L_s) = \frac{\sum_{q,s} \binom{n_{qs}}{2} - [\sum_q \binom{n_q}{2} \sum_s \binom{n_s}{2}]}{\frac{1}{2} [\sum_q \binom{n_q}{2} + \sum_s \binom{n_s}{2}] - [\sum_q \binom{n_q}{2} \sum_s \binom{n_s}{2}]} / \binom{n}{2} \quad (13)$$

337 where  $n_e$  and  $n_t$  are the single cell numbers in cluster  $q$  and  $s$ , respectively;  $n_{qs}$  is the number of single cells shared  
338 between clusters  $q$  and  $s$ ; and  $n$  is the total number of single cells. ARI ranges from 0 to 1, where a higher value  
339 represents a higher level of similarity between the query and subject clusters.

340

## 341 Acknowledgements

342 This research was supported by the National Institute of Health grant R01 HL129132 (awarded to YL).

343

## 344 Author Contributions

345 Conceptualization: Yuchen Yang, Gang Li, Yun Li.

346 Formal analysis: Yuchen Yang, Gang Li, Huijun Qian.  
347 Funding acquisition: Yun Li.  
348 Methodology: Yuchen Yang, Gang Li.  
349 Software: Yuchen Yang, Gang Li.  
350 Supervision: Yun Li.  
351 Visualization: Yuchen Yang, Gang Li, Huijun Qian.  
352 Writing – original draft: Yuchen Yang, Gang Li.  
353 Writing – review & editing: Yun Li, Kirk C. Wilhelmsen, Yin Shen.

354

## 355 REFERENCES

- 356 1. Rozenblatt-Rosen O, Stubbington MJ, Regev A, Teichmann SA (2017) The human cell atlas: from vision to  
357 reality. *Nature News* 550: 451.
- 358 2. Stegle O, Teichmann SA, Marioni JC (2015) Computational and analytical challenges in single-cell  
359 transcriptomics. *Nature Reviews Genetics* 16: 133.
- 360 3. Chen M, Zhou X (2017) Controlling for confounding effects in single cell RNA sequencing studies using both  
361 control and target genes. *Scientific reports* 7: 13587.
- 362 4. Stuart T, Satija R (2019) Integrative single-cell analysis. *Nature Reviews Genetics*: 1.
- 363 5. Smyth GK (2005) Limma: linear models for microarray data. *Bioinformatics and computational biology solutions*  
364 using R and Bioconductor: Springer. pp. 397-420.
- 365 6. Johnson WE, Li C, Rabinovic A (2007) Adjusting batch effects in microarray expression data using empirical  
366 Bayes methods. *Biostatistics* 8: 118-127.
- 367 7. Leek JT (2014) Svaseq: removing batch effects and other unwanted noise from sequencing data. *Nucleic acids*  
368 *research* 42: e161-e161.
- 369 8. Haghverdi L, Lun AT, Morgan MD, Marioni JC (2018) Batch effects in single-cell RNA-sequencing data are  
370 corrected by matching mutual nearest neighbors. *Nature biotechnology* 36: 421.
- 371 9. Van Der Maaten L (2014) Accelerating t-SNE using tree-based algorithms. *The Journal of Machine Learning*  
372 *Research* 15: 3221-3245.
- 373 10. Nestorowa S, Hamey FK, Sala BP, Diamanti E, Shepherd M, et al. (2016) A single-cell resolution map of mouse  
374 hematopoietic stem and progenitor cell differentiation. *Blood* 128: e20-e31.
- 375 11. Spitzer MH, Gherardini PF, Fragiadakis GK, Bhattacharya N, Yuan RT, et al. (2015) An interactive reference  
376 framework for modeling a dynamic immune system. *Science* 349: 1259425.
- 377 12. Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, et al. (2019) Comprehensive Integration of Single-  
378 Cell Data. *Cell* 177: 1888-1902 e1821.
- 379 13. Duò A, Robinson MD, Sonesson C (2018) A systematic performance evaluation of clustering methods for single-  
380 cell RNA-seq data. *F1000Research* 7.
- 381 14. Kiselev VY, Andrews TS, Hemberg M (2019) Challenges in unsupervised clustering of single-cell RNA-seq  
382 data. *Nature Reviews Genetics*: 1.
- 383 15. Zhu L, Lei J, Klei L, Devlin B, Roeder K (2019) Semisoft clustering of single-cell data. *Proceedings of the*  
384 *National Academy of Sciences* 116: 466-471.
- 385 16. Sun Z, Chen L, Xin H, Jiang Y, Huang Q, et al. (2019) A Bayesian mixture model for clustering droplet-based  
386 single-cell transcriptomic data from population studies. *Nature communications* 10: 1649.
- 387 17. Butler A, Hoffman P, Smibert P, Papalexi E, Satija R (2018) Integrating single-cell transcriptomic data across  
388 different conditions, technologies, and species. *Nature biotechnology* 36: 411.

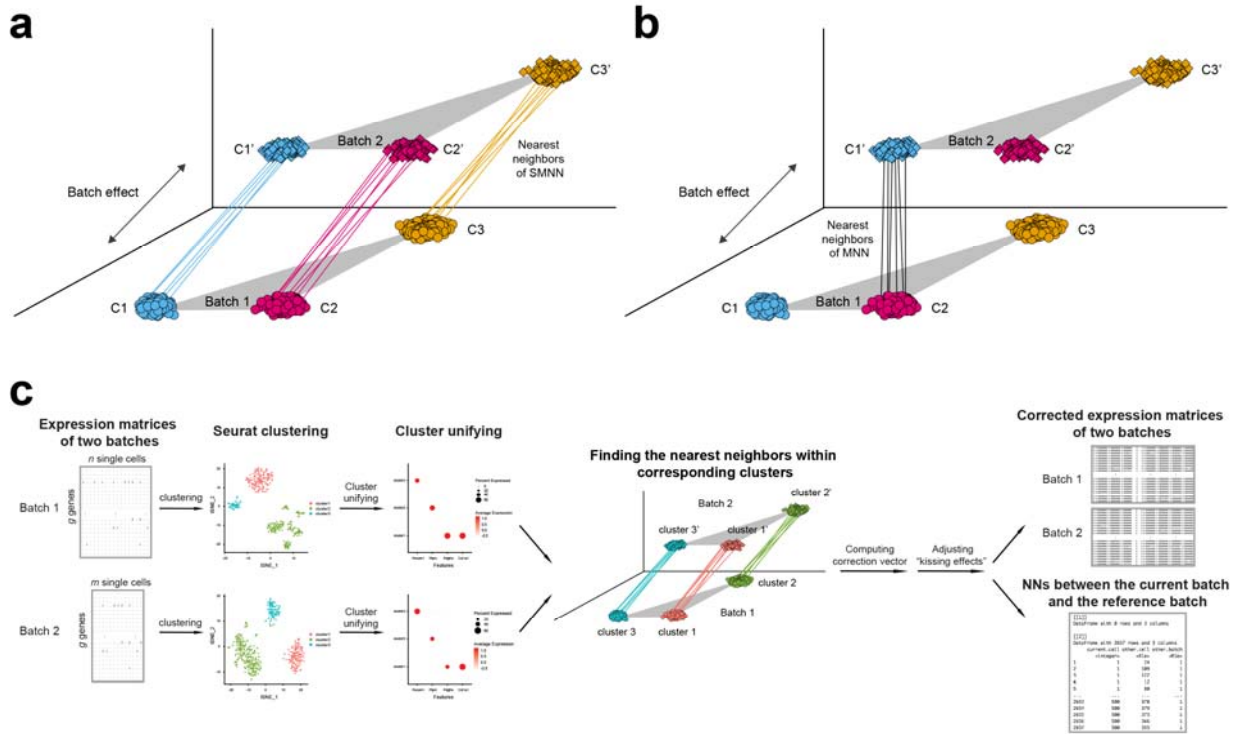
- 389 18. Yang Y, Huh R, Culpepper HW, Lin Y, Love MI, et al. (2019) SAFE-clustering: Single-cell aggregated (from  
390 ensemble) clustering for single-cell RNA-seq data. *Bioinformatics* 35: 1269-1277.
- 391 19. Van Loan CF, Golub GH (1983) *Matrix computations*: Johns Hopkins University Press.
- 392 20. Lieu YK, Reddy EP (2012) Impaired adult myeloid progenitor CMP and GMP cell function in conditional c-  
393 myb-knockout mice. *Cell Cycle* 11: 3504-3512.
- 394 21. Paul F, Arkin Ya, Giladi A, Jaitin DA, Kenigsberg E, et al. (2015) Transcriptional heterogeneity and lineage  
395 commitment in myeloid progenitors. *Cell* 163: 1663-1677.
- 396 22. Maaten Lvd, Hinton G (2008) Visualizing data using t-SNE. *Journal of machine learning research* 9: 2579-2605.
- 397 23. Yu G, Wang L-G, Han Y, He Q-Y (2012) clusterProfiler: an R package for comparing biological themes among  
398 gene clusters. *OmicS: a journal of integrative biology* 16: 284-287.
- 399 24. Grün D, Muraro MJ, Boisset J-C, Wiebrands K, Lyubimova A, et al. (2016) De novo prediction of stem cell  
400 identity using single-cell transcriptome data. *Cell stem cell* 19: 266-277.
- 401 25. Muraro MJ, Dharmadhikari G, Grün D, Groen N, Dielen T, et al. (2016) A single-cell transcriptome atlas of the  
402 human pancreas. *Cell systems* 3: 385-394. e383.
- 403 26. Zheng GX, Terry JM, Belgrader P, Ryvkin P, Bent ZW, et al. (2017) Massively parallel digital transcriptional  
404 profiling of single cells. *Nature communications* 8: 14049.
- 405 27. Hubert L, Arabie P (1985) Comparing partitions. *Journal of classification* 2: 193-218.
- 406

407



408 **Figure Legend**

409



410

411 **Fig 1. Overview of SMNN.** Schematics for detecting mutual nearest neighbors between two batches under a non-

412 orthogonal scenario **(a)** in SMNN; and **(b)** in MNN. **(c)** Workflow of SMNN. Single cell clustering is first

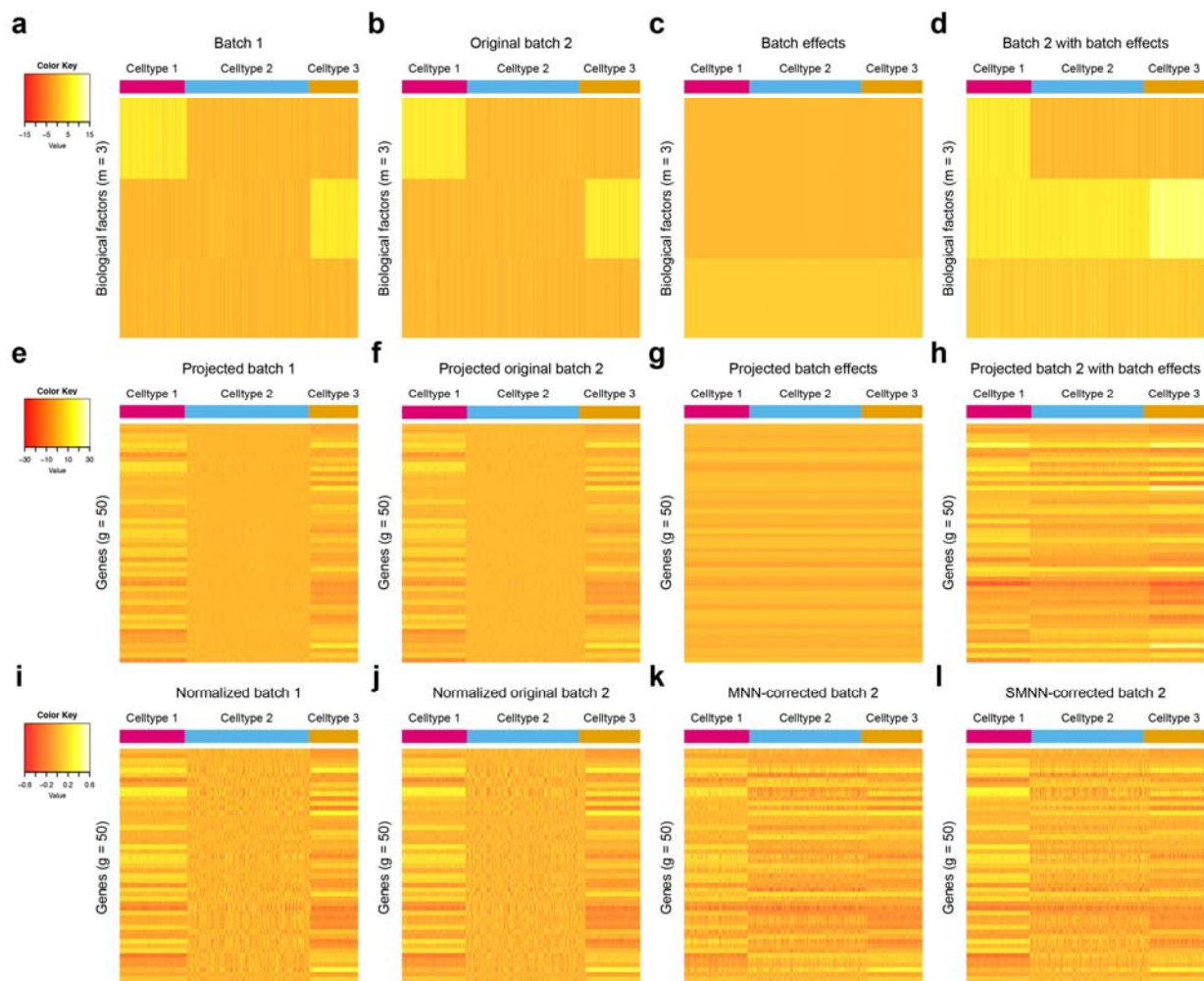
413 performed within each batch using Seurat; and then SMNN takes user-specified marker gene information for each

414 cell type to match clusters/cell types across batches. With the clustering and cluster-specific marker gene

415 information, SMNN searches mutual nearest neighbors within each cell type and performs batch effect correction

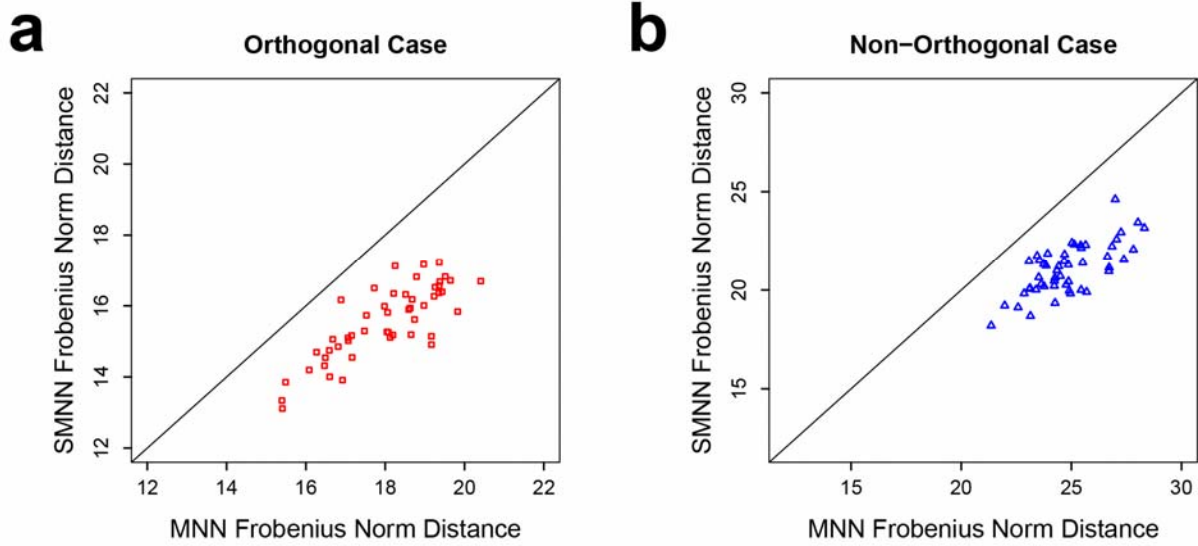
416 accordingly.

417



418  
 419 **Fig 2. Heatmap of gene expression matrices for simulated data under non-orthogonal scenario.** (a), (b), (c) and  
 420 (d) show the 3-dimensional biological space with rows of each heatmap representing biological factors and columns  
 421 corresponding to single cells. (e), (f), (g) and (h) show the high dimensional gene expression profiles with rows  
 422 corresponding to genes and columns again representing single cells. (a), (e) and (i) correspond to the batch 1, and  
 423 (b), (f) and (j) correspond to batch 2. (c) and (g) provide a visualization for the direction of batch effects in low-  
 424 dimensional biological space and high-dimension gene expression spaces, respectively. (d) and (h), sum of (b) and (c)  
 425 and sum of (f) and (g) respectively, are “observed” data for cells in batch 2 in low and high dimensional space  
 426 respectively. (i) and (j) are the cosine-normalized data for batch 1 and original batch 2. Note “original” is in the  
 427 sense that no batch effects have been introduced to the data yet. (k) and (l) are the MNN and SMNN corrected  
 428 results, respectively.

429

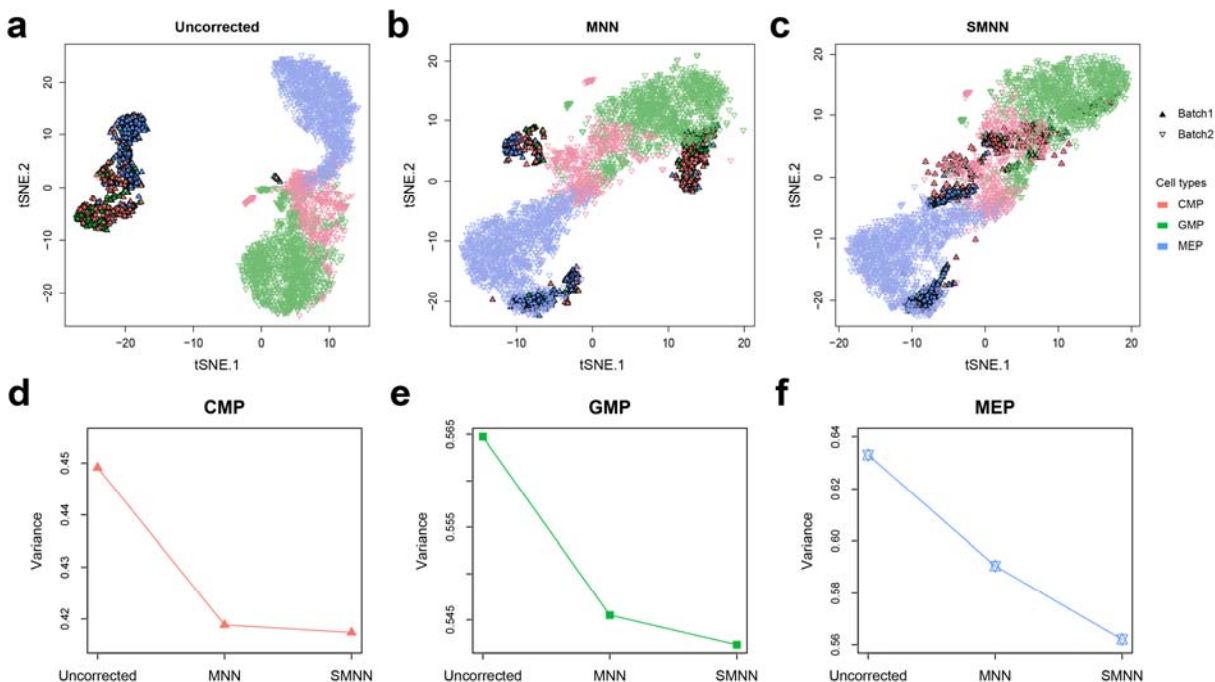


430

431 **Fig 3. Frobenius norm distance between two batches after SMNN and MNN correction in simulation data**

432 **under orthogonal (left) and non-orthogonal scenarios (right).**

433



434

435 **Fig 4. Performance comparison between SMNN and MNN in two hematopoietic datasets.** (a-c) t-SNE plots for

436 two hematopoietic datasets before and after correction with SMNN and MNN. Solid and inverted triangle represent

437 the first and second batch, respectively; and different cell types are shown in different colors. (d-f) Variance

438 comparisons for the three different cell types: CMP (d), GMP (e) and MEP (f), in merged data by pooling batch 1

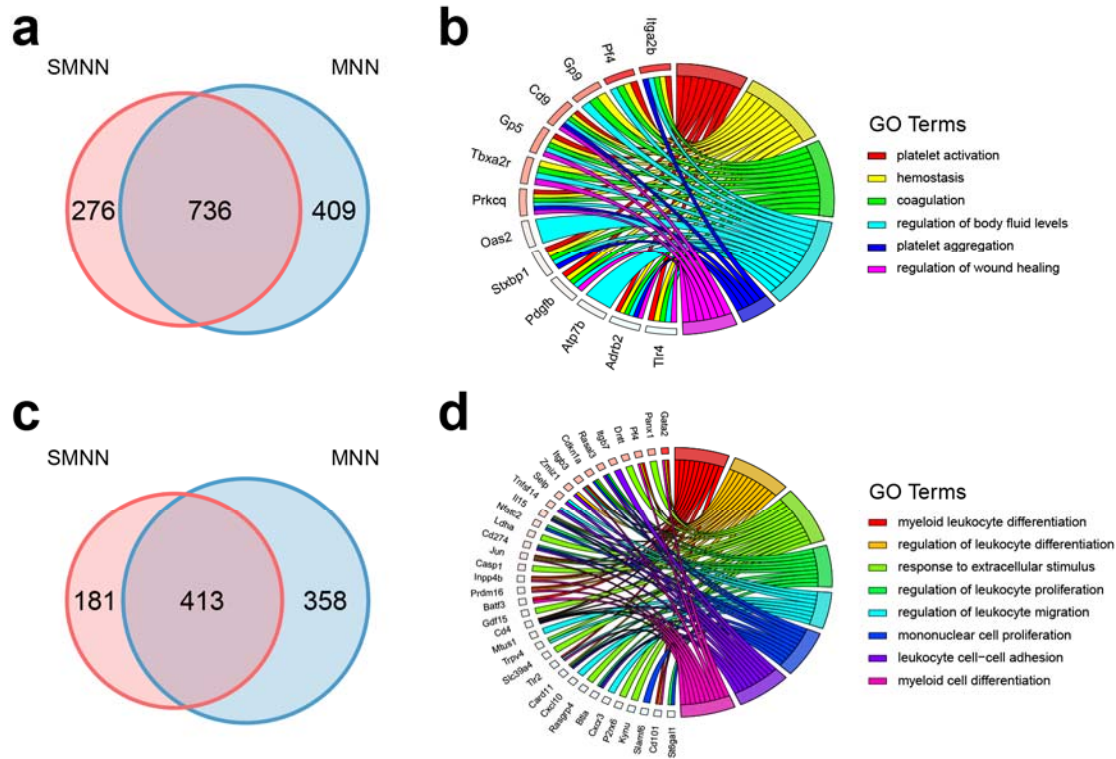
439 with different versions of batch 2. Specifically, we show the following three versions of batch 2 data: original

440 observed (uncorrected), MNN-corrected (MNN) and SMNN corrected (SMNN). The SMNN corrected version

441 resulted in variances slightly (for CMP and GMP cells) or substantially (for MEP cells) smaller than those from the

442 MNN corrected version, suggesting improved mixing of cells across batches.

443



444

445 **Fig 5. Comparison of differentially expressed genes (DEGs), identified in the merged dataset by pooling batch**

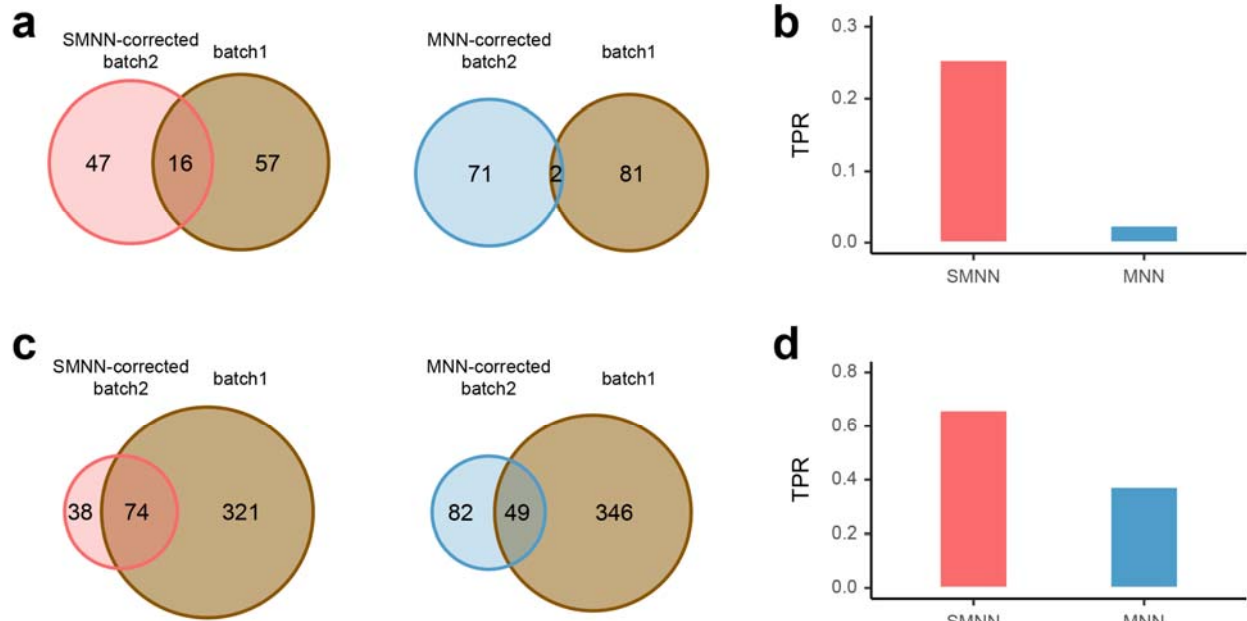
446 **1 data with batch 2 data after SMNN and MNN correction. (a) Overlap of DEGs up-regulated in CMP over**

447 **GMP after SMNN and MNN correction. (b) Feature enriched GO terms and the corresponding DEGs up-regulated**

448 **in CMP over GMP. (c) Overlap of DEGs up-regulated in CMP over MEP after SMNN and MNN correction. (d)**

449 **Feature enriched GO terms and the corresponding DEGs up-regulated in CMP over MEP.**

450



451

452 **Fig 6. Reproducibility of DEGs (between CMP and GMP), identified in uncorrected batch 1 and in SMNN or**

453 **MNN-corrected batch 2. (a)** Reproducibility of DEGs up-regulated in CMP over GMP, detected in batch 1, versus

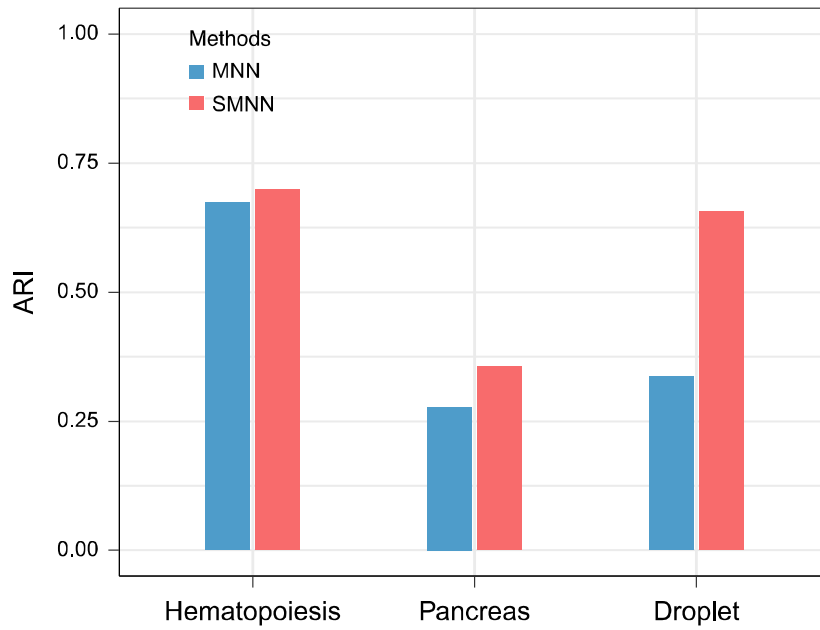
454 SMNN (left) or MNN-corrected (right) batch 2. **(b)** True positive rate (TPR) of the DEGs (between CMP and GMP)

455 identified in batch 2 after SMNN and MNN correction. **(c)** Reproducibility of DEGs up-regulated in GMP over

456 CMP, identified in the uncorrected batch 1, and in SMNN (left) or MNN-corrected (right) batch 2. **(d)** TPR of the

457 DEGs up-regulated in GMP over CMP identified in batch 2 after SMNN and MNN correction.

458



459

460 **Fig 7. Clustering accuracy in three datasets after batch effect correction.** Adjusted Rand Index (ARI) is

461 employed to measure the similarity between clustering results before and after batch effect correction.