

1 **Data Matrix Normalization and Merging Strategies Minimize**
2 **Batch-specific Systemic Variation in scRNA-Seq Data**

3
4 Benjamin R. Babcock¹, Astrid Kusters¹, Junkai Yang¹, Mackenzie L. White¹, Eliver E. B. Ghosn^{1,2,*}

5
6 ¹Department of Medicine, Division of Immunology, Lowance Center for Human Immunology, Emory
7 University School of Medicine, Atlanta, GA 30322, USA

8 ²Emory Vaccine Center, Yerkes National Primate Research Center, Emory University School of Medicine,
9 Atlanta, GA 30322, USA

10
11 *Corresponding author:

12 Eliver E.B. Ghosn

13 E-mail: eliver.ghosn@emory.edu

16 **Abstract**

17 Single-cell RNA sequencing (scRNA-seq) can reveal accurate and sensitive RNA abundance in a single
18 sample, but robust integration of multiple samples remains challenging. Large-scale scRNA-seq data
19 generated by different workflows or laboratories can contain batch-specific systemic variation. Such
20 variation challenges data integration by confounding sample-specific biology with undesirable batch-specific
21 systemic effects. Therefore, there is a need for guidance in selecting computational and experimental
22 approaches to minimize batch-specific impacts on data interpretation and a need to empirically evaluate the
23 sources of systemic variation in a given dataset. To uncover the contributions of experimental variables to
24 systemic variation, we intentionally perturb four potential sources of batch-effect in five human peripheral
25 blood samples. We investigate sequencing replicate, sequencing depth, sample replicate, and the effects of
26 pooling libraries for concurrent sequencing. To quantify the downstream effects of these variables on data
27 interpretation, we introduced a new scoring metric, the Cell Misclassification Statistic (CMS), which
28 identifies losses to cell type fidelity that occur when merging datasets of different batches. CMS reveals an
29 undesirable overcorrection by popular batch-effect correction and data integration methods. We show that
30 optimizing gene expression matrix normalization and merging can reduce the need for batch-effect
31 correction and minimize the risk of overcorrecting true biological differences between samples.

32

33 Introduction

34 Recent advances in throughput and commercial availability of single-cell RNA-sequencing (scRNA-
35 seq) technology have increased accessibility and led to widespread adoption of this technology by the
36 scientific community (Supplementary Fig. 1). A single study may encompass thousands of cells from
37 multiple samples, often spanning time points and conditions, resulting in large and heterogenous scRNA-
38 seq datasets^{1,2}. This dramatic shift is empowering scientists to massively profile multiple samples in parallel
39 at extremely high resolution. As experiments grow more complicated, the need arises to align and co-
40 analyze ever larger and more diverse outputs of single-cell workflows.

41 Minute and often uncontrollable technical variations in sample collection and data processing can
42 manifest as noticeable effects which confound the interpretation of data³. This systemic variation, referred to
43 commonly as “batch-effect,” can pose an obstacle to data interpretation by confounding biologically-derived
44 variation (desirable) with technically-derived variation (undesirable). An inability to discern the source of a
45 particular signal can even lead to over-interpretation of data, in that systemic variation arising from technical
46 differences may be interpreted as a biologically driven phenotypic difference. A typical case in which batch-
47 effect confounds data interpretation will present as an over-merging or under-merging of cell types.
48 Uncorrected batch-effect can cause similar cell populations between samples to appear divergent. In the
49 inverse case, batch-effect can cause two biologically distinct populations to appear as one due to a shared
50 technical signal. Both the prevalence and persistence of batch-specific signals have been highlighted by
51 prior work, as well as the spectrum of methods existing to correct and remove them^{4,5}. However, there
52 remains an insufficient understanding of how experimental design and data analysis approaches play a role
53 in producing batch-effect or identifying batch-effect when present in a sample. An understanding of the
54 source of batch-effect and the informed selection of tools to identify batch-effect has the potential to alter
55 the outcomes and conclusions of scRNA-seq studies.

56 Technical sources of variation most apparently manifest themselves in the Principal Component
57 Analysis (PCA) matrix, representing a shared low-dimensional space. Therefore, isolation and removal of
58 these effects in PCA dimensions are critical, as PCAs are the foundation used to produce cell cluster

59 assignments and UMAP visualizations. Towards this aim, many methods have been applied to isolate and
60 remove batch-effect from scRNA-seq data⁶⁻¹¹. These methods try to merge biologically similar populations
61 into a shared low-dimensional space while disregarding the influence of undesirable signals. A commonly
62 shared assumption of current methods is that batch-specific, technically-derived signals are contained
63 *within* the sample, while true biologically-derived signals are shared *between* samples. However, current
64 methods are mostly agnostic to the fundamental sources of systemic variation and the underlying biological
65 heterogeneity contained within each sample.

66 Here, we present a novel approach to validating batch-correction methods by demonstrating
67 experimental variables which contribute most to systemic variation. To assess the degree of batch-effect,
68 we introduce a biologically-grounded metric, the Cell Misclassification Statistic (CMS). While most current
69 scoring systems are agnostic to cell identity, the CMS directly grounds itself in the cell-type classification of
70 every single cell and is therefore uniquely able to quantify the loss of biological information during sample
71 merging. Using CMS to quantify batch-associated systemic variation, we show that sequencing replicates
72 and sequencing depth contribute only minimally to batch-effect and that pooling samples together for
73 sequencing does not meaningfully improve the measured or observed batch-effect. Instead, we find that
74 sample donor, along with the microfluidic encapsulation and library preparation steps, represent the main
75 source of batch-associated variation. We test three popular batch-correction algorithms (Harmony⁸,
76 LIGER¹¹, and Seurat V3¹⁰), which have been previously scored as the best performing⁵. Our CMS scoring,
77 which accounts for cell identity as a biological feature, revealed misclassifications of major cell lineages in
78 all three commonly used batch-correction methods. We further applied CMS in a supervised approach to
79 reveal that selecting a proper dataset normalization and merging strategy can perform comparably to
80 popular batch-correction algorithms. Furthermore, we revealed that much of the batch-effect present is due
81 to low expression levels of broadly expressed genes, which can be minimized by selecting a proper
82 normalization and dataset merging strategy. Our unique and biologically-grounded approach allows for
83 effective data integration in a carefully supervised workflow without the need for corrective algorithms.

84 **Results**

85 **A Cell Misclassification Statistic quantifies data integration and batch-correction fidelity**

86 We generated seven single-cell RNA-seq datasets from five individual donors of human peripheral
87 blood mononuclear cells (PBMCs) (Fig. 1a). To quantify the contributions of known variables to batch effect
88 and evaluate the performance of data integration and batch correction methods, we developed a cell
89 misclassification statistic (CMS) metric (Fig. 1b). CMS is grounded in the premise that if different systems
90 identify a single cell as different biological types, both cannot be correct. To calculate a CMS, we first gather
91 the cell-type identities obtained by classifying cells from only a single sample, meaning that all cells were
92 processed under identical conditions and are influenced identically by potential technical sources of
93 variance, if at all. Then, to calculate a CMS score, we compare the cell-type classifications of each sample
94 individually to cell-type classifications after multiple sample merging. We measure the fraction of cells that
95 have changed classification and generate a statistic, such that a CMS score of 0 means that no cells
96 changed classification/cell-type identity after merging datasets, while a CMS of 0.2 means a
97 misclassification of 20% of cells. A higher CMS score will result when cell barcodes change cell-type identity
98 after sample merging and indicates over-correction of the sample. By relying on invariable biological
99 principles (a single, non-doublet cell barcode must hold only one cell classification), we directly interpret
100 CMS scores as a measure of biological signal loss during data integration.

101 To produce the initial cell type classifications required for CMS, we analyzed seven datasets
102 generated from five PBMC samples. Each dataset contained two matrices holding independent modalities
103 of data: gene expression (GEX) and antibody-derived tag (ADT) cell-surface protein expression. First, we
104 confirmed the cell types present in each single sample by removing an aliquot of cells and performing flow
105 cytometry staining and analysis by a gold-standard “gating strategy” approach (representative data shown
106 in Supplementary Fig. 2). In parallel, we processed the GEX sequencing data to generate clusters
107 according to the Seurat V3 workflow, following default settings¹⁰. We next assigned each cluster to a major
108 lineage using both ADT and GEX markers based on a gating strategy similar to flow cytometry, as we have
109 previously reported¹² (representative data shown in Supplementary Fig. 3). With cell-type classifications

110 independently established for each sample, we can assess the impacts of batch-effect on cell-type
111 classification in merged datasets.

112 CMS measures consistency in cell-type classification, which is only one of two metrics we used to
113 validate data integration. Another important aspect of data integration and batch correction methods is
114 generating UMAP embeddings for data visualization. To assess the performance of UMAP embeddings for
115 visualization of integrated datasets, we employed a modified Local Inverse Simpson's Index (LISI)^{8,13}, which
116 we used specifically to measure the final UMAP integration, or integration LISI (iLISI) (Fig. 1b). In our
117 system, an iLISI score of 1 represents a UMAP completely segregated by sample ID, while an iLISI of 0
118 represents a perfectly integrated UMAP. Therefore, the best-integrated data will have both CMS and iLISI
119 scores approaching 0, while the most segregated data will return both CMS and iLISI scores approaching 1
120 (Fig. 1b). As described below, we applied the CMS and iLISI methods to compare the fidelity of the various
121 data-integration/batch-correction methods we evaluated in this study.

122

123 **Batch-associated systemic variation is observed when sample replicates and sequencing depth are** 124 **perturbed, but not sequencing replicate**

125 To examine the contributions of experimental variables to batch-associated data effects, we first
126 describe four unique steps of a typical scRNA-seq workflow in which sample-specific variation may be
127 introduced:

- 128 1. Sample donor
- 129 2. Sample and library preparation
- 130 3. Library sequencing
- 131 4. Data analysis

132 The sample donor (i) represents the baseline variation intrinsic to each human subject who donated blood
133 for this study. Sample preparation (ii) includes the processes of extracting cells from donor tissue and
134 preparing a single-cell suspension (the sample) (Fig. 1a). Together, sample and library preparation (ii)
135 encompasses all steps spanning microfluidic encapsulation of donor cells through to the generation of a

136 barcoded cDNA library (Fig. 1a). (Note: As sample and library preparation will largely be tissue and
137 platform-specific and will not typically be varied for a single-donor aliquot of cells, we consider them here as
138 processes intrinsically linked to sample donor and not as independent sources of batch-effect). Library
139 sequencing (iii) consists of converting the cDNA molecules (i.e., barcoded libraries) into aligned reads, and
140 finally, an expression matrix of gene counts per cell (Fig. 1a). Data analysis (iv) constitutes the last unit of
141 the workflow. During data analysis, we make interpretations on the expression matrix by assigning cell type
142 IDs to cell barcodes and clusters, visualizing cells by UMAP embeddings, and detecting differential gene
143 expression (DGE) (Fig. 1a). To evaluate the contributions of each of these steps (i-iv) to batch-effect, we
144 knowingly introduced differences in key variables to sets of human PBMC samples and evaluated the
145 downstream batch-effect observed. Finally, we directly evaluated whether pooling libraries for concurrent
146 sequencing could reduce sample-associated batch-effect compared to the same libraries sequenced
147 independently.

148 Sample donor: We assessed the systemic batch variation in identically-prepared PBMC samples from three
149 healthy adult donors (samples 1-3, Supplementary Table 1). We processed samples simultaneously for
150 PBMC isolation, single-cell encapsulation, and barcoded cDNA library generation. Following library
151 preparation, we pooled samples for simultaneous sequencing. Depth of sequence (reads per cell) was not
152 significantly different between samples (Wilcoxon rank-sum test, $p = 0.3487, 0.2445, 0.8471$ for samples 1,
153 2, 3, respectively), and we proceeded with 16,946 cell barcodes after concatenating the GEX matrices.
154 Next, we assigned each cluster to a cell type using both GEX and ADT data (Supplementary Fig. 3). After
155 sample merging, the three PBMC samples were mostly segregated by donor when visualized on a UMAP
156 (Fig. 2a). The failure to effectively integrate the three samples, as readily apparent in the UMAP embedding,
157 was confirmed by our modified iLISI scoring (iLISI = 0.861, Fig. 2a). We then applied CMS scoring and
158 revealed approximately ten percent cell-type misclassification after sample merging (CMS = 0.099, Figs. 2a
159 and Supplementary Fig. 4). Hence, we demonstrate that sample-specific variations cause a loss of
160 biological signal in the integrated data analysis and contribute to undesirable batch-associated data effects.

161 Library Sequencing Replication: To demonstrate the contributions of library sequencing to systemic
162 variation, we aliquoted a single library in two and performed duplicate sequencing under similar conditions
163 (samples 4-A and 4-B) (Fig. 2b). We sequenced each aliquot in a separate flow cell, on different days, with
164 similar target sequence depth (~40,000 reads/cell) (Fig. 2b). After merging, the dataset contained
165 approximately 13,000 cells, representing a shared set of approximately 6,500 barcodes, duplicated. We
166 found that 6,392, or greater than 99 percent, of cell barcodes, were shared between both sequencing
167 replicates, indicating minimal new cells were identified by additional sequencing. As above, we followed a
168 default Seurat workflow to generate UMAP visualizations and Louvain clusters (Fig. 2c) and used the GEX
169 and ADT data to classify clusters by cell type. We constructed a table of the proportion of each sample that
170 contributed to each cell type and cluster (Supplementary Table 2). As both libraries are derived from two
171 aliquots of the same cDNA pool, we would expect the proportions of cell types to be equal by replicate,
172 which they were (Supplementary Table 2; chi-square test, $p = 1$). Scoring the replicates by iLISI confirmed a
173 homogenous distribution of replicates in the UMAP (iLISI = 0.030) while CMS revealed only 3.2% of cells
174 changed classification (CMS = 0.032, Supplementary Fig. 4), which represents fewer cells than the
175 expected cell-cell doublet rate for this technology. Hence, we find that variation introduced by library
176 sequencing was a non-significant contributor to batch-effect by all measures.

177 Library Sequencing Depth: Next, we quantified the effect of sequencing depth by intentionally altering the
178 number of reads between replicates. For this comparison, we constructed two identical aliquots of a PBMC
179 library containing approximately 4,000 cells, sequenced to different depths (samples 5-A and 5-B) (Fig. 2b).
180 Sample 5-B contained three-fold more total reads than sample 5-A (170M reads for 5-A vs. 560M reads for
181 5-B). Notably, we sequenced samples 5-A and 5-B on the same Illumina flow cell as samples 4-A and 4-B,
182 respectively, from the experiment above (Fig. 2b). Having already established that library sequencing
183 replication (i.e., sample 4) produced minimal batch-effect (Fig. 2c), we interpret that any differences
184 observed between samples 5-A and 5-B can be attributed to sequence depth. We recovered similar
185 numbers of cells per replicate, resulting in averages of 45,160 reads/cell in sample 5-A and 138,274

186 reads/cell in 5-B, with 94% overlap in cell barcode sequences (3,758 of 4,007). Of the 249 non-shared cell
187 barcodes, all were present only in 5-B, the sample with greater sequencing depth (i.e., higher reads/cell). As
188 before, we analyzed the gene expression matrices following the Seurat workflow under default settings and
189 assigned cell types using a GEX and ADT-based gating approach (Fig. 2d). Surprisingly, in this comparison
190 cell-type compositions were biased by sequencing replicate (Supplementary Table 2; chi-square, $p <$
191 2.2×10^{-16}) while iLISI and CMS scores remained favorable (iLISI = 0.028, CMS = 0.033). One notable cell
192 type, erythrocytes, stood out as batch-biased: 0.3% of total cells in sample 5-A were classified as
193 erythrocytes, compared to 5% in 5-B. Indeed, removing the erythrocyte-classified population resulted in a
194 homogenous composition of cell types by sample (Supplementary Table 2; chi-square, $p = 1$). Removing
195 erythrocytes also improved the CMS score from 0.033 to 0.017, still well below the expected cell-cell
196 doublet rate.

197 We hypothesized that erythrocyte-classified barcodes appeared only in the greater read depth
198 sample because they contain fewer unique molecular identifiers (UMIs), requiring more reads to capture
199 their relatively rare transcripts. We confirmed that erythrocytes did contain fewer unique mRNA molecules
200 (Supplementary Fig. 5, Wilcoxon rank-sum test, $p < 2.2 \times 10^{-16}$). In addition, we observed that hematopoietic
201 stem and progenitor cells (HSPCs) co-clustered with erythrocytes, which may explain the few high-UMI dots
202 present in Supplementary Fig. 5. We also confirmed that sample 5-B (the higher read-depth sample)
203 contained more UMIs per cell than sample 5-A (Supplementary Fig. 5, Wilcoxon rank-sum test, $p =$
204 1.482×10^{-10}), and that sample 5-B also contained more unique genes per cell (Supplementary Fig. 5,
205 Wilcoxon rank-sum test, $p < 2.2 \times 10^{-16}$). As a standard part of the 10X Genomics Cell Ranger workflow (i.e.,
206 the sequencing read alignment and demultiplexing steps), cell barcodes with low-UMI counts are excluded
207 as potential empty droplets containing ambient mRNA/noise. To confirm that the erythrocyte cell barcodes
208 unique to sample 5-B were also present in the aliquot sequenced for 5-A but instead had been artificially
209 excluded as low-UMI barcodes, we investigated the unfiltered sequencing matrices of sample 5-A. We
210 discovered 100% of cell barcodes specific to 5-B were contained within the unfiltered 5-A data matrix but
211 were excluded as part of the Cell Ranger quality control steps. Only in sample 5-B, the high-depth replicate,

212 were enough UMIs sequenced to distinguish erythrocytes from ambient noise. Hence, our results
213 demonstrate that imbalanced sequencing depth may result in cell type biases, contributing somewhat to
214 batch-associated data effects.

215 Pooled Library Sequencing: It is widely assumed that technical effects can be minimized by pooling libraries
216 for sequencing in a single batch. However, we demonstrate above that duplicated sequencing of identical
217 library aliquots can yield highly similar results with minimal batch-specific variation (Fig. 2c). To directly
218 evaluate the benefits of sequencing libraries together in a single pool, compared to sequencing un-pooled
219 libraries independently, we again analyzed the sequenced replicates of PBMC samples 4 and 5 (samples 4-
220 A, 4-B, and 5-A). We do not expect samples 4 and 5 to be identical, as they originate from different sample
221 donors. However, we will directly contrast the analysis of merged samples 4-A/5-A, which we sequenced in
222 a single pool, against the analysis of merged samples 4-B/5-A, which we sequenced in different pools (see
223 experimental design, Fig. 3a). If indeed pooled sequencing alleviates batch-effect, we would expect to
224 identify less batch-specific variation (lower iLISI and CMS scores) when merging samples that were pooled
225 for sequencing together (samples 4-A vs. 5-A) as compared to un-pooled samples sequenced separately
226 (samples 4-B vs. 5-A). Notably, we sequenced all samples for these comparisons to a highly similar read
227 depth (~40,000 reads/cell). Surprisingly, we find that pooling libraries for sequencing yields nearly identical
228 results as sequencing unpooled libraries independently (Figs. 3b-c). Cell misclassification rates rose
229 compared to the previously presented experiments of duplicated sequencing of identical libraries (Fig. 2d)
230 but remained constant between pooled and un-pooled comparisons (CMS = 0.070 and 0.070, respectively).
231 The UMAP homogeneity similarly remained near-constant between pooled and un-pooled sequencing
232 batches (iLISI = 0.279, 0.295, respectively; Figs. 3b-c). Our results assert that pooling libraries for
233 sequencing neither reduced nor contributed to a major source of batch-associated data effects in these
234 samples.

235

236 **Commonly used batch-effect correction and data integration methods imperfectly resolve batch-** 237 **associated systemic variation**

238 We examined the extent to which commonly cited data integration algorithms can minimize the
239 batch-specific systemic variation we quantified in the prior experiments. Above, we showed that the sample
240 donor variable significantly confounded our ability to correctly assign cell types and resulted in a sample-
241 biased UMAP from a merged dataset (samples 1-3; Fig. 2a, iLISI = 0.861). Here, we attempt to improve the
242 low-dimensional embeddings (i.e., UMAP and clusters) and the robustness of cell-type classifications by
243 applying commonly-used batch-effect correction algorithms. First, we selected three of the most cited
244 packages for integrating gene expression data: Harmony⁸, LIGER¹¹, and Seurat V3 (SCTransform and data
245 anchoring)^{6,10}, which have been established as the best performing by prior study⁵. Next, we followed the
246 default workflow settings recommended by the respective publications and applied each method to our
247 PBMC samples 1-3. Finally, we generated UMAP embeddings and cluster/cell-type assignments to assess
248 each corrected dataset (Fig. 4).

249 We first evaluated Harmony, implemented through the Seurat function “RunHarmony”^{8,14}. Although
250 homogeneity of the UMAP is greatly improved (Fig. 4a; iLISI improved from 0.861 to 0.183), CMS scoring of
251 the harmonized sample revealed a greater loss of cell-type fidelity (Fig. 4a; CMS rose from 0.099 to 0.154).
252 Interestingly, the increased proportion of cell-type misclassification (from 9.9% in the uncorrected sample to
253 15.4% in the harmonized sample) comes largely from over-merging CD4⁺ and CD8⁺ T cells into mixed-
254 lineage clusters (Figs. 4a and Supplementary Fig. 4). This highlights the potential to over-homogenize data
255 while ignoring specific cell type-exclusive signals.

256 We next tested LIGER, as implemented through the R package “rliger” and the Seurat functions
257 “RunOptimizeALS” and “RunQuantileNorm”^{11,15,16}. We generated a UMAP from the LIGER integrative non-
258 negative matrix factorization (iNMF) components (Fig. 4b). LIGER produced a homogeneously distributed
259 UMAP (Fig. 4b; iLISI = 0.132) but severely over-merged clusters, resulting in a major loss of cell-type
260 information for nearly one-quarter of all cells, as measured by a severely increased CMS score (Fig. 4b;
261 CMS rose from 0.099 to 0.257). As with Harmony, the most affected cell types were the CD4⁺ and CD8⁺ T

262 cells. LIGER could not differentiate those two distinct subsets of T cells (Supplementary Fig. 4) and incurred
263 a misclassification cost of greater than 25% of cells (i.e., CMS = 0.257).

264 Lastly, we evaluated the performance of Seurat V3, as implemented through the Seurat
265 SCTransform and data anchoring workflows^{6,17}. Data integration by Seurat V3 resulted in improved
266 homogeneity of UMAP embeddings (Fig. 4c; iLISI improved from 0.861 to 0.191), accompanied by an
267 increase in CMS scores compared to uncorrected data (Fig. 4c; CMS rose from 0.099 to 0.182). Again, as
268 with both Harmony and Liger, the Seurat data integration methods were unable to effectively segregate
269 subsets of CD4⁺ and CD8⁺ T cells (Supplementary Fig. 4). We conclude that similar yet transcriptionally
270 distinct types of PBMCs pose a problem for data integration methods as all three failed to resolve these
271 specific T-cell subsets. While we focus our subsequent efforts on detailing the batch-effect impacting T-cell
272 subsets, it is important to note that these effects were global and impacted other major lineages, including
273 NK cells, B cells, and monocytes (Supplementary Fig. 4). Crucially, each method generated a different set
274 of clusters when applied to the same merged dataset (Fig. 4), highlighting a potential to misinterpret results
275 from a single method, especially when relying on clusters as a meaningful descriptor of biological status.
276 Notably, by applying the CMS score, we are uniquely able to reveal and quantify the hazards that dataset
277 integration can pose to data interpretation.

278

279 **Optimized dataset normalization and scaling before integration resolves batch-effect without the** 280 **need for batch-correction methods**

281 We showed that batch-specific variations negatively impact UMAP embeddings (iLISI scores) and
282 cell-type classifications (CMS scores) in merged samples, and yet commonly-used batch-correction
283 methods imperfectly resolve these effects (Figs. 2 and 4). Hence, we sought to identify and optimize the
284 specific data processing steps which can influence downstream UMAP embeddings and clustering. We
285 focused on the steps of data normalization and data scaling, which are required to produce comparable
286 transcript counts between cells and across genes, and must be applied prior to any PCA and downstream
287 UMAP and clustering. To quantify the effect of data normalization and scaling on resolving batch-associated

288 systemic variation, we applied two unique data normalization and scaling methods with two different data
289 pooling workflows to generate a total of four datasets for comparison (Fig. 5).

290 We first evaluated the common normalization/scaling method of log-transformation and gene
291 centering, as implemented by Seurat methods `NormalizeData` and `ScaleData`¹⁰. Because each cell has a
292 different number of UMI sequenced, `NormalizeData` divides gene count values by the total number of read
293 counts per cell and multiplies by a scaling factor (10,000 by default). The result is a scaling of each cell to a
294 total of 10k UMIs to avoid the effect of a different sequencing depth across cell types in a sample.
295 `NormalizeData` then adds a pseudocount of 1 (to avoid transcript zero-values) and takes the natural log of
296 each count. In this way, we normalize cells for per-cell sequencing depth, fostering more similar cell-cell
297 comparisons. The data requires further scaling, however, to stabilize the relationship between gene
298 expression level and variance. `ScaleData` employs a simple gene-level centering and scaling, meaning that
299 each gene will be mean-centered to zero and expression values scaled by the standard deviation. The
300 resulting scaled values (a z-score) are clipped to a maximum (default of 10) to reduce the effect of outlier
301 high-variance genes expressed by a minority subset of cells.

302 The second method evaluated is the `SCTransform` method included in Seurat V3¹⁰. Briefly,
303 `SCTransform` models UMI counts using a regularized negative binomial model to remove the cell-cell
304 variation caused by differing sequencing depth between cell barcodes. In accomplishing this, `SCTransform`
305 pools genes with similar abundance to obtain stable parameter estimates, preventing the overfitting caused
306 by a global scaling model. In this way, `SCTransform` simultaneously corrects for influences of both total UMI
307 and mean expression on the gene variance.

308 We applied the above two normalization/scaling methods to PBMC samples 1-3, which had the
309 greatest batch-effect in our prior experiments (Fig. 2a), and which were ineffectively integrated by common
310 batch-effect correction algorithms (Fig. 4). Importantly, each dataset was processed by each method, either
311 log-normalization and scaling (Fig. 5a) or `SCTransform` (Fig. 5b), both prior to data merging or after data
312 merging as a single unified/concatenated dataset. We then assessed the differences in the final UMAP
313 visualization and cluster compositions and generated iLISI and CMS scores for each method (Figs. 5 and

314 Supplementary Fig. 4). SCTransform normalization produced undesirable results, showing greater numbers
315 of cell-type misclassification (high CMS scores) and sample-stratified UMAP embeddings (high iLISI
316 scores), irrespective of whether it was performed prior to or after dataset merging (Fig. 5b; CMS = 0.202
317 and iLISI = 0.416 prior to merging; CMS = 0.204 and iLISI = 0.559 after merging). We identified one
318 method, log-normalization with mean-centering performed on each sample independently and prior to data
319 merging, that produced the best homogenized UMAP (Fig. 5a; iLISI = 0.211). However, nearly one in every
320 six cells was misclassified after dataset merging (Fig. 5a; CMS = 0.159). In contrast, the same
321 normalization/scaling method performed after dataset merging showed a highly stratified UMAP embedding,
322 even though it had the best CMS score (Fig. 5a; iLISI = 0.861, CMS = 0.099). The simple default methods
323 of log-normalization with mean-centering performed prior to dataset merging surprised us by their ability to
324 nearly match the performance of dedicated batch-effect correction algorithms (Figs. 4 and 5a).

325 Taken together, our results implicate data normalization and scaling as an effective method to de-
326 emphasize the systemic variation present in the sample-specific data matrix, performing similarly to the
327 Harmony and Seurat V3 batch-effect correction methods. However, even the best normalization
328 approaches led to high levels of cell misclassification, which carry a dangerous potential for
329 misinterpretation of data and false conclusions.

330

331 **Low levels of highly variable transcripts are associated with batch-effect**

332 We show above that normalization and merging strategies can have dramatic impacts on the
333 observed levels of batch-effect. Here, we attempt to isolate the gene transcripts, which may be responsible
334 for the batch effects observed when the integrated samples are normalized differently.

335 Zero-count gene transcripts can indicate either the true absence of a transcript, or that the
336 expression level of a given transcript is very low and may not be captured by the assay, resulting in what is
337 known as gene “dropout” events. Both true zero-count transcripts and gene dropouts account for a large
338 proportion of the gene expression matrix and can contribute significantly to batch-effect^{18,19}. One key
339 difference between general normalization methods commonly used for bulk RNA-seq and those specifically

340 developed for scRNA-seq is the ability to cope with excessive zeros in the scRNA-seq data matrix¹⁸. Hence,
341 we reasoned that low expression level transcripts (fewer than 5 UMI) may be susceptible to gene dropout
342 (i.e., zero-count transcripts) and that proper normalization may minimize the effect of these low expression
343 level genes.

344 To identify precisely how normalization methods could improve cell-type classification, we first
345 isolated the misclassified cells of sample 1, the individual sample with the greatest total cell-type
346 misclassification (i.e., highest CMS score) of the PBMCS samples 1-3 (Figs. 2a, 4, and 5). We found that a
347 specific group of cells in sample 1, the CD8⁺ T cells, frequently changed their initial classification when the
348 cell-typing workflow to classify clusters (Supplementary Fig. 3) was repeated/validated after sample
349 merging. This group of CD8⁺ T cells was prone to inappropriately co-cluster with NK cells, along with
350 another group of CD8⁺ T cells which were co-clustered with CD4⁺ T cells and regulatory T cells (Fig. 6a).
351 Notably, this observation was not exclusive to a single computational workflow (Supplementary Fig. 4) and
352 was repeated across all normalization methods tested (Fig. 5).

353 We next sought to identify the specific gene transcripts which can distinguish the misclustered CD8+
354 T cells of sample 1 from the CD8⁺ T cells of samples 2 and 3. We directly performed differential gene
355 expression (DGE) analyses between the misclassified CD8⁺ T cells of sample 1 vs. the correctly classified
356 CD8⁺ T cell cluster using a likelihood ratio test for differential expression, as implemented by Seurat
357 “FindMarkers” function^{10,20}, and identified a set of significant differentially-expressed genes (Fig. 6b). Next,
358 to restrict this list to genes that could directly influence the UMAP and clusters, we retained only genes
359 detected as “highly variable genes” (HVGs). Only HVGs are included as input to the PCA, which is used as
360 the input for the downstream UMAP and clustering. Finally, to assess whether the genes in our list could
361 cause batch effect, we excluded the genes that were in the HVGs of the dataset that showed the least batch
362 effect (Fig. 5a, left panel). The final list is a limited set of genes that are differentially expressed between the
363 CD8⁺ T cells of sample 1 and the samples 2/3, are capable of influencing PCA (and downstream UMAP and
364 clustering), and are present in the HVGs of datasets which show strong batch effect. This final list of 26
365 genes is presented in Fig. 6b.

366 For 23 out of 26 genes present in this final list, the median transcript count (of cells with nonzero
367 transcripts) was two or fewer copies (i.e., ≤ 2 UMI). The more highly expressed genes were all ribosomal
368 subunits (*RPS26*, *RPS4Y1*, *RPS4X*), which had median transcript counts of 8, 3, and 18 UMIs, respectively
369 (Fig. 6b). Remarkably, we show that specific removal of these genes from the HVGs of samples 1-3
370 resulted in a marked improvement in data integration without the need for any other optimization (Fig. 6c).
371 iLISI scores improved from 0.861 to 0.313, and CMS scores approached the scores achieved by the best
372 batch-correction and normalization approaches (CMS = 0.175). Therefore, our results show that the
373 selection of normalization/scaling and sample merging workflow plays an important role to either exacerbate
374 or minimize batch-associated systemic variation by properly controlling contributions of low-expression
375 genes to the total variance of the sample. In conclusion, here, we experimentally demonstrated the various
376 sources of batch-associated effects and developed a new scoring system that takes into account cell-type
377 identity as a key biological feature in data integration efficiency. Taken together, these results can inform
378 the optimal experimental design, and data integration approaches.

379

380

381 Discussion

382 As accessibility to scRNA-seq dramatically increases, so too do the challenges of integrating and co-
383 analyzing diverse datasets. One challenge to integrating multiple scRNA-seq samples is the influence of
384 batch-specific technical variance on UMAP visualization and clustering. Here, we apply our novel Cell
385 Misclassification Statistic (CMS) score alongside a modified Local Inverse Simpson's Index (iLISI) scoring to
386 reveal previously unclarified sources of these batch-specific effects and the potential dangers of using
387 batch-correction algorithms. We conclude that batch-effect can be partially mitigated by supervised
388 optimization of normalization and scaling methods, which work by minimizing the influence of low-
389 expression gene transcripts.

390 Existing batch-correction scoring metrics such as Local Inverse Simpson's Index (LISI)^{8,13} or
391 Average Silhouette Width (ASW)²¹ quantify the mixing of merged samples in either the integrated clusters or
392 in the final integrated UMAP. However, these and other current batch-correction approaches remain
393 agnostic to the biology of each cell, such that potential cell-type misclassification after data integration is not
394 considered in the analysis. Yet, as we show here, cell-type misclassification is a pervasive phenomenon in
395 data integration and can dramatically influence data interpretation. Therefore, we contend that in addition to
396 measuring optimal mixing of integrated samples (e.g., LISI or ASW), the potential for cell-type
397 misclassification post sample merging must be taken into account when evaluating the fidelity of batch-
398 correction and data integration methods.

399 In this study, we developed the CMS score to directly measure cell-type misclassification that often
400 occurs following batch correction and data integration. CMS differs from existing approaches such as
401 LISI^{8,13} or ASW²¹ in that CMS quantifies batch-integration based on a known biological classification of cells,
402 rather than by simply measuring the mixing of different samples. CMS is uniquely sensitive to the incorrect
403 merging of dissimilar cell types caused by the over-homogenization of datasets. To complement these
404 strengths, we combine CMS with a modified LISI score, or the integration LISI (iLISI), to measure the
405 efficiency of the UMAP integration. iLISI is uniquely sensitive to detecting the under-merging of similar
406 samples, a phenomenon that can result in a UMAP that is segregated by sample instead of cell type. We

407 find that the biologically-grounded measures produced by our CMS scoring prove to be a robust indicator of
408 cell-type fidelity post data integration, while our modified iLISI scoring is an excellent benchmark of UMAP
409 integration. We demonstrate the ability of this combined CMS and iLISI approach to measuring the most
410 impactful contributors to batch-effect between samples.

411 Using CMS and iLISI, we reported the ability of common experimental variables to influence the
412 biological interpretation of clusters. In isolation, replicating the library sequencing did not cause any
413 significant loss in cell-type fidelity, while sequencing depth and sample donor both did. However, the effects
414 of sequencing depth were contained entirely within a single, defined cell type while sample donor effects
415 (which include the entangled effects of microfluidic encapsulation and library preparation) were widespread
416 across all cells surveyed. While we identified this phenomenon on PBMCs here, we would expect similar
417 results in other heterogeneous tissues, both with regards to sequencing depth and the effects of sample
418 donor. Pooling PBMC libraries for sequencing together or sequencing libraries independently did not
419 change the observed batch effect. This suggests that pooling libraries for sequencing together may not
420 supply a significant benefit to experimental design or effectively prevent batch-associated systematic
421 variation. We emphasize the critical role of CMS scores in quantifying cell-type fidelity, which helped us to
422 conclude that sample donor is the major contributing variable to batch-effect in the PBMCs evaluated.

423 While PBMC populations may subtly differ between healthy adult donors, we should be able to pool
424 PBMCs from different donors for the purposes of constructing a reference sample (i.e., a healthy control
425 sample or a cell atlas). To integrate PBMCs from different donors, we applied three popular and often cited
426 batch-correction algorithms⁵ (Harmony⁸, LIGER¹¹, and Seurat V3¹⁰). Although, as expected, we find that all
427 three methods effectively integrated UMAP embeddings, we show that explicit batch-correction may be
428 unnecessary and, in some cases, even harmful. By selecting a suitable normalization and merging strategy
429 through a CMS-guided approach, we produced cell-type integration results comparable to the best batch-
430 correction algorithms. We endorse a supervised approach (such as our CMS-guided optimization) over
431 selecting a batch-correction method for the following reasons: while batch-correction methods do not
432 directly adjust the gene expression matrix, they do affect the dimensional reduction matrix and therefore

433 cluster assignments, which are used to infer cell-type identity. Most significantly, each batch correction
434 method evaluated here produced subtly different cluster assignments from the same data, which could
435 provoke a dangerous potential to misclassify cells and misinterpret results. Therefore, we caution against
436 the application of batch-effect correction methods as a “fix” for systemic batch-effect and place emphasis on
437 the proper, supervised selection of appropriate normalization methods. The choice of the best approach
438 may be aided by a measure such as CMS, or similar, which directly measures the impact of any
439 computational approach on biologically assigned cell-type classifications. By directly quantifying the level of
440 biological signal loss, we can reveal systemic variation that may otherwise go undetected.

441 To complete our study, we investigated the most symptoms of such systemic variations on the gene
442 expression matrix. We identified the cells most sensitive to misclassification and highlighted a pattern of
443 gene expression which distinguished them from the similar cell type, but correctly classified cells, of the
444 other samples. The differentially expressed genes between the misclassified and correctly classified cells
445 reveal a broad pattern of expression with a low total expression level. We find that, depending on the
446 normalization/merging strategy chosen, these genes may be excluded from the highly variable gene (HVG)
447 list. Strikingly, we show that simple exclusion of these genes significantly improves batch-effect but is
448 insufficient to eliminate all batch-effect in an improperly normalized/merged sample. This suggests that
449 more sources of variation may be present, which cannot be corrected by gene curation alone. We conclude
450 that removing problematic genes *post-hoc* is insufficient and that optimizing the normalization/merging
451 strategy is the best approach to consistently reduce sample-specific variance and highlight biologically
452 relevant signals.

453 Together, our experiments describe the variables which contribute most to batch-effect (sample
454 donor, sequencing depth) and those which do not significantly contribute (sequencing replicate, sequencing
455 pool). We apply a combined approach of CMS and iLISI scoring to reveal that batch correction algorithms
456 pose a risk to over-merge diverse cell types if not properly supervised, while data normalization offers a
457 simple strategy for minimizing the influence of batch-specific, low expression gene transcripts to the sample.
458 We recognize that while we identify the best strategy to minimize batch-effect in the PBMC samples tested,

159 the optimal strategy may vary by tissue or even by dataset. Rather than endorse a single approach to
160 remove batch-effect, we offer a novel tool, CMS scoring, that, in combination with iLISI can assess cell
161 mismatching and UMAP integration and aid in choosing the correct computational methods for any dataset.

462 **Materials and Methods**

463 **Sample preparation and single-cell encapsulation**

464 Peripheral venous blood was collected from healthy volunteer donors through Emory University's
465 Children's Clinical and Translational Discovery Core. Peripheral Blood Mononuclear Cells (PBMCs) were
466 isolated using the Direct Human PBMC Isolation Kit (StemCell Technologies) according to the
467 manufacturer's protocol and washed/resuspended in custom RPMI-1640 deficient in biotin, L-glutamine,
468 phenol red, riboflavin, and sodium bicarbonate (defRPMI), and containing 3% newborn calf serum and
469 Benzonase. Cell number and viability were assessed using ViaStain™ Acridine Orange/Propidium Iodide
470 (AOPI, Nexcelom) on a Cellometer K2 cell counter (Nexcelom), following manufacturer recommendations. A
471 maximum of 1×10^6 cells/donor were stained by oligo-barcoded antibodies (TotalSeq-A or TotalSeq-C;
472 Biolegend, Expedeon) for 30 min on ice, followed by two washes in defRPMI + 0.04% BSA. Cells were
473 resuspended at 1200-1500 cells/ul in defRPMI + 0.04% BSA, passed through a 20uM filter, and counted
474 prior to loading onto a Chromium Controller (10X Genomics). For generating single-cell RNA-seq libraries,
475 cells were loaded to target encapsulation of six thousand cells.

477 **Library generation and sequencing**

478 Single-cell RNA-seq gene expression and ADT libraries were generated following the
479 manufacturer's instructions using Chromium Single Cell 3'Library & Gel Bead Kit v2, with ADT library
480 generation according to (Nat Methods 14, 865–868 (2017) for PBMC samples 1-3 and Chromium Single
481 Cell 5'Library & Gel Bead Kit v1 with feature barcoding (10X Genomics) for PBMC samples 4 and 5. Gene
482 expression libraries were sequenced to an average depth of 84,098 reads per cell using the Novaseq 6000
483 platform (Illumina). ADT libraries were sequenced to a target depth of one hundred reads per antibody, per
484 cell, on the Next-seq platform (Illumina). A separate aliquot of cells for each donor was stained with a
485 fluorescent antibody cocktail (Supplementary Table 3) for 30 minutes on ice, followed by washing with
486 defRPMI/3%NBCS + Benzonase (Sigma). Cells were fixed for 30 minutes at RT with FACS lysis solution

487 (BD Biosciences), followed by a final wash and analysis on a BD LSRII or FACSymphony A5 cytometer.
488 Data were analyzed with FlowJo V10. The gating strategy for excluding debris, doublets, and non-viable
489 cells is shown in Supplementary Fig. 3.

490

491 **Gene expression (GEX) data processing**

492 Raw sequence data (FASTQ files) was processed in a Linux environment using CellRanger V3 (10X
493 Genomics, version details in Supplementary Table 1) to generate a digital expression matrix. Specifically,
494 splicing-aware aligner STAR was used for sequence alignment to GRCh38 (Ensembl 93). Viable cell
495 barcodes were found automatically by CellRanger. Digital expression matrices were exported for further
496 analyses. Further data analysis was performed in an R environment (Version 3.6.1, CRAN) using the Seurat
497 toolkit (Version 3.2.2, Satija Lab) following previously published workflows²². Following Seurat standard
498 recommendations, data were first filtered for quality using specific QC criteria (maximum mitochondrial
499 content) to limit the analysis to cells with transcriptomes that were not apoptotic. Cell barcodes with
500 mitochondrial transcripts (>5 standard deviations above the median level of mitochondrial transcripts) were
501 suspected to be dying cells and excluded. On average, 99.3% of putative cells met QC criteria and were
502 included in the analysis (range = 99.2 - 99.4).

503 From the five human PBMC donors (seven libraries), we included a total of 37,442 cells, with a
504 range of 3,758-6,374 cells per donor and a median of 5,729 cells per sample. Genes were selected for
505 differential expression across the sample using Seurat's highly variable gene selection tool,
506 "FindVariableFeatures." A setting was chosen to select the 2,000 most variable features identified per
507 sample. Principle Components Analysis (PCA) was used to reduce the dimensionality of the gene
508 expression matrix. A singular value decomposition (SVD) PCA was performed on the subset of highly
509 variable genes. To identify an appropriate number of PCs, we employed a z-scoring method. We run a
510 complete PCA reduction and z-score the contribution of each PC to the total variance. PCs with $z \geq 1$
511 were considered significant and used further in the analysis. The SVD PCA returns the right singular values,

512 representing the embeddings of each cell in PC space, and left singular values, representing the loadings
513 (weights) of each gene in the PC space. Cell embeddings (right singular values) were weighted by the
514 variance of each PC.

516 **Antibody-derived tag (ADT) data processing**

517 ADT data was co-analyzed with GEX data as raw sequence data (FASTQ files) in a Linux
518 environment using Cell Ranger V3 (10X Genomics, version details in Supplementary Table 1). ADT data
519 was aligned directly to a feature reference file containing the sequence of each barcode mapped to the
520 corresponding antibody. Digital expression matrices for ADT protein expression were exported for further
521 analyses. ADT count matrices were normalized and denoised using the R package “dsb” version 0.1.0²³.

523 **Data embedding and cell clustering**

524 For 2-dimensional visualization, Uniform Manifold Approximation and Projection (UMAP) reduction
525 was performed on the PCA matrix (<https://github.com/lmcinnes/umap>). In parallel, independent of the
526 UMAP coordinates, Louvain–Jaccard clustering was performed on the PC-space. This “bottom-up”
527 clustering method employs a stochastic shared-nearest-neighbor (SNN) approach, in which cells are
528 grouped according to their neighbors in PC space. The nearness of two cells is weighted by the Jaccard
529 index or the degree of sharing between the lists of each cell’s nearest neighbors. The algorithm will build
530 small groups of cells and attempt to iteratively merge them into clusters until the modularity is maximized.
531 We found that a resolution of 0.8 was most proper for building biologically meaningful clusters.

533 **Cell-type identification and classification**

534 Following PCA (performed on the GEX matrix of highly variable genes), we generated clusters by
535 Louvain–Jaccard clustering with a resolution of 0.8, as implemented by the Seurat function “FindClusters”.
536 Using detection of lineage-specific ADT and GEX markers, we assigned each cluster to a single major

537 lineage (markers and representative data depicted in Supplementary Fig. 2). When indicated in the text,
538 differentially expressed genes were determined using the likelihood-ratio test for single-cell gene
539 expression, as implemented in the “bimod” method of the Seurat function “FindMarkers”^{10,20}.

540

541 **Cell Misclassification Statistic**

542 To generate the Cell Misclassification Statistic (CMS) score, we first assign cell type labels to
543 clusters in a merged, multi-sample dataset. This can be carried out either by following a gating strategy and
544 assigning cell types to entire clusters which fall within each gate or by transferring cluster labels from a
545 previously gated dataset containing shared cells. The transferred labels can be used to aid cell-type
546 identification but should always be confirmed by marker expression following the overall gating strategy.
547 After cell types have been recorded for the merged samples (target dataset), we match a vector of cell
548 identities from the same sample, processed by a different workflow (reference dataset). We then compare
549 the target cell IDs to a vector of the reference cell type IDs. The CMS score generated is the sum of
550 matching values (total number of cell type ID matches) between the two cell-typing replicates, divided by the
551 total number of cells. We then invert the score by subtracting it from one so that the result is a measure of
552 the fraction of cells that change cell type ID between data analysis workflows. CMS can be directly
553 interpreted as the percentage (expressed as a decimal) of cells which have been misclassified by one of the
554 workflows. Where reported, the CMS is the mean score of all samples present.

555

556 **Modified Local Inverse Simpson’s Index (LISI):**

557 Our modified LISI score, the integration LISI or iLISI, is an effective indicator of sample mixing in the
558 UMAP. A LISI score represents, for a given cell, the number of neighboring cells which need to be sampled
559 before the same identity class is sampled (in this case, sample ID). In other words, LISI effectively counts
560 how many identity classes are represented in the local neighborhood of each cell¹³. If we merge three
561 samples, then the maximum LISI score will be close to three, and the minimum close to one. However, the
562 actual maximum score may differ by UMAP. To account for differing numbers of samples and to place LISI

563 on a similar scale as CMS, we first take the median LISI score for a merged set of samples using the r
564 package “LISI”¹³ with sample ID used as the cell label. We subtract the minimum LISI score (one) from the
565 median LISI and divide the result by the maximum LISI (the number of samples). The result is an iLISI score
566 on a scale of zero to one, where a score of one is ideal integration and zero is a perfectly segregated
567 sample. We then inverted these scores by subtracting them from 1, such that the best integrated UMAP will
568 achieve a score approaching 0. This was done to align the values with CMS so that higher scores mean
569 better integration across both measurements. $iLISI = 1 - (LISI - 1) / (\text{maximum score} - 1)$. Where reported,
570 the iLISI is the mean score of all samples present.

571

572 **Data availability**

573 Gene expression and Antibody Derived Tag matrices have been deposited on the Gene Expression
574 Omnibus (GEO).

575

576 **Code availability**

577 Full data analysis workflow and R scripts are made available at github.com/Ghosn-Lab/BBabcock

578

579

580 References

- 581 1 Yu, P. & Lin, W. Single-cell Transcriptome Study as Big Data. *Genomics Proteomics Bioinformatics* **14**,
582 21-30, doi:10.1016/j.gpb.2016.01.005 (2016).
- 583 2 Luecken, M. D. & Theis, F. J. Current best practices in single-cell RNA-seq analysis: a tutorial. *Mol*
584 *Syst Biol* **15**, e8746, doi:10.15252/msb.20188746 (2019).
- 585 3 Aliverti, E. *et al.* Projected t-SNE for batch correction. *Bioinformatics* **36**, 3522-3527,
586 doi:10.1093/bioinformatics/btaa189 (2020).
- 587 4 Chen, W. *et al.* A multicenter study benchmarking single-cell RNA sequencing technologies using
588 reference samples. *Nat Biotechnol*, doi:10.1038/s41587-020-00748-9 (2020).
- 589 5 Tran, H. T. N. *et al.* A benchmark of batch-effect correction methods for single-cell RNA sequencing
590 data. *Genome Biol* **21**, 12, doi:10.1186/s13059-019-1850-9 (2020).
- 591 6 Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-seq data
592 using regularized negative binomial regression. *Genome Biol* **20**, 296, doi:10.1186/s13059-019-
593 1874-1 (2019).
- 594 7 Hie, B., Bryson, B. & Berger, B. Efficient integration of heterogeneous single-cell transcriptomes
595 using Scanorama. *Nat Biotechnol* **37**, 685-691, doi:10.1038/s41587-019-0113-3 (2019).
- 596 8 Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat*
597 *Methods* **16**, 1289-1296, doi:10.1038/s41592-019-0619-0 (2019).
- 598 9 Stein-O'Brien, G. L. *et al.* Decomposing Cell Identity for Transfer Learning across Cellular
599 Measurements, Platforms, Tissues, and Species. *Cell Syst* **8**, 395-411 e398,
600 doi:10.1016/j.cels.2019.04.004 (2019).
- 601 10 Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902 e1821,
602 doi:10.1016/j.cell.2019.05.031 (2019).
- 603 11 Welch, J. D. *et al.* Single-Cell Multi-omic Integration Compares and Contrasts Features of Brain Cell
604 Identity. *Cell* **177**, 1873-1887 e1817, doi:10.1016/j.cell.2019.05.006 (2019).
- 605 12 Xu, C. *et al.* Comprehensive multi-omics single-cell data integration reveals greater heterogeneity in
606 the human immune system. *bioRxiv* (2021).
- 607 13 Korsunsky, I. *LISI*, <<https://github.com/immunogenomics/LISI>> (2019).
- 608 14 Korsunsky, I. *How to use Harmony with Seurat V3*,
609 <<https://github.com/immunogenomics/harmony/blob/master/docs/SeuratV3>> (2019).
- 610 15 Welch, J. D. *LIGER*, <<https://github.com/welch-lab/liger>> (2021).
- 611 16 Butler, A. *Integrating Seurat objects using LIGER*, <[https://github.com/satijalab/seurat-](https://github.com/satijalab/seurat-wrappers/blob/master/docs/liger.md)
612 [wrappers/blob/master/docs/liger.md](https://github.com/satijalab/seurat-wrappers/blob/master/docs/liger.md)> (2021).
- 613 17 Satija, R. *Integration and Label Transfer: SCTransform Vignette*,
614 <<https://satijalab.org/seurat/archive/v3.0/integration.html>> (2019).
- 615 18 Lytal, N., Ran, D. & An, L. Normalization Methods on Single-Cell RNA-seq Data: An Empirical Survey.
616 *Front Genet* **11**, 41, doi:10.3389/fgene.2020.00041 (2020).
- 617 19 Hou, W., Ji, Z., Ji, H. & Hicks, S. C. A systematic evaluation of single-cell RNA-sequencing imputation
618 methods. *Genome Biol* **21**, 218, doi:10.1186/s13059-020-02132-x (2020).
- 619 20 McDavid, A. *et al.* Data exploration, quality control and testing in single-cell qPCR-based gene
620 expression experiments. *Bioinformatics* **29**, 461-467, doi:10.1093/bioinformatics/bts714 (2013).

- 521 21 Rousseeuw, P. J. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis.
522 *Journal of Computational and Applied Mathematics* **20**, 53-65, doi:10.1016/0377-0427(87)90125-7
523 (1987).
- 524 22 Ocasio, J. *et al.* scRNA-seq in medulloblastoma shows cellular heterogeneity and lineage expansion
525 support resistance to SHH inhibitor therapy. *Nat Commun* **10**, 5829, doi:10.1038/s41467-019-
526 13657-6 (2019).
- 527 23 Mulè, M. P., Martins, A. J. & Tsang, J. S. Normalizing and denoising protein expression data from
528 droplet-based single cell profiling. *bioRxiv*, doi:<https://doi.org/10.1101/2020.02.24.963603> (2021).
529

530

531 **Acknowledgments**

532 This study was supported in part by Georgia Clinical and Translational Science Alliance (CTSA)
533 through the National Center for Advancing Translational Sciences of the National Institutes of Health under
534 Award number NIH UL1TR002378; Pediatric Research Alliance, Center for Transplantation, and Immune-
535 Mediated Disorders (Children's Healthcare of Atlanta); and Lowance Center for Human Immunology. We
536 thank Sachin Kumar (Emory University) for helpful conversations. We thank Emory University's Children's
537 Clinical and Translational Discovery Core (CCTDC) for providing peripheral blood samples from healthy
538 donors. Flow cytometry data were collected at the Emory's Pediatrics/Winship Flow Cytometry Core
539 (access supported in part by Children's Healthcare of Atlanta). Single-cell libraries were sequenced at the
540 Emory Integrated Genomics Core (EIGC), which is subsidized by the Emory University School of Medicine
541 and is one of the Emory Integrated Core Facilities, and at PerkinElmer Genomics Inc.

542

543 **Author Contributions**

544 EEBG and BRB conceived the study. BRB performed all computational analyses and developed
545 the CMS scoring method under the supervision of EEBG. BRB, EEBG and AK wrote, reviewed and
546 edited the manuscript. AK performed all tissue processing and flow cytometry, and generated the
547 scRNA-seq libraries. JY performed all data pre-processing workflows, including GEX and ADT
548 alignment. MLW performed flow cytometry analyses. All authors read and approved the final
549 manuscript.

550

551 **Competing Interest**

552 The authors declare no competing interests.

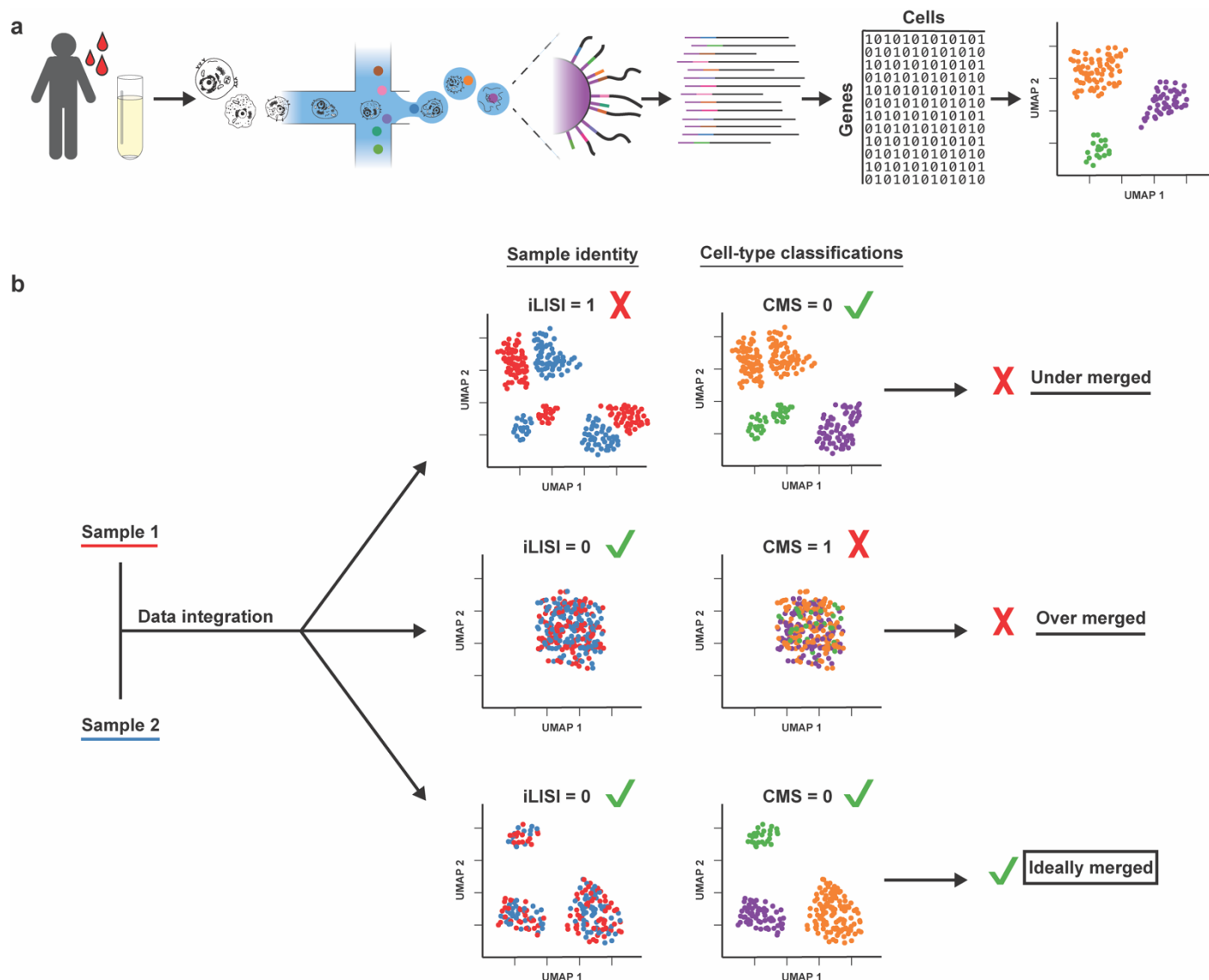


Fig. 1: Schematic of analysis workflow to quantify data integration and batch-correction fidelity.

a, Tissue (blood) is extracted from the sample donor, washed, and PBMCs isolated. PBMCs are encapsulated in microfluidic droplets, along with a barcode-bearing bead. Cells are lysed in the droplets, and mRNA is captured on the bead, resulting in a barcoded cDNA library. Libraries are sequenced to generate a GEX matrix (cells x genes) containing transcript counts. Cells are analyzed and plotted by UMAP and clustered according to transcript similarity, after which cell types are classified. **b**, Cell Misclassification Statistic (CMS) and modified integration LISI (iLISI) metrics provide measures of cell-type fidelity and UMAP mixing, respectively. CMS and iLISI can assist in optimizing a workflow for the best-mixed UMAP while preserving the cell-type-specific signals in the data.

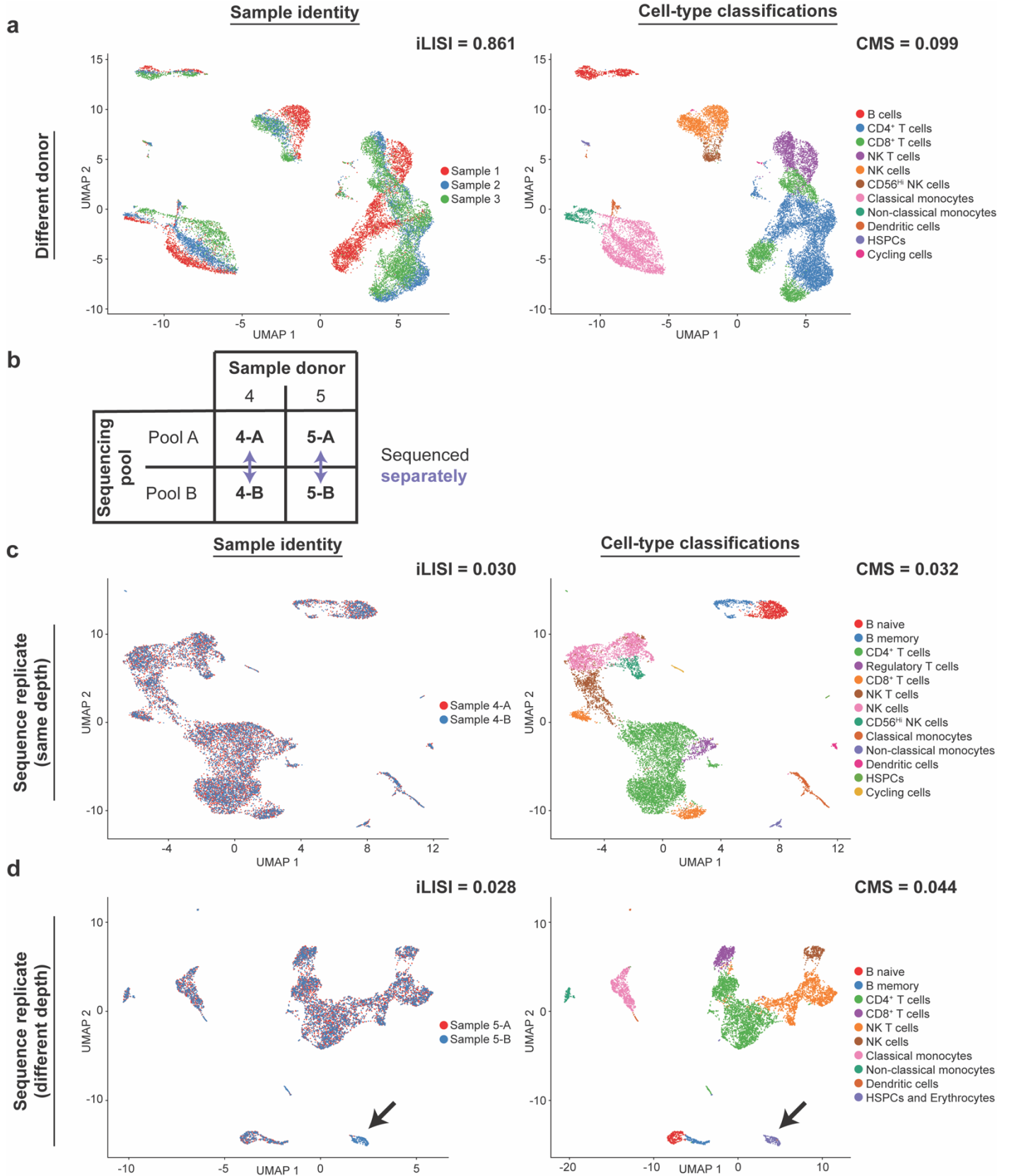


Fig. 2: Batch effects are generated by sample donor and sequencing depth, but not sequencing replicates alone. **a**, PBMCs from three different donors, processed and sequenced simultaneously and mixed for co-analysis, produce a favorable cell-type classification fidelity (CMS = 0.099) but an unfavorable (i.e., poorly mixed) UMAP (iLISI = 0.861). **b**, A schematic of sample sequencing strategy. Libraries from two donors (samples 4 and 5) were sequenced twice. **c**, A library generated from a single PBMC sample, divided and sequenced twice to similar read depth (average reads per cell), produces a well-integrated UMAP and high-fidelity cell classifications (iLISI = 0.030, CMS = 0.032). **d**, A library generated from a single PBMC sample, divided and sequenced twice with a three-fold difference in read depth produces favorable UMAP and cell-type fidelity metrics (iLISI = 0.028, CMS = 0.044). However, there are local sample-specific UMAP islands (black arrow) and misclassified cell types that reveal sequencing-depth-specific batch effects.

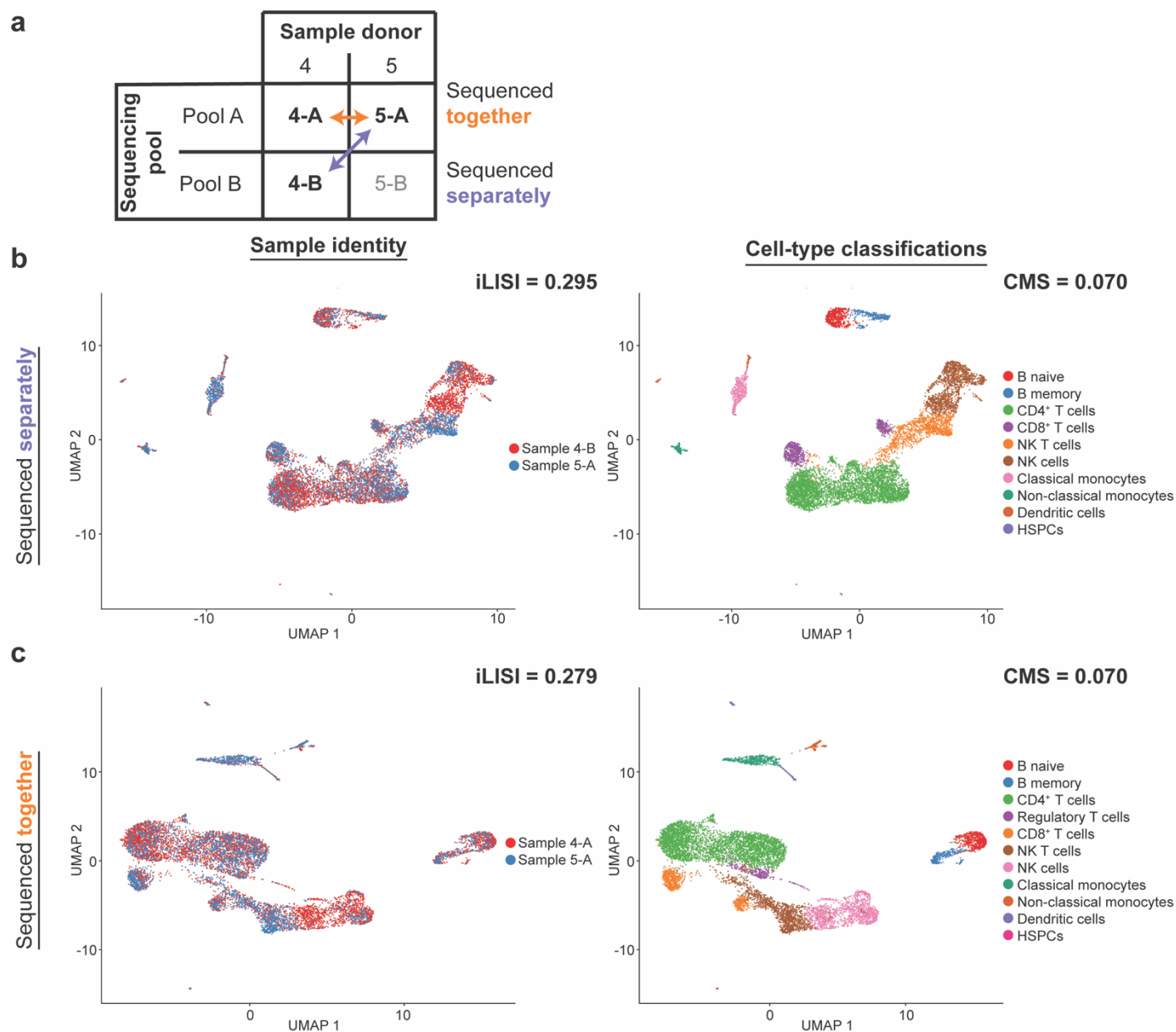


Fig. 3: Pooling samples for sequencing does not appreciably improve the measured batch effect. a, A schematic of pooling and comparison strategy between samples 4 and 5, sequenced across two independent pools (-A and -B). The same sample (5-A) was co-analyzed with either a sample from the same sequencing pool (4-A) or a different pool (4-B). **b,** Analysis of two PBMC samples from different sequencing pools results in a poorly mixed UMAP and the misclassification of 7% of cells (iLISI = 0.295, CMS = 0.070). **c,** Pooling libraries to sequence PBMC samples in a single pool, followed by mixing data for co-analysis, does not improve batch effects and, similarly, results in a poorly mixed UMAP and an equivalent 7% cell misclassification rate (iLISI = 0.279, CMS = 0.070).

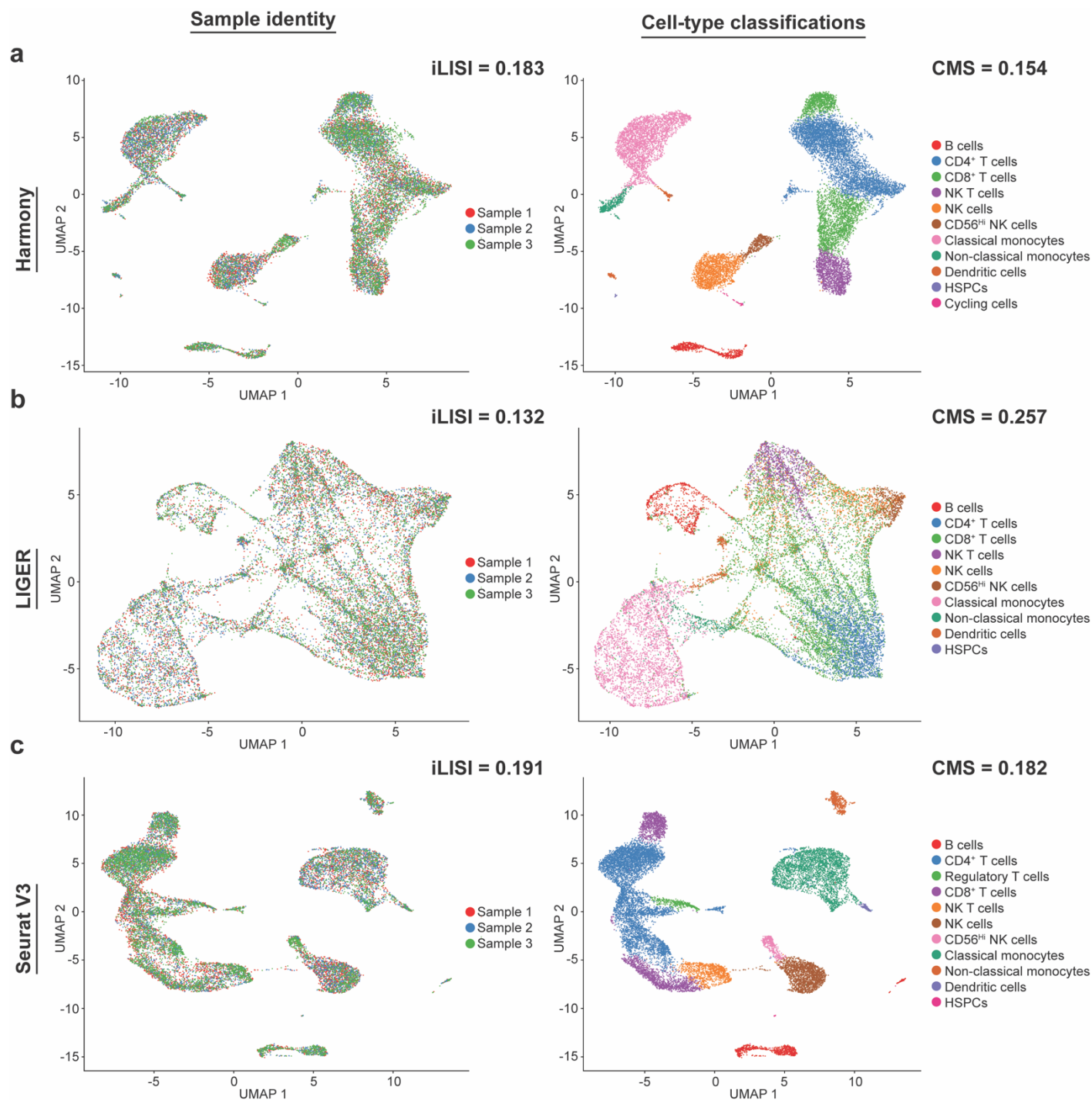


Fig. 4: Common batch-effect correction methods imperfectly resolve batch effect. **a**, Harmony correction merges samples to produce an integrated UMAP, but at 15.4% loss in cell-type classification fidelity (iLISI = 0.183, CMS = 0.154). **b**, LIGER correction produces an integrated UMAP, but at a 25.7% loss in cell-type classification fidelity (iLISI = 0.132, CMS = 0.257). **c**, Seurat V3 correction results in an integrated UMAP but at a cost of an 18.2% loss in cell-type classification fidelity (iLISI = 0.191, CMS = 0.182).

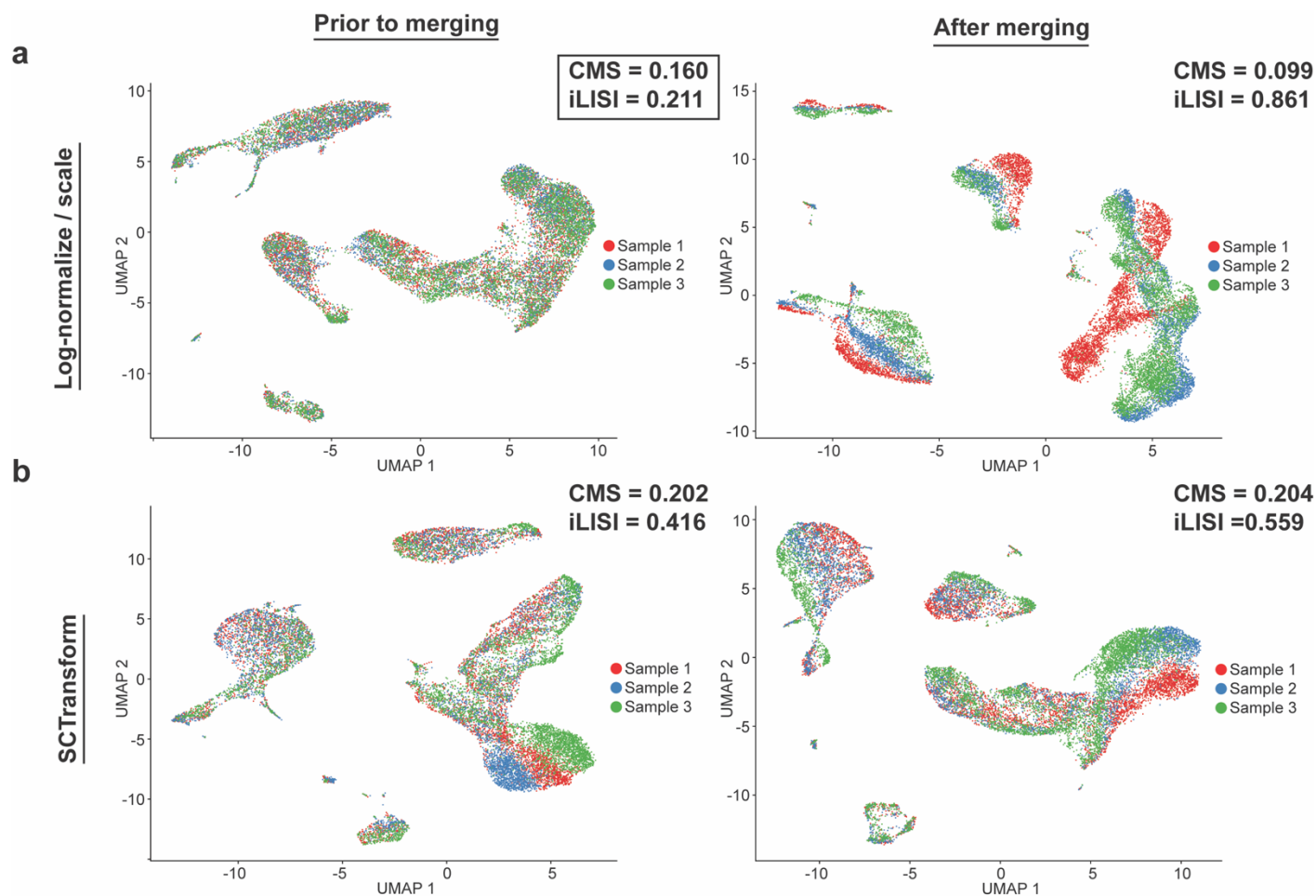
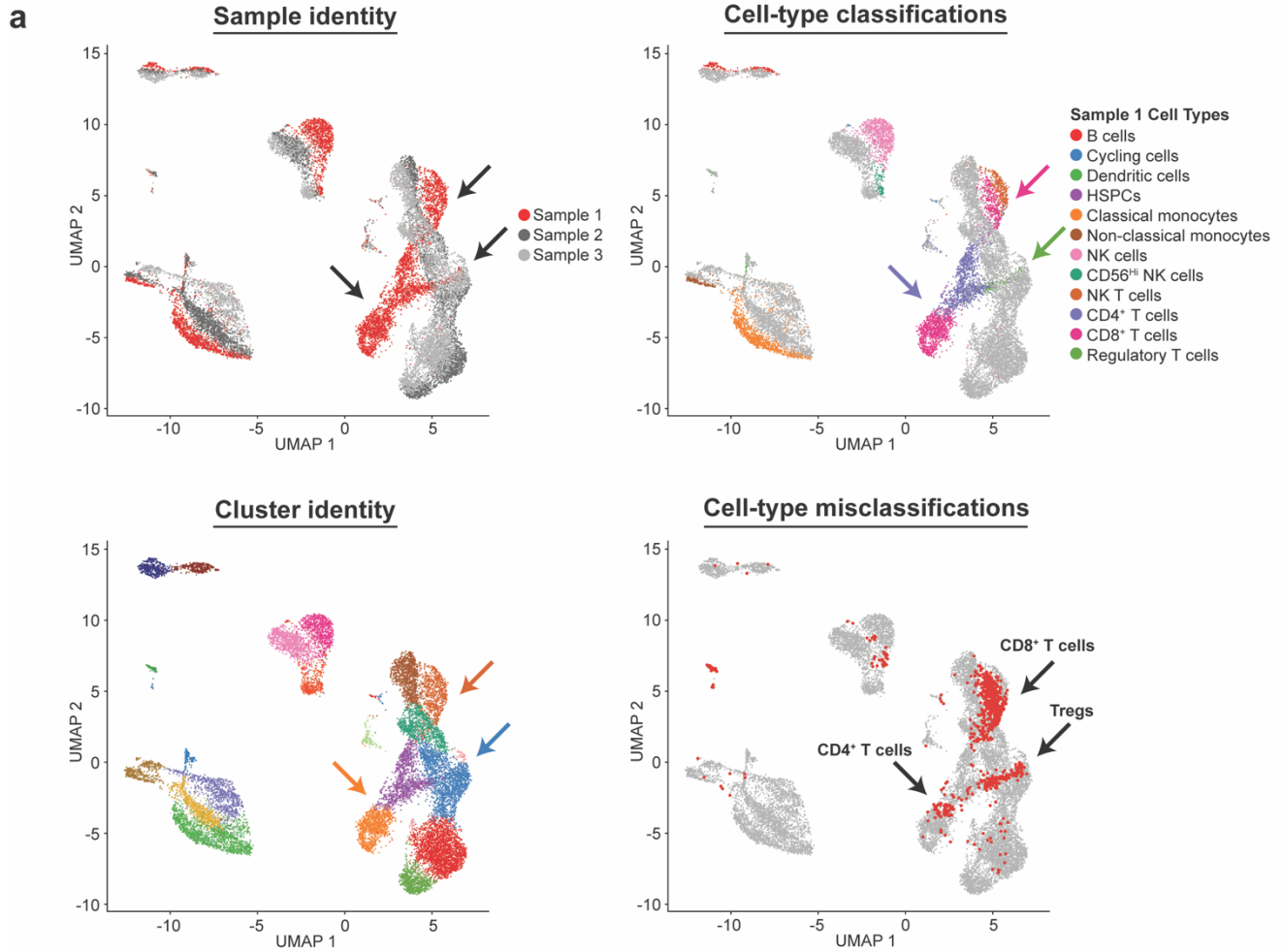


Fig. 5: Data normalization and merging strategies differentially impact the measured batch effect. a, UMAPs, colored by sample, when the same normalization methods (log-normalization and data scaling) are performed on each sample either prior to sample merging (i.e., normalized separately) or after sample merging (normalized together). **b,** UMAPs, colored by sample, when the SCTransform normalization method is performed before or after sample merging. For this data set, log-normalization/scaling performed on each sample separately, prior to sample merging, generated the best combination of both integration scores, CMS and iLISI (panel a, left).



b

Gene	RPS26	RPS4Y1	XIST	RPS4X	EIF1AY	SLC4A10	SRGN	MX1	HLA-DQA2	CTSA	FAM118A	RUNX3	NOSIP
Fraction of cells expressing	90.2%	27.5%	36.1%	98.1%	6.3%	1.2%	64.1%	13.8%	7.2%	13.2%	15.8%	23.4%	49.4%
Median count (of positive cells)	8	3	1	18	1	1	2	1	1	1	1	1	2
Gene	IFI44	LINC00685	SUB1	PYHIN1	CEP78	PLP2	IDH2	MTRNR2L8	SIT1	JAKMIP1	LY6E	PTRHD1	H1FX
Fraction of eells expressing	9.5%	9.6%	70.8%	18.3%	9.6%	36.1%	20.8%	13.8%	17.4%	3.0%	59.6%	16.5%	25.9%
Median count (of positive cells)	1	1	2	1	1	1	1	1	1	1	2	1	1

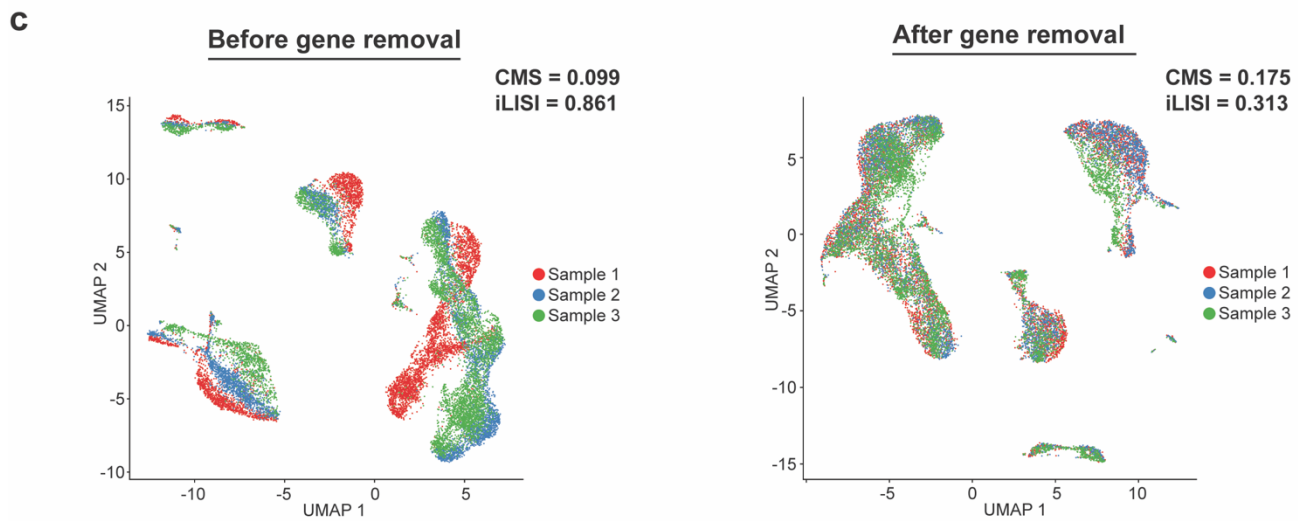
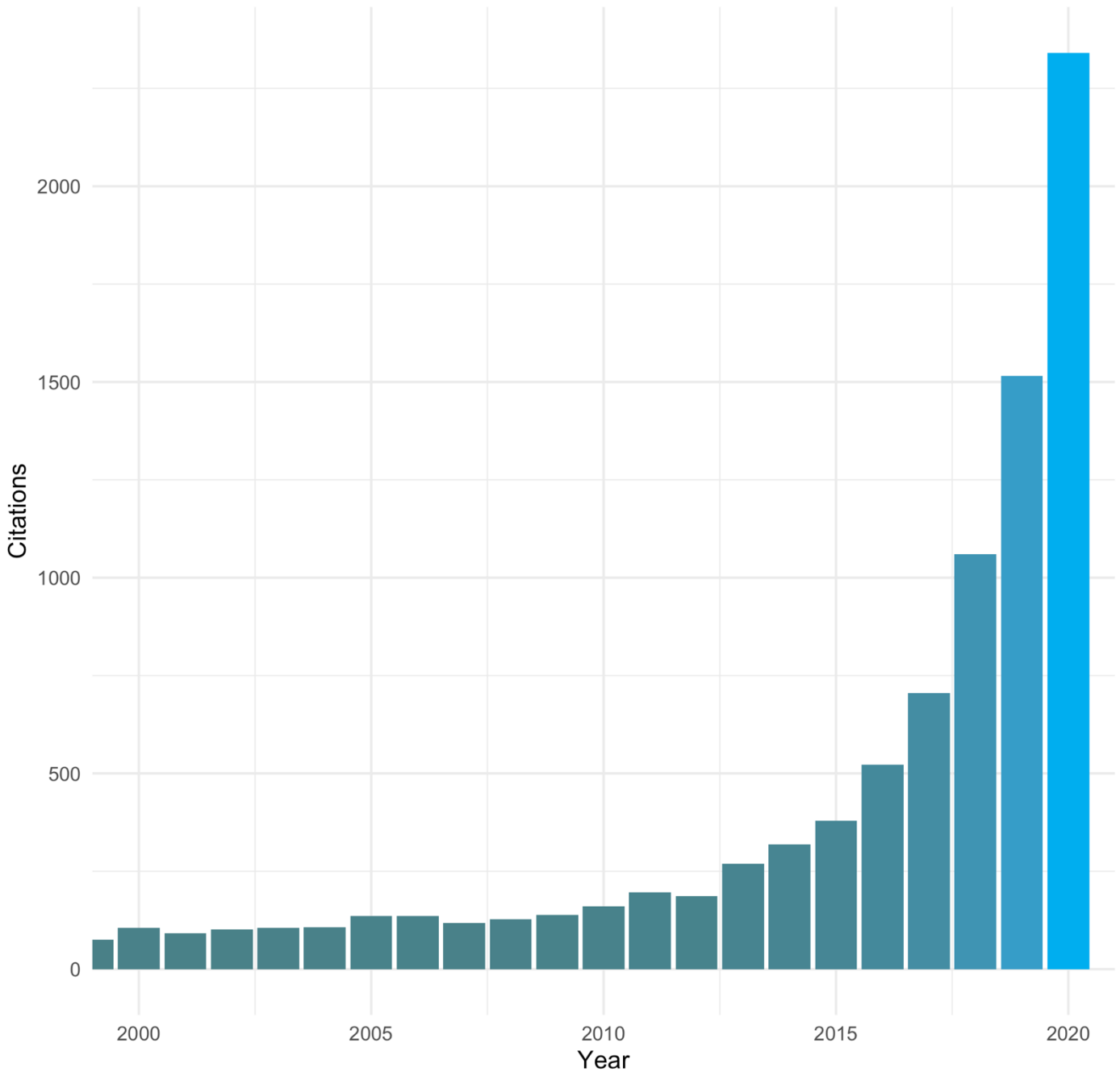
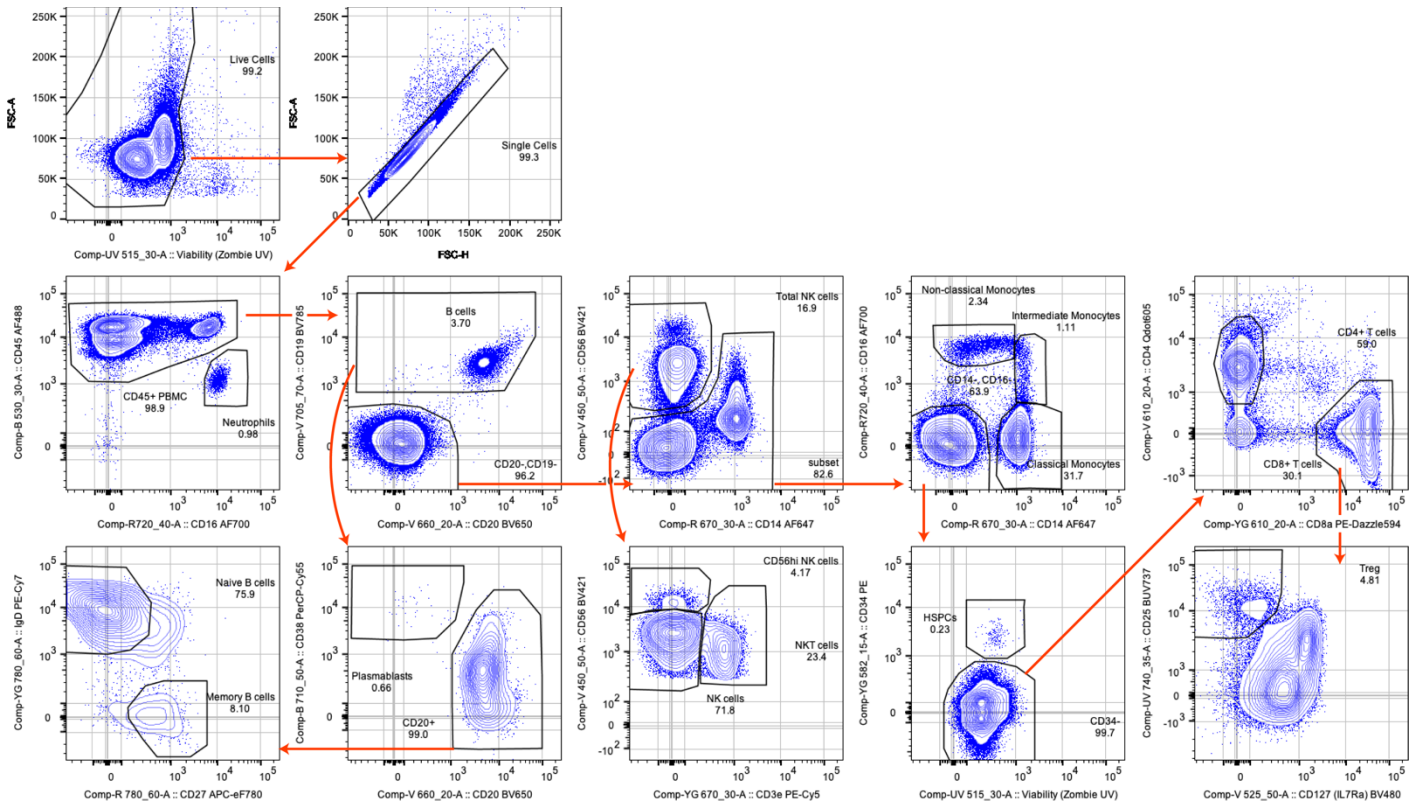


Fig. 6: Low expression levels of highly-variable genes can play a role in batch effect, which can be amplified by the data normalization and merging strategies. **a**, Misclassified cells are not distributed at random, and specific analysis of PBMC sample 1 reveals a loss of cell-type identity resulting from the over-merging of T cell subsets into clusters containing a majority of a different subset (misclassified T cells highlighted by arrows and red points). **b**, Selected genes identified as differentially expressed between the incorrectly classified CD8⁺ T-cells of sample 1, marked by an arrow in (a), compared to the correctly classified CD8⁺ T-cell cluster. Genes were restricted to those contained within the HVGs of sample 1, but not the HVGs of samples 2 or 3. Note that genes showing low expression levels (low median count) and broad expression (higher fraction of cells) are susceptible to gene dropout. **c**, Removal of selected genes (b) lessens but does not completely remove the batch effect from the sample (iLISI changes from 0.861 to 0.313, CMS changes from 0.099 to 0.175).

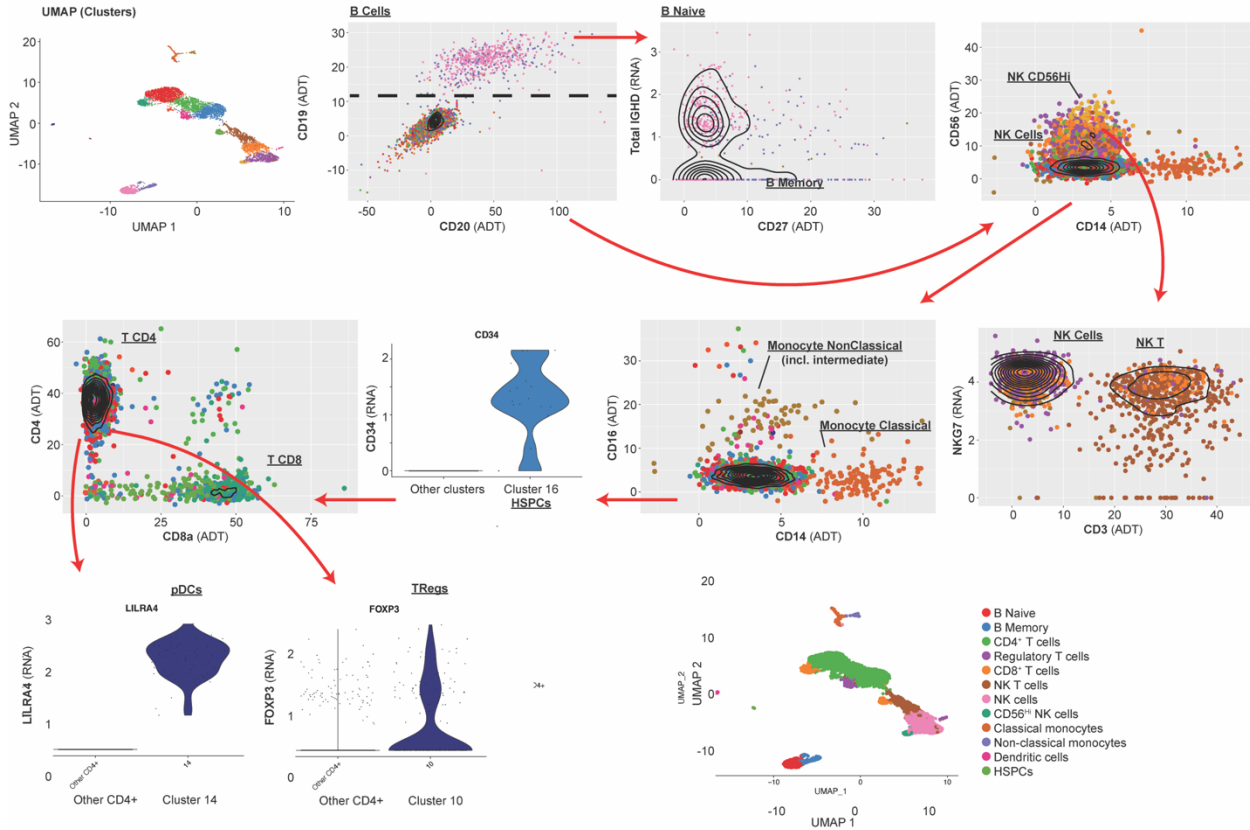
"Single Cell RNA" Pubmed Results by Year, 1999-2020



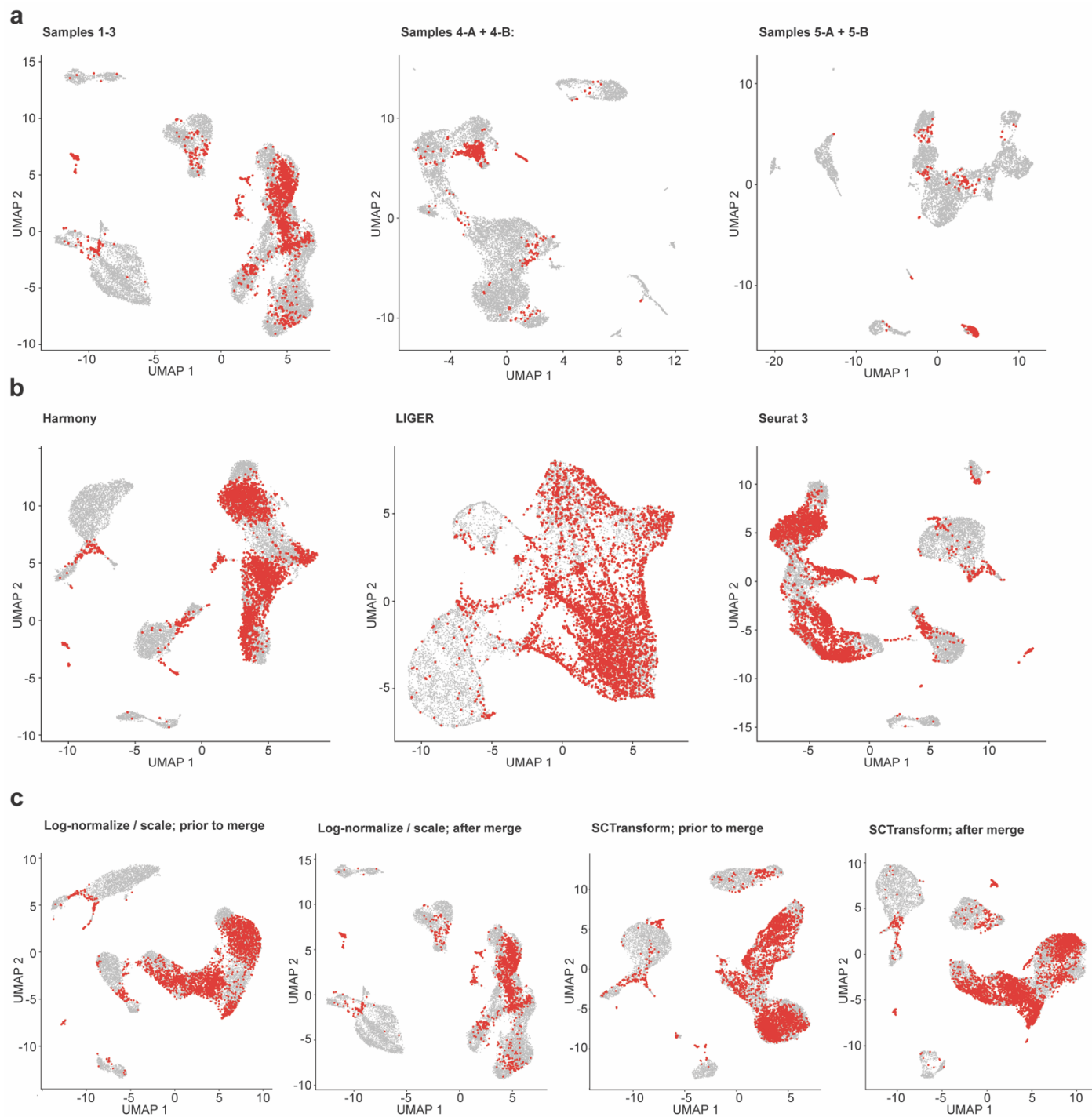
Supplementary Figure 1. PubMed results for “Single Cell RNA” by year. Increasing trend of publications citing single-cell RNA-seq data in recent decades.



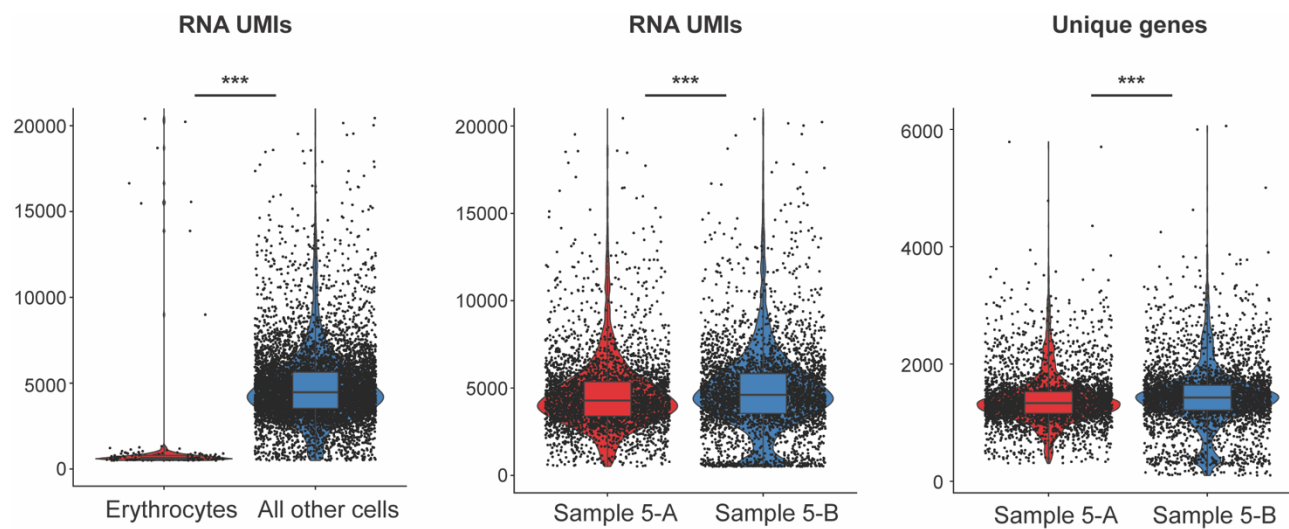
Supplementary Figure 2. High-dimensional flow cytometry gating strategy to identify major immune lineages in the PBMC samples. Representative gating strategy generated from sample 3, showing markers and gates used to identify major immune cell lineages using an 18-parameter flow cytometry panel.



Supplementary Figure 3. Gating strategy to identify major immune lineages using multi-omics single-cell sequencing data. Representative gating strategy generated from sample 4-A, showing markers and gates used to assign entire clusters to major immune cell lineages based on gene expression (GEX matrix) and cell-surface proteins (ADT matrix).



Supplementary Figure 4. Distribution of misclassified cells after data merging. Red points represent cells that change classification after data merging, as generated by the CMS scoring method. UMAPs and data generated for Figs. 2, 4, and 5.



Supplementary Figure 5. Comparison of sample 5 high- and low-depth sequencing replicates.

Comparison of key metrics (number of UMI and number of unique genes) distinguishing high and low read-depth (reads per cell) sequencing replicates of sample 5. Also shown are the number of UMI for sample 5 erythrocytes, as compared to all non-erythrocyte cells. *** $p < 1 \times 10^{-9}$, Wilcoxon rank-sum test.

Supplementary Table 1. Donor demographics and library sequencing details

Donor ID	Sample ID	Source Tissue	Sex	Age	Isolation method	10X Chemistry	ADT version	Sequence Platform	Flow cell	CellRanger Version	Reference genome
Donor 1	PBMC Sample 1	PBMC	M	37	PBMC isolation SCT kit	3' v2	TotalSeq-A	Novaseq	S4-300	3.0.1	GRCh38 (Ensembl 93)
Donor 2	PBMC Sample 2	PBMC	F	41	PBMC isolation SCT kit	3' v2	TotalSeq-A	Novaseq	S4-300	3.0.1	GRCh38 (Ensembl 93)
Donor 3	PBMC Sample 3	PBMC	F	43	PBMC isolation SCT kit	3' v2	TotalSeq-A	Novaseq	S4-300	3.0.1	GRCh38 (Ensembl 93)
Donor 4	PBMC Sample 4-A	PBMC	M	46	PBMC isolation SCT kit	5' v1	TotalSeq-C	Novaseq	S4-300	3.1.0	GRCh38 (Ensembl 93)
Donor 4	PBMC Sample 4-B	PBMC	M	46	PBMC isolation SCT kit	5' v1	TotalSeq-C	Novaseq	S4-300	3.1.0	GRCh38 (Ensembl 93)
Donor 5	PBMC Sample 5-A	PBMC	F	36	PBMC isolation SCT kit	5' v1	TotalSeq-C	Novaseq	S4-300	3.1.0	GRCh38 (Ensembl 93)
Donor 5	PBMC Sample 5-B	PBMC	F	36	PBMC isolation SCT kit	5' v1	TotalSeq-C	Novaseq	S4-300	3.1.0	GRCh38 (Ensembl 93)

Reads	Cells
582,860,987	5778
582,860,987	5178
582,860,987	6121
582,860,987	6392
582,860,987	6423
582,860,987	3783
582,860,987	4034

Supplementary Table 2. Cell-type and cluster composition (fraction of total sample)

Cell-type composition by samples							
Cell type	B Naive	B Memory	T CD4	TReg	T CD8	NK_T	NK
PBMC Sample 4-A	0.067748543	0.030880731	0.472664251	0.02851741	0.062391681	0.086394675	0.162596502
PBMC Sample 4-B	0.06761845	0.031220584	0.473956699	0.028082837	0.062598055	0.084719172	0.1625353
PBMC Sample 5-A	0.041522491	0.02874634	0.410700027	---	0.076124567	0.26031408	0.047112057
PBMC Sample 5-B	0.038682306	0.027202396	0.395557774	---	0.072373347	0.246568505	0.044172698
PBMC Sample 5-A (Erythrocytes Removed)	0.041655541	0.028838451	0.412016021	---	0.076368491	0.261148198	0.047263017
PBMC Sample 5-B (Erythrocytes Removed)	0.040714473	0.028631468	0.416338324	---	0.076175466	0.259521933	0.046493302
Cluster composition by samples							
Cluster	0	1	2	3	4	5	6
PBMC Sample 4-A	0.20529384	0.108870332	0.081770915	0.080825587	0.075783835	0.076098944	0.074996061
PBMC Sample 4-B	0.203953561	0.108095388	0.080953875	0.081581425	0.076090367	0.075462818	0.07530593
PBMC Sample 5-A	0.118179398	0.102475379	0.097684323	0.085174341	0.083311153	0.083311153	0.076124567
PBMC Sample 5-B	0.10706264	0.093835787	0.095333167	0.079111555	0.080359371	0.079361118	0.072373347

NK CD56HI	Monocyte_C1 assical	Monocyte_No nClassical	Dendritic_Cel ls	HSPCs	Cycling_Cells	Erythrocytes		
0.033401607	0.031668505	0.010241059	0.008192847	0.002835986	0.003466205	---		
0.033260119	0.031377471	0.010197678	0.008158142	0.002823972	0.003451522	---		
---	0.098749002	0.022624434	0.010912962	---	---	0.003194038		
---	0.093835787	0.021462441	0.010232094	---	---	0.049912653		
---	0.099065421	0.022696929	0.010947931	---	---	---		
---	0.098765432	0.022589966	0.010769635	---	---	---		
7	8	9	10	11	12	13	14	15
0.067748543	0.037655585	0.033401607	0.031668505	0.02851741	0.025366315	0.024736096	0.010241059	0.00929573
0.06761845	0.03796674	0.033260119	0.031377471	0.028082837	0.025729526	0.024631315	0.010197678	0.009256354
0.068937982	0.065211605	0.047112057	0.04285334	0.041522491	0.02874634	0.003194038	0.022624434	0.015437849
0.066633392	0.060144747	0.044172698	0.039430996	0.038682306	0.027202396	0.049912653	0.021462441	0.014474669

Supplementary Table 3. Key reagents and resources

Antigen	Clone	Vendor	Cat. no.	10X Chemistry compatibility
CD34 (gp105-120)	581	Biolegend	343537	3'
CD24 (Ly-52)	ML5	Biolegend	311137	3'
CD138 (Syndecan-1)	MI15	Biolegend	356533	3'
CD3 (T3, CD3ε)	UCHT1	Biolegend	300475	3'
Annexin V	Annexin	Biolegend	94700	3'
CD16 (FcγRIII)	3G8	Biolegend	302061	3'
IgD	IA6-2	Biolegend	348243	3'
IgM	MHM-88	Biolegend	314541	3'
CD56 (NCAM-1)	5.1H11	Biolegend	362557	3'
CD127 (IL-7Ra)	A019D5	Biolegend	351352	3'
CD38 (T10)	HIT2	Biolegend	303541	3'
CD20 (B1)	2H7	Biolegend	302359	3'
CD27 (S152)	O323	Biolegend	302847	3'
CD14 (LPS receptor)	63D3	Biolegend	367131	3'
CD19 (B4)	HIB19	Biolegend	302259	3'
CD5 (Ly-1)	UCHT2	Biolegend	300635	3'
CD43 (Ly-48)	CD43-10G7	Biolegend	343209	3'
CD25	BC96	Biolegend	302643	3'
CD4	SK3	Biolegend	344649	3'
CD8	SK1	Biolegend	344751	3'
CCR7 (CD197)	G043H7	Biolegend	353251	5'
CD11b	ICRF44	Biolegend	301359	5'
CD127 (IL-7Ra)	A019D5	Biolegend	351356	5'
CD14 (LPS receptor)	M5E2	Biolegend	301859	5'
CD16 (FcγRIII)	3G8	Biolegend	302065	5'
CD19 (B4)	HIB19	Biolegend	302265	5'
CD20 (B1)	2H7	Biolegend	302363	5'
CD21 (C3dR)	Bu32	Biolegend	354923	5'
CD25	BC96	Biolegend	302649	5'
CD27 (S152)	O323	Biolegend	302853	5'
CD3 (T3, CD3ε)	UCHT1	Biolegend	300479	5'
CD38 (T10)	HIT2	Biolegend	303543	5'
CD4	RPA-T4	Biolegend	300567	5'
CD45	2D1	Biolegend	368545	5'
CD45RA	HI100	Biolegend	304163	5'
CD56	QA17A16	Biolegend	392425	5'

CD69	FN50	Biolegend	310951	5'
CD80	2D10	Biolegend	305243	5'
CD86	IT2.2	Biolegend	305447	5'
CD8a	RPA-T8	Biolegend	301071	5'
CD95 (APO-1)	DX2	Biolegend	305651	5'
HLA-DR	L243	Biolegend	307663	5'
PDL1 (CD274)	29E.2A3	Biolegend	329751	5'
IgM	MHM-88	Biolegend	314547	5'
<i>Custom oligo-conjugated in house</i>				
Antigen	Clone	Vendor	Cat. no.	
CD11c (α X integrin)	S-HCL-3	Biolegend	371502	5'
CD133 (Prominin-1)	clone 7	Biolegend	372802	5'
CD138 (Syndecan-1)	MI15	Biolegend	356502	5'
CD24 (Ly-52)	ML5	Biolegend	311102	5'
CD34 (gp105-120)	581	Biolegend	343602	5'
CD41 (gpIIb)	HIP8	Biolegend	303702	5'
CD43 (Ly-48)	CD43-10G7	Biolegend	343202	5'
CD45R (B220)	RA3-6B2	Biolegend	103202	5'
CD5 (Ly-1)	UCHT2	Biolegend	300602	5'
IgD	IA6-2	Biolegend	348212	5'
CD20 (B1)	2H7	Biolegend	302302	5'