

# 1 The differential impacts of dataset imbalance in

## 2 single-cell data integration

**3** Hassaan Maan<sup>1,2,3</sup>, Lin Zhang<sup>3</sup>, Chengxin Yu<sup>4,6</sup>,  
 Michael Geuenich<sup>4,6</sup>, Kieran R Campbell<sup>†4,5,6</sup>, & Bo Wang<sup>†2,3,7,8</sup>

<sup>4</sup> <sup>1</sup>Department of Medical Biophysics, University of Toronto, Toronto, ON,  
<sup>5</sup> Canada

**6** *<sup>2</sup>Vector Institute, Toronto, ON, Canada*

<sup>7</sup> *<sup>3</sup>Peter Munk Cardiac Centre, University Health Network, Toronto, ON, Canada*

**8**      <sup>4</sup>Department of Molecular Genetics, University of Toronto, Toronto, ON,  
**9**      Canada

<sup>5</sup>Department of Statistical Sciences, University of Toronto, Toronto, ON, Canada

<sup>12</sup> <sup>6</sup>Lunenfeld-Tanenbaum Research Institute, Toronto, ON, Canada

<sup>7</sup>Department of Computer Science, University of Toronto, Toronto,

<sup>8</sup>Department of Laboratory Medicine and Pathobiology, University of Toronto,

15                   Toronto, ON, Canada

<sup>†</sup> Corresponding authors: [kierancampbell@lunenfeld.ca](mailto:kierancampbell@lunenfeld.ca), [bo.wang@uhmresearch.ca](mailto:bo.wang@uhmresearch.ca)

17

## Abstract

18 Single-cell transcriptomic data measured across distinct samples has led to a  
19 surge in computational methods for data integration. Few studies have ex-  
20 plicitly examined the common case of cell-type imbalance between datasets to  
21 be integrated, and none have characterized its impact on downstream analy-  
22 ses. To address this gap, we developed the *IniQuitate* pipeline for assessing the  
23 stability of single-cell RNA sequencing (scRNA-seq) integration results after  
24 perturbing the degree of imbalance between datasets. Through benchmarking  
25 5 state-of-the-art scRNA-seq integration techniques in 1600 perturbed integra-  
26 tion scenarios for a multi-sample peripheral blood mononuclear cell (PBMC)  
27 dataset, our results indicate that sample imbalance has significant impacts on  
28 downstream analyses and the biological interpretation of integration results. We  
29 observed significant variation in clustering, cell-type classification, marker gene-  
30 based annotation, and query-to-reference mapping in imbalanced settings. Two  
31 key factors were found to lead to quantitation differences after scRNA-seq inte-  
32 gration - the cell-type imbalance within and between samples (*relative cell-type*  
33 *support*) and the relatedness of cell-types across samples (*minimum cell-type*  
34 *center distance*). To account for evaluation gaps in imbalanced contexts, we  
35 developed novel clustering metrics robust to sample imbalance, including the  
36 balanced Adjusted Rand Index (bARI) and balanced Adjusted Mutual Infor-  
37 mation (bAMI). Our analysis quantifies biologically-relevant effects of dataset  
38 imbalance in integration scenarios and introduces guidelines and novel metrics  
39 for integration of disparate datasets. The *IniQuitate* pipeline and balanced clus-  
40 tering metrics are available at <https://github.com/hsmaan/IniQuitate> and  
41 <https://github.com/hsmaan/balanced-clustering>, respectively.

## 42 Introduction

43 Single-cell sequencing technologies developed in the past decade have led to  
44 breakthrough discoveries due to the high resolution that they offer in deter-  
45 mining biological heterogeneity [1–3]. A major challenge associated with the  
46 analysis of high throughput sequencing data is that of accounting for batch  
47 effects, which are technical artifacts caused by factors such as differences in se-  
48 quencing protocols, experimental reagents, and ambient conditions that lead to  
49 quantification changes that are not biologically driven [4]. Batch effects can lead  
50 to major discrepancies in comparisons of similar experimental groups that can  
51 easily be misinterpreted as biological signal [5]. The amount of mRNA captured  
52 and reads sequenced per cell in single-cell RNA sequencing (scRNA-seq) assays  
53 is very low compared to their bulk counterparts, leading to measurements that  
54 tend to be sparse and noisy [6, 7]. These factors, combined with measurements  
55 often conducted across separate experimental groups without balanced designs  
56 [7], leads to a higher susceptibility of scRNA-seq data to batch effects. Methods  
57 for removing batch effects from bulk RNA sequencing data have demonstrated  
58 poor performance in single-cell settings due to invalid assumptions of shared  
59 populations and linear application of technical effects [8]. To account for this  
60 gap in methodology, batch correction/integration techniques have been devel-  
61 oped specifically for scRNA-seq data [8].

62 Current single-cell integration methods underperform in settings where datasets  
63 are imbalanced based on cell-types [9]. More specifically, this form of imbalance  
64 is dictated by differences in the cell-types present, number of cells per cell-type,  
65 and cell-type proportions across samples [9, 10]. Imbalanced datasets occur  
66 in many integration contexts, including developmental and cancer biology. In  
67 developmental data, it is unlikely that cell populations and proportions will be  
68 shared across samples from different developmental time-points due to factors  
69 such as depletion of stem-like progenitors and differentiation [11]. In tumor  
70 samples, both clonal and subclonal heterogeneity can be present, as well as  
71 different levels of immune and stromal cell infiltration, both within and across  
72 samples [12]. Therefore, as imbalanced contexts can be common in single-cell  
73 data analysis, integration methods and analysis pipelines must be able to ex-  
74 plicitly address these imbalances or integration results may lead to inaccurate  
75 biological conclusions.

76 In comprehensive single-cell integration benchmarking studies by Tran et

al. and Luecken et al. [9, 13], scRNA-seq integration methods were found to perform poorly in terms of both batch-correction and cell-type identity conservation metrics, particularly in large and imbalanced datasets. Ming et al. [10] highlighted dataset imbalance limitations through simulation studies for balanced and imbalanced cell-type compositions in scRNA-seq integration settings, and demonstrated that cell-type proportion imbalance leads to skewed distributions in standardized gene expression values between datasets. This drives major changes in the dimensionality reduction step in scRNA-seq analysis, and subsequently leads to inaccurate integration results [10]. Currently, no existing study has quantified the effects of dataset imbalance on both integration results and downstream biological conclusions. This aspect is highly relevant, as mechanisms to account for dataset imbalance do not readily exist in frequently utilized integration techniques [9, 13].

Here, we present an extensive analysis of the effects of dataset imbalance on scRNA-seq data integration. We begin by examining two balanced scRNA-seq batches of human peripheral blood mononuclear cell (PBMC) data [9, 14, 15] as a controlled setting. To determine the effects of dataset imbalance on integration results and downstream analyses, we perform 1600 perturbation experiments using the *Iniquitate* pipeline that involve control, downsampling, and ablation simulations in a cell-type-specific manner with replicates. Downstream analyses tested include unsupervised clustering [8], differential expression to determine marker genes [8], nearest-neighbor-based cell-type classification [16], and query-to-reference cell-type annotation [17]. To extend the analyses to more complex settings, we analyze datasets with prevalent imbalance, including imbalanced PBMC datasets [18], longitudinal mouse hindbrain developmental data [19], and pancreatic ductal adenocarcinoma (PDAC) samples from different patients [20]. Our analyses reveals that dataset imbalance has cell-type-specific effects on integration performance, as well as the downstream results, and that these effects are largely method-agnostic. We further define two key aspects of multi-sample single-cell data that act in concert to affect downstream results - *relative cell-type support* and *minimum cell-type center distance*. To address limitations with respect to dataset imbalance in benchmarking single-cell integration, we reformulate current integration metrics to consider imbalance explicitly. Finally, we provide a series of guidelines and recommendations to help minimize and mitigate the impacts of dataset imbalance in scRNA-seq integration settings.

112

## Results

113

### Development of a comprehensive perturbation pipeline to determine the impacts of imbalance in scRNA-seq integra- tion

114

115

To assess the impacts of dataset imbalance in scRNA-seq integration, we developed a pipeline termed *Iniquitate*, that quantifies imbalance prevalent in datasets using global and per-cell-type statistics, determines the differences in these quantities between samples/batches, and tests the effects of down-sampling perturbations on integration and downstream analysis results (Figure 1A). Datasets utilized were annotated by experts in their respective studies, with the exception of the PDAC data which was re-annotated to better identify malignant cells (Online Methods). We tested five state-of-the-art scRNA-seq integration methods, including BBKNN [21], Harmony [22], Scanorama [23], scVI [24] and Seurat [25]. A uniform integration pipeline embedded within *Iniquitate* was utilized to make comparisons between methods and across datasets comparable, with some noted exceptions (Online Methods). We measured cell-type heterogeneity conservation and batch effect correction for each technique across datasets and perturbations using the Adjusted Rand Index (ARI) [26], Adjusted Mutual Information (AMI) [27], Homogeneity Score [28], and Completeness Score [28] (Online Methods).

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

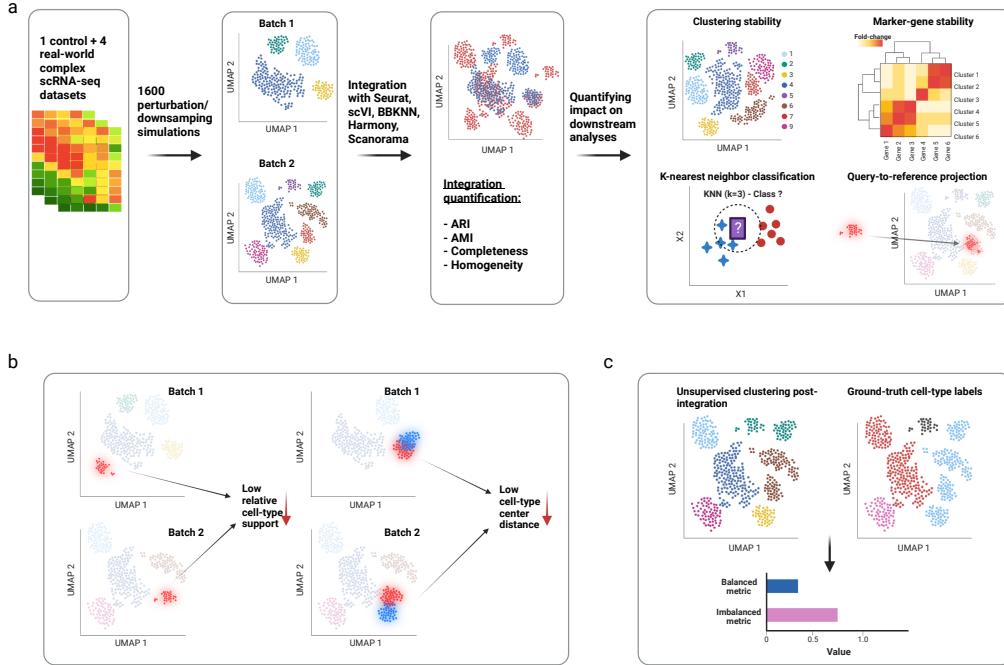
To determine the impacts of dataset imbalance on downstream analyses, we analyzed post-integration impacts on unsupervised clustering [8, 29], cell-type classification [16], differential gene expression [8] [30], and query-to-reference annotation [17] results (Figure 1A). Clustering impacts were assessed based on changes in the number of clusters post-integration using unsupervised clustering (Figure 1A). To assess the impacts on cell-type classification, a nearest-neighbor cell-type classifier was trained on the post-integration embeddings and tested on a holdout set (Figure 1A). Differential gene expression results and variation was assessed using a global importance metric of the ranking changes of marker genes specific to each cell-type analyzed, before and after perturbing the dataset balance (Figure 1A). Query-to-reference annotation was done using the Seurat 4.0 method [31], which projects each batch to be integrated onto a reference scRNA-seq dataset, and accuracy of annotation was utilized as an endpoint (Figure 1A). The details of the evaluations utilized for all of the

146 downstream analyses and parameters of the integration pipeline are outlined in  
147 Online Methods.

148 Through perturbation and cell-type-specific analysis of both balanced and  
149 complex imbalanced scRNA-seq datasets, we determined that cell-type imbal-  
150 ance affects the scores of typical integration metrics in a cell-type and method-  
151 specific manner. Further, we discovered that cell-type imbalance in datasets  
152 to be integrated can lead to significant deviations in the results of downstream  
153 analyses. After investigating factors of imbalance that can quantifiably lead to  
154 distinct downstream results, we found that the distance between cell-types in the  
155 embedding space (*minimum cell-type center distance*) and imbalance between  
156 cell-types (*relative cell-type support*) to be the most relevant and predictive in  
157 this regard (Figure 1B). Finally, we determined that typical clustering metrics  
158 utilized in benchmarking single-cell integration techniques, such as ARI and  
159 AMI, are inadequate in imbalanced scenarios as they weigh the more prevalent  
160 cell-types disproportionality compared to rare cell-types. Therefore, we de-  
161 velop and introduce novel balanced clustering metrics, including the *Balanced*  
162 *Adjusted Rand Index* (bARI), *Balanced Adjusted Mutual Information* (bAMI),  
163 *Balanced Homogeneity Score*, and *Balanced Completeness Score* (Figure 1C).  
164 The balanced metrics reweigh the base scores such that each ground-truth cell-  
165 type's contribution to the score is considered equally.

166 **I. Perturbation-induced imbalance in a PBMC cohort in-**  
167 **icates cell-type-specific effects on integration results**

168 The ideal test case for assessing impacts of dataset imbalance should begin with  
169 a balanced dataset as a baseline, and thus we analyzed a peripheral blood mono-  
170 nuclear (PBMC) cohort of two batches/samples processed independently from  
171 two different healthy donors [9, 14, 15]. We downsampled each batch to have  
172 6 major cell-types and an equal number of cells within each cell-type (400 cells  
173 for each cell-type) (Figure 2A) (Online Methods). The cell-types were selected  
174 such that they are equivalent between the batches. Therefore, the cell-types  
175 present, number of cells per cell-type, and cell-type proportions between the  
176 batches are equal and the integration scenario is balanced (Figure 2A, Figure  
177 2B). A batch effect is prevalent between the samples (Figure 2B), which is ex-  
178 pected as they were processed at different centers using different technologies  
179 (10x 3' vs 5' protocols - Online Methods) [14, 15]. In this balanced setup, we

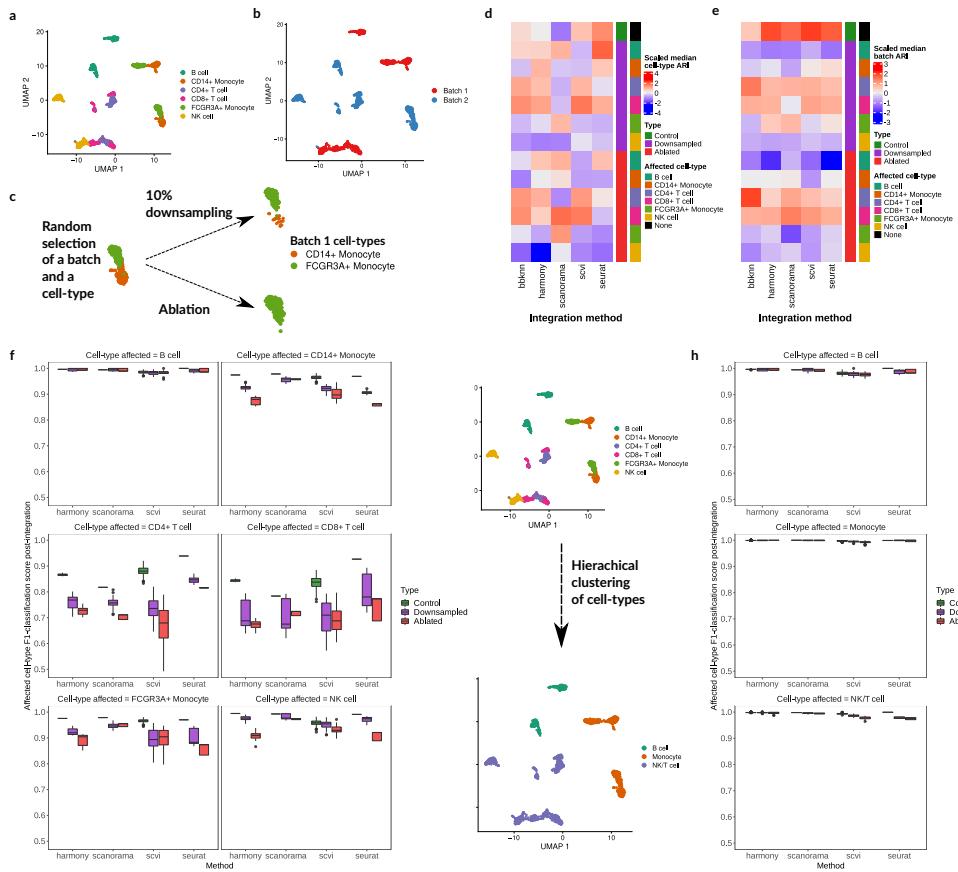


**Figure 1: Overview of the Iniquitate pipeline and analysis results.** (a) To determine the effects of dataset imbalance in scRNA-seq integration, 1 controlled balanced PBMC dataset and 4 complex datasets with imbalance already present were integrated using current state-of-the-art scRNA-seq integration techniques. A total of 1600 perturbation experiments involving downsampling on the controlled dataset were performed and the effects of imbalance on integration results as well as downstream analyses (clustering, differential gene expression, cell-type classification, query-to-reference prediction) were quantified. (b) In complex datasets, results in the controlled setting were verified, and two key data characteristics were found to contribute to altered downstream results in imbalanced settings - relative cell-type support and minimum cell-type center distance. (c) To account for imbalanced scRNA-seq integration scenarios in evaluation and benchmarking, typically utilized metrics and scores were reformulated to reweigh disproportionate cell-types, which includes the Balanced Adjusted Rand Index (bARI), Balanced Adjusted Mutual Information (bAMI), Balanced Homogeneity Score, and Balanced Completeness Score.

180 aimed to assess how the integration results of the two PBMC batches, from a  
181 typical integration metric standpoint as well as their impacts on downstream  
182 analyses, varied between the control balanced data and perturbation-induced  
183 imbalanced data. For each perturbation, we randomly selected one of the two  
184 batches and one cell-type within the selected batch to either downsample to 10%  
185 of the original population or ablate/remove completely from the selected batch  
186 (Figure 2C). These perturbations were repeated 400 times for both downsam-  
187 pling and ablation of a random batch/cell-type. Control experiments with no  
188 perturbations to the balanced data were repeated 800 times, resulting in 1600  
189 integration experiments where each integration technique was tested (Online  
190 Methods).

191 To determine how cell-type-specific changes in dataset balance affected typi-  
192 cal integration metrics, we examined the  $\text{ARI}_{\text{cell-type}}$  and  $(1 - \text{ARI}_{\text{batch}})$  scores [9]  
193 for each run and method independently.  $\text{ARI}_{\text{cell-type}}$  represents conserved het-  
194 erogeneity of annotated cell-types post-integration, and  $(1 - \text{ARI}_{\text{batch}})$  represents  
195 the degree to which the two batches being integrated overlap post-integration  
196 [9]. As variation between methods was not the main objective of the analysis,  
197 the two scores were Z-score normalized for each method across the perturbation  
198 experiments and the median value was utilized due to the presence of replicates  
199 (Online Methods). Neither the scaled median  $\text{ARI}_{\text{cell-type}}$  or scaled median  $(1 -$   
200  $\text{ARI}_{\text{batch}})$  indicated distinct patterns for the perturbation experiments (Figures  
201 2D, 2E). In fact, there seemed to be a high degree of method-specific variation  
202 in these results, making the interpretation challenging. In terms of the median  
203  $(1 - \text{ARI}_{\text{batch}})$  scores, for 4 out of 5 methods the top score occurred in the  
204 control setup which shows that imbalance leads to worsening performance in terms  
205 of batch-mixing (Figure 2E). The results for cell-type heterogeneity were even  
206 less clear and indicated differences based on both the method utilized and cell-  
207 type downsampled/ablated. Overall, the results did not contain clear patterns  
208 and point to the fact that global clustering metrics do not account for dataset  
209 balance and may not be adequate for assessing performance in scenarios with  
210 imbalanced datasets and rare cell-types.

211 To overcome this limitation of global metrics, we examined integration per-  
212 formance at a cell type-specific level through a k-nearest-neighbor (KNN) clas-  
213 sifier [16, 32] that was trained on 70% of the post-integration embeddings from  
214 each method independently and the remaining 30% was used as a test-set for  
215 cell-type classification. The train/test split was stratified by cell-type label,



**Figure 2: Perturbation analysis of controlled PBMC dataset and effects on cell-type-specific integration.** (a), (b) The cell-type and batch representations of the balanced two-batch PBMC dataset.(c) The perturbation setup for the balanced PBMC data - in each iteration, one batch and one cell-type is randomly selected, and the cell-type is randomly either downsampled to 10% of its original number or ablated. Control experiments are also performed where no downsampling occurs. (d), (e) Z-score normalized median  $ARI_{cell\_type}$  (cell-type integration accuracy) (d) and median  $(1-ARI_{batch})$  (batch mixing) (e) results across experiment type (control, cell-type downsampling, cell-type ablation), specific-cell-type downsampled, and integration method utilized. (f) KNN-classification within the integrated embedding space in control, downsampling and ablation replicates and across methods. The F1-scores are indicated for the same cell-type that was downsampled. (g) Hierarchical clustering of similar cell-types in the balanced two-batch PBMC data. (h) Cell-type-specific integration results using a KNN-classifier after hierarchical clustering across perturbation experiments with the same setup as (f). The cell-types here are based on the label after hierarchical clustering from (g).

such that an equal proportion of cell-types occurred in both subsets, allowing for comparison of classification at the cell-type level (Online Methods). Overall, the classification results provide evidence for cell-type-specific effects of dataset imbalance, as downsampling a specific cell-type led to a statistically significant decrease in the KNN classification F1-score [33] for the same cell-type post-integration, based on an analysis-of-variance (ANOVA) model [34] (ANOVA *p*-value << 0.05, F-statistic = 1304.96, Supplementary Figure S1) (Figure 2F). This result is method agnostic as the ANOVA test factored in method utilized and cell-type downsampled (Online Methods). The only cell-type that exhibited stability were B cells (Figure 2F - standard deviation of median F1-score across methods and experiment types < 0.01). Comparing the cell-type-specific results with the global ARI metrics, we found weak correlation across all methods (Supplementary Figures S2, S3, Spearman's  $\rho \leq 0.4$  across methods and metrics). The uniformly worsening F1 classification scores for the majority of cell-types being perturbed, when compared with the global ARI metrics, shows that global metrics may not adequately capture the integration performance in imbalanced settings. Instead, cell-type-specific metrics such as the KNN-classification score can capture more granular information.

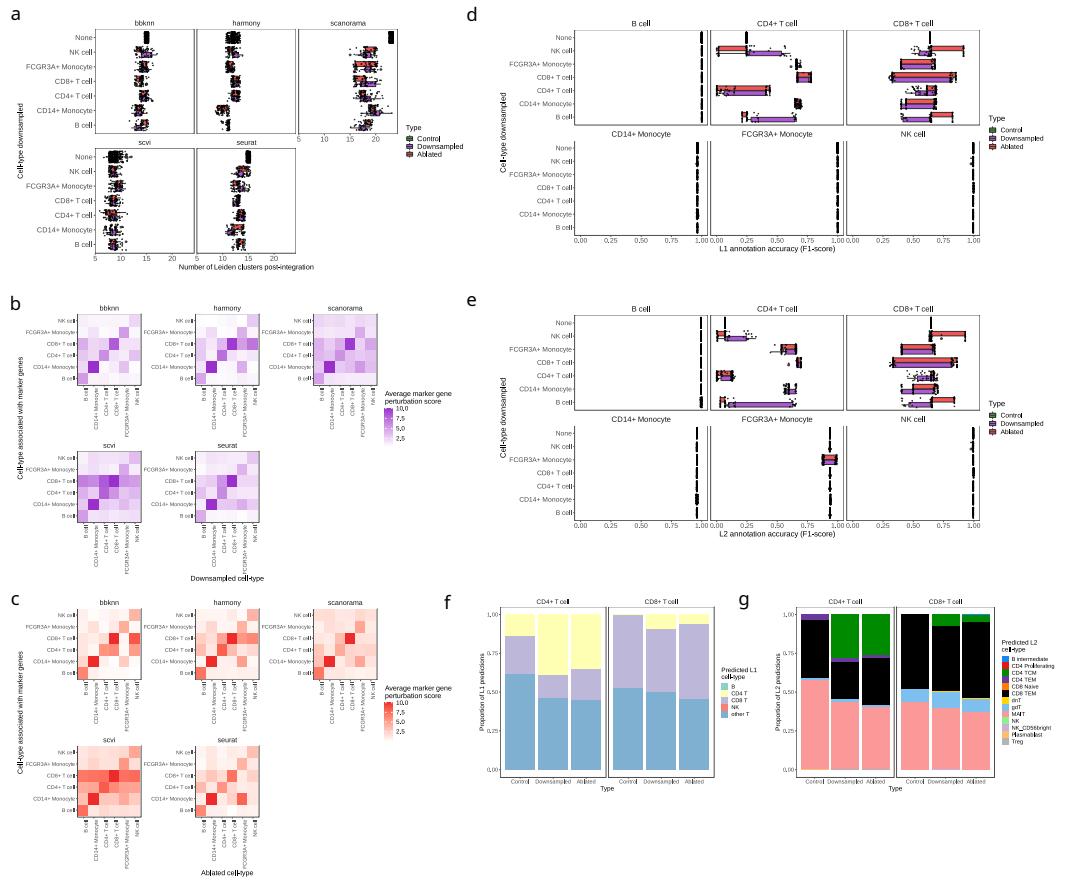
We hypothesized that the integration performance for B cells was unaffected in the perturbation experiments because they are highly distinct from the other cell-types. The two monocyte subsets (CD14+ Monocytes and FCG3RA+ Monocytes) are transcriptionally similar, and the two T-cell subsets (CD4+ T cells and CD8+ T cells) and NK cells are also very similar (Figure 2A, Supplementary Figure S13). As a test, we performed hierarchical clustering of the cell-types into three higher-level subsets - B cells, Monocytes, and NK/T cells (Figure 2G) (Online Methods). As expected, downsampling these subsets did not result in worsening performance to the same degree as the base cell-types (Figure 2F, Figure 2H) (ANOVA F-statistic = 374.46 (hierarchical) << 1304.96 (base), Supplementary Figure S1) (Online Methods). This initial result on the balanced PBMC cohort indicates that the relative transcriptomic similarity of cell-types can drive cell-type-specific performance of integration techniques when considering differing levels of dataset imbalance.

248           **II. Biological interpretation of integration results is con-**  
249           **tingent on relative cell-type proportions between batches**

250           To further analyze the impact of the perturbation experiments on the balanced  
251           PBMC cohort, we quantified the effects of imbalance on downstream analyses  
252           typically performed after integration, including unsupervised clustering, differ-  
253           ential gene expression/marker gene selection, and query-to-reference annotation  
254           (Figure 1A). As we observed significant impacts on KNN-based cell-type classi-  
255           fication in the same setting, it is likely that the impacts of imbalance on integra-  
256           tion may also affect other aspects of single-cell analysis. Therefore, we utilized  
257           the same perturbation setup and downsampling experiments in the balanced  
258           PBMC cohort to analyze these effects.

259           **Stability of unsupervised clustering of samples post-integration**

260           We observed a significant variation in the inferred number of clusters after  
261           integration across all tested methods due to perturbation of cell-type balance  
262           (ANOVA  $p << 0.05$ , F-statistic = 990.79) (Figure 3A). After integration in  
263           both balanced and perturbed simulations, clustering was performed using the  
264           Leiden clustering algorithm with a fixed resolution (Online Methods). Although  
265           all methods indicated at least some degree of variation in the number of clus-  
266           ters between control and downsampled/ablation experiments, there were also  
267           method-dependent effects present (Figure 3A). For instance, while Harmony  
268           exhibited variation in the number of clusters regardless of cell-type downsam-  
269           pled, ablation of CD14+ Monocytes specifically led to a much smaller number  
270           of clusters overall post-integration (Figure 3A). A similar effect was observed for  
271           Seurat and BBKKN, while Scanorama's post-integration clusters diverged most  
272           from the control experiments when ablating CD4+ and CD8+ T cells (Figure  
273           3A). scVI's post-integration clustering results were relatively more stable after  
274           perturbation (Figure 3A). There was variation observed for the control experi-  
275           ments across methods as well, but clear deviation after perturbation was present  
276           in all tested methods. This result indicates that differing levels of imbalance  
277           can cause significant deviations in cluster number, even though the number  
278           of clusters should be stable as the number of cell-types across all batches re-  
279           mains the same in both perturbed and unperturbed experiments. As cell-types  
280           are typically annotated using unsupervised clustering and subsequent marker  
281           gene analysis [8, 17], varying degrees of dataset imbalance can lead to distinct  
282           biological conclusions.



**Figure 3: Quantification of the effects of perturbation-induced dataset imbalance on downstream analyses.** (a) After integration of the PBMC balanced dataset in different perturbation scenarios (type) and based on the cell-type downsampled, the number of unsupervised clusters from the results of each method based on Leiden clustering across replicates. (b) The average marker gene ranking change in differential gene expression (average marker gene perturbation score) for cell-types downsampled and marker gene sets of specific cell-types, across methods. The rankings are averaged across replicates for the ‘downsampled’ experiment type. (c) The average change in marker gene ranking in differential gene expression averaged across replicates for the ‘ablation’ experiment type. (d), (e) The cell-type-specific L1 annotation (coarse-grained) (d) and L2 annotation (fine-grained) (e) accuracy scores across replicates for query-to-reference results for individual batches based on experiment type (control, downsampling, ablation) and cell-type downsampled. (f), (g) The L1 predictions (f) and L2 predictions (g) by proportion across experiment types and replicates for CD4+ T cells and CD8+ T cells.

283 An important caveat to this result is that a reduced number of total cells  
284 (through downsampling or ablation) might lead to less clusters in general at a  
285 fixed resolution due to less overall heterogeneity in the data. We argue that this  
286 is not a major limitation due to two factors - (1) This case is still reflective of the  
287 effects of perturbing the cell-type balance, even if the effects are uniform across  
288 cell-types downsampled, (2) We did not observe a uniform reduction in cluster  
289 number based on the integration methods utilized. For example, scVI's results  
290 for cluster number were fairly stable after downsampling or ablation, while the  
291 results from Scanorama indicated a drastic reduction after perturbation (Figure  
292 3A). Moreover, within each method, the results for reduction in cluster number  
293 were not uniform based on the cell-type that was downsampled (Figure 3A).  
294 Therefore, although we are limited in evaluation at a fixed clustering resolution,  
295 the results nevertheless show that the perturbation setup can lead to cell-type  
296 and method-specific effects that can potentially alter further analyses.

## 297 Differential gene expression and marker gene stability

298 Frequently, the next step after integration and unsupervised clustering in  
299 a scRNA-seq analysis workflow is differential gene expression analysis [8, 35].  
300 Typically, a series of one-versus-all differential expression experiments, using  
301 statistical tests such as the non-parametric Wilcoxon Rank-Sum Test or more  
302 RNA-seq specific techniques such as DESeq2, is done for each cluster to deter-  
303 mine the top ranking “marker genes” specific to all clusters [8, 17, 35]. These  
304 marker genes are indicative of cell-type identity for each cluster and are used  
305 to annotate clusters into putative cell-types [8, 17, 35]. One way to assess  
306 marker gene stability before and after perturbation is to constrain the number  
307 of clusters to be equivalent across simulations, but this would be unrealistic  
308 as variation in cluster number in both control and perturbed experiments was  
309 observed across methods (Figure 3A). As the *ranking* of marker genes is typi-  
310 cally utilized to annotate clusters from scRNA-seq data [8, 35], we considered  
311 deviation in ranking for genes with known cell-type associations to be an im-  
312 portant end-point. Using the unintegrated data separately for each batch, we  
313 determined the top 10 marker genes for each cell-type, and assessed the stability  
314 of their ranking before and after perturbation (Online Methods). Changes in  
315 ranking for marker genes across replicates for a given subset of experiments were  
316 defined as the *marker gene perturbation score*, indicating the standard deviation  
317 of the rank (Online Methods). In the case of examining all marker genes for a  
318 given cell-type, the standard deviation of ranking of all of the marker genes was

319 averaged, and this is indicated as the *average marker gene perturbation score*  
320 (Online Methods).

321 For the majority of marker genes, we observed deviations in ranking after  
322 downsampling and ablation, with many diverging as much as 10 ranks (Sup-  
323 plementary Figure S4 - Marker gene perturbation score) which could lead to  
324 significant changes in biological interpretation of results if the top 10 marker  
325 genes are used as a heuristic for annotation. An ANOVA test factoring in the  
326 specific marker gene, method, and downsampled cell-types indicated that per-  
327 turbation led to statistically significant changes in ranking (ANOVA  $p << 0.05$ ,  
328 F-statistic = 48.99 - highest of all factors) (Online Methods). There was strong  
329 correlation in marker gene perturbation across methods, with the exception of  
330 scVI, which exhibited stronger deviations in rankings for some marker genes  
331 (Supplementary Figure S4).

332 Next, we examined whether downsampling or ablation of a specific cell-type  
333 will change the ranking of marker genes for the same cell-type, and we observed  
334 that this was the case across all methods (Figure 3B, Figure 3C). The strongest  
335 ranking change of marker genes occurred after downsampling or ablation of  
336 CD8+ T cells and CD14+ Monocytes (Figure 3B, Figure 3C). As these two  
337 cell-types are highly similar to CD4+ T cells and FCGR3A+ Monocytes re-  
338 spectively, downsampling likely induces a collapse of cells in the downsampled  
339 cell-types into clusters corresponding to their neighboring cell-types. This con-  
340 tributes to significant deviations in marker gene ranking and possible changes in  
341 biological interpretation in both the downsampling and ablation experiments.  
342 Significant changes in marker gene ranks were also observed for cell-types that  
343 were not downsampled or ablated, such as in NK cells, which were pronounced  
344 for Harmony and scVI results (Figure 3B, Figure 3C). Once again, this is likely  
345 due to mixing of cell-types within clusters after an imbalance is introduced,  
346 as NK cells are very transcriptionally similar to CD4+ and CD8+ T cell sub-  
347 sets. Similarity of cell-types and effects in integration are investigated further  
348 in Results III. and IV.

349 To more definitively determine whether or not these perturbations in marker  
350 gene rankings could change the biological conclusions of an analysis, we per-  
351 formed a case-study with clusters that contained a majority of CD4+ and CD8+  
352 T cells after integration using Seurat (Supplementary Figure S16). Considering  
353 a permissive list of the top 50 marker genes and using canonical markers for  
354 CD4+ and CD8+ T cells, we observed that the fraction of clusters annotated

355 as either CD4+ or CD8+ T do change after downsampling/ablation induced  
356 imbalance is introduced (Supplementary Figure S16).

357 **Query-to-reference projection and cell-type annotation**

358 With the increasing availability of public scRNA-seq datasets with high  
359 quality annotations, query-to-reference annotation has become a major appli-  
360 cation for scRNA-seq data integration [17]. However, the accuracy of annotation  
361 depends on the quality of the integrated space. To examine the effects of im-  
362 balance in this setting, we utilized the Seurat 4.0 query-to-reference annotation  
363 pipeline and a large-scale multi-modal PBMC dataset of 211 000 cells as a ref-  
364 erence [31]. In the Seurat 4.0 pipeline, each batch (query) is projected to the  
365 reference dataset, such that the integration is performed individually for each  
366 batch [31] (Online Methods). In this setup, the effects of inter-batch imbalances  
367 are not relevant, but only the imbalance relative to each query batch and the  
368 reference dataset. In this setting, the perturbations were done for the query  
369 batches (balanced PBMC 2 batch data), and the reference was static (Online  
370 Methods). We assessed the accuracy of query-to-reference projection through a  
371 “fuzzy-matching” of cell-type labels between the balanced PBMC batches and  
372 multi-modal PBMC reference from Seurat [9, 31] (Online Methods).

373 The majority of cell-types were stable across control and downsampling/ablation  
374 experiments with near perfect scores. However, the two T-cell subsets had vary-  
375 ing performance to a high degree, regardless of which cell-type was downsam-  
376 pled or ablated (Figure 3D, Figure 3E). This result is indicative of the fact  
377 that the imbalance between the projected batch (which was perturbed) and the  
378 reference dataset (held constant across all experiments) is driving variance in  
379 integration and subsequent annotation results. This highlights a similar prob-  
380 lem concomitant with previous results, in that perturbing the degree of balance  
381 for transcriptionally similar cell-types can lead to biologically distinct results  
382 compared to the balanced scenario. In this case, the CD4+ T cell and CD8+  
383 T cell populations are transcriptionally similar, and a trade-off in their anno-  
384 tation performance can be observed in the control unperturbed data (Figure  
385 3D, Figure 3E). After perturbing the degree of balance within a given batch,  
386 the trade-off point is moved in favor of either subset (Figure 3D, Figure 3E).  
387 Further, the result highlights that perturbation of dataset balance can affect  
388 downstream results even when there is a degree of imbalance already present  
389 between integrated datasets, which was the case between the query and refer-  
390 ence data in these experiments (Supplementary Table 3, Supplementary Table

391 9).

392 Examining the cell-type annotations more closely at two levels of resolution,  
393 we observed that both the CD4+ and CD8+ T cells were largely mis-annotated  
394 as Mucosal Associated Invariant T-Cells (MAIT) (Figure 3F, Figure 3G). After  
395 downsampling or ablation of a given cell-type and subsequent analysis of anno-  
396 tation accuracy of the same cell-type, we find that CD4+ T cells were annotated  
397 more accurately, while CD8+ T cells were further mis-annotated, compared to  
398 their respective control scores (Figure 3F). The transcriptional similarity be-  
399 tween not just the CD4+/CD8+ subsets, but the many subsets that fall under  
400 “other T”, is a challenging problem for integration and subsequent label-transfer  
401 [36]. This challenge is potentially exacerbated when imbalance is present, as  
402 indicated by the perturbation experiments and their effects on the annotation  
403 results.

404 Overall, cell-type imbalance affected all three major aspects of downstream  
405 analysis that were tested, and we observed strong evidence of impact on biolog-  
406 ical interpretation of the results. This observation is likely even more relevant  
407 in complex datasets, as the balanced PBMC cohort is not representative of the  
408 ever-increasing throughput of current scRNA-seq protocols [37]. The limita-  
409 tions of the reference dataset utilized may also be a major source of variation  
410 in the query-to-reference integration results. It may be the case that a more  
411 suitable reference may not lead to high variance in the results of the two T-cell  
412 subsets, however assessment and selection of reference datasets is outside the  
413 scope of this study and the multi-modal PBMC reference used is one of the  
414 most comprehensive single-cell references to date.

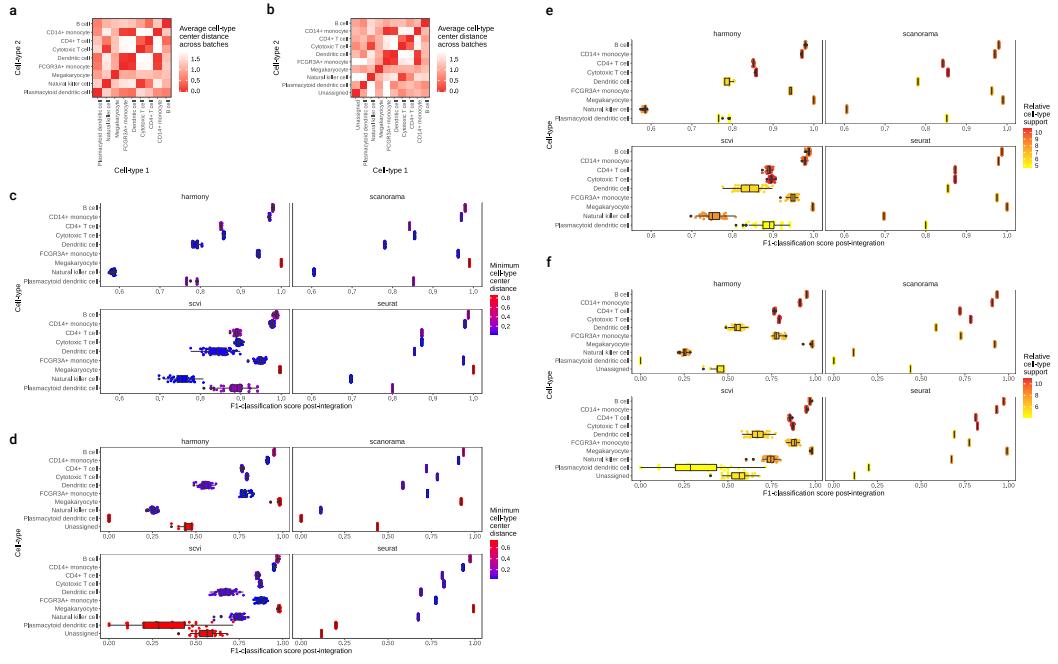
### 415 III. Analysis of imbalanced complex datasets reveals key 416 metrics for stability of integration results

417 While perturbation experiments of the balanced two batch PBMC cohort re-  
418 vealed the effects of dataset imbalance on integration and downstream analyses  
419 in a controlled setting, current scRNA-seq datasets typically involve a much  
420 larger number of cells and cell-types captured [37]. Therefore, we examined the  
421 effects of dataset imbalance when integrating complex datasets with multiple  
422 samples that are not perturbed, but already have inherent cell-type imbalance  
423 between samples. To this end, we analyzed an imbalanced 2 batch PBMC co-

424 hort [18], imbalanced 4 batch PBMC cohort [18], imbalanced cohort of 6 batches  
425 of mouse hindbrain developmental data [19], and an imbalanced heterogeneous  
426 tumor cohort of 8 batches of pancreatic ductal adenocarcinoma (PDAC) data  
427 [20] (Online Methods). No downampling was done in these experiments and  
428 we aimed to analyze the effects of integration on cell-types that are imbalanced  
429 with respect to others, both within and across batches.

430 As determined in the analysis using the balanced PBMC data and pertur-  
431 bations, transcriptomic similarity (relatedness) of cell-types and cell-type im-  
432 balance are two important factors that impact downstream results. We sought  
433 to observe if these properties also led to differences in integration performance  
434 per-cell-type in more complex datasets without perturbations. We formalized  
435 these two properties as the *relative cell-type support* and *minimum cell-type*  
436 *center distance*, quantifying the degree of imbalance and relatedness to other  
437 cell-types (Online Methods). The *relative cell-type support* is defined as the  
438 number of cells specific to a cell-type present across all batches, and the *mini-*  
439 *mum cell-type center distance* considers the average distance across all batches  
440 between cell-types in a principal component analysis (PCA) dimensionality re-  
441 duction representation space and selects the distance of the closest neighboring  
442 cell-type for each cell-type (Online Methods).

443 To correlate these properties with integration performance, we used the same  
444 KNN-classification setup as before to determine performance on a per-cell-type  
445 basis (Online Methods). We started by analyzing the cell-type center distances  
446 on the imbalanced PBMC 2 and 4 batch datasets. Examining the average cell-  
447 type center distances - which were averaged across all batches if the cell-types  
448 were present in more than one batch - there is a clear pattern evident between  
449 cell-types that were observed before in the 2 batch balanced PBMC data (Figure  
450 4A, Figure 4B). NK cells have small relative distance between the T cell  
451 subsets, while dendritic cells share transcriptional similarity with the monocyte  
452 subsets (Figure 4A, Figure 4B). B cells and megakaryocytes have the great-  
453 est distance between the rest of the cell-type centers (Figure 4A, Figure 4B),  
454 and thus we expect these cell-types to have strong performance in integration  
455 which was previously observed for B cells. This pattern did hold when examin-  
456 ing integration performance through KNN-classification for both datasets on a  
457 per cell-type basis compared to their *minimum cell-type center distance* across  
458 batches, as Megakaryocytes and B cells had strong performance regardless of  
459 integration technique utilized (Figure 4C, Figure 4D).



**Figure 4: Factors in imbalanced complex datasets predictive of altered integration and downstream results.** (a) The average cell-type center distance across cell-types in the imbalanced PBMC 2 batch dataset. For each batch, the distance from the centers of cell-type clusters in principal component analysis (PCA) reduction space are calculated, and the relative distances between cell-types are determined and averaged across batches. (b) The average cell-type center distance across cell-types in the imbalanced PBMC 4 batch dataset. (c), (d) Comparison of F1-classification accuracy of each cell-type in the imbalanced PBMC 2 batch dataset (c) and imbalanced PBMC 4 batch dataset (d), specific to method and across replicates, compared with the *minimum cell-type center distance* value. (e), (f) Comparison of F1-classification accuracy of each cell-type in the imbalanced PBMC 2 batch dataset (e) and imbalanced PBMC 4 batch dataset (f), across methods and replicates, compared with the *relative cell-type support* value. The *relative cell-type support* is based on the number of cells in the integrated embedding space present for each cell-type.

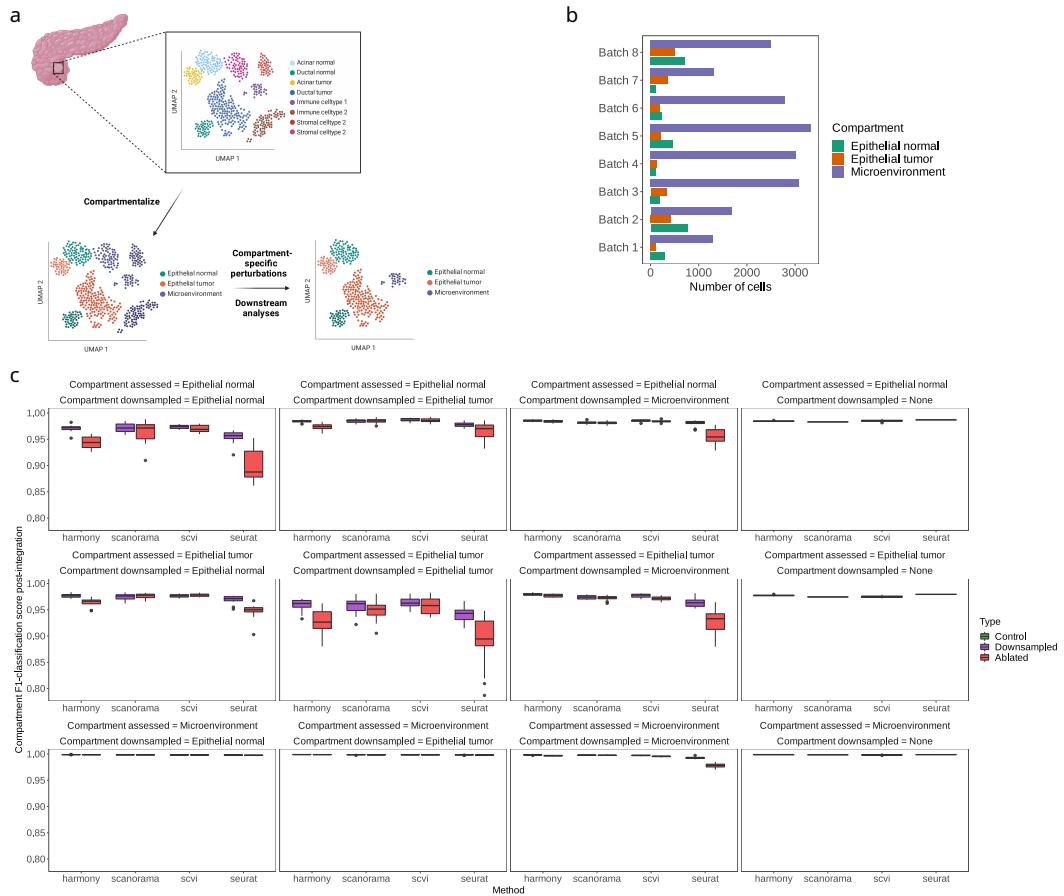
460        However, the results were not straightforward for other cell-types. Exam-  
461        ining the plasmacytoid dendritic cells, we expected strong performance due to  
462        their high relative distance between other cell-types, but this was not the case  
463        for both the 2 and 4 batch datasets (Figure 4C, Figure 4D). Although these cells  
464        have a large relative distance between other cell-types, they occur in a much  
465        smaller number compared to others (Figure 4E, Figure 4F). We quantified *rel-*  
466        *ative cell-type support* as the log-transformation of the total number of cells  
467        for each cell-type across batches (Online Methods), and plasmacytoid dendritic  
468        cells have the lowest value across cell-types within these two datasets (Figure  
469        4E, Figure 4F). This result indicates that *minimum cell-type center distance* is  
470        necessary but not sufficient to explain variation in integration results on a per-  
471        cell-type basis. Overall, higher *relative cell-type support* does seem to lead to  
472        higher performance in integration for some cell-types, such as the CD14+ and  
473        CD16+ Monocyte subsets, but is also not sufficient for higher performance, as  
474        NK cells perform poorly across integration techniques due to having a low *min-*  
475        *imum cell-type center distance* and overlap with the T-cell subsets despite not  
476        having low *relative cell-type support*. Examining the results across both met-  
477        rics and performing an ANOVA test to determine variance in scores explained  
478        by the metric, we did find statistically significant associations for both *mini-*  
479        *mum cell-type center distance* and *relative cell-type support* [ANOVA *p*-value  
480        << 0.05 for both metrics across all datasets, with the exception of the mouse  
481        hindbrain 6 batch dataset, for which *relative cell-type support* is non-significant  
482        (Supplementary Table 1) (Online Methods)].

483        Analysis of the 6 batch mouse hindbrain developmental data and 8 batch  
484        PDAC datasets indicated similar results, albeit much less easily interpretable  
485        due to the presence of a very large number of cell-types and more batches (Sup-  
486        plementary Figures S5-S8). Cell-types in close proximity within an embedding  
487        space have an interpretable explanation for poor integration performance, as  
488        they may collapse and become merged with their overlapping counterpart in  
489        the integration step. This was observed in Results sections I and III. Low  
490        cell-type support leads to less data for a given cell-type that an integration  
491        method/model can utilize, and therefore models may not be able to learn the  
492        correct embedding for these cell-types.

493           **IV. Perturbation analysis in PDAC samples reveals tumor  
494           compartment-specific effects of dataset imbalance**

495           To further analyze the effects of dataset imbalance in realistic scenarios, we  
496           considered the pancreatic ductal adenocarcinoma (PDAC) dataset of 8 batches  
497           comprising tumor samples across 8 different biopsies [20]. One major challenge  
498           in the analysis of PDAC data is accurate annotation of tumor cells, and be-  
499           ing able to separate these from normal non-cancerous epithelial cells [38, 39].  
500           As both acinar and ductal epithelial cells have been proposed as cell of origin  
501           candidates in PDAC across numerous studies [40, 41], reliably classifying tu-  
502           mor cells from these normal epithelial cell-types in scRNA-seq data remains a  
503           major computational challenge. Given this difficulty, we sought to determine  
504           if different levels of imbalance between epithelial normal and epithelial tumor  
505           compartments can influence the accuracy of PDAC tumor tissue integration  
506           and subsequent classification of tumor cells. As scRNA-seq data from tumor  
507           tissue is often integrated across multiple biopsy sites, patients, and cohorts [42],  
508           the ability to reliably quantify tumor cells is imperative to the biological valid-  
509           ity of subsequent downstream analyses. We preprocessed and annotated tumor  
510           cells in the PDAC samples through integration, unsupervised clustering, and  
511           marker gene-based annotation (Online Methods). In setting up the perturba-  
512           tion experiments, we grouped epithelial normal cells (acinar and ductal) into an  
513           “epithelial normal” compartment, tumor cells into “epithelial tumor” compart-  
514           ment and the remaining microenvironment cells into the “microenvironment”  
515           compartment (Figure 5A) (Online Methods). Overall, the microenvironment  
516           heavily outnumbered the epithelial tumor and epithelial normal populations  
517           (Figure 5B), which is reflective of the low tumor purity typical of PDAC biopsy  
518           samples [20]. Perturbation experiments included downsampling or ablation of  
519           a randomly selected compartment within 4 randomly selected batches out of  
520           8 (Figure 5A, Figure 5B). We also performed replicates for control, downsam-  
521           pling, and ablation, for a total of 200 simulations for integration of all 8 batches  
522           (Online Methods).

523           Examining the KNN-classification scores on a per-compartment assessed,  
524           per-compartment downsampled basis indicated that downsampling or ablating  
525           the microenvironment compartment leads to stable compartment classification  
526           across all methods, with a slight decrease in performance observed for Seu-  
527           rat in the epithelial normal and tumor compartments (Figure 5C). This result  
528           is concordant with previous analysis indicating that proximity is a key factor



**Figure 5: Compartment-wise perturbation experiments for 8 batches of PDAC samples.** (a) Overview of the experimental setup. To determine the effects of dataset imbalance across tumor compartments, various microenvironment tumor were collapsed into the ‘microenvironment’ compartment, normal ductal and acinar cells into the ‘epithelial normal’ compartment, and malignant ductal and acinar cells into the ‘epithelial tumor’ compartment. The perturbation experiments involved the sample downsampling (10% of a compartment) and ablation (complete removal of a compartment) setup for 4/8 randomly selected batches. Note that all batches are integrated at once using each method. (b) Number of cells in each compartment after cell-type collapse, across batches/biopsy samples in the PDAC data. (c) F1-classification score for KNN classification post-integration, specific to each compartment when compared with the compartment that was downsampled or ablated, across replicates and methods utilized for integration.

529 that dictates the degree to which perturbations in cell-type balance can af-  
530 fect integration results. The *minimum cell-type center distance* in the PDAC  
531 data shows that acinar and ductal cells, which comprise the epithelial normal  
532 and epithelial tumor populations, are two of the most distant cell-types from  
533 others in the data (Supplementary Figure S8, Supplementary Figure S15)).  
534 Similar to the discrepancy between  $\text{ARI}_{\text{cell-type}}$  and  $(1 - \text{ARI}_{\text{batch}})$  observed in  
535 the integration of the balanced and imbalanced PBMC datasets, we observed  
536 that higher  $\text{ARI}_{\text{compartment}}$  based on downsampling of the microenvironment  
537 did not lead to higher batch mixing scores (Supplementary Figure S9, Supple-  
538 mentary Figure S10). In fact, we observed the opposite effect almost uniformly  
539 across all methods, as downsampling the epithelial compartments decreased the  
540  $\text{ARI}_{\text{compartment}}$  and increased  $(1 - \text{ARI}_{\text{batch}})$  (Supplementary Figure S9, Supple-  
541 mentary Figure S10). Downsampling the microenvironment had the opposite  
542 effect. As the microenvironment is quite large, downsampling likely leads to  
543 decreases in batch mixing scores because these metrics are driven by more  
544 prevalent compartments/cell-types. Epithelial cells and their tumor/normal  
545 dichotomy is of strong interest in analyzing PDAC data, and therefore batch  
546 mixing is likely a poor quantifier of integration performance and the biological  
547 validity and utility of the results. This result also reiterates the limitation of  
548 global clustering metrics that do not take into account less prevalent cell-types  
549 and their overall difficulty in interpretation, as the increased performance in  
550  $\text{ARI}_{\text{compartment}}$  after downsampling the microenvironment was not concordant  
551 with the KNN-classification results (Supplementary Figure S9, Figure 5C).

552 Tumor and normal epithelial compartment KNN-classification scores wors-  
553 ened as either compartment was downsampled (Figure 5C). More specifically,  
554 downsampling either the epithelial normal or epithelial tumor compartments led  
555 to the greatest decrease in the integration performance of the same compartment  
556 through the KNN-classification setup (Figure 5C) (ANOVA F-statistic<sub>Normal epithelial</sub>,  
557 F-statistic<sub>Tumor epithelial</sub> > F-statistic<sub>Microenvironment</sub> - Supplementary Figure S11)  
558 (Online Methods). This indicates that relative proportions of tumor to normal  
559 cells can lead to differing results in integration performance in this setting,  
560 which is reflective of the earlier observations in the balanced PBMC dataset  
561 with highly similar populations, such as the NK and T cell subsets. Overall,  
562 these results demonstrate that the degree of imbalance between the similar com-  
563 partments across tumor tissue cohorts can significantly affect the downstream  
564 results and possibly subsequent analyses. This result is not specific to PDAC  
565 data, as tumor samples across cancer types share typical characteristics, but

566 may not have the same compartments or compartment proportions [43]. We  
567 formalize recommendations for the integration of highly imbalanced datasets in  
568 Section VI.

569 **V. Balanced clustering metrics accurately benchmark im-**  
570 **balanced integration**

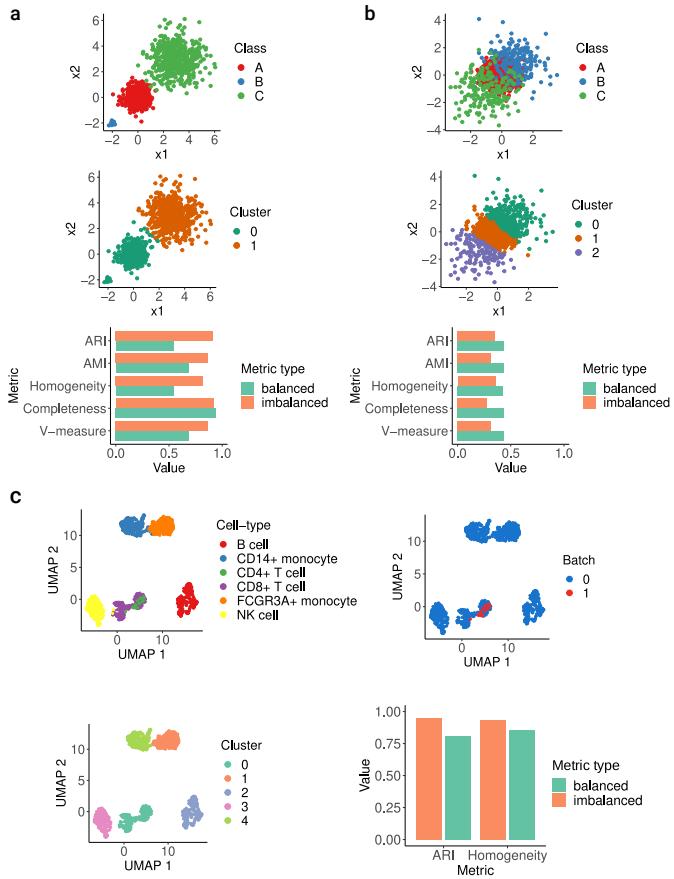
571 Through extensive analysis of both simulated balanced and real-world imbal-  
572 anced scRNA-seq datasets, we have shown that clustering metrics commonly  
573 used in quantifying scRNA-seq integration may be insufficient in imbalanced  
574 contexts. Metrics such as the AMI and ARI are agnostic to information on  
575 label proportions [26, 27] and are thus inadequate for assessing integration per-  
576 formance in imbalanced datasets, which is a common case in single-cell integra-  
577 tion. To overcome limitations of routinely used metrics, we developed balanced  
578 versions of these scores, including the *Balanced Adjusted Rand Index* (bARI),  
579 *Balanced Adjusted Mutual Information* (bAMI), *Balanced Homogeneity*, and  
580 *Balanced Completeness*. Combining *Balanced Homogeneity* and *Balanced Com-*  
581 *pleteness* also allows us to attain the *Balanced V-measure* [28]. These metrics  
582 are robust to dataset imbalance and allow for more nuanced comparisons of  
583 integration results in the aforementioned cases, as they weigh each cell-type  
584 present equally and are not driven by cell-types present in high proportions  
585 (Online Methods).

586 We first demonstrated the utility of the proposed balanced clustering metrics  
587 on simulated data. In the first scenario, we examined a dataset with 3 classes  
588 that are incorrectly clustered into 2 instances using K-means clustering [32]  
589 (Figure 6A) (Online Methods). This scenario can occur in single-cell settings  
590 when a cell-type is highly related to a neighboring cell-type and unsupervised  
591 clustering leads to a collapse of both into the same cluster. As expected, the  
592 base/imbalanced metrics all overestimated the clustering accuracy as they do  
593 not weigh the smaller class (B) as much as the larger classes when assessing the  
594 incorrect assignment (Figure 6A). However, the balanced metrics account for  
595 the incorrect clustering of class B, and indicated worse performance with the  
596 exception of the Balanced Completeness measure. This is an expected result as  
597 Completeness only measures whether all members of a given class are members  
598 of the same cluster [28].

599 In a second scenario, we sought to check whether the balanced metrics can  
600 return higher performance than their base counterparts in the appropriate sce-  
601 nario, and to do so we simulated a dataset where a larger class (A) partially  
602 overlaps two smaller classes (B, C) (Figure 6B) (Online Methods). K-means  
603 clustering with a preset cluster number of 3 slices the larger class in a manner  
604 that the two smaller classes are mostly assigned to the correct cluster/label while  
605 the larger class is split between the 3 clusters (Figure 6B) (Online Methods). In  
606 this setting, the base metrics penalized the results based on the prevalence of  
607 the larger class (A) and the associated mis-clustering, but the balanced metrics  
608 took into account the strong performance on the smaller classes (B, C) and  
609 returned higher scores.

610 These two simulated scenarios demonstrate that balanced metrics can reveal  
611 information not present in typical global clustering scores and benchmark results  
612 in a manner that takes into account class imbalance.

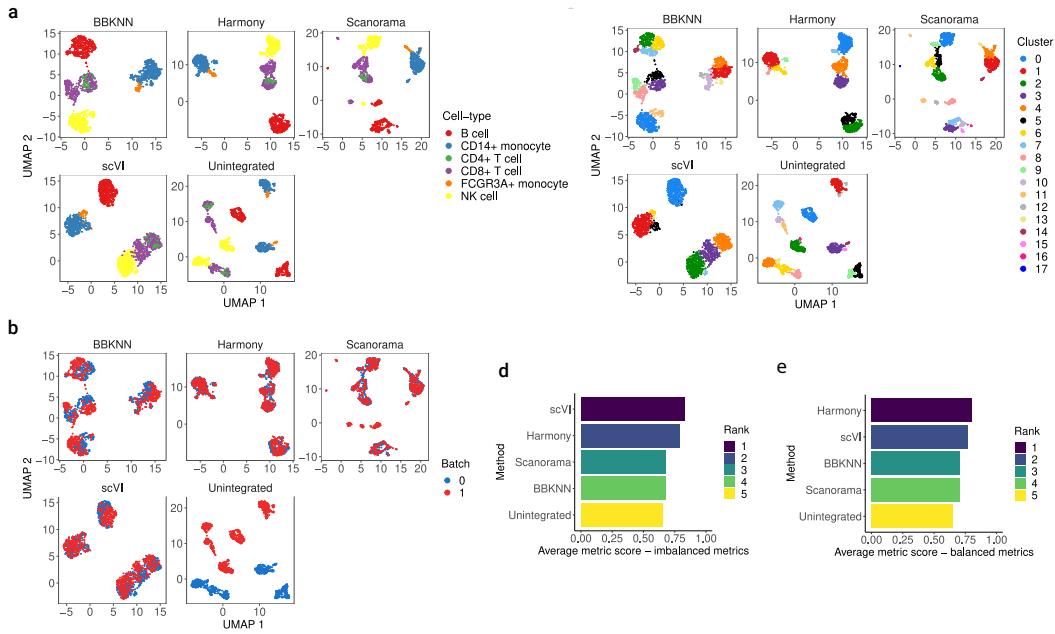
613 We further assessed the applicability of the balanced metrics in single-cell  
614 data by considering the balanced PBMC cohort of two batches (Results I.).  
615 The first test assessed whether or not the first simulated case holds in single-  
616 cell data, as we downsampled CD4+ T cells in one batch in a manner where they  
617 overlapped with CD8+ T cells after clustering (Figure 6C) (Online Methods).  
618 Comparing the balanced and base ARI and Homogeneity scores, we found that  
619 the balanced scores did in fact decrease by a significant margin (Figure 6C).  
620 This is because the balanced metrics are considering the mis-clustering of the  
621 CD4+ T cells in a manner that is weighted equally to the correct clustering  
622 of the other cell-types, even though these cells are present in a much smaller  
623 number overall. To determine if the balanced metrics can change the results  
624 of a benchmarking analysis, we downsampled CD4+ T cells and FCGR3A+  
625 monocytes from one batch in the balanced 2 batch PBMC dataset and per-  
626 formed integration using BBKNN, Harmony, Scanorama, and scVI (Figure 7A,  
627 Figure 7B) (Online Methods). After integration, Leiden graph-based unsuper-  
628 vised clustering [29] was performed and the results of clustering were compared  
629 with ground-truth labels using the base and balanced metrics by considering  
630 an average of their scores across ARI, AMI, Homogeneity, Completeness, and  
631 V-measure values (Figure 7C) (Online Methods). Examining the integration  
632 results, most methods mixed the CD4+ T cells with the CD8+ T Cells and the  
633 FCGR3A+ Monocytes with the CD14+ Monocytes to differing extents, while  
634 also having varying success in integrating the two batches overall (Figure 7A,



**Figure 6: Demonstration of balanced clustering metrics on simulated data and scenarios.** (a) Simulated data of 3 well separated imbalanced isotropic Gaussian classes with imbalance that are incorrectly clustered into two clusters that collapses the smaller class (B) with another. The concordance of the class labels with the clustering result for the base (imbalanced) and balanced ARI, AMI, Homogeneity, Completeness and V-measure for this result are indicated. (b) Simulated data of 3 imbalanced and overlapping isotropic Gaussian classes that are clustered into 3 clusters that mix the larger class (middle - A) with the smaller classes (B, C) and concordance of the class labels with the clustering result for the base (imbalanced) and balanced metrics. (c) Constructed scenario with balanced two batch PBMC single-cell data where a very small subset of CD4+ T cells (10% of original proportion) present in only one batch are incorrectly clustered with CD8+ T cells after integration. In this scenario, the concordance of the unsupervised clustering labels and the ground-truth cell-type labels are indicated for both the base (imbalanced) and balanced ARI and Homogeneity scores.

635       Figure 7C). When using the base metrics and their averaged scores, scVI ranked  
636       the highest and BBKNN ranked the worst. Surprisingly, the base metric scores  
637       for Scanorama and BBKNN, which ranked the worst in this subset, were al-  
638       most the same as using the unintegrated embedding (Figure 7D). Scanorama  
639       and BBKNN have shown strong performance with low variance results for all  
640       of our previous analyses and performed well in comprehensive benchmarking  
641       studies [9, 13], which is not in concordance with this result. When analyz-  
642       ing the result with the balanced metrics, the rankings changed significantly, as  
643       Harmony became the top performer (switched with scVI) and BBKNN now  
644       performed better than Scanorama (Figure 7E). Of particular note is the fact  
645       that there is a larger separation in scores between the unintegrated embedding  
646       and the results of BBKNN and Scanorama using the balanced metrics, and this  
647       result is more valid as we expect the integration methods to perform signifi-  
648       cantly better than an unintegrated baseline. This ranking shift occurred while  
649       the magnitude of the overall scores did not diverge significantly.

650       Lastly, we reexamined the initial uninformative results obtained using the  
651       base ARI<sub>cell-type</sub> scores for the perturbation experiments on the balanced 2 batch  
652       PBMC dataset (Figure 2D). Utilizing the balanced ARI (bARI) instead of the  
653       base metric for calculating ARI<sub>cell-type</sub>, the results indicated more clear/distinct  
654       patterns that reflected both the *relative cell-type support* and *minimum cell-type*  
655       *center distance* properties (Supplementary Figure S12). Specifically, downsam-  
656       pling/ablating B cells did not lead to decreases in the balanced ARI<sub>cell-type</sub>  
657       across all methods, which is in line with the cell-type center distance property  
658       as the B cells are distant from all other cell-types in this dataset (Supplementary  
659       Figure S13). Further, there is a clear pattern of worsening performance in the  
660       balanced ARI<sub>cell-type</sub> scores when the CD4+ or CD8+ T cells are downsampled  
661       or ablated, which is concordance with both the expected results in terms of  
662       *minimum cell-type center distance* (Supplementary Figure S13) and the KNN-  
663       classification results (Figure 2F). Similar results hold for downsampling or abla-  
664       tion of the monocyte subsets, although the scores are more method-specific and  
665       performance decreases are less pronounced (Supplementary Figure S12). Over-  
666       all, the balanced clustering metrics can capture nuances in the data related to  
667       cell-type imbalance in a manner that the base metrics cannot. The balanced  
668       metrics are also more concordant with cell-type specific results, such as the  
669       KNN F1-score, as they weigh classes equally when considering performance.



**Figure 7: Benchmarking single-cell data integration using balanced clustering metrics.** (a) Cell-type values for the balanced two batch PBMC data with FCG3RA+ Monocytes and CD4+ T cells downsampled to 10% of their original proportion in one batch, after integration with the tested methods as well as an unintegrated representation. (b) Batch values for the integrated and unintegrated downsampled two batch PBMC data. (c) Unsupervised clustering results for Leiden clustering in the embedding space of the integrated and unintegrated results for the downsampled two batch PBMC data. (d), (e) Scoring and ranking of integration results, when considering concordance of the unsupervised clustering labels and ground-truth cell-type labels for each integration method and the unintegrated subset, using the average results of the base (imbalanced) clustering metrics (d) (ARI, AMI, Completeness, Homogeneity) and average of the balanced clustering metrics (e) (bARI, bAMI, Balanced Completeness, Balanced Homogeneity).

670

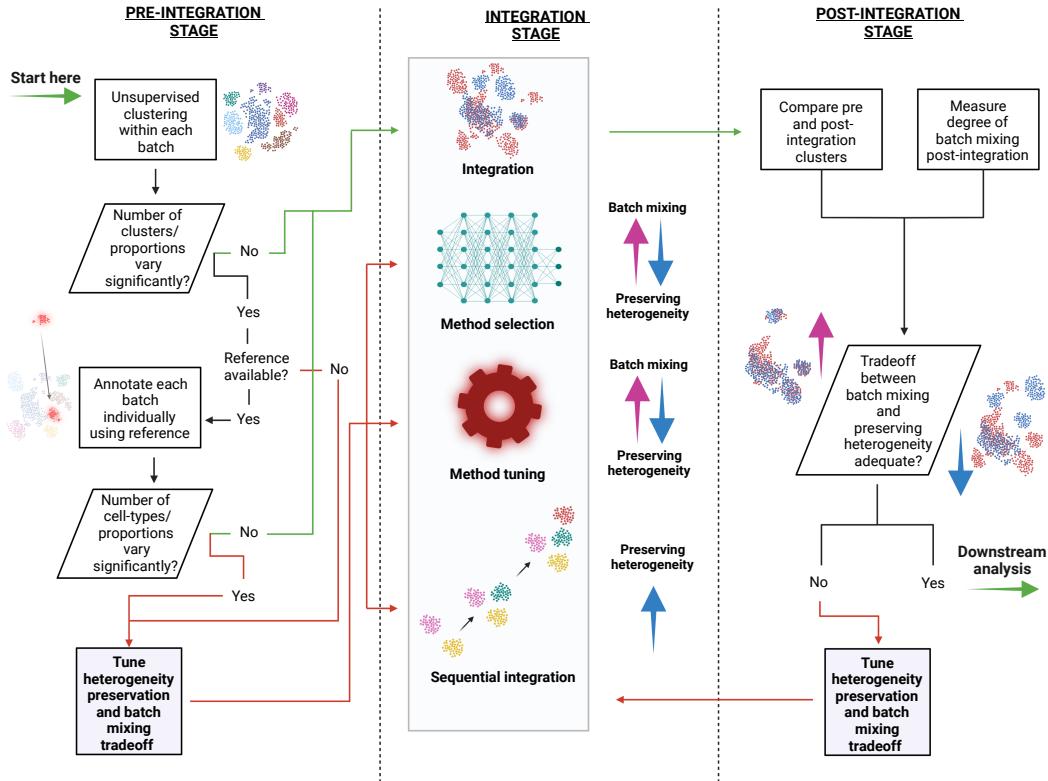
## VI. Guidelines for imbalanced single-cell data integration

671

To aid in the integration of imbalanced datasets, we introduce general guidelines for users of integration techniques (Figure 8, Supplementary Table 2). We note that these guidelines are not meant to be strict rules, but rather suggestive in nature, as scRNA-seq and multi-modal single-cell data from different samples can have very different properties even after taking imbalance into account [4]. The guidelines are method agnostic, as our analysis revealed that all frequently utilized techniques in scRNA-seq integration are susceptible to the outlined effects of dataset imbalance (Results sections I-IV). An important aspect to consider when utilizing these guidelines is prior knowledge on potential disparity in the datasets can help guide the degree of desired batch mixing. For instance, in analyzing heterogeneous tumor samples from distinct patients with disparate cell-types and proportions, biological heterogeneity conservation is likely to be poor if batch-mixing is prioritized in integration [8]. However, this may be a desired result if the end analysis goal is only to assess common variation between the tumor samples and perform downstream analyses such as differential abundance of shared cell-types [8]. Judging the degree of desired batch mixing is often very difficult in practice [8]. As such, we emphasize an iterative process where imbalance, degree of batch-correction, and conservation of biological heterogeneity are assessed at multiple steps in the scRNA-seq integration pipeline (Figure 8).

691

Overall, potential imbalance within datasets to be integrated can be assessed based on pre-integration tests using unsupervised clustering and/or query-to-reference annotation (Figure 8, Supplementary Table 2). The latter would yield a more accurate representation of potential imbalance, but can only be used when a reference dataset is available for the given tissue samples. Unsupervised clustering can be used in any situation as it does not require cell-type labels but can be a noisy readout as clustering is highly sensitive to the technique, parameters, and the underlying data distribution [44]. If either of these outlined pre-integration tests reveals disparity in the datasets, the integration step itself can be altered by: (i) picking an integration method that is suitable for preserving biological heterogeneity - Luecken et al. [13] provide an extensive overview of scRNA-seq and multi-modal integration techniques in this respect and provide selection criteria, (ii) tuning the integration method itself to better preserve biological heterogeneity across datasets - the availability of such parameters will vary based on method [13], and (iii) performing sequential in-



**Figure 8: Guidelines for single-cell integration in imbalanced settings.** A stepwise procedure is outlined, starting with diagnostic tests in the pre-integration stage that dictate whether or not to tune integration methods at the integration stage or perform further steps in the pre-integration stage. After integration, the trade off between batch-mixing and conservation of biological heterogeneity can also be diagnosed, and if determined inadequate, further tuning at the pre-integration and integration stages can be done. Complete details as well as examples of implementations for each recommendation are given in Supplementary Table 2.

706 tegration if the datasets are known or suspected to have temporal structure  
707 [45] (e.g. developmental data) (Figure 8). There is also the possibility of inte-  
708 grating only shared putative cell-types between batches if a reference dataset  
709 is available, as this would better ensure imbalance is minimized in the integra-  
710 tion step (Supplementary Table 2). After the integration step, post-integration  
711 techniques to assess preservation of biological heterogeneity and degree of batch  
712 mixing can be used to determine the current balance between the two desired  
713 outcomes [8], and integration and pre-integration steps can be further tuned  
714 to strike the desired balance (Figure 8). A complete description with specifics  
715 and code implementations, in the R and python programming languages, of the  
716 outlined recommendations are indicated in Supplementary Table 2.

## 717 Discussion

718 In this work, we thoroughly analyzed the effects of dataset imbalance in scRNA-  
719 seq integration scenarios, and its impacts on downstream analyses and over-  
720 all biological conclusions. When the level of imbalance between batches was  
721 perturbed, we observed varying degrees of effects on unsupervised clustering,  
722 neighbor-based cell-type annotation, differential gene expression analysis, and  
723 query-to-reference annotation. More importantly, these effects were not method-  
724 specific, and thus have implications for single-cell data integration overall, where  
725 biological conclusions plausible under one scenario may not be concordant if the  
726 pre-integration data distribution is different due to many possible underlying  
727 factors of variation. These results have significant ramifications for single-cell  
728 data integration, as most datasets being integrated will likely not have a high  
729 degree of shared variation with the increasing complexity of the tissues be-  
730 ing analyzed and higher throughput of current scRNA-seq and multi-modal se-  
731 quencing protocols [37, 46]. We further examined these results on more complex  
732 data and concluded that the potential of dataset imbalance to affect integra-  
733 tion results can be summarized by two key metrics - *relative cell-type support*  
734 and *minimum cell-type center distance*. To aid the integration and subsequent  
735 downstream analyses in scenarios with imbalanced datasets, we introduce sev-  
736 eral guidelines pre-integration, at the integration step, and post-integration, as  
737 well as balanced clustering metrics for more accurate assessment and bench-  
738 marking in such cases.

739        Although single-cell data integration is ubiquitous in current computational  
740        analysis pipelines for both scRNA-seq and multi-modal single-cell sequencing  
741        data, analyzing the nuanced properties and behavior of integration techniques  
742        on different datasets has lagged in lieu of performance-based studies. Extensive  
743        benchmarking studies have been performed for scRNA-seq integration, but these  
744        analyses have largely focused on performance in specific settings, as determined  
745        by batch-mixing and conservation of biological heterogeneity [9, 47]. Some  
746        studies have raised specific concerns towards the impacts dataset imbalance  
747        can have on integration and a few methods have been developed specifically to  
748        address this challenge [10, 36, 48–50], but an extensive analysis on downstream  
749        effects had yet to be performed. Further understanding of the properties of both  
750        the pre-integration and post-integration representation spaces will likely shed  
751        light on gaps in performance between different techniques and the situational  
752        trade offs between batch mixing and conservation of biological heterogeneity.  
753        For example, although anchor-based techniques are used to link both scRNA-seq  
754        and multi-modal datasets in integration [45], the conditions that lead to false-  
755        positive and false-negative (missing) anchors between batches and multi-modal  
756        samples have not been extensively characterized. Such analyses will further the  
757        understanding of the limiting conditions of single-cell data integration, and lead  
758        to better tools, guidelines, and a sounder foundation for downstream analyses  
759        and inference of biological phenomena.

760        In benchmarking single-cell integration techniques, often standardized datasets  
761        do not contain high degrees of cell-type imbalance across batches that may be  
762        encountered in common real-world scenarios such as temporal integration [9,  
763        13, 45]. Therefore, a more principled approach to benchmarking may involve  
764        a stronger focus on these cases and non-trivial datasets such as tumor samples  
765        from multiple-patients and cohorts. The trade-off between batch mixing and  
766        biological heterogeneity conservation is an important research direction, as con-  
767        serving biological signal can be much more complex than what is indicated by  
768        clustering metrics post-integration, particularly if integration is being done on  
769        the entire count matrix and not within an embedding space [13]. In our anal-  
770        ysis, we introduced four novel balanced clustering metrics that can be utilized  
771        to better benchmark integration techniques in imbalanced scenarios. These  
772        metrics are used to analyze clustering results post-integration, but more salient  
773        scores for preserving biological signal such as the conservation of highly-variable  
774        genes introduced by Lueken et al. [13] will also allow for a complete picture  
775        of the potential downstream impacts of integration. As our understanding of

776 the limitations of current integration methods evolves, we envision more comprehensive guidelines that incorporate our analysis on situational integration  
777 setups and utilization of different methods in the correct contexts, as opposed  
778 to a single method for every integration scenario. Our findings and guidelines  
779 can be extended to multi-modal analysis of disjoint samples (e.g. scRNA-seq  
780 and scATAC-seq of similar but distinct tissue samples), but the finer details  
781 of the impacts of imbalance on both joint and separately profiled multi-modal  
782 integration and subsequent analysis remains unknown. An important future  
783 research direction in multi-modal integration is better understanding of integration  
784 results at both the technique and data level, as comprehensive benchmarks  
785 specifically focused on multi-modal data have yet to be completed.  
786

787 Our analysis is limited by the extent of datasets analyzed and methods  
788 tested. We sought to identify the effects of downsampling in a highly controlled  
789 scenario where imbalance was not already present, which was the balanced  
790 PBMC 2 batch dataset, but extrapolating the results of the downsampling  
791 experiments to more complex cases was not straightforward. Thus, only the  
792 cell-type-specific effects in already imbalanced datasets with no perturbations  
793 were examined for complex cases, although we did analyze perturbation of com-  
794 plex and imbalanced PDAC samples. Further, although we included frequently  
795 utilized and best performing scRNA-seq integration techniques based on pre-  
796 vious benchmarking studies [9, 13], we did not include recent methods that  
797 focus specifically on preserving biological heterogeneity when differing cell pop-  
798 ulations are present between samples, such as CIDER [49]. There have also  
799 been strides in this direction in the multi-omic integration space, with tech-  
800 niques such as SCOTv2 [50]. As the aim of this analysis was not to determine  
801 the best performing method, but to analyze the impacts of imbalance on in-  
802 tegration results with frequently utilized methods, we deemed this omission  
803 to be acceptable. However, future method-based benchmarking studies should  
804 feature techniques that have sought to explicitly address the issue of dataset  
805 imbalance and several datasets with a high degree of imbalance present. Lastly,  
806 this analysis focused on scRNA-seq integration and did not incorporate multi-  
807 modal datasets and techniques, and although extrapolation may be possible,  
808 this must be confirmed by future work addressing integration when imbalance  
809 across jointly and separately profiled multi-modal datasets is present.

810        **Online Methods**

811        **1 Dataset preprocessing**

812        **1.1 Preprocessing and normalization**

813        All datasets utilized in the study were preprocessed using a uniform pipeline.  
814        Datasets were only further processed if it was clear that no filtering was done  
815        on the raw scRNA-seq data, such as removal of low quality cells and genes [8].  
816        If it was indicated that no filtering was done, quality-control (QC) metrics were  
817        calculated using the *Scuttle* R package (version 1.4.0) [51], including cells with  
818        the most genes having low counts, cells with a high percentage of mitochondrial  
819        genome content, and cells with a low library size (total number of reads overall).  
820        An approach recommended by Amezquita et al. [8] was taken, and the cells and  
821        features with values more than 3 median absolute deviations (MADs) for two  
822        out of three criteria were filtered out. As normalization and log-transformation  
823        need to be tuned specific to the method being utilized and are done in the  
824        integration pipeline where necessary, these were not done in the preprocessing  
825        steps.

826        Datasets were split and saved as individual batches in *h5ad* format, and the  
827        *scanpy* library (version 1.8.2) [52] was used for all further downstream processing  
828        within the integration pipeline, including total-count per cell normalization,  
829        log1p transformation, and highly variable gene selection [52]. These steps were  
830        carried out uniformly for each method tested, with the exception of scVI, as the  
831        technique must utilize the raw counts [24]. Therefore, total-count normalization  
832        and log1p transformation were not done for scVI.

833        All scRNA-seq datasets utilized were preprocessed in this manner, including  
834        the balanced 2 batch PBMC dataset [15][14][9], imbalanced 2 batch PBMC  
835        dataset [18], 4 batch PBMC dataset [18], 6 batch mouse hindbrain development  
836        dataset [19] and 8 batch pancreatic ductal adenocarcinoma (PDAC) dataset  
837        [20].

838        The **ground-truth annotations for cell-types** in each dataset across  
839        batches were determined specific to the annotation protocol followed by each  
840        original study, with the exception of the 8 batch PDAC data, which was re-  
841        annotated (see 1.3).

842        **1.2 Setting up the PBMC control dataset**

843        For testing in a scenario where the cell-types and cell-type proportions are  
844        perfectly balanced between batches, and subsequent perturbation experiments,  
845        a dataset that was preprocessed by Tran et al.[9] comprising of two batches  
846        of peripheral blood mononuclear cells (PBMCs) sequenced using two variants  
847        of 10x genomics protocols - 5' versus 3' end. As these technologies capture  
848        different regions of mRNA, there is an expected batch effect present. To create  
849        a balanced dataset, the two batches were downsampled for cell-types that had at  
850        least 200 cells in each batch - leaving B cells, CD14+ Monocytes, CD4+ T cells,  
851        CD8+ T cells, FCGR3A+ Monocytes, and Natural Killer (NK) cells. Within  
852        each batch, these remaining cell-types were randomly downsampled to 200 cells,  
853        leading to a perfectly balanced control setup for perturbation experiments.

854        **1.3 Setting up the pancreatic ductal adenocarcinoma (PDAC)  
855        dataset**

856        The pancreatic ductal adenocarcinoma dataset was taken from the Peng et al.  
857        [20] multi-patient study, which comprised of 23 samples. For this data, custom  
858        annotations of tumor cells were done in the following manner:

859        Cells from different samples were integrated with Harmony [22] and clustered  
860        with Seurat [25]. A cell type label was then assigned to each Seurat cluster,  
861        based on the expression of specific marker genes for each cell type (Supplemen-  
862        tary Table 8). To identify tumour cells, all epithelial cells including those from  
863        normal tissue were clustered again using Seurat. All cells that clustered with  
864        the normal samples were assigned as 'Epithelial normal', while all others were  
865        assigned as 'Epithelial tumor'.

866        After annotation of the epithelial normal and tumor cells, the rest of the cells

867 were collapsed into the 'Microenvironment' compartment. Ductal and acinar  
868 cells that did not fall into the epithelial normal or epithelial tumor popula-  
869 tions were removed, as these were likely mis-annotated. Batches/samples were  
870 filtered based on the presence of at least 50 cells in each of the three com-  
871 partments (Epithelial normal, Epithelial tumor, Microenvironment). This left  
872 8 batches/samples, which were utilized in subsequent experiments.

873 **2 scRNA-seq integration methods and**  
874 **parameters**

875 Five state of the art scRNA-seq methods were utilized, based on their perfor-  
876 mance in previous benchmarking papers [9] [13], including BBKNN (version  
877 1.5.1) [21], Harmony (python implementation - version 0.0.5) [22], scVI (scvi-  
878 tools version 0.14.4) [24], Scanorama (version 1.7.1) [23], and Seurat (version  
879 4.0.6) [25]. LIGER (version 0.5.0) [53] was also originally tested, but did not in-  
880 dicate strong performance and resulted in a high degree of variability due to the  
881 removal of seeding in different steps. Therefore, the results from LIGER were  
882 omitted from the main findings, as the high variance of results even within the  
883 control experiments did not allow for a statistically sound comparison between  
884 control and perturbation groups.

885 Because the perturbation experiments were carried out in replicates, to get a  
886 more clear sense of variability within replicates, seeding mechanisms within each  
887 method were removed. This included removing any calls in the method source-  
888 code to R-based seeding for Seurat, and any calls to seeding from the following  
889 libraries for BBKNN, Harmony, Scanorama, and scVI: random, numpy, torch.  
890 This led to a more true estimation of the variability in performance of each  
891 method, as well as a more reliable estimation of the effects of perturbation  
892 because the variability can no longer be simply attributed to variability in the  
893 method which has been accounted for.

894 With the exception of scVI, each method utilized the same processing pipeline  
895 for the data, where scanpy's functions [52] were utilized in the following manner:

- 896 1. Count normalization for each cell to total value of  $1 * 10^4$

- 897        2. Transformation of counts using the  $\log(1 + x)$  function
- 898        3. Highly variable gene selection using the 'seurat' method for 2500 genes
- 899        4. Principal component analysis (PCA) reduction to top 20 principal com-
- 900        ponents that explain the highest variance (PCs) for counts
- 901        5. **Integration at this step for Harmony, Scanorama, and Seurat**
- 902        6. Creating neighborhood graph using the embedding with 20 dimensions
- 903        for embedding and 15 nearest-neighbors - **integration at this step for**
- 904        **Scanorama (replaces neighborhood graph step in scanpy pipeline)**
- 905        7. Leiden clustering on the integrated neighborhood graph using scanpy's
- 906        default parameters
- 907        8. Uniform Manifold Approximation and Projection (UMAP) on the inte-
- 908        grated neighborhood graph using scanpy's default parameters

909        The only exception to this setup was scVI, which requires raw scRNA-seq  
910        expression counts [24], and utilized the entire set of genes for each dataset and  
911        the raw counts. For scVI, steps 1-4 outlined are omitted, and it simply integrates  
912        the raw data and returned a 20 dimensional embedding, which replaces the PCs.  
913        The rest of the steps (6-8) are the same.

914        BBKNN performs integration on the embedding neighborhood representa-  
915        tion [21], and as a result, many of the downstream analyses that required embed-  
916        dings did not have data for BBKNN as it was untestable. Default parameters  
917        were utilized for all methods to ensure fairness in across-method comparisons,  
918        as well as comparisons before and after perturbations.

### 919        3 Perturbation experiments

920        Perturbation experiments were carried out in two settings - the balanced 2  
921        batch PBMC data, and the pancreatic ductal adenocarcinoma data. In both  
922        instances, batches are randomly selected to be perturbed, as well as given  
923        cell-types/compartments. There are three types of perturbation experiments  
924        performed - control, downsampling, and ablation. Control experiments don't

925 downsample any data but allow for replicates of integration runs across meth-  
926 ods to get a sense of intra and inter-method variance on the data without  
927 perturbation. Downsampling experiments involve randomly selecting cells of  
928 a selected cell-type across the indicated number of batches, and downsam-  
929 pling to 10% of the original cell-type population. Ablation experiments in-  
930 volve completely removing selected cell-types from the indicated number of  
931 batches. Randomness of selection for the batches, cell-types, and cells within  
932 indicated cell-type are ensured through randomly generated numbers for each  
933 perturbation simulation/run. To determine the effects of perturbation, results  
934 from the control experiments are compared with results from downsampling  
935 and ablation experiments, across all methods and selected datasets. The code  
936 for the experimental setup, as well as the Iniquitate pipeline, are available at  
937 <https://github.com/hsmaan/Iniquitate>.

### 938 3.1 Balanced 2 batch PBMC data

939 Within the balanced 2 batch PBMC dataset, perturbation experiments were  
940 performed for one of two batches (randomly selected) at a time, and for one  
941 cell-type at a time. 400 replicates were done for the control experiments, and  
942 200 replicates were done for the downsampling and ablation experiments, ensur-  
943 ing that both batches ( $n=2$ ) and each cell-type ( $n=6$ ) is sampled repeatedly and  
944 method performance variance within control experiments is taken into account  
945 adequately. Within the **hierarchical setup**, where similar cell-types were hi-  
946 erarchically clustered into 3 groups (B cell, Monocyte, NK/T cell), the same  
947 number of replicates were done for for the control, downsampling and ablation  
948 experiments.

### 949 3.2 8 batch PDAC data

950 For the 8 batch pancreatic ductal adenocarcinoma dataset, where cells were  
951 grouped into three major compartments, 4 batches were randomly selected for  
952 downsampling or ablation, and one compartment was downsampled or ablated  
953 within one replicate. In total, 100 control replicates were performed, and 50  
954 downsampling and ablation replicates, as the number of compartments is small  
955 ( $n=3$ ) and there will be adequate sampling and repetition within 50 runs.

956        **4 Benchmarking integration performance - PBMC**  
957        **2 batch control dataset**

958        Performance in integration and downstream tasks was assessed using the in-  
959        tegrated embeddings and Leiden clustering [29] results from each integration  
960        technique. After integration at either the embedding or neighborhood calcu-  
961        lation stage through the scanpy library, Leiden clustering was used. Default  
962        values were used for the embedding and clustering steps in the scanpy library  
963        [52]. Only BBKNN did not result in embeddings to be utilized as it performs in-  
964        tegration at the neighborhood, and therefore was not included in the *K-nearest*  
965        *neighbors classification* experiments as these relied on integrated embeddings.

966        **4.1 Quantifying cell-type conservation and batch-mixing**  
967        **with clustering metrics**

968        Four metrics were calculated for all integration experiments, including pertur-  
969        bations and replicates, - Adjusted Rand Index (ARI) [26], Adjusted Mutual  
970        Information (AMI) [27], Completeness [28], and Homogeneity [28]. The *sklearn*  
971        (version  $\geq 1.0.1$ ) implementation of these metrics was utilized. Details of these  
972        metrics can be found in [the scikit learn documentation](#) [32]. These values were  
973        calculated by comparing the known annotated labels with the cluster labels  
974        obtained after integration for each technique. Both cell-type and (1 - batch)  
975        values were calculated for each metric, where cell-type metrics compared the  
976        known cell-type annotations with the cluster labels to determine how well the  
977        integrated embeddings corresponded to known cell-type labels, and the (1 -  
978        batch) values used the batch annotations and the cluster labels to determine  
979        how well the different batches co-aggregate in the embeddings. The assumption  
980        of the latter is that integration should lead to strong batch mixing, and the  
981        shadow of the value is used (1 - batch) to reflect this desired property. The  
982        **median** value across all replicates for a given combination of {method, ex-  
983        periment type, downsampled cell-type} was determined. For the main analysis,  
984        the cell-type and (1 - ARI<sub>batch</sub>) was utilized, but values for all metrics were  
985        calculated (Supplementary table 9).

986        **4.1.1 Z-score normalization of ARI metrics**

987        As the focus of the analysis was not to assess inter-method variation, but de-  
988        termine intra-method variation based on the properties of perturbations versus  
989        the control experiments, the median values for cell-type and  $(1 - \text{ARI}_{batch})$ , for  
990        all combinations of {method, experiment type, downsampled cell-type}, were  
991        Z-score normalized. E.g. for cell-type ARI values for a specific subset:

$$\frac{\text{Median ARI}_{\{method_x, type_y, cell-type_z\}} - \mu(\text{Median ARI}_{\{method, type, cell-type\}})}{\sigma(\text{Median ARI}_{\{method, type, cell-type\}})} \quad (1)$$

992        The exact same procedure is followed for the  $(1 - \text{ARI}_{batch})$  values.

993        **4.2 Downstream analysis - unsupervised clustering**

994        In this downstream analysis test post-integration, the number of unsupervised  
995        Leiden clusters are determined and compared between the different perturba-  
996        tion experiments and control groups. The same indicated setup is used, and the  
997        number of clusters are determined using the default parameters for Leiden clus-  
998        tering in the scanpy library [52]. As each method resulted in different Leiden  
999        clusters, these were analyzed independently and intra-method and experiment-  
1000      type comparisons were performed.

1001        **4.3 Downstream analysis - k-nearest neighbor (KNN)  
1002        classification**

1003        The goal of this downstream analysis test was to determine the performance  
1004        of integration techniques at a per-cell-type level before and after perturbation.  
1005        After obtaining the integrated embeddings for all methods, with the exception  
1006        of BBKNN, a KNN classifier is trained on a 70/30 training/test split of the  
1007        integrated embeddings to predict the cell-type labels of the test data. Stratified  
1008        sampling was used for the split to ensure that all classes were represented  
1009        in the same proportions between train and test sets. The *sklearn* (version  $\geq$

1010        1.0.1) library was used for the data preparation, test/train split, stratified sam-  
1011        pling, KNN-classifier training and prediction [32]. The explicit formulation for  
1012        prediction of a class on a test data point  $x_i$  is:

$$\text{class } x_i = \max_{y \in \mathbb{Y}} \sum_{\substack{y_i \in \mathbb{N}_{x_i} \\ y_i = y}}^k \delta(y_i, y) \quad (2)$$

$$\delta(y_i, y) = \begin{cases} 1, & \text{if } y_i = y \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

1013        Where  $k$  indicates the number of neighbors used in the classifier, which was  
1014        set to 15 for all runs. As different runs could possible lead to different test/train  
1015        splits of the integrated embedding, a seed was used to ensure that the same split  
1016        occurs across all experiments. This also ensured that each method was tested  
1017        on the same split of the data. Using the results of the predictions, the cell-  
1018        type-specific precision, recall, and F1 scores were determined, and the F1-score  
1019        specific to each cell-type was used as key metric. These metrics was calculated  
1020        Primarily, cases were examined where a specific cell-type was downsampled or  
1021        ablated, and the effects of performance on the same cell-type based on the  
1022        KNN-classification F1-score was analyzed. The form of the score is given by  
1023        [33]:

$$\text{F1 score per cell-type} = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} = \frac{TP}{TP + \frac{1}{2}(FP + FN)} \quad (4)$$

1024        Where  $TP$  is the number of true positive calls,  $FP$  is the number of false  
1025        positive calls, and  $FN$  is the number of false negative calls, all on a per-cell-type  
1026        basis.

#### 1027        4.4 Downstream analysis - marker gene ranking

1028        Differential gene expression (DGE) analysis is typically performed after clus-  
1029        tering (or clustering an integrated representation of many batches/samples) to

1030 determine marker genes specific to each cluster that are then used to annotate  
1031 cells within those clusters [8]. The goal of this analysis was to determine to  
1032 what extent can the results of DGE be altered after perturbation of balanced  
1033 data.

1034 First, the top 10 marker genes corresponding to each cell-type were de-  
1035 termined in each batch for each dataset (e.g. 2 batch PBMC dataset) using  
1036 the Wilcoxon Rank-Sum Test and the scanpy package (sc.tl.rank\_genes\_groups)  
1037 [52]. To ensure selection of relevant markers, ribosomal and mitochondrial genes  
1038 were removed from the pool of tested genes. After this, to obtain a consensus on  
1039 the marker genes across batches, the union of markers for each cell-type across  
1040 batches was determined and duplicate gene calls across batches were dropped.  
1041 This will lead to an uneven number of markers for some cell-types if completely  
1042 distinct sets are called from different batches, but leads to a more complete set  
1043 as the integrated space across batches is being analyzed. This set of markers  
1044 for each cell-type in a dataset was deemed the **master marker list**.

1045 From here, after each the integration step using each method in each control  
1046 or perturbation simulation for a given dataset, unsupervised clustering using the  
1047 Leiden clustering algorithm with the scanpy default parameters was done for  
1048 the integrated embedding and DGE using the Wilcoxon Rank-Sum Test was  
1049 performed for each of the obtained unsupervised clusters [52]. A challenge here  
1050 is that *there is no correspondence between the unsupervised clusters obtained*  
1051 *in this integrated embedding and the cell-types used for determining the master*  
1052 *marker list*. However, a way to get around this is to do DGE for each cluster  
1053 and check the **maximum ranking** of a given marker gene across all clusters.  
1054 Ranking is defined by how significant the DGE *p*-value is for a given gene, where  
1055 the highest rank is the most statistically significant differentially expressed gene  
1056 for a given cluster. If a cluster still corresponds to a given cell-type (which is  
1057 the central assumption in unsupervised integration), then that cluster should  
1058 return a high ranking for a given marker gene corresponding to that cell-type  
1059 in DGE. Therefore, for the markers in the master marker list, we can analyze  
1060 their **maximum ranking** across unsupervised clusters in the integrated space  
1061 - to see if biological information specific to that cell-type and its markers is still  
1062 being retained after integration.

1063 This is precisely the operation carried out, and the change in ranking for all  
1064 of the marker genes corresponding to the different known cell-types in datasets  
1065 quantified by their standard deviation in a given subset of experiments (control,

1066 downsampling, or ablation) **change in maximum ranking**. This change in  
1067 maximum ranking within an experiment group (e.g. the control group) was  
1068 indicated as the **marker gene perturbation score**:

$$\text{Marker gene perturbation score} = \sigma(\max \text{ marker gene ranking}) \quad (5)$$

1069 If this value is being averaged over many genes (e.g. for a cell-type), this is  
1070 indicated as the **average marker gene perturbation score**. If there are  $m$   
1071 marker genes for a given cell-type:

$$\text{Avg. marker perturbation score} = \frac{1}{m} \sum_{i=1}^m \text{Marker } i \text{ perturbation score} \quad (6)$$

1072 **4.4.1 Case study - CD4/CD8 T cell assignment based on marker**  
1073 **genes**

1074 To determine if the changes in marker gene ranking that were observed could  
1075 realistically influence the results of a single-cell analysis, the same marker gene  
1076 perturbation set-up was utilized. In this case however, each of the unsupervised  
1077 clusters after integration were annotated as specific cell-types based on a major-  
1078 ity of their cells present (e.g. Cluster 1 - $\downarrow$  Majority CD4+ T cells - $\downarrow$  CD4+ T).  
1079 In this setup, only clusters that contained a majority of either CD4+ or CD8+  
1080 T cells were kept. For simplifying the case study, only integration results from  
1081 the Seurat method were utilized.

1082 After integration and selection of clusters with a majority of CD4+ and  
1083 CD8+ T cells, differential expression analysis was performed as previously indi-  
1084 cated, and a permissive threshold of the top 50 marker genes was used to select  
1085 markers for the CD4+ and CD8+ majority clusters. From here, each of the  
1086 CD4+ and CD8+ T cell majority clusters were predicted to be either CD4+  
1087 or CD8+ based on the presence of canonical marker genes: IL7R for CD4+ T  
1088 cells and CD8A for CD8+ T cells [25].

1089 Examining the top 50 marker genes for each cluster, the rules for predicting

1090 the cell-types each of the CD4+ and CD8+ T majority clusters comprised of  
1091 were the following:

```
1092 if IL7R and CD8A present then
1093   if Rank(IL7R) > Rank (CD8A) then
1094     Annotate as CD4+ T
1095   else if Rank(CD8A) > Rank (IL7R) then
1096     Annotate as CD8+ T
1097   end if
1098   else if IL7R present then
1099     Annotate as CD4+ T
1100   else if CD8A present then
1101     Annotate as CD8+ T
1102   else
1103     Annotate as Undefined
1104   end if
```

1105 From here, the fraction of unsupervised clusters that contained a majority
1106 of CD4+ and CD8+ T cells were predicted for their cell-types based on differen-
1107 tial expression in control and perturbation (downsampling and ablation) experi-
1108 ments, including replicates. Only downsampling and ablation experiments that
1109 affected CD4+ and CD8+ T cells were analyzed, as downsampling/ablating
1110 these were found to most likely affect the marker gene rankings of either cell-
1111 type.

## 1112 4.5 Downstream analysis - query-to-reference annota- 1113 tion

1114 To test the robustness of query-to-reference annotation techniques across vary-
1115 ing degrees of unshared variation, the Seurat 4.0 multi-modal projection tech-
1116 nique was utilized [31]. Although the control PBMC 2 batch dataset has only
1117 scRNA-seq information, a multi-modal reference can still be utilized as is, as
1118 only the RNA-seq modality will be integrated. The reference dataset utilized is
1119 from Hao et al., and the same parameters indicated in the vignette were utilized
1120 [31].

1121        It's important to note that **integration was not performed before projection** using the Seurat 4.0 method. Instead, **each batch/sample is individually projected/integrated to the reference dataset and annotated**,  
1122        as per the guidelines for Seurat 4.0 [31]. Therefore, there are no method-specific  
1123        comparisons to be made in this analysis.  
1124

1126        As the annotations in the reference will not exactly match the annotations  
1127        from the PBMC 2 batch data (mostly due to a higher degree of granularity  
1128        and different naming conventions) [9] [31], a scoring guide was created to de-  
1129        termine if the annotation correctly matches the ground-truth cell-type label  
1130        by using "fuzzy-matching" of ground-truth cell-type labels from the PBMC 2  
1131        batch dataset and the labels in the reference data. The following table summa-  
1132        rizes the guide for the PBMC 2 batch data, and acceptable annotations for L1  
1133        (coarse-grained label from Hao et al. [31]) and L2 (fine-grained label from Hao  
1134        et al. [31]):

Ground-truth label	Acceptable L1 reference	Acceptable L2 reference
CD4 T cell	CD4 T	CD4 TCM, CD4 Naive, CD4 CTL, CD4 Proliferating, CD4 TEM
CD8 T cell	CD8 T	CD8 Naive, CD8 TEM, CD8 TCM, CD8 Proliferating
Monocyte_CD14	Mono	CD14 Mono
Monocyte_FCGR3A	Mono	CD16 Mono
NK cell	NK	NK, NK Proliferating, NK_CD56bright

1135        Using this annotation guide, the annotation accuracy as determined by the  
1136        F1 score (4.3) was determined for each experiment and experiment type (control,  
1137        downsampling, ablation). This value was calculated using both the L1 and L2  
1138        annotations, and the annotations for each cell in each experiment were saved.

1139        **5 Complex imbalanced dataset analysis**

1140        After quantifying the effects of unshared variation in the control 2 batch PBMC  
1141        dataset through perturbation experiments, complex datasets that are multi-  
1142        batch and already imbalanced were analyzed, including: imbalanced 2 batch  
1143        PBMC dataset, batch PBMC dataset, 6 batch mouse hindbrain development  
1144        dataset, and 8 batch pancreatic ductal adenocarcinoma dataset.

1145        **5.1 Cell-type center distance**

1146        To determine the distance between cell-types in the embedding space utilized  
1147        for integration, across all batches to be integrated, the following preprocessing  
1148        steps were performed on the raw data for each batch in a dataset:

- 1149        1. Count normalization for each cell to total value of  $1 \times 10^4$
- 1150        2. Transformation of counts using the  $\log(1 + x)$  function
- 1151        3. Highly variable gene selection using the 'seurat' method for 2500 genes
- 1152        4. PCA to top 20 (PCs) for on the counts data

1153        After obtaining the PCs for a given dataset, the ground-truth cell-type labels  
1154        are used to determine the cell-type center distance between all cell-types in  
1155        the data in a pairwise manner. The cell-type center distance is defined as the  
1156        weighted cosine distance between the center (average) of the PCA representation  
1157        for each cell-type in a given batch.

1158        For each batch  $b$ , and cell-type  $a$  with  $n$  cells and a PCA reduction of the  
1159        data:

$$PC_{a_b} \in \mathbb{R}^{n \times 20} \quad (7)$$

$$\text{cell-type } a_b \text{ center} = \frac{1}{n} \sum_{i=1}^n PC_{a_b i} \in \mathbb{R}^{1 \times 20} \quad (8)$$

1160

Then for quantifying the distance between cell-types  $a$  and  $c$  in batch  $b$ :

Let  $v \in \mathbb{R}^{1 \times 20}$  be the variance explained by each of the top 20 PCs

Let  $CC_{ab}$  be the cell-type center for cell-type  $a$  in batch  $b$

Let  $CC_{cb}$  be the cell-type center for cell-type  $c$  in batch  $b$

$$\text{Cell-type center distance } ac_b = 1 - \frac{(CC_{ab} \circ v) \cdot (CC_{cb} \circ v)}{\| (CC_{ab} \circ v) \| \| (CC_{cb} \circ v) \|} \quad (9)$$

1161

Where  $CC_{ab} \circ v$  is the element-wise rescaling of the cell-type center of  $a$  based on the variance explained by the PCs.

1163

The rationale behind a reweighted cosine distance is that the distance itself between cell-types should be scaled according to the variance explained by each PC because the distance is being calculated in the joint PCA reduction of all cells, and not every PC axis will have equal contribution for the variance explained.

1168

We can take the average of this cell-type center distance across  $p$  batches:

$$\text{Avg. cell-type center distance} = \frac{1}{p} \sum_{b=1}^p \text{Cell-type center distance } ac_b \quad (10)$$

1169

From here, the minimum cell-type center distance, or the distance corresponding to the cell-type closest to cell-type  $a$  is simply the minimum value across all batches ( $p$  total). Assume there are  $k$  total cell-types across all batches and missing cell-type pairs (e.g. cell-type present in batch 1 and not batch 2) have an imputed maximum cosine distance of 1. Values between the same cell-types are also imputed as 1. Then using the tensor of cell-type center distances across batches  $D$ :

$$D \in \mathbb{R}^{k \times k \times b} \quad (11)$$

**Minimum cell-type center distance**  $a = \min(D_{a,:,:})$  (12)

Where  $D_a$  is the subset of the first axis for cell-type  $a$ . The cell-type and batch corresponding to this value can also be found through the *argmin*.

The minimum distance in any batch is taken instead of averaging distances across all batches because this minimum distance will correspond to the most 'haphazard' scenario for a given batch being integrated. There are two scenarios possible here:

1. The similarity between cell-types is largely similar across batches, and the minimum value will correspond roughly to the average
2. The similarity between cell-types can be very different across batches, due to scenarios/factors such as developmental data or treatment-effects

The first case is most readily applicable to the PBMC datasets, but the second scenario may be more applicable to the PDAC and hindbrain developmental data. However, even in these cases, taking the minimum may lead to a better approximation of proximity affecting integration results because it will factor in the worst possible scenario (across batches) for a given cell-type.

## 5.2 Cell-type support

The cell-type support (or relative cell-type support) was simply the log2-transformation of the number of cells for each cell-type  $a$  across all batches  $b$ :

$$\text{Relative cell-type support } a = \log_2\left(\sum_{b=1}^p \text{Cell-type } a_{batch\ b}\right) \quad (13)$$

1194

## 6 Statistical testing

1195

### 6.1 One-way ANOVA tests

1196

To determine statistical significance for the effects of perturbations, the following generic one-way analysis-of-variance (ANOVA) setup was utilized [34]:

1197

$$\text{response} \sim x_0 + x_1 + x_2 + \dots + x_m + \text{type} \quad (14)$$

$$\mathcal{H}_0 : \text{response} = x_0 + x_1 + x_2 + \dots + x_m \quad (15)$$

1198

Where  $\mathcal{H}_0$  is the null hypothesis, *response* can be an endpoint of interest in the analysis (e.g. number of clusters post integration),  $x_0$  is a constant (intercept/bias),  $x_1, \dots, x_m$  are factors we'd like to control before testing significant with respect to perturbations (e.g. method, cell-type that was downsampled), and *type* is a binary covariate indicating the experiment type that was done:

1199

1200

1201

1202

$$\text{type} = \begin{cases} 1, & \text{if } y_i = \text{downsampling, ablation} \\ 0, & \text{if } y_i = \text{control} \end{cases} \quad (16)$$

1203

After accounting for the various factors we'd like to control ( $x_1, \dots, x_m$ ), we can assess the statistical significance of perturbation of unshared variation (*type*) with respect to the *response* covariate through the **ANOVA F-statistic** and **p-value** associated with the *type* covariate.

1204

1205

1206

In situations where significance is achieved across various groups due to factors such as intra- and inter-method variance, the magnitude of the F-statistic is compared.

1207

1208

1209

1210        **6.2 Control PBMC 2 batch dataset**

1211        **6.2.1 KNN classification per cell-type**

1212        For assessing the effects of perturbation on the F1 classification scores post-  
1213        integration on a per-cell-type level, the following ANOVA (6.1) setup was uti-  
1214        lized:

$$F1 \text{ classification score} \sim x_0 + \text{method} + \text{downsampled cell-type} + \text{type} \quad (17)$$

$$\mathcal{H}_0 : F1 \text{ classification score} = x_0 + \text{method} + \text{downsampled cell-type} \quad (18)$$

1215        The F1-classification scores here are across all cell-types in the integrated  
1216        dataset. **The cell-type being analyzed (for the F1 classification score**  
1217        **in each instance) is equivalent to the downsampled cell-type in each**  
1218        **sample included in the test.**

1219        **6.2.2 Unsupervised clustering**

1220        For comparing the significance of perturbation on the number of unsuper-  
1221        vised clusters obtained post-integration using Leiden clustering, the following  
1222        ANOVA (6.1) setup was utilized:

$$n \text{ clusters} \sim x_0 + \text{method} + \text{downsampled cell-type} + \text{type} \quad (19)$$

$$\mathcal{H}_0 : n \text{ clusters} = x_0 + \text{method} + \text{downsampled cell-type} \quad (20)$$

1223        **6.2.3 marker gene ranking**

1224        To test the statistical significance of perturbations for each marker gene ana-  
1225        lyzed (4.4), the following ANOVA setup was used for each marker gene  $g$ :

Marker  $g$  max rank  $\sim x_0 + \text{method} + \text{downsampled cell-type} + \text{type}$  (21)

$$\mathcal{H}_0 : \text{Marker } g \text{ max rank} = x_0 + \text{method} + \text{downsampled cell-type} \quad (22)$$

1226 Then, to test the overall effects on marker gene ranking, considering all  
1227 marker genes at once, the following test was done:

Marker max rank  $\sim x_0 + \text{gene} + \text{method} + \text{downsampled cell-type} + \text{type}$  (23)

$$\mathcal{H}_0 : \text{Marker max rank} = x_0 + \text{gene} + \text{method} + \text{downsampled cell-type} \quad (24)$$

1228 **6.3 Complex imbalanced datasets**

1229 **6.3.1 Cell-type support and cell-type center distance**

1230 To determine if the two key metrics that were determined in the complex dataset  
1231 analysis - **relative cell-type support** (5.2) and **cell-type center distance**  
1232 (5.1) - are in fact predictive of integration performance, the following ANOVA  
1233 setups were used where the F1-classification score for each experiment, cell and  
1234 associated ground-truth cell-type was tested:

F1 classification score  $\sim x_0 + \text{method} + \text{minimum cell-type center distance}$  (25)

$$\mathcal{H}_0 : \text{F1 classification score} = x_0 + \text{method} \quad (26)$$

F1 classification score  $\sim x_0 + \text{method} + \text{relative cell-type support}$  (27)

$$\mathcal{H}_0 : \text{F1 classification score} = x_0 + \text{method} \quad (28)$$

1235 Where:

1236                  *minimum cell-type center distance*  $\in \mathbb{R}_{0,1}^+$                   (29)

1237                  *relative cell-type support*  $\in \mathbb{N}$                   (30)

1236                  The *cell-type analyzed* was not included as a factor to control, because the  
1237                  *minimum cell-type center distance* and *relative cell-type support* metrics were  
1238                  calculated on a **per cell-type basis**. Therefore, these metrics are perfectly  
1239                  collinear with cell-type, and this would absorb the residuals that would be  
1240                  picked up by the key metrics.

1241                  **6.3.2 PDAC perturbation analysis**

1242                  Perturbations were performed for the compartmentalized PDAC data (1.3 and  
1243                  3.2) to determine the effects of downsampling/ablation on the classification  
1244                  scores of all compartments. Here, the following ANOVA setup was used to  
1245                  determine the effects on F1-scores **for a specific compartment** based on  
1246                  **downsampling of the same compartment** for each compartment  $c$ :

1247                   $F1 \text{ classification score } c \sim x_0 + \text{method} + \text{type}$                   (31)

1248                   $\mathcal{H}_0 : F1 \text{ classification score } c = x_0 + \text{method}$                   (32)

1249                  These results were analyzed independently and jointly for all compartments  
downsampled, where joint-comparison included comparison of F-values for per-  
turbation in each setup.

1250                  **7 Balanced clustering scores**

1251                  None of the utilized clustering metrics, which in this analysis and other inte-  
1252                  gration benchmarking/methods papers are used to compare the concordance  
1253                  of ground-truth cell-type labels and unsupervised clusters attained in an em-  
1254                  bedding, factor in class balance. The metrics utilized include: the Adjusted

1255 Rand Index (ARI), Adjusted Mutual Information (AMI), Homogeneity Score,  
1256 and Completeness Score. The implementation details of these metrics can be  
1257 found in the scikit learn documentation [32].

1258 Strictly speaking, the **Homogeneity Score and Completeness Scores**  
1259 **are not metrics, because they are not symmetric**. However, this sym-  
1260 metry is not necessary in the case of single-cell benchmarking, and the general  
1261 case of comparing clustering labels with ground-truth annotations, because one  
1262 set of labels is known to be ground-truth. In fact, balancing the ARI and AMI  
1263 will break their symmetry as well.

1264 To introduce the procedure behind reweighing these metrics, we'll begin with  
1265 the balanced ARI. Then we'll extrapolate this procedure to the entropy-based  
1266 metrics/scores (AMI, Homogeneity, and Completeness), as this extrapolation  
1267 only involves a slight modification to the ARI procedure for these scores.

1268 Code notebooks on implementing the balanced clustering scores with usage  
1269 demonstrations and relevant examples are available at <https://github.com/>  
1270 [hsmaan/balanced-clustering/tree/main/notebooks](https://github.com/hsmaan/balanced-clustering/tree/main/notebooks).

## 1271 7.1 The *Balanced Adjusted Rand Index*

### 1272 7.1.1 The Rand Index and Adjusted Rand Index

1273 For a set of  $n$  objects,  $S = \{O_1, O_2, O_3, \dots, O_n\}$ , the goal of clustering is to  
1274 partition these objects into meaningful subsets, which we can call partitioning  
1275  $V$ . Assuming we have access to either ground-truth labels or clusters from  
1276 another technique, which we can denote partitioning  $U$ . Both  $U$  and  $V$  contain  
1277 subsets, which we call either classes or clusters:  $U = \{u_1, u_2, \dots, u_R\}$  and  $V =$   
1278  $\{v_1, v_2, \dots, v_C\}$ . These clustering results are subject to the following constraints  
1279 to be valid for calculating the Rand Index:

- 1280 1. All  $n$  objects within the set  $S$  must be within sets  $U$  and  $V$ :

$$U_{i=1}^R u_i = U_{j=1}^C v_j = S \quad (33)$$

- 1280 2. No element from set  $S$  can belong in to more than one subset in either  $U$

or  $V$

$$1 \leq i \neq i' \leq R \quad (34)$$

$$1 \leq j \neq j' \leq C \quad (35)$$

$$u_i \cap u_{i'} = \emptyset = v_j \cap v_{j'} \quad (36)$$

1281 To quantify the overlap between the partitions  $U$  and  $V$  (e.g. in determining  
 1282 overlap between a set of ground-truth labels and the results of clustering), we  
 1283 can start by creating a contingency table which indicates the overlap:

	$v_1$	$v_2$	$v_3$	$\dots$	$v_c$	
$u_1$	$t_{11}$	$t_{12}$	$t_{13}$	$\dots$	$t_{1C}$	$t_{1..}$
$u_2$	$t_{21}$	$t_{22}$	$\dots$	$\dots$	$\dots$	$\dots$
$u_3$	$t_{31}$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$
$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$
$u_R$	$t_{R1}$	$\dots$	$\dots$	$\dots$	$t_{RC}$	$t_{R..}$
	$t_{..1}$	$\dots$	$\dots$	$\dots$	$t_{..C}$	$t_{..}$

1284 Each element of this table indicates overlapping elements. E.g.  $t_{11}$  indicates  
 1285 the number of samples that have the label  $v_1$  in  $V$  and  $u_1$  in  $U$ . The total  
 1286 number of values in the matrix is  $\binom{n}{2}$  if  $n$  objects/samples are present. As we  
 1287 now have a table/matrix that represents the overlap of assignments to subsets  
 1288 in  $U$  and  $V$  for  $n$  objects, we can determine the concordance of partitions  $U$   
 1289 and  $V$  for these objects using the Rand Index:

$$a = \sum_{r=1}^R \sum_{c=1}^C \binom{t_{rc}}{2}, \binom{x}{2} = 0 \text{ if } x = 0 \quad (37)$$

$$b = \left[ \sum_{r=1}^R \binom{t_{r.}}{2} \right] - a \quad (38)$$

$$c = \left[ \sum_{c=1}^C \binom{t_{.c}}{2} \right] - a \quad (39)$$

$$d = \binom{n}{2} - a - b - c \quad (40)$$

$$\textbf{Rand Index (RI)} = \frac{a + d}{a + b + c + d} \quad (41)$$

1290            Intuitively, the Rand Index aims to calculate how many pairs are concordantly in the same subsets in  $V$  and  $U$  ( $a$ ), how many pairs are concordantly in different subsets in  $V$  and  $U$ , and how many are discordant (in the same group in one partition and otherwise in the other). It's important to note that **pairs here refer to all combinations of two different objects, not the same object being considered in the two partitions.**

1296            Although the Rand Index is normalized (lower bound = 0, upper bound = 1), it is not adjusted for chance clustering. A correction can be made [26] to the  
 1297            RI formula that takes into account the **expected value of the RI for two**  
 1298            **partitions of the objects  $U$  and  $V$** , denoted by the Adjusted Rand Index  
 1299            (ARI) [26]:

$$\textbf{ARI} = \frac{\binom{n}{2}(a + d) - [(a + b)(a + c) + (c + d)(b + d)]}{\binom{n}{2}^2 - [(a + b)(a + c) + (c + d)(b + d)]} \quad (42)$$

1301            With this correction, the ARI is a metric (symmetric, positive-definite, and  
 1302            the triangle inequality) [27] that has the properties of **normalization** and  
 1303            **expectation** [26].

1304

### 7.1.2 Balancing the ARI

1305

Rebalancing the ARI (as well as the other entropy-based scores/metrics) will amount to **rescaling the total number of values in the subsets of the partition we consider ground-truth**. In this case, **assume  $U$  is the partition with the ground-truth information**. We want each subset from  $U$  to have an equal contribution to the ARI value - this is concomitant with each class from the ground-truth data (which we have assumed to be  $U$ ) having an equal weighing in the calculating of the score. This can be done in the following step-wise manner:

1306

- 1307 1. Determine contribution of each subset of  $U$  ( $t_1, t_2, \dots, t_R$ ) to score through  
1308 mean of marginals from contingency table:

$$(t_{1.}, t_{2.}, t_{3.}, \dots, t_{R.}) \quad (43)$$

1309

- 1310 2. Get the mean contributions of all subsets:

$$C = \frac{1}{R} \sum_{i=1}^R t_{i.} \quad (44)$$

- 1311 3. For each subset ( $t_i$ ) of  $U$ , normalize the contribution to be equal to the  
1312 mean using a scaling factor:

$$S_i = \frac{C}{t_{i.}} \quad (45)$$

$$\forall t_i, (t_{i1}, t_{i2}, \dots, t_{iC}) = S_i * (t_{i1}, t_{i2}, \dots, t_{iC}) \quad (46)$$

1313

After these steps, we've essentially rescaled the contingency table such that the contribution from each subset in  $U$  will be considered equally in calculations using the table results. To calculate the *Balanced Adjusted Rand Index*, we can apply the 7.1.2 normalization procedure and use the same ARI formula as before (42):

$$\text{Balanced ARI} = \frac{\binom{n}{2}(a+d) - [(a+b)(a+c) + (c+d)(b+d)]}{\binom{n}{2}^2 - [(a+b)(a+c) + (c+d)(b+d)]} \quad (47)$$

1321 Examining the values needed to calculate the RI and ARI (37), we can see  
1322 that this normalization procedure will effectively rescale the calculations for  $a$ ,  
1323  $b$ , and  $c$ . This procedure does this while still retaining the total counts ( $n$ ),  
1324 such that the calculation for  $d$  will be unaffected. Because the RI and ARI  
1325 calculations simply depend on these values in the contingency table that can  
1326 be calculated independently, the application of this normalization procedure is  
1327 straightforward and does not require any further steps.

## 1328 **7.2 *Balancing entropy-based scores***

### 1329 **7.2.1 Mutual information**

1330 Central to the Adjusted Mutual Information (AMI), Homogeneity, and Com-  
1331 pleteness scores is the calculation of mutual information between partitions  $U$   
1332 and  $V$  [27] [28]. For the contingency table previously defined in 47, the mutual  
1333 information between these two partitions is equal to the following [27]:

$$I(U, V) = \sum_{r=1}^R \sum_{c=1}^C \frac{t_{rc}}{n} \log \frac{t_{rc}/n}{t_r t_c/n^2} \quad (48)$$

1334 We'll follow the same normalization procedure that we did in 7.1.2, as we  
1335 are starting from the same contingency table of overlapping objects in subsets  
1336 of partitions  $U$  and  $V$ . From here, we can calculate the mutual information  
1337 value and proceed with the rest of the calculations for the entropy-based scores.

### 1338 **7.2.2 Entropy**

1339 Aside from mutual information, the other important factor that is used by all  
1340 of the entropy-based scores is the calculation of the entropy of the labelling -  
1341 i.e. how ordered/disordered are the objects in partitions  $U$  and  $V$ . This can  
1342 also be calculated from the contingency table from 47 in the following manner  
1343 [27]:

$$H(U) = - \sum_{r=1}^R \frac{t_{r.}}{n} \log \frac{t_{r.}}{n} \quad (49)$$

$$H(V) = - \sum_{c=1}^C \frac{t_{.c}}{n} \log \frac{t_{.c}}{n} \quad (50)$$

1344 The Homogeneity and Completeness scores also require the conditional en-  
 1345 tropy formulation [28] [27]:

$$H(U|V) = - \sum_{r=1}^R \sum_{c=1}^C \frac{t_{rc}}{n} \log \frac{t_{rc}/n}{t_{.c}/n} \quad (51)$$

$$H(V|U) = - \sum_{r=1}^R \sum_{c=1}^C \frac{t_{rc}}{n} \log \frac{t_{rc}/n}{t_{r.}/n} \quad (52)$$

### 1346 7.2.3 Balanced entropy-based scores

1347 The calculation of the entropy and mutual information can proceed as-is after  
 1348 the normalization procedure from 7.1.2, and this will balance the contributions  
 1349 from a presumed ground-truth partition  $U$  in calculating the entropy and mutual  
 1350 information. From here the Balanced Adjusted Mutual Information, Balanced  
 1351 Homogeneity, and Balanced Completeness scores can be calculated using these  
 1352 two values, the rescaled contingency matrix after 7.1.2, and the base formulas  
 1353 for these scores [27] [28]:

$$\text{Balanced AMI} = \frac{I(U, V) - \mathbb{E}[I(U, V)]}{\frac{1}{2}[H(U) + H(V)] - \mathbb{E}[I(U, V)]} \quad (53)$$

$$\text{Balanced Homogeneity} = 1 - \frac{H(U|V)}{H(U)} \quad (54)$$

$$\text{Balanced Completeness} = 1 - \frac{H(V|U)}{H(V)} \quad (55)$$

1354        The V-measure and Balanced V-measure are simply the harmonic mean of  
1355        the Completeness and Homogeneity scores [28]:

$$\text{Balanced V-measure} = \frac{(2 \times \text{Bal. Homog.} \times \text{Bal. Compl.})}{(\text{Bal. Homog.} + \text{Bal. Compl.})} \quad (56)$$

1356        **7.3 Balanced clustering evaluations**

1357        The following section details the evaluations that were utilized for the balanced  
1358        clustering metric analysis. Seeding was set for all of these cases to ensure repro-  
1359        ducibility of the simulations, downsampling, and integration methods (where  
1360        possible).

1361        **7.3.1 3 imbalanced well-separated classes, 2 clusters**

1362        In this scenario, 3 well separated but imbalanced classes were utilized and a  
1363        mis-clustering of the smaller class was done with k-means clustering with k=2.  
1364        This data was simulated using 2D Gaussian densities with the following values  
1365        for each class:

- 1366            • Class A  $\sim N(0, 0.5)$  - 500 samples  
1367            • Class B  $\sim N(-2, 0.1)$  - 20 samples  
1368            • Class C  $\sim N(3, 1)$  - 500 samples

1369        K-means clustering with k=2 led to class B overlapping with class A in the  
1370        clustering result.

1371        The balanced and imbalanced metrics were compared when calculating the  
1372        concordance of the ground-truth labels (class labels) and k-means clustering  
1373        labels.

1374       **7.3.2 3 imbalanced overlapping classes, 3 clusters**

1375       In this case, 3 classes that are overlapping and imbalanced (2 smaller classes on  
1376       edges of larger class) were analyzed, and k-means clustering with k=3 was done  
1377       and the result correctly clustered most of the samples from the smaller classes,  
1378       but due to slicing of the larger class present because of overlap, mis-clustered a  
1379       large number of majority class samples.

1380       This data was simulated using 2D Gaussian densities with the following  
1381       values for each class:

- 1382       • Class A  $\sim N(0, 0.5)$  - 1500 samples (larger class)  
1383       • Class B  $\sim N(1, 1)$  - 200 samples  
1384       • Class C  $\sim N(-1, 1)$  - 200 samples

1385       **7.3.3 Balanced 2 batch PBMC - co-clustered CD4+ T cells and**  
1386       **CD8+ T cells**

1387       The balanced 2 batch PBMC dataset was utilized here (1.2). Batch 1 was kept  
1388       as is, and batch 2 had all of the cells ablated except for CD4+ T cells, which  
1389       were downsampled to 10% of their original proportion.

1390       The default Leiden clustering resolution of 1 in the scanpy implementation  
1391       was changed to 0.1, as this value perfectly clusters all of the cell-types with the  
1392       exception of the CD4+ T cells, which get collapsed into a cluster with CD8+  
1393       T cells, simulating a case where a smaller cell-type is co-clustered with a larger  
1394       cell-type.

1395       The resultant embedding with no integration was utilized, and the ground-  
1396       truth cell-type labels and unsupervised clustering labels were used to compare  
1397       the balanced and imbalanced/vanilla scores - where the ARI and Homogeneity  
1398       scores were shown.

1399           **7.3.4 Balanced 2 batch PBMC data - downsampled CD4+ T cells**  
1400           **and FCGR3A+ monocytes**

1401       In this evaluation, the 2 batch balanced PBMC dataset was once again utilized.  
1402       For the two batches, each one had either the CD4+ T cells or FCGR3A+ mono-  
1403       cytes downsampled to 10% of their original population, creating an imbalanced  
1404       scenario specific to these two cell-types.

1405       After this, integration was done using BBKNN, Harmony, Scanorama, and  
1406       scVI. The same integration pipeline from 2 was utilized. An 'unintegrated' con-  
1407       trol subset was used, where the pipeline from 2 was followed without integration  
1408       with any method.

1409       From here, the average value of the balanced and imbalanced metrics was  
1410       used for comparison. e.g.:

1411        $\text{Avg imbalanced} = \frac{1}{5} \sum(\text{ARI}, \text{AMI}, \text{Homog.}, \text{Completeness}, V - \text{measure}).$

1412           **8 Code and data availability**

1413       The python package for implementing the balanced clustering metrics can be  
1414       found here:

1415       <https://github.com/hsmaan/balanced-clustering>

1416       All of the code necessary to reproduce the results of the Iniquitate pipeline  
1417       are available at:

1418       <https://github.com/hsmaan/Iniquitate>

1419       The datasets utilized in this study, which are associated with the various  
1420       configurations used in the Iniquitate GitHub repository, can be all found here:

1421       [https://drive.google.com/file/d/102ntQuclUzQILRxMVXo1-yQR43t97Q3r/  
view?usp=sharing](https://drive.google.com/file/d/102ntQuclUzQILRxMVXo1-yQR43t97Q3r/view?usp=sharing)

1423       This directory is in the exact necessary structure needed to run the Iniqui-

1424       tate pipeline, and can be copied into the cloned GitHub repository for Iniqui-  
1425       tate under **Iniquitate/resources**. Instructions are also given in the Iniquitate  
1426       GitHub link.

1427

## Bibliography

- 1428 1. Argelaguet, R. *et al.* Multi-omics profiling of mouse gastrulation at single-  
1429 cell resolution. *Nature* **576** (2019).
- 1430 2. Chiou, J. *et al.* Interpreting type 1 diabetes risk with genetics and single-  
1431 cell epigenomics. *Nature* **594**, 398–402 (2021).
- 1432 3. Pijuan-Sala, B. *et al.* A single-cell molecular map of mouse gastrulation  
1433 and early organogenesis. *Nature* **566**, 490–495 (Feb. 2019).
- 1434 4. Lähnemann, D. *et al.* *Eleven grand challenges in single-cell data science*  
1435 **1**, 1–35 (Genome Biology, 2020).
- 1436 5. Leek, J. T. *et al.* Tackling the widespread and critical impact of batch  
1437 effects in high-throughput data. *Nature Reviews Genetics* **11**, 733–739  
1438 (2010).
- 1439 6. Hou, W., Ji, Z., Ji, H. & Hicks, S. C. A systematic evaluation of single-cell  
1440 RNA-sequencing imputation methods. *Genome Biology* **21**, 1–30 (2020).
- 1441 7. Hicks, S. C., Townes, F. W., Teng, M. & Irizarry, R. A. Missing data and  
1442 technical variability in single-cell RNA-sequencing experiments. *Biostatistics*  
1443 **19**, 562–578 (2018).
- 1444 8. Amezquita, R. A. *et al.* Orchestrating single-cell analysis with Bioconduc-  
1445 tor. *Nature Methods* **17**, 137–145 (2020).
- 1446 9. Tran, H. T. N. *et al.* A benchmark of batch-effect correction methods for  
1447 single-cell RNA sequencing data. *Genome Biology* **21**, 1–32 (2020).
- 1448 10. Ming, J. *et al.* FIRM: Flexible integration of single-cell RNA-sequencing  
1449 data for large-scale multi-tissue cell atlas datasets. *Briefings in Bioinfor-  
1450 matics*, 1–14 (2022).
- 1451 11. He, P. *et al.* *The changing mouse embryo transcriptome at whole tissue*  
1452 *and single-cell resolution* **7818**, 760–767 (Springer US, 2020).
- 1453 12. Zhang, Y. *et al.* Single-cell RNA sequencing in cancer research. *Journal of*  
1454 *Experimental and Clinical Cancer Research* **40**, 1–17 (2021).
- 1455 13. Luecken, M. D. *et al.* *Supplementary Material - Benchmarking atlas-level*  
1456 *data integration in single-cell genomics.* (2021).
- 1457 14. Zheng, G. X. *et al.* Massively parallel digital transcriptional profiling of  
1458 single cells. *Nature Communications* **8** (2017).

- 1459 15. Genomics, 1. 8k PBMCs from a Healthy Donor, Single Cell Gene Expression Dataset by Cell Ranger 2.1.0 (2019).
- 1460
- 1461 16. Abdelaal, T. *et al.* A comparison of automatic cell identification methods for single-cell RNA sequencing data. *Genome Biology* **20**, 1–19 (2019).
- 1462
- 1463 17. Clarke, Z. A. *et al.* Tutorial: guidelines for annotating single-cell transcriptomic maps using automated and manual methods. *Nature Protocols* **16**, 2749–2764 (2021).
- 1464
- 1465
- 1466 18. Ding, J. *et al.* Systematic comparison of single-cell and single-nucleus RNA-sequencing methods. *Nature Biotechnology* **38**, 737–746 (2020).
- 1467
- 1468 19. Vladoiu, M. C. *et al.* Childhood cerebellar tumours mirror conserved fetal transcriptional programs. *Nature* **572**, 67–73 (2019).
- 1469
- 1470 20. Peng, J. *et al.* Single-cell RNA-seq highlights intra-tumoral heterogeneity and malignant progression in pancreatic ductal adenocarcinoma. *Cell Research* **29**, 725–738 (2019).
- 1471
- 1472
- 1473 21. Polański, K. *et al.* BBKNN: Fast batch alignment of single cell transcriptomes. *Bioinformatics* **36**, 964–965 (2020).
- 1474
- 1475 22. Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony. *Nature Methods* **16**, 1289–1296 (2019).
- 1476
- 1477 23. Hie, B., Bryson, B. & Berger, B. Efficient integration of heterogeneous single-cell transcriptomes using Scanorama. *Nature Biotechnology* **37**, 685–691 (2019).
- 1478
- 1479
- 1480 24. Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for single-cell transcriptomics. *Nature Methods* **15**, 1053–1058 (2018).
- 1481
- 1482
- 1483 25. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888–1902.e21 (2019).
- 1484
- 1485 26. Hubert, L. & Arabie, P. Comparing partitions. *Journal of Classification* **2**, 193–218 (1985).
- 1486
- 1487 27. Vinh, N. X., Epps, J. & Bailey, J. Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance. *Journal of Machine Learning Research* **11**, 2837–2854 (2010).
- 1488
- 1489

- 1490 28. Rosenberg, A. & Hirschberg, J. V-Measure: A conditional entropy-based  
1491 external cluster evaluation measure. *EMNLP-CoNLL 2007 - Proceedings of*  
1492 *the 2007 Joint Conference on Empirical Methods in Natural Language Pro-*  
1493 *cessing and Computational Natural Language Learning*, 410–420 (2007).
- 1494 29. Traag, V. A., Waltman, L. & van Eck, N. J. From Louvain to Leiden: guar-  
1495 anteeing well-connected communities. *Scientific Reports* **9**, 1–12 (2019).
- 1496 30. Wang, T., Li, B., Nelson, C. E. & Nabavi, S. Comparative analysis of  
1497 differential gene expression analysis tools for single-cell RNA sequencing  
1498 data. *BMC Bioinformatics* **20**, 1–16 (2019).
- 1499 31. Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**,  
1500 3573–3587.e29 (2021).
- 1501 32. Buitinck, L. *et al.* API design for machine learning software: experiences  
1502 from the scikit-learn project, 108–122 (2013).
- 1503 33. Goutte, C. & Gaussier, E. A Probabilistic Interpretation of Precision, Re-  
1504 call and F-Score, with Implication for Evaluation. *Lecture Notes in Com-*  
1505 *puter Science* 345–359 (2005).
- 1506 34. Winer, B. J., Brown, D. R. & Michels, K. M. *Statistical principles in*  
1507 *experimental design* 3rd ed. (McGraw-Hill, New York, 1991).
- 1508 35. Luecken, M. D. & Theis, F. J. Current best practices in single-cell RNA-seq  
1509 analysis: a tutorial. *Molecular Systems Biology* **15** (2019).
- 1510 36. Andreatta, M. & Carmona, S. J. STACAS: Sub-Type Anchor Correction  
1511 for Alignment in Seurat to integrate single-cell RNA-seq data. *Bioinfor-*  
1512 *matics* **37**, 882–884 (2021).
- 1513 37. Svensson, V., da Veiga Beltrame, E. & Pachter, L. A curated database  
1514 reveals trends in single-cell transcriptomics. *Database* **2020**, 1–7 (2020).
- 1515 38. Dohmen, J. *et al.* Identifying tumor cells at the single-cell level using ma-  
1516 chine learning. *Genome Biology* **23**, 1–23 (2022).
- 1517 39. Trinh, M. K. *et al.* Precise identification of cancer cells from allelic imbal-  
1518 ances in single cell transcriptomes. eng. *Communications biology* **5**, 884  
1519 (Sept. 2022).
- 1520 40. Xu, Y., Liu, J., Nipper, M. & Wang, P. Ductal vs. acinar? Recent insights  
1521 into identifying cell lineage of pancreatic ductal adenocarcinoma. *Annals*  
1522 *of Pancreatic Cancer* **2**, 1–12 (2019).

- 1523 41. Backx, E. *et al.* On the Origin of Pancreatic Cancer: Molecular Tumor  
1524 Subtypes in Perspective of Exocrine Cell Plasticity. *Cmgh* **13**, 1243–1253  
1525 (2022).
- 1526 42. Richards, L. M. *et al.* A comparison of data integration methods for single-  
1527 cell RNA sequencing of cancer samples. *bioRxiv*, 2021.08.04.453579 (2021).
- 1528 43. Aran, D., Sirota, M. & Butte, A. J. Systematic pan-cancer analysis of  
1529 tumour purity. *Nature Communications* **6**, 1–12 (2015).
- 1530 44. Krzak, M., Raykov, Y., Boukouvalas, A., Cutillo, L. & Angelini, C. Bench-  
1531 mark and Parameter Sensitivity Analysis of Single-Cell RNA Sequencing  
1532 Clustering Methods. *Frontiers in Genetics* **10**, 1–19 (2019).
- 1533 45. Argelaguet, R., Cuomo, A. S., Stegle, O. & Marioni, J. C. Computational  
1534 principles and challenges in single-cell data integration. *Nature Biotech-*  
1535 *nology* **39**, 1202–1215 (2021).
- 1536 46. Ogbeide, S., Giannese, F., Mincarelli, L. & Macaulay, I. C. Into the multi-  
1537 verse: advances in single-cell multiomic profiling. *Trends in Genetics* **38**,  
1538 831–843 (2022).
- 1539 47. Luecken, M. D. *et al.* Benchmarking atlas-level data integration in single-  
1540 cell genomics. *Nature Methods* **19**, 41–50 (2022).
- 1541 48. Johansen, N. & Quon, G. ScAlign: A tool for alignment, integration, and  
1542 rare cell identification from scRNA-seq data. *Genome Biology* **20**, 1–21  
1543 (2019).
- 1544 49. Hu, Z., Ahmed, A. A. & Yau, C. CIDER: an interpretable meta-clustering  
1545 framework for single-cell RNA-seq data integration and evaluation. *Genome*  
1546 *Biology* **22**, 337 (Dec. 2021).
- 1547 50. Demetçi, P., Santorella, R., Sandstede, B. & Singh, R. Unsupervised In-  
1548 tegration of Single-Cell Multi-omics Datasets with Disproportionate Cell-  
1549 Type Representation. *Lecture Notes in Computer Science* **13278 LNBI**  
1550 (ed Pe'er, I.) 3–19 (2022).
- 1551 51. McCarthy, D. J., Campbell, K. R., Lun, A. T. & Wills, Q. F. Scater: Pre-  
1552 processing, quality control, normalization and visualization of single-cell  
1553 RNA-seq data in R. *Bioinformatics* **33**, 1179–1186 (2017).
- 1554 52. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell  
1555 gene expression data analysis. *Genome Biology* **19**, 15 (Dec. 2018).
- 1556 53. Welch, J. D. *et al.* Single-Cell Multi-omic Integration Compares and Con-  
1557 trasts Features of Brain Cell Identity. *Cell* **177**, 1873–1887.e17 (2019).

## 1558        9 Appendix A: Python and R library/package 1559        version numbers

1560        The following two environments (9.1, 9.2) were used in the benchmarking and  
1561        analysis phases, where all integration experiments and downstream analysis  
1562        tests were done with the pipeline environment (9.1), and all results analysis  
1563        and plotting was done with the analysis environment (9.2). The only exception  
1564        were the balanced metric analyses and tests (7), which utilized the analysis  
1565        environment (9.2) for generation and testing of the various scenarios outlined.

1566        Configurations for these environments are also available at <https://github.com/hsmaan/Iniquitate/tree/main/workflow/envs>.

1568        The library for the balanced metrics (7) was developed independently, and  
1569        all of the information on dependency versions is available at <https://github.com/hsmaan/balanced-clustering>.

### 1571        9.1 Iniquitate pipeline environment

- 1572            • python>=3.7,<=3.10
- 1573            • numpy>=1.19.0
- 1574            • pandas>=1.2.0
- 1575            • scipy>=1.5.0
- 1576            • leidenalg>=0.8.0
- 1577            • umap-learn>=0.5.0
- 1578            • mnnpy>=0.1.9
- 1579            • scikit-learn>=1.0.1
- 1580            • scanpy=1.8.2
- 1581            • anndata=0.8.0

- 1582           • faiss-cpu>=1.7.0
- 1583           • pytorch=1.10.1
- 1584           • torchmetrics<=0.6.0
- 1585           • cudatoolkit=10.2
- 1586           • scvi-tools=0.14.4
- 1587           • bbknn=1.5.1
- 1588           • harmonypy=0.0.5
- 1589           • scanorama=1.7.1
- 1590           • r-base>=4.0.0
- 1591           • r-liger=0.5.0
- 1592           • r-seurat=4.0.6
- 1593           • r-seuratdisk>=0.0.9
- 1594           • r-data.table>=1.14.0
- 1595           • r-reticulate=1.24
- 1596           • cython>=0.29.25
- 1597           • r-rann=2.6.1

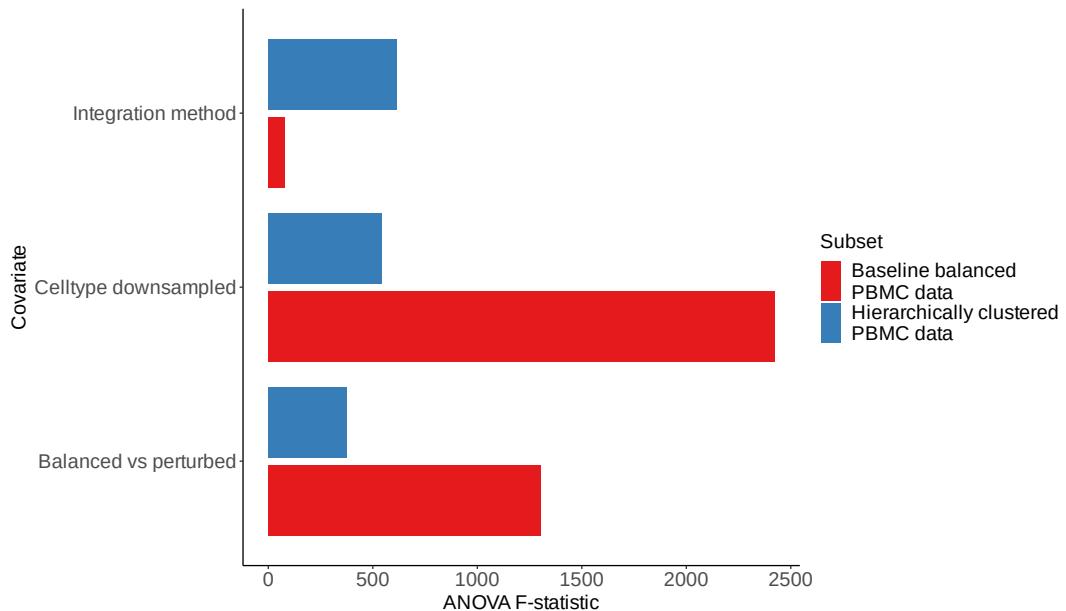
- ## 1598        9.2 Analysis scripts environment
- 1599           • python>=3.7,<=3.10
  - 1600           • numpy>=1.19.0
  - 1601           • pandas>=1.2.0
  - 1602           • scipy>=1.5.0
  - 1603           • seaborn>=0.11.2

- 1604           ● plotnine>=0.8.0
- 1605           ● leidenalg>=0.8.0
- 1606           ● umap-learn>=0.5.0
- 1607           ● scikit-learn>=1.0.1
- 1608           ● scanpy=1.8.2
- 1609           ● anndata>=0.7.5
- 1610           ● ipykernel>=6.4.0
- 1611           ● jupyterlab>=3.2.9
- 1612           ● notebook>=6.4.2
- 1613           ● scvi-tools=0.14.4
- 1614           ● pytorch=1.10.1
- 1615           ● torchmetrics<=0.6.0
- 1616           ● cudatoolkit=10.2
- 1617           ● bbknn=1.5.1
- 1618           ● harmonypy=0.0.5
- 1619           ● scanorama=1.7.1
- 1620           ● r-base>=4.0.5
- 1621           ● r-seurat>=4.0.5
- 1622           ● r-data.table>=1.14.0
- 1623           ● r-ggplot2>=3.3.0
- 1624           ● r-tidyverse>=1.2.1
- 1625           ● r-reshape2>=1.4.3
- 1626           ● r-data.table>=1.14.0

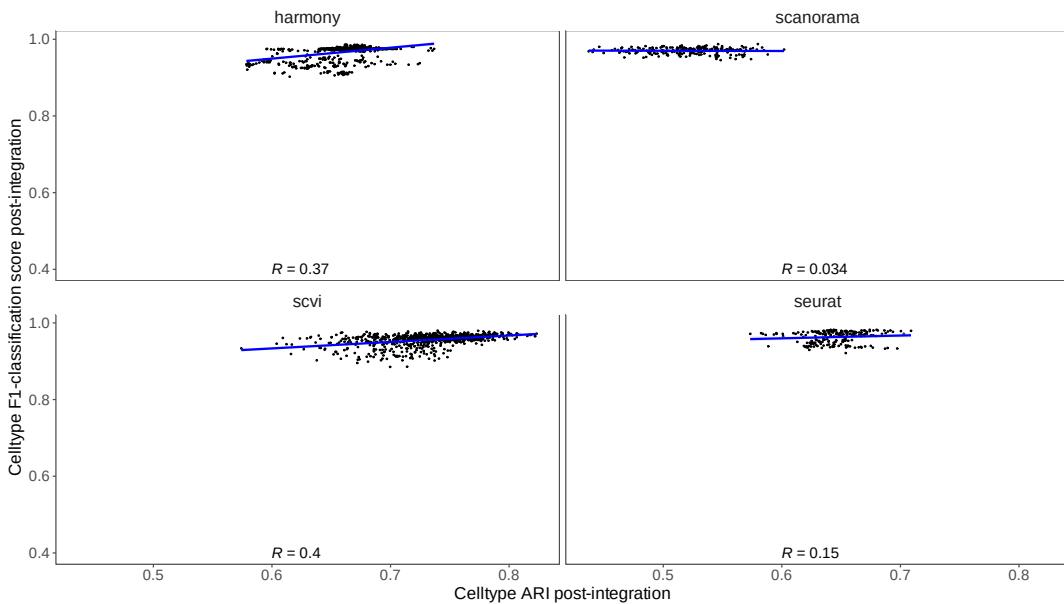
- 1627           ● r-ggthemes>=4.2.0
- 1628           ● r-ggextra>=0.8.0
- 1629           ● r-dotwhisker>=0.7.4
- 1630           ● r-seuratdisk>=0.0.9019
- 1631           ● r-deldir>=1.0.2
- 1632           ● r-ggpubr>=0.4.0
- 1633           ● r-cowplot>=1.1.1
- 1634           ● r-ggrepel>=0.9.1
- 1635           ● r-rcolorbrewer>=1.1
- 1636           ● r-ggbump>=0.1.0
- 1637           ● bioconductor-complexheatmap<=2.9.0
- 1638           ● r-venndiagram>=1.7.1
- 1639           ● r-multipanelfigure>=2.1.2
- 1640           ● r-gridextra>=2.3
- 1641           ● r-cairo>=1.5
- 1642           ● r-lemon>=0.4.5
- 1643           ● r-networkd3>=0.4
- 1644           ● r-emt>=1.2
- 1645           ● cython>=0.29.25

1646

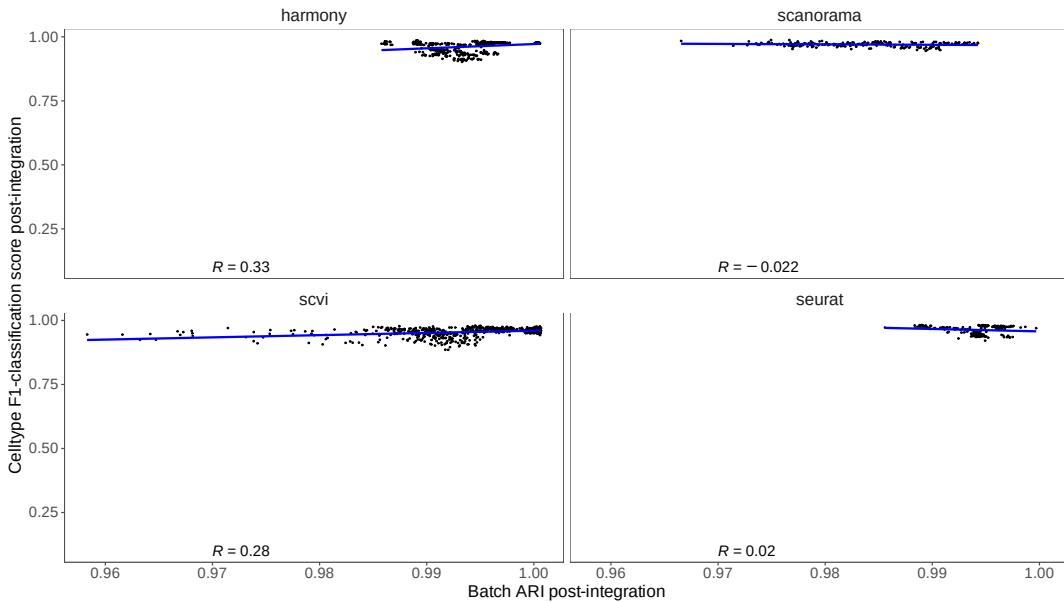
## Supplementary Figures



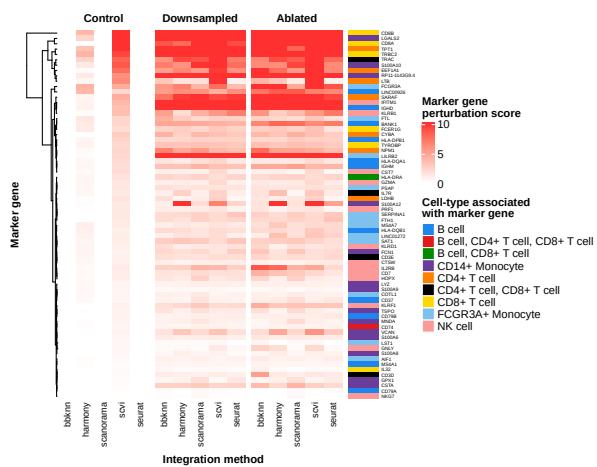
**Figure S1: ANOVA F-statistic values for cell-type specific KNN classification in the baseline and hierarchical 2 batch balanced PBMC data.** The ANOVA F-statistic values, indicating the ratio of variation between between sample means and variation within the samples themselves, for the covariates used in the KNN-classification task ANOVA for the 2 batch PBMC balanced dataset (Online Methods). F-statistics are shown for integration method (first covariate in model), cell-type that was downsampled (second covariate in model), and which type of experiment was performed (control balanced vs. perturbed - last covariate in model). The F-statistics are compared between the baseline setup (6 cell-types initially utilized) and the hierarchical setup after merging closely related cell-types.



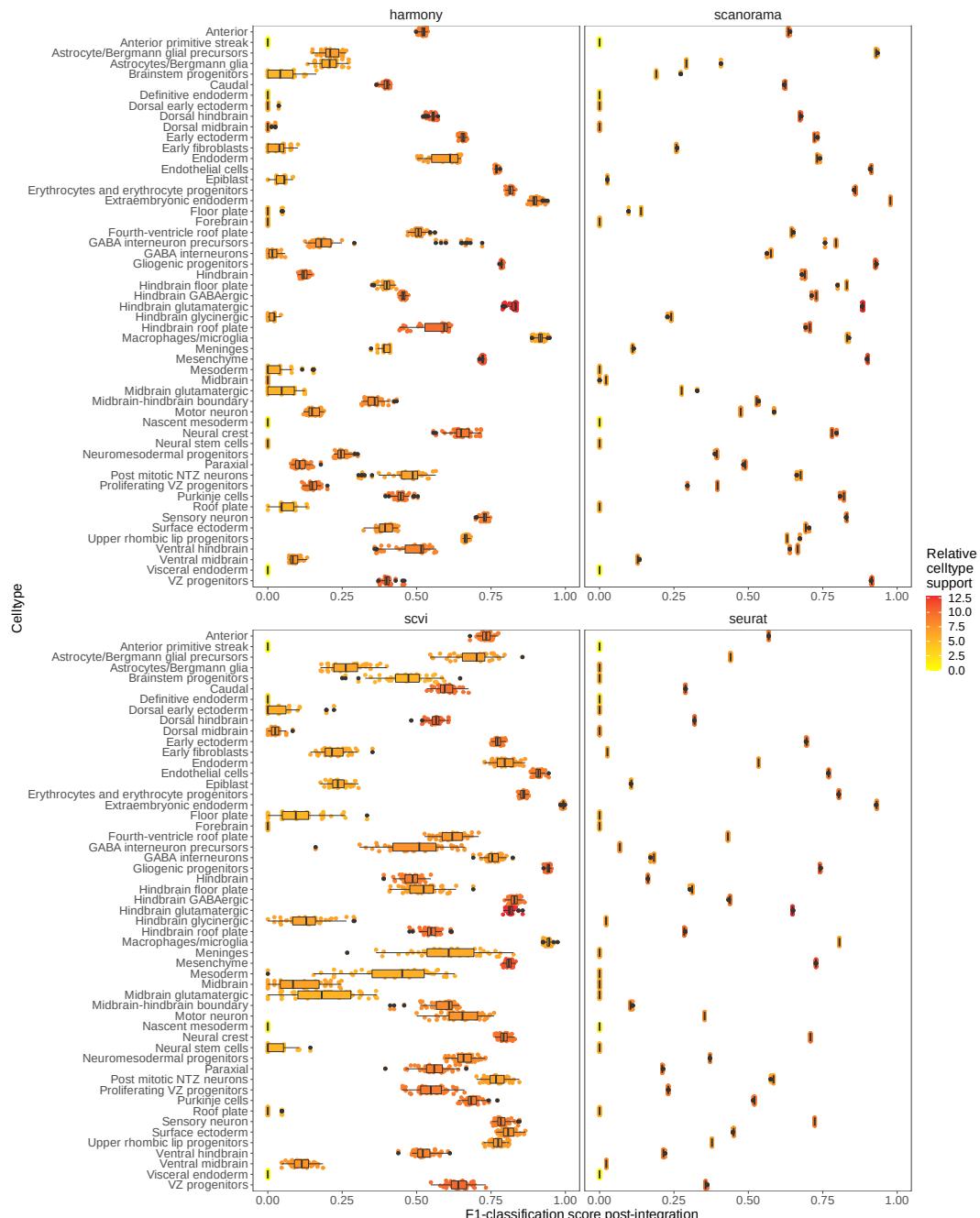
**Figure S2: Correlation between cell-type-specific F1-classification scores and cell-type  $ARI_{cell-type}$  in balanced 2 batch PBMC data.** All experiments (control, downsampling, and ablation) are indicated for the baseline 2 batch PBMC data. For perturbation experiments, values are subset for only where the cell-type being classified (F1-classification score) is equivalent to the cell-type that was down-sampled. The median cell-type classification F1-score across all cell-types is shown, grouped by method and experiment type, for direct comparison with the ARI values which are calculated per replicate. The Spearman correlation value between the scores is indicated.



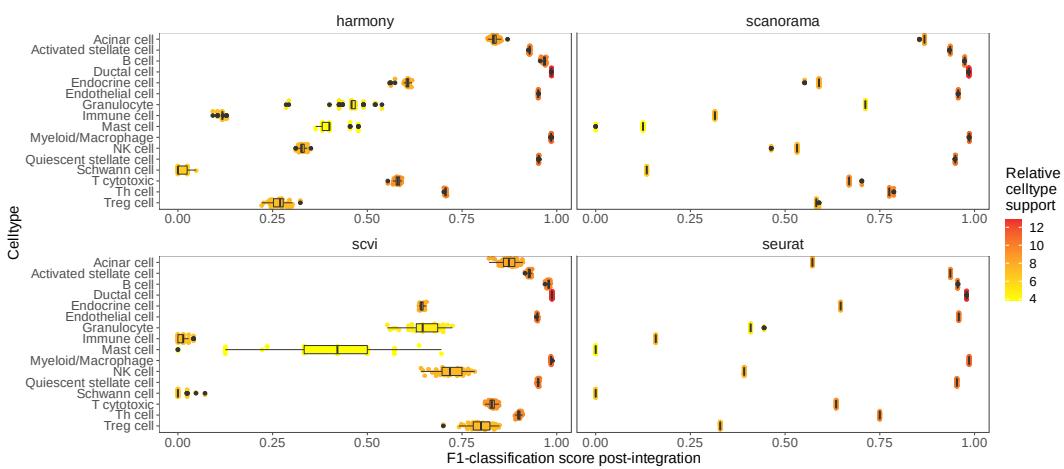
**Figure S3: Correlation between cell-type-specific F1-classification scores and  $(1 - \text{ARI}_{\text{batch}})$  values in balanced 2 batch PBMC data.** All experiments (control, downsampling, and ablation) are indicated for the baseline 2 batch PBMC data. For perturbation experiments, values are subset for only where the cell-type being classified (F1-classification score) is equivalent to the cell-type that was down-sampled. The median cell-type classification F1-score across all cell-types is shown, grouped by method and experiment type, for direct comparison with the ARI values which are calculated per replicate. The Spearman correlation value between the scores is indicated.



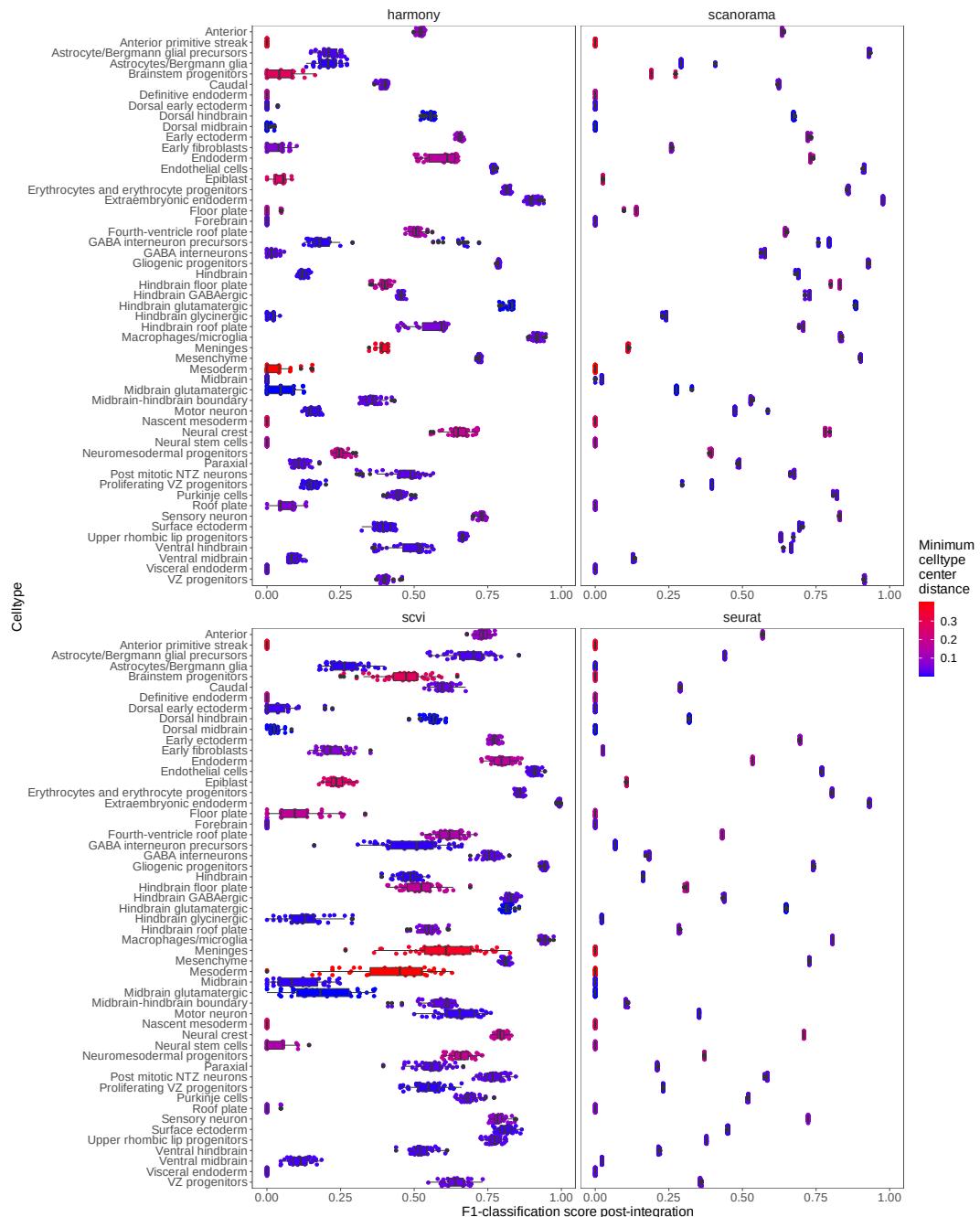
**Figure S4: marker gene perturbation scores across all marker genes for the cell-types in the balanced 2 batch PBMC dataset.** Marker genes were determined through differential gene expression analysis within each batch (Online Methods), and their perturbation score, indicating change in maximum ranking across unsupervised clusters post-integration are shown across control, downsampling, and ablation experiments (Online Methods). Note that downsampling and ablation (perturbation) experiments are not subset here for the marker gene being analyzed and its associated cell-type (e.g. maximum-rank change for B-cell markers in only runs where B-cells are downsampled).



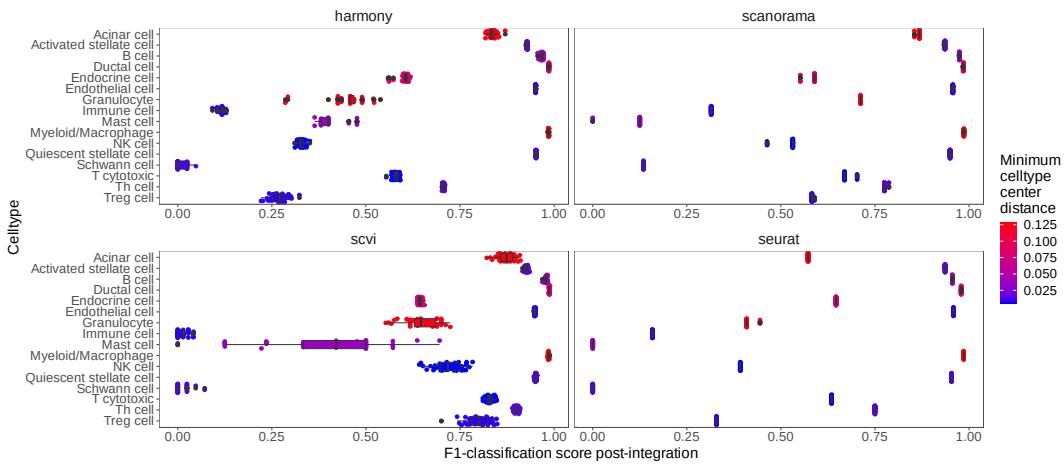
**Figure S5: Comparison of F1-classification accuracy and relative cell-type support of each cell-type in the imbalanced 6 batch mouse hindbrain development dataset.** The relative cell-type support is based on the number of cells in the integrated embedding space present for each cell-type (Online Methods).



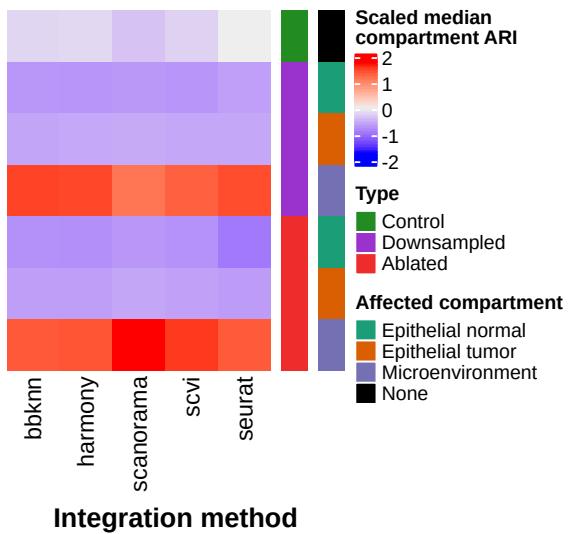
**Figure S6: Comparison of F1-classification accuracy and relative cell-type support of each cell-type in the imbalanced 8 batch PDAC dataset.** The relative cell-type support is based on the number of cells in the integrated embedding space present for each cell-type (Online Methods).



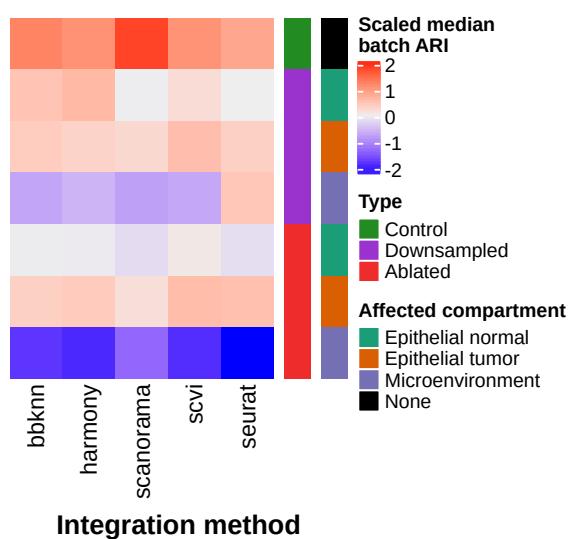
**Figure S7: Comparison of F1-classification accuracy and minimum cell-type center distance of each cell-type in the imbalanced 6 batch mouse hindbrain development dataset.** The minimum cell-type center distance value indicates how close is the closest other cell-type across batches in PCA space (Online Methods).



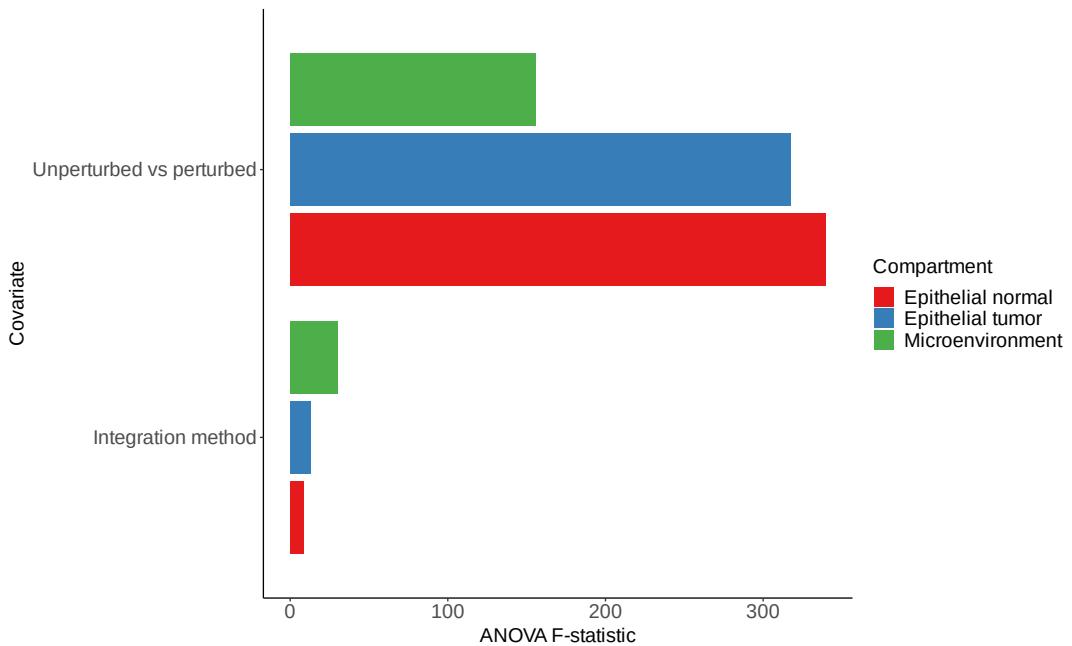
**Figure S8: Comparison of F1-classification accuracy and minimum cell-type center distance of each cell-type in the imbalanced 8 batch PDAC dataset.** The minimum cell-type center distance value indicates how close is the closest other cell-type across batches in PCA space (Online Methods).



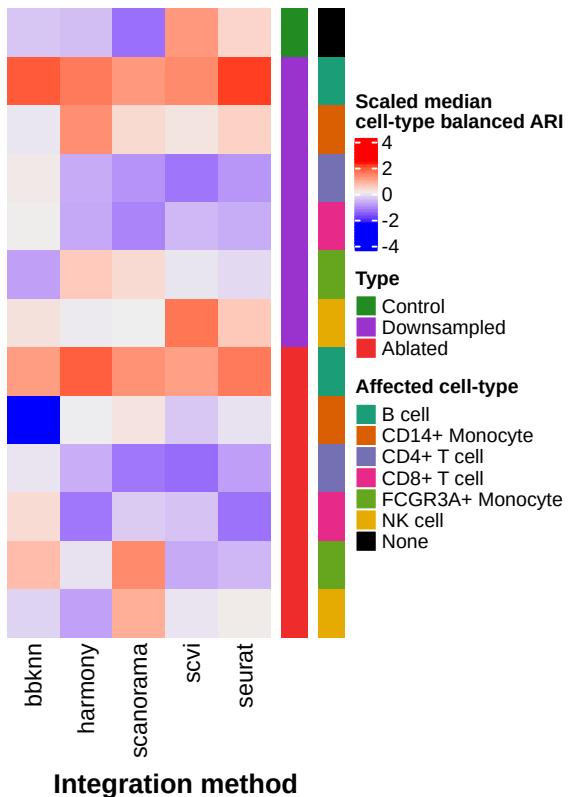
**Figure S9: Comparison of compartment heterogeneity conservation ARI results across PDAC data perturbation experiments.** Z-score normalized median  $ARI_{compartment}$  (compartment integration accuracy) results across experiment type (control, compartment downsampling, compartment ablation), specific-compartment downsampled, and integration method utilized.



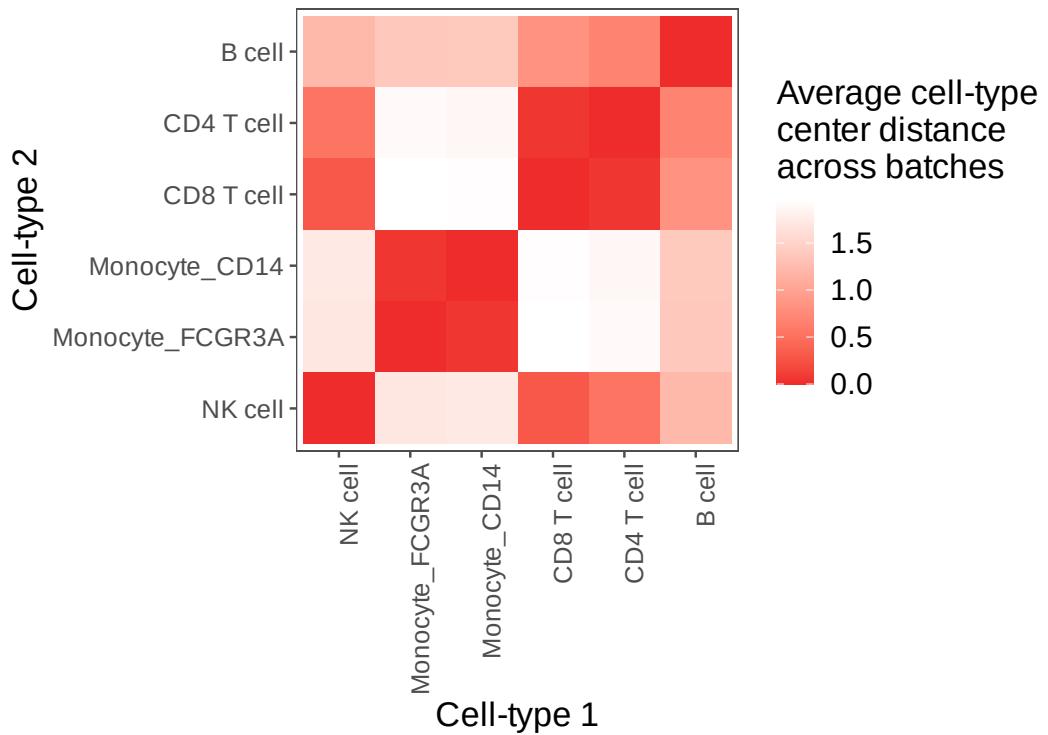
**Figure S10: Comparison of batch-mixing ARI results across PDAC data perturbation experiments.** Z-score normalized median ( $1 - \text{ARI}_{\text{batch}}$ ) (batch mixing) results across experiment type, compartment downsampled, and integration method.



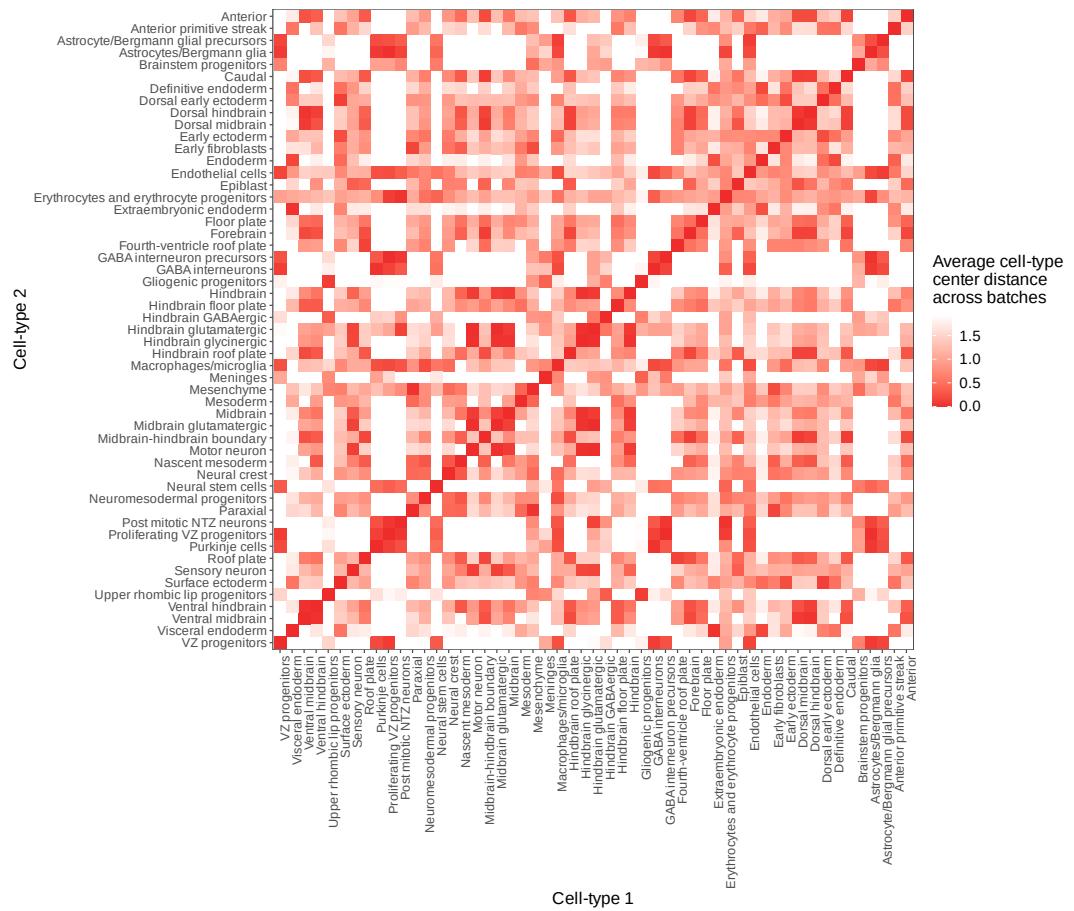
**Figure S11: ANOVA F-statistic values for compartment specific KNN-classification in the 8 batch compartmentalized PDAC data.** The ANOVA F-statistic values, indicating the ratio of variation between sample means and variation within the samples themselves, for KNN-classification of individual compartments before and after perturbation (Online Methods). F-statistics are shown for integration method (first covariate in model), and which type of experiment was performed (control vs. perturbed - last covariate in model). The ANOVA tests were performed individually for each compartment, and the compartment-specific F-statistics are shown. Note that the perturbations here are specific to the compartment being analyzed (e.g. microenvironment subset will only contain perturbations that targeted the microenvironment) (Online Methods).



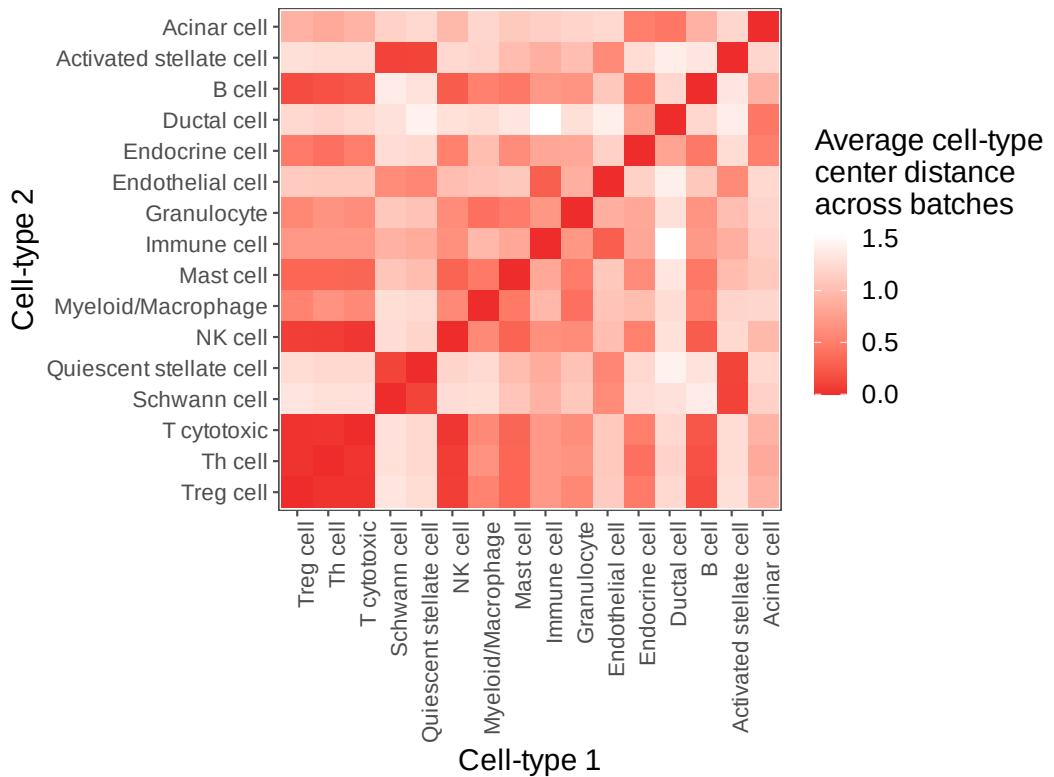
**Figure S12: Comparison of cell-type heterogeneity conservation ARI results across the PBMC 2 batch perturbation experiments, using the balanced ARI (bARI) score.** Z-score normalized median  $ARI_{cell-type}$  (cell-type integration accuracy) results across experiment type (control, compartment down-sampling, compartment ablation), specific-cell-type downsampled, and integration method utilized, using the bARI instead of the base ARI metric.



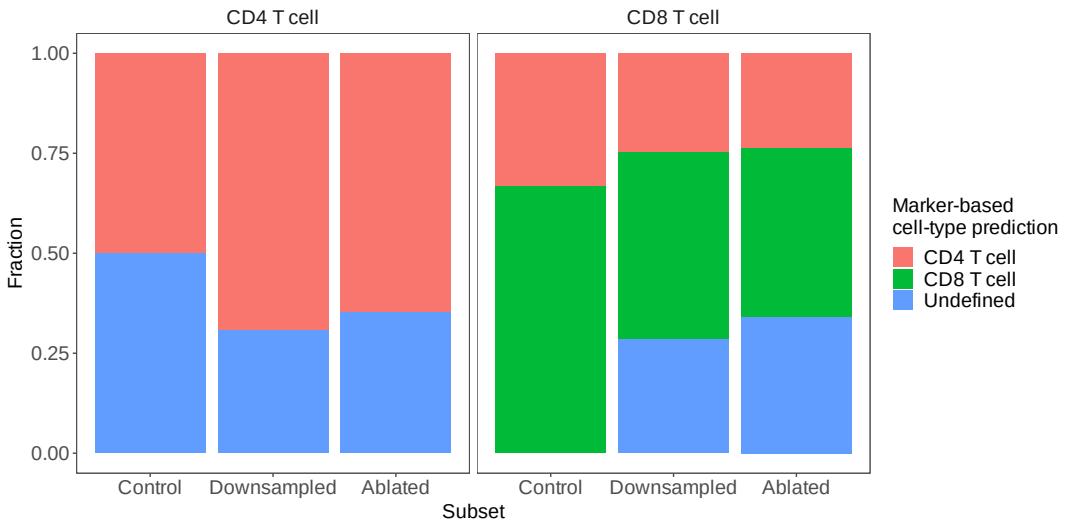
**Figure S13: Average cell-type center distance between cell-types in the balanced PBMC 2 batch dataset.** For each batch, the distance from the centers of cell-type clusters in principal component analysis (PCA) reduction space are calculated, and the relative distances between cell-types are determined and averaged across batches (Online Methods).



**Figure S14: Average cell-type center distance between cell-types in the 6 batch mouse hindbrain development dataset.** For each batch, the distance from the centers of cell-type clusters in principal component analysis (PCA) reduction space are calculated, and the relative distances between cell-types are determined and averaged across batches (Online Methods).



**Figure S15: Average cell-type center distance between cell-types in the 8 batch pancreatic ductal adenocarcinoma (PDAC) dataset.** For each batch, the distance from the centers of cell-type clusters in principal component analysis (PCA) reduction space are calculated, and the relative distances between cell-types are determined and averaged across batches (Online Methods).



**Figure S16: Predicted cell-types for CD4/CD8 T-cell majority clusters in the balanced PBMC 2 batch data, based on marker gene differential expression.** In this setup, the top 50 marker genes were analyzed based on differential expression for unsupervised clusters from the Seurat integration method across experimental subsets (Control, Downsampling, Ablation). The Downsampling and Ablation subsets here contain only instances where CD4+ and CD8+ T cells were affected. Only clusters that contained a majority of cells (based on ground-truth annotations) of CD4+ or CD8+ T were kept. The canonical marker genes for CD4+ T cells (IL7R) and CD8+ T cells (CD8A) were used to predict the cell-type for each cluster based on their relative ranking in the top 50 marker genes for the given clusters (Details in Online Methods: Downstream analysis - marker gene ranking - Case study - CD4/CD8 T cell assignment based on marker genes). The fraction of clusters that contain a majority of CD4+ or CD8+ T cells, and their predicted cell-type based on the aforementioned marker gene setup are indicated, across experimental subsets.