1  **Benchmarking algorithms for joint integration of unpaired and paired single-cell**

2  **RNA-seq and ATAC-seq data**

3

4  Michelle Y. Y. Lee[1,3], Klaus H. Kaestner[1], Mingyao Li[2]

5

6  1. Department of Genetics, University of Pennsylvania, Philadelphia, PA 19104, USA

7  2. Department of Biostatistics, Epidemiology and Informatics, Perelman School of Medicine,

8     University of Pennsylvania, Philadelphia, PA, USA

9  3. Graduate Group in Genomics and Computational Biology, University of Pennsylvania

10    Perelman School of Medicine, Philadelphia, Philadelphia, PA, 19104, US

11  Address correspondence to: kaestner@pennmedicine.upenn.edu or

12  mingyao@pennmedicine.upenn.edu

# Abstract

Single-cell RNA-sequencing (scRNA-seq) measures gene expression in single cells, while single-nucleus ATAC-sequencing (snATAC-seq) enables the quantification of chromatin accessibility in single nuclei. These two data types provide complementary information for deciphering cell types/states. However, when analyzed individually, scRNA-seq and snATAC-seq data often produce conflicting results regarding cell type/state assignment. In addition, there is a loss of power as the two modalities reflect the same underlying cell types/states. Recently, it has become possible to measure both gene expression and chromatin accessibility from the same nucleus. Such paired data make it possible to directly model the relationships between the two modalities. However, given the availability of the vast amount of single-modality data, it is desirable to integrate the paired and unpaired single-modality data to gain a comprehensive view of the cellular complexity. Here, we benchmarked the performance of seven existing single-cell multi-omic data integration methods. Specifically, we evaluated whether these methods are able to uncover peak-gene associations from single-modality data, and to what extent the multiome data can provide additional guidance for the analysis of the existing single-modality data. Our results indicate that multiome data are helpful for annotating single-modality data, but the number of cells in the multiome data is critical to ensure a good cell type annotation. Additionally, when generating a multiome dataset, the number of cells is more important than sequencing depth for cell type annotation. Lastly, Seurat v4 is the best at integrating scRNA-seq, snATAC-seq, and multiome data even in the presence of complex batch effects.

# Background

Over the past ten years, hundreds of single-cell RNA-seq (scRNA-seq) (for transcript abundance in single cells) or single-nucleus ATAC-seq (snATAC-seq) (for chromatin accessibility in single nuclei) have been produced by laboratories worldwide, leading to the discovery of new cell types and regulatory circuits. In addition, by applying single-cell assays to two-state models such as the comparison between control and mutant tissues, changes in gene expression or chromatin accessibility caused by a gene mutation could be analyzed at the cell type-specific level easily for the first time. Unfortunately, each single-modality dataset measures either the gene expression or the chromatin accessibility of a given cell. Although the two datasets are generated from the same cell population, they measure different cells. Most of the time, the two experimental modalities result in the identification of similar cell types, as the promoters of highly expressed genes used to define cell types at the transcript levels are frequently also identified as highly accessible by the ATAC-seq modality. However, there are situations in which the two profiles are discordant. In these situations, simultaneous, joint profiling of gene expression and chromatin accessibility is paramount for resolving inconsistency and revealing novel cell types and states that show modality-specific features. Moreover, the joint profiling of gene expression and chromatin accessibility of the same exact cells offers the most direct link between *cis*-regulatory elements and their target genes [1].

Recently, the simultaneous determination of both transcript levels and chromatin state in the same nucleus has become possible, using so-called "multi-omics" approaches. An example is the 10x Genomics single cell Multiome ATAC + gene expression technology [2]. Multi-omics datasets are clearly superior at refining cell types and revealing gene regulatory networks [1]. However, it is not practical to repeat all prior studies of interest performed using the single-modality assays with the multiome approaches, as frequently precious samples are either no longer available or funding is limited. Therefore, it is highly desirable to integrate pre-existing single-modality scRNA-seq and snATAC-seq datasets with multiome data generated subsequently using the newer technology to achieve more accurate cell type annotations.

Several methodologies have been developed for multi-omic data integration. Here, we refer to multi-omic integration as the integration of RNA-seq and ATAC-seq profiles measured in single cells, either with or without the guidance of multiome data. These methods attempt to align cells profiled by separate technologies and project them into one common low-dimensional

69    space to ensure consistent cell type calling. However, we still lack an objective evaluation of

70    whether the addition of the multiome data improves the annotation of single-modality datasets.

71    Furthermore, some of the methods try to impute the missing modality for the single-modality

72    datasets and identify peak-gene pairs using these 'pseudo-paired' datasets. Thus, it is still

73    uncertain if the imputed missing modality can truly provide additional biological insights to the

74    same degree as provided by the experimentally produced multiome datasets. Finally, given the

75    availability of many methods for multi-omic data integration, at present, we do not know which

76    method performs the best when integrating all three data types.

77

78        The current multi-omic integration methods can be divided into two categories. Methods

79    in the first category perform multi-omic integration using only the single-modality datasets,

80    aiming to find a mapping between gene expression profiles and chromatin accessibility states to

81    create an aligned space that explains both modalities; we call these approaches 'unpaired

82    integration'. Representative methods in this category include Seurat version 3 (Seurat v3) [3],

83    which performs canonical correlation analysis (CCA) to align experimentally measured gene

84    expression with pseudo-gene expression obtained from chromatin accessibility. One example of

85    pseudo-gene expression is the gene activity score, calculated by summing up peak counts

86    within the gene body plus 2kb upstream in the ATAC-seq data. LIGER [4] also uses the gene

87    expression and gene activity score to obtain shared features between the two modalities and

88    then derives a low-dimensional embedding through a non-negative matrix factorization

89    approach. FigR [5] aligns the snATAC-seq and scRNA-seq data using a CCA-based approach.

90    In addition, it provides matching of snATAC-seq and scRNA-seq cells, which enables the

91    identification of *cis*-regulatory elements similar to paired multiome data. BindSC [6] goes beyond

92    the simple construction of gene activity scores. Instead, bindSC uses a bi-directional CCA to

93    empirically construct a cell-by-gene matrix for the snATAC-seq cells that preserve its similarity

94    with the ATAC-seq input and simultaneously maximizes the correlation with the scRNA-seq

95    matrix it is being integrated with.

96

97        Methods in the second category encompass more recent approaches that incorporate

98    information from multiome cells and integrate all three data types for a more comprehensive

99    exploration of cellular identities; we term these approaches 'multiome-guided integration'.

100    Representative methods in this category include Seurat version 4 (Seurat v4) [7], an approach

101    that first learns a low-dimensional representation of the cells profiled by the multiome

102    methodology using both the RNA-seq and ATAC-seq profiles by weighted nearest neighbors

103    (WNN) analysis [7]. Subsequently, the two single-modality datasets are projected onto the WNN

104    embedding space in a supervised manner. MultiVI [8] and Cobolt [9] use a deep-learning

105    approach called 'variational autoencoder' to embed all three data types. Both methods employ

106    the encoder-decoder system to learn a low-dimensional representation of the data. Specifically,

107    two encoders and two decoders are set up, one for each modality. However, there are different

108    model choices. MultiVI assumes a negative binomial distribution for the RNA-seq data and a

109    Bernoulli distribution for the ATAC-seq data, while Cobolt assumes a Multivariate Normal

110    distribution for both modalities. Furthermore, the two methods integrate the modality-specific

111    representation for the paired cells differently. MultiVI first aligns the two embeddings through a

112    symmetric Kullback-Leibler (KL) divergence loss and then obtains an average of the two

113    embeddings. On the other hand, Cobolt simply multiplies the two embeddings to represent the

114    paired cells, while the representation of the unpaired cells is first generated by the

115    corresponding encoder and refined using a linear transformation to ensure enough similarity

116    between the RNA-seq derived embedding and the ATAC-seq derived embedding.

117

118        All methods described above aim to project cells from different data types into one

119    shared space to facilitate the identification of cell types through clustering. Nevertheless, a

120    common goal for studies profiling chromatin accessibility and gene expression at the single-cell

121    level is to understand cell type-specific *cis*-regulatory logic. Since the two single-modality

122    datasets are generated from different cells in a given population, albeit representing the same

123    cell types, the single-modality datasets cannot be naïvely combined to test for association

124    between chromatin accessibility and gene expressions. Therefore, multiple efforts have

125    attempted to impute the missing modality for the single-modality datasets, aiming to

126    computationally generate paired profiles similar to those measured experimentally by the

127    multiome technology. Some methods mentioned above, e.g., Seurat v3, FigR, bindSC, Seurat

128    v4, and MultiVI, are capable of this task. However, an objective evaluation of how reliable the *in-*

129    *silico* imputed profiles are compared to what is directly measured by the paired multiome

130    technologies is still lacking. Therefore, we aimed to conduct an extensive benchmarking

131    analysis to evaluate the above-mentioned methods by addressing two important questions. First,

132    do multiome data help the integration of single-modality datasets? Second, what is the best

133    computational method for the integration of scRNA-seq, snATAC-seq, and multiome data?

134

135    **Results**

**Overview of the benchmarking scheme and evaluation strategies**

The overall workflow of our benchmarking evaluations is summarized in Figure 1. Figure 1A illustrates our approach to evaluate whether multiome data integration can improve the value of single-modality datasets, while Figure 1B outlines how we assess the effectiveness of each integration methods, at various conditions of the multiome dataset. To answer the proposed questions, we simulated situations where all three data types are available by using two publicly available multiome datasets [10, 11]. The first multiome dataset [10] profiled 10,085 peripheral blood mononuclear cells (PMBCs) and represents a simple biological system, because PBMCs can be easily divided into seven well-separated cell types (Supplementary Figure 1A). The second dataset profiled bone marrow mononuclear cells (BMMC) [11], an example of highly complex cell populations. BMMCs are closely related to each other transcriptionally, and contain, for example, myeloid progenitors and their closely related descendants, CD16+ and CD14+ monocytes (Supplementary Figure 1B). The individual BMMC cell types are therefore much harder to separate compared to the PBMC populations, thus allowing us to thoroughly evaluate the performance of each method in both simple and complex biological systems. Moreover, the BMMC dataset is composed of samples generated from four research sites and nine donors [12], which enables the analysis of batch effects and technical replicates.

We evaluated four popular unpaired integration methods (Seurat v4, LIGER, FigR, bindSC), and three multiome-guided integration methods (Seurat v4, MultiVI, and Cobolt). To account for the increased power resulting from the larger absolute number of cells employed during the integration process by the multiome-guided methods, we created another scenario termed 'unpaired (multiome-split)' in which the RNA-seq and ATAC-seq data from the multiome samples were treated as independent datasets and appended to the single-modality datasets. This category again includes the four unpaired-integration methods, the only difference being that the single-modality datasets now include additional single-modality cells that were converted from the multiome cells.

To evaluate the performance of each method for cell type identification, we performed Louvain clustering [13] on the integrated embedding. For methods capable of missing modality imputation, we imputed gene expression using snATAC-seq profiles. We then evaluated the integration results in four aspects as shown in Figure 1A. Specifically, we evaluated cell type annotation accuracy using two metrics: Adjusted Rand Index (ARI) [14] and Normalized Mutual Information (NMI) [15]. Both metrics range from 0 to 1, with 1 being the best. The detailed

6

170 approach is described in the Methods section. The accuracy of cell type annotation depends on

171 the number of cell clusters identified; therefore, an additional way to measure data integration

172 quality is via the accuracy of cell type separation. Using the ground-truth annotation, we

173 evaluated how well cells of different identities are separated, using a cell type specific average

174 silhouette width (ASW) [16] and a cell type Local Inverse Simpson's Index (cLISI) [17, 18].

175 Furthermore, because the three data types could have technology-specific differences, we used

176 a batch ASW [16] and the k-nearest neighbor batch effect test (kBET) [16] to measure batch-

177 mixing of the integrated results. These four measurements were normalized to be in the range

178 of 0 and 1 in which 1 is the best result, namely high separation between cell types and complete

179 mixing of data batches.

180

181 We also evaluated the quality of 'peak to gene pair' predictions by assessing the

182 accuracy of assigning an ATAC-seq peak to a specific gene. Using the measured ATAC-seq

183 and imputed RNA-seq data, we computed the percentage of significant peak-gene pairs

184 recovered as compared to a ground truth obtained using all cells in the multiome dataset. To

185 penalize for the presence of false positives reported by the data integration methods, we also

186 calculated an F1 score [15], which normalized the absolute percent recovery of the true peak-

187 gene pairs by the occurrence of false positive and false negative relationships.

188

189 **Do Multiome data improve the annotation of single-modality datasets?**

190

191 **PBMC**

192 To answer if multiome data improve the analysis of single-modality datasets (scRNA-seq and

193 scATAC-seq), we first simulated the situation with 1,000 scRNA-seq cells and 1,000 snATAC-

194 seq cells based on the PBMC data. These single-modality cells were integrated using each of

195 the four unpaired integration methods. To evaluate if multiome data improve the analysis of

196 single-modality datasets, we considered the situation where we have a multiome dataset,

197 potentially with different numbers of cells (e.g., 1000, 3000, or 8000). These multiome data were

198 integrated with the single-modality datasets using the multiome-guided methods. However,

199 because the number of cells used during clustering and gene expression imputation impacts the

200 clustering accuracy and peak-gene association identification, we ran the unpaired integration

201 methods again, this time treating the multiome dataset as single-modality cells and adding them

202 to the existing single-modality data. Here, any increase in performance is solely caused by the

203 increase in cell number; the results from these evaluations are labeled as the 'unpaired

7

204    (multiome-split)' category. For each simulation, we randomly drew the cells from the 10,085

205    PBMC dataset and each condition was repeated five times. The parameters used for this

206    simulation are summarized in Figure 2A.

207

208         The evaluation result for each method is summarized in Figure 2B. Without the

209    incorporation of multiome data, the cell type annotation accuracy was already good, being 0.81

210    in ARI and 0.81 in NMI when integrating the unpaired data using the bindSC program (Figure

211    2B). Surprisingly, in the presence of 1,000 multiome cells, the multiome-guided approaches

212    performed worse than simply integrating the data from the 2,000 single-modality cells by

213    themselves (Figure 2B). This unexpected result is likely caused by the fact that 1,000 multiome

214    cells alone do not achieve good cell type separation, which is a critical requirement for the

215    multiome-guided methods to succeed. However, when we used 3,000 or 8,000 multiome cells,

216    Seurat v4, one of the multiome-guided methods, achieved the best results in terms of cell type

217    annotation (Figure 2B). Furthermore, when comparing the multiome-guided results with

218    unpaired (multiome-split) results, the performance of Seurat v4 remained higher when there are

219    3,000 or 8,000 cells (Figure 2B). Thus, our findings indicate that the multiome data can improve

220    cell type annotation of the single modality datasets, provided that there is a sufficient number of

221    multiome cells available.

222

223         Next, we evaluated the performance of each method in predicting peak-gene pairs.

224    Peak-gene pairs are calculated using 1,000 measured chromatin accessibility profiles and the

225    corresponding 1,000 imputed gene expression profiles. Here, we compared predicted peak-

226    gene pairs to the ground-truth list calculated using multiome cells in the full PBMC data. Seurat

227    v3 performed very well at recovering the absolute number of peak-gene pairs, and the

228    incorporation of data from multiome cells through splitting only marginally increased the

229    performance (Figure 2B). BindSC had a slightly better F1 score than Seurat v3, meaning that

230    the Seurat v3 results contained more false positives (Figure 2B). For the multiome-guided

231    methods, the more multiome cells available during gene expression imputation resulted in

232    higher peak-gene pair recovery (Figure 2B). Nevertheless, the incorporation of data from

233    multiome cells using the multiome-guided methods did not perform better than the unpaired

234    methods, with the exception that the F1 score was higher in MultiVI (Figure 2B).

235

236         The number of cells used for predicting peak-gene pairs influences the accuracy. To

237    give a general idea of how well the predicted gene expression profiles are, we compared the

238   peak-gene pair identification result to the one obtained using the real paired profiles. We

239   included a red dashed line in Figure 2B to indicate the percentage of peak-gene pair recovery

240   and F1 score calculated using the measured, paired gene expression and chromatin

241   accessibility profiles of the 1,000 cells being evaluated, instead of the gene expression profile

242   imputed from chromatin accessibility. What's surprising is that the in-silico prediction profile from

243   Seurat v3 revealed a higher percentage of recovered peak-gene pairs and a better F1 score

244   than the measured paired gene expression and chromatin accessibility profile from 1,000 cells.

245   This is likely due to the dropout issue common to single-cell assays and the predicted RNA

246   profile can borrow information from similar cells, thus recovering the trend better. However, we

247   also note that the predicted profiles only recovered less than 45% of the ground-truth list

248   calculated using the full PBMC data with 10,412 cells. Although the predicted profiles are better

249   than the measured gene expression profiles, it is only recovering a small percentage of peak-

250   gene pairs revealed by the experimentally generated multiome dataset.

251

252   **BMMC**

253   Having evaluated the various data integration platforms with the PBMC data, which represent a

254   low-complexity situation with clearly defined major cell types, we next sought to determine how

255   the different methodologies perform when analyzing data from highly complex cell populations,

256   as is the case for bone marrow mononuclear cells (BMMC). Here, to avoid complexity caused

257   by batch differences, we only used 6,740 multiome cells from one sample (site 1 donor 2). We

258   again started with 1,000 scRNA-seq and 1,000 snATAC-seq cells, and then tested the result

259   when incorporating 1000, 2000, and 4000 multiome cells, composed of 21 cell types (Figure 2A).

260   In this biological system, we found that including a larger number of multiome cells improved

261   cell type annotation, with Seurat v3 performing the best among the unpaired (multiome-split)

262   methods (Figure 2C). Among the multiome-guided methods, Seurat v4 achieved the best

263   performance when the input data included 4,000 multiome cells. Remarkably, when we

264   employed data from only 1,000 or 2,000 multiome cells, all multiome-guided methods performed

265   worse than when inputting the multiome data as two separate, unpaired modalities (Figure 2C).

266   A similar trend was observed in the peak-gene pair prediction (Figure 2C). The likely reason

267   causing the poor performance of the multiome-guided methods is the limited quality of multiome

268   data and the high complexity of the biological system being profiled. Note that peak-gene

269   prediction recovery and F1 score obtained via the unpaired Seurat v3 algorithm are still higher

270   than the association calculated from the observed multiome profile indicated by the red dash

271   line in Figure 2C.

272

### **Comparison of run time and visualization of integration**

274  Another important issue to consider when comparing various computational approaches is the

275  computation time needed to complete a given task. All methods were run with 8 CPU cores and

276  64GB of RAM. Figure 2D shows the runtime, measured in seconds. Unpaired methods all have

277  similar runtimes, and the increase in the unpaired (multiome-split) category was due to the

278  incorporation of the additional data from multiome experiments. Importantly, the multiome-

279  guided methods vary greatly in runtime and thus costs. Cobolt was the fastest method, but

280  unfortunately, it exhibited comparatively low clustering accuracy and peak-gene recovery.

281  Seurat v4 had a shorter runtime than the unpaired (multiome-split) methods, while MultiVI took

282  the longest to complete the assigned tasks, due to its use of variational autoencoder.

283

284  To visually examine the integration results, we generated UMAP plots using the

285  integrated latent embedding and colored the cells by the ground-truth annotation, the predicted

286  identity, and the dataset origin (Figure 2E). We showed the best-performing results from both

287  the unpaired (multiome-split) and multiome-guided categories for each of the PBMC and BMMC

288  simulations. Additional evaluation on cell type separation and batch-mixing are shown in

289  Supplementary Figure 2. Most metrics show method-specific values, meaning the rankings of

290  methods do not change across different numbers of multiome cells. Among the unpaired

291  methods, Seurat v3 is the best at separating cell types in the integrated space, but it has the

292  worst batch mixing result. On the other hand, FigR shows the opposite trend; it ranked the

293  highest for batch mixing, but the lowest for cell type separation. Among the multiome-guided

294  integration methods, MultiVI mixes the batches better while Seurat v4 often results in a higher

295  cell type silhouette score, especially when there is a greater number of multiome cells. We also

296  evaluated the integration results visually, through examining UMAP projection of the integration

297  results as shown in Supplementary Figure 3 for the PBMC simulations and Supplementary

298  Figure 4 for the BMMC simulations. Visually, we do not see drastic differences between

299  methods and there are no methods showing particularly poor cell type separation or batch

300  mixing result. Therefore, we conclude that the incorporation of multiome cells improves cell type

301  annotation when there are enough cells to resolve the cell type heterogeneity in the multiome

302  dataset alone.

303

304

305  **How to spend your sequencing dollars: more cells or increased sequencing depth?**

306

307    Experimentalists are commonly constrained by budget limitations and need to consider whether

308    sequencing a larger number of cells at low depth or a smaller number of cells at high depth is

309    the more productive approach.  To answer this question, we evaluated how the sequencing

310    depth of the multiome approach influences the integration result. Since we know that including

311    multiome data improves cell type annotation for the single-modality datasets, for this analysis,

312    we aimed to evaluate the cell type annotation accuracy of the three data types together. Table 1

313    shows the sequencing depth of the original multiome samples. To simulate data with lower

314    depths, we down-sampled the reads for both RNA and ATAC profiles to 25%, 50%, 75% of the

315    original data (Figure 3A) and compared these results to the original samples. We performed this

316    experiment on both the PBMC dataset (Figure 3B) and the BMMC dataset (Figure 3C). For the

317    PBMC study, the increase in sequencing depths resulted in an increase in cell type annotation

318    accuracy for all methods, with Seurat v4 achieving the highest ARI and NMI among all methods

319    for 75% and 100% depth. In contrast, when we used the BMMC data set as the input, we noted

320    that when including only 2,000 multiome cells, regardless of sequencing depth, the unpaired

321    method (Seurat v3) performed the best. However, when we included 4,000 cells in the BMMC

322    multiome sample, 50% of read depth was sufficient for Seurat v4 to annotate the cell types most

323    accurately. These conflicting results prompted us to ask whether sequencing depth is less

324    important than cell number.

325

326        To answer this question, we designed another simulation. Given a fixed cost for

327    1,000,000 RNA-seq reads and 4,000,000 ATAC-seq reads, we used either 400 cells with 100%

328    of the depth (see Table 1), or 10% of the reads for 4,000 cells. Next, we analyzed the datasets

329    using Seurat v3 and Seurat v4, the best-performing method in each category based on Figure

330    3C. For cell type annotation accuracy, the sequencing depth curve plateaued sooner than the

331    number of cells curve. For Seurat v4, the ARI and NMI did not increase much beyond 60%

332    sequencing depth, while both scores increased consistently as the number of cells increases.

333    Comparing Seurat v3 with Seurat v4, we noted that Seurat v4 performed better when there was

334    30% sequencing depth given 4,000 cells or 2,600 cells given 100% depth. Therefore, for the

335    accuracy of cell type annotation for integrated data, having more cells is more important than

336    having a higher sequencing depth. Importantly, once a sufficient number of cells has been

337    profiled to capture the complexity of a given sample, the multiome-guided methods, specifically

338    Seurat v4, are the best. Our analysis also demonstrated that the 'sufficient' number of cells

339    depends on the complexity of the biological system in question. For PBMC, we see that if the

340    goal is to detect seven distinct cell types, 2,000 cells is already enough. However, for BMMC

341    with its more complex cell type composition at least 2,600 cells are needed.

342

343    In addition to the cell type annotation accuracy, we also evaluated recovery of peak-

344    gene association for the 1,000 single-modality ATAC-seq cells when incorporating mulitome

345    samples generated at ten different depths and numbers of cells. We see that Seurat v3 is

346    consistently better than Seurat v4 (Figure 3D). Moreover, the number of cells and sequencing

347    depth did not affect the percentage of peak-gene pair recovery nor the F1 score. This is likely

348    because Seurat v3 predicts RNA expression using a nearest neighbor approach on the

349    integrated space, and the software had enough cells in the scRNA-seq dataset for the prediction,

350    thus changes in the multiome data did not affect the result.

351

352    Next, we evaluated cell type separation and batch mixing results as summarized in

353    Supplementary Figure 5. Most metrics increased slightly as sequencing depth increased, but

354    the ranking of methods is similar as described before. Overall, Seurat v4 shows the best

355    separation of cell types in the integrated space, but the mixing of batches is the worst, across

356    sequencing depths. A UMAP projection of each method under each simulated scenario is

357    shown in Supplementary Figures 6-8 for visual comparison.

358

359    Overall, we conclude that the number of cells in the multiome data is more critical than

360    sequencing depth for annotating cell types in the integrated data. On the other hand, treating

361    multiome data as unpaired single-modality datasets recovers peak-gene pairs at a higher

362    accuracy.

363

364    **Which method is the best at removing batch effects?**

365

366    It is common that scRNA-seq and snATAC-seq data are generated by different labs or from

367    different individuals than the multiome data. Therefore, another key characteristic for integration

368    methods is whether they can integrate samples displaying batch effects. To answer this

369    question, we leveraged the complex batch structure present in the BMMC dataset. Figure 4A

370    shows the technical batch or biological batch structure we aimed to evaluate, with the multiome

371    cells coming from a different research site, or a different donor. Figure 4B shows the results of

372    cell type annotation accuracy for unpaired integration methods and the multiome-guided

373    methods. We again saw increasing cell type annotation accuracy as the number of multiome

374  cells increased. With 3,000 or more multiome cells, Seurat v4 again was the best-performing

375  method. Seurat v4 is a supervised approach, meaning that the multiome sample serves as a

376  reference to which the single-modality datasets are mapped to. Figure 4B shows that although

377  the multiome sample has strong batch effects (Supplementary Figure 9), the supervised

378  mapping approach resulted in the most accurate cell type annotation. Additional integration

379  results are shown in Supplementary Figure 10 and the UMAP projections are shown in

380  Supplementary Figures 11-12.

381

382  To further challenge all methods in the situation of complex mixtures of samples, we

383  considered a situation where the multiome sample includes cells from a mixture of two donors,

384  and the scRNA-seq and snATAC-seq data come from the same or different research sites. Due

385  to batch effects in the multiome samples, we added one more category called 'Seurat v4

386  integrate', in which the integration of samples was first done on each modality separately, then

387  two modalities were joined using the Seurat v4 weighted nearest neighbor approach, and lastly

388  combined with the single modality dataset (see more in Supplementary methods). Figure 4D

389  shows that in the case of low batch effects between the two donors, Seurat v4 and 'Seurat v4

390  integrate' performed similarly well at annotating cell types. However, in the presence of stronger

391  batch effects, 'Seurat v4 integrate' outperformed all other methods for cell type annotation, with

392  much higher cell type separation as measured in cell type average silhouette width (ASW)

393  (Supplementary Figure 13). From the UMAP projection in Supplementary Figure 14, we see that

394  'Seurat v4 integrate' mixes cells from the two multiome samples much better than Seurat v4.

395  Therefore, when the multiome data include two donors with strong batch effects, integration

396  across the batches is required before mapping the single-modality datasets.

397

398

# Discussion

400  In summary, we evaluated seven multi-omic integration methods under three realistic scenarios.

401  Firstly, we showed that the incorporation of multiome data improves the cell type annotation

402  accuracy of scRNA-seq and snATAC-seq data when there are sufficient number of cells in the

403  multiome data to reveal cell type identities. Secondly, we showed that the number of cells in the

404  multiome data plays a more important role than sequencing depth per cell for cell type

405  annotation accuracy. Thus, when generating a multiome dataset with a fixed budget, a better

406  strategy is to profile more cells so that rare cell types can be captured. Lastly, when the three

13

407 datasets to be integrated are confounded by batch effects, Seurat v4 resulted in the best cell

408 type annotation accuracy.

409

410   In all evaluations, Seurat v4 demonstrated superior performance at resolving cell type

411 heterogeneity when data from many multiome-profiled cells are available. This makes sense as

412 Seurat v4 is a supervised approach in which single-modality cells are merely projected to the

413 integrated space learned from the multiome dataset. Therefore, when the multiome data have

414 an insufficient number of cells to reveal accurate cell types, the integration will lead to poor

415 annotation accuracy. The other two multiome-guided methods, e.g., Cobolt and MultiVI, claim to

416 be able to make use of all three data types. The hope is that the single-modality cells can help

417 the clustering when multiome cells are small. However, as shown in Figures 2 and 3, both

418 Cobolt and MultiVI performed worse than the unpaired integration methods that do not leverage

419 the paired relationship of the multiome data. Therefore, when the multiome dataset has a small

420 number of cells, it is better to treat the multiome cells as unpaired and append them to the

421 single-modality datasets for the integration of three datasets.

422

423   There are several limitations of this study. Firstly, our simulations represent the most

424 ideal situation, where the single-modality cells are generated from the exact same dataset as

425 the multiome cells. In reality, the single-modality and the multiome data are generated from

426 different experimental kits that could have slight differences since the multiome workflow is

427 optimized to capture both gene expression and chromatin accessibility. Moreover, the gene

428 expression captured through the multiome workflow is in fact measuring mRNA in individual

429 nuclei, while scRNA-seq captures mRNA in whole cells. Slight differences between snRNA-seq

430 and scRNA-seq datasets have been reported [19]. Lastly, the PBMC dataset did not have

431 expert-annotated cell type labels. We followed a tutorial by Seurat v4 to obtain annotations [20],

432 thus the evaluation of PBMC-simulated scenarios might favor Seurat v4. However, the BMMC

433 data were manually annotated by experts and Seurat v4 still showed outstanding performance

434 in evaluations based on this dataset.

435

436   Secondly, although Seurat v4 was the best at annotating cell types, it performed worse

437 than unpaired integration methods at recovering peak-gene associations. Furthermore, even the

438 best method only revealed 45% of peak-gene pairs detected in the paired multiome dataset,

439 and many of the detected pairs are false positives. Moreover, we did not explore the possibility

440 of imputing chromatin accessibility from scRNA-seq or appending imputed profile with observed

441    multiome sample. To truly integrate the three data types and understand the underlying *cis-*

442    regulatory logic, one would hope to impute the missing modality for both the scRNA-seq and

443    snATAC-seq data, and then append the imputed profiles with the multiome dataset to identify

444    peak-gene pairs with the largest number of cells. Therefore, additional work needs to be done to

445    evaluate the performance of different methods in jointly integrating the imputed single-modality

446    datasets with the multiome samples for downstream analyses.

447

448

449    # Conclusions

450    Our benchmarking evaluations showed that multiome data are helpful for annotating single-

451    modality data. The number of cells in the multiome data is critical to ensure a good cell type

452    annotation after integration and the exact number of cells depends on the complexity of the

453    biological system. When generating a multiome dataset, the number of cells is more important

454    than sequencing depth for cell type annotation. Lastly, Seurat v4 is the best at integrating

455    scRNA-seq, snATAC-seq, and multiome data even in the presence of complex batch effects.

## Methods

**Datasets**

**Peripheral blood mononuclear cell (PBMC) dataset**

This dataset was generated using the 10x Genomics Single Cell Multiome ATAC + Gene Expression kit [10]. The PBMC dataset with granulocytes removed was downloaded from the 10x Genomics website, which included 11,909 cells. The dataset was processed and annotated into 30 cell types following the Seurat tutorial [7, 20]. We grouped similar cell types and refined the annotations into 9 broad cell types (similar to the level 1 categories from the Azimuth database [3]): B-cells ('B'), CD4 T cells ('CD4 T'), CD8 Naïve T cells ('CD8 Naïve'), CD8 Effector T cells ('CD8 TEM'), Dendritic cells ('DC'), Monocytes ('Mono'), Nature killer cell ('NK'), other T cell ('other_T'), and other cell categories ('other'). The ATAC-seq profile released on 10x Genomics website was counting the Tn5 insertion events in each genomic region. Here, we retabulated the cell-peak matrix by the number of reads overlapping each genomic region, using the Signac's FeatureMatrix function [21]. We used the peak-based counting result as input for the peak-gene pair identification (described below) and subsequent simulations. The list of peak-gene pairs identified using all cells in the multiome dataset (10,412 cells) is treated as the ground truth when calculating percentage of peak-gene pair recovery or F1 score. 'Other_T' and 'other' cells were excluded from the data simulation due to their extensive separation in the UMAP embedding. After removal of cells, there are 10,085 cells used for simulation.

**Bone marrow mononuclear cells (BMMC) dataset**

This dataset was generated as part of the "Open Problems in Single-cell Analysis" competition [12]. BMMC cells from nine healthy donors were profiled at four different research sites using the 10x Multiome ATAC + Gene Expression kit. The dataset was analyzed by Lance and colleagues [12], who annotated the cells into 22 cell types. The values in the cell-peak matrix of the ATAC-seq data was also the insertion-based counting, so we again converted it into peak-based counting as mentioned above. Data simulations related to Figures 2 and 3 were performed using cells from the site 1 donor 2 (S1D2) BMMC sample. This sample contains 6,740 cells, annotated into 21 cell types. The peak-gene pair prediction accuracies shown in Figures 2 and 3 were calculated by comparing the result to a ground-truth list generated with the S1D2 sample. To simulate technical batch and biological batch effects (Figure 4), we used cells

16

489    generated at research site 1 or from donor 1, which includes a total of 29,486 cells, composed

490    of 21 cell types (Supplementary Figure 1B).

491

492    **Evaluation metrics**

493    Annotation accuracy

494    Each integration method returns an integrated latent embedding matrix for cells. Louvain

495    clustering was performed to identify k clusters, in which k is the number of cell types in the

496    ground-truth annotation. To evaluate annotation accuracy, Adjusted Rand Index (ARI) [14] and

497    Normalized Mutual Information (NMI) [15] from the Scib package (v1.0.2) [18] were calculated to

498    compare the predicted cluster labels with the ground truth. Specifically, ARI compares every

499    pair of cells in the dataset and calculates a similarity measurement by considering the number

500    of cell pairs that are in the same cluster in both annotation results, versus the number of cell

501    pairs showing discordant annotations. This metric is then adjusted by chance, as there will be a

502    non-zero similarity between the two clustering results just due to random permutation of labels.

503    The resulting metric ranges from 0 to 1 in which 1 means perfect matching between the two

504    results while 0 means random labeling of cells. NMI is another measurement commonly used

505    for comparison of two clustering results. NMI measures if knowing one label provides

506    information about the other label. If the two lists are highly correlated, then it has high mutual

507    information. NMI is then normalized by a factor to control for differences due to the number of

508    clusters in each set of labels.

509

510    Cell type separation

511    We evaluated the separation of clusters and the tightness of cells in the integrated latent space

512    derived from each method. We calculated cell type-specific average silhouette width (ASW) [18],

513    using the ground-truth annotation and the joint embeddings. The resulting score is between 0

514    and 1 in which 1 means small intra-cluster distance and high inter-cluster distance. We also

515    calculated a cell type Local Inverse Simpson's Index (cLISI) [18], which is an adaptation of LISI

516    previously used to quantify the degree of batch effects [17]. Here, cLISI was calculated using

517    the ground-truth labels again in which it evaluates how many cells need to be drawn from a

518    cell's neighborhood to draw a second cell of the same type. The score is normalized again so

519    that 1 means good local neighborhood preservation of the same cell type while 0 is otherwise.

520

521    Batch mixing

17

522    To evaluate batch mixing, two metrics were employed. A batch ASW score was used to

523    evaluate the within-batch distance and the across-batch distance [18]. The score was rescaled

524    so that 0 is the worst and 1 is the best separation. To evaluate the local neighborhood accuracy,

525    k-nearest neighbor batch effect test (kBET) was also performed [16]. Specifically, kBET

526    measures the difference between observed batch frequency in the k-nearest neighbors

527    compared to an expected frequency based on the number of cells in each batch. The value is

528    rescaled to 0 and 1 in which 1 represents the optimal mixing of cells from different batches in

529    which cells in the neighborhood are highly similar to the expected frequency.

530

531    <u>Peak-gene pair recovery</u>

532    To identify correlated peak-gene pairs, we used the methodology introduced in the SHARE-seq

533    paper [1]. Specifically, a Pearson correlation is calculated between the raw accessibility count of

534    every peak and the normalized UMI count of every gene if the peak is within 50,000 base pairs

535    from the transcription start site (TSS) of the gene. The null distribution of correlation coefficients

536    was then generated through selecting 100 peaks that have similar GC content, length, and

537    accessibility as the target peak, and calculating correlation of the background peaks and the

538    target gene. A one-sided t-test was used to calculate a p-value for every peak-gene pair by

539    comparing to the background peaks and the peak-gene pairs with p-value less than 0.05 and z-

540    score greater than 0.05 identified as significant peak-gene pairs. Associated peak-gene pairs

541    were identified using all cells from each dataset. To evaluate the performance of each method

542    at imputing gene expression from snATAC-seq data, a peak-gene association was calculated in

543    the same manner using the raw cell-peak count of the unpaired ATAC data and the predicted

544    gene expression generated by the evaluated methods. To evaluate the *in silico* imputed gene

545    expression results, we calculated the percentage of peak-gene pairs recovered using the

546    imputed gene expression and the observed snATAC-seq peak counts. To account for false

547    negative results, we calculated an F1 score. Thus, the peak-gene pair percent recovery and the

548    F1 score were used to evaluate each method that can impute missing gene expression.

549

550    **Evaluation scenarios**

551    We simulated three scenarios to evaluate the performance of each method. For each scenario,

552    we simulated five independent replicates. Details regarding how each method was implemented

553    are described in the Supplementary Methods.

554

555    **Scenario 1: evaluating the effect of multiome cells on single-modality integration.**

556

Data simulation

558 In this task, we first defined the number of cells to be drawn for each data type with an example

559 shown in Figure 2A. Then, we randomly selected cells from the ground-truth multiome dataset

560 according to the desired number of cells for each data type. For scRNA-seq, we kept the gene

561 expression matrix; for snATAC-seq, we kept the cell-by-peak matrix and the fragment file; lastly,

562 for the multiome sample, we kept all three data files. The cells were sampled without

563 replacement.

564

Evaluated methods

566 We first ran the four unpaired integration methods (Seurat v3, LIGER, FigR, and bindSC) to

567 integrate the simulated scRNA-seq and snATAC-seq datasets and the results were summarized

568 under the 'Unpaired' categories.  To make use of the multiome data, we ran the four methods

569 again, with the multiome cells treated as unpaired. Specifically, the RNA profile from the

570 multiome cells was appended to the scRNA-seq dataset, and the ATAC-seq profile was

571 appended to the snATAC-seq dataset. The results from this category were summarized under

572 'Unpaired (multiome-split)'. Lastly, we ran the multiome-guided methods with the scRNA-seq,

573 snATAC-seq, and multiome datasets as input.

574

Evaluations

576 To evaluate if the presence of multiome cells improves the integration of single-modality

577 datasets, we evaluated the annotation accuracy, peak-gene pair recovery, cell type separation,

578 and batch mixing of the scRNA-seq and snATAC-seq cells.

579

580 **Scenario 2: evaluating the impact of sequencing depth in multiome cells on multi-omic**

581 **data integration.**

582

Data simulation

584 For this task, we first defined the number of cells in each data type as well as the percentage of

585 original depth the multiome cells will be down-sampled to; an example is shown in Figure 3A.

586 We first generated the three data types according to the number of cells defined. Then, we

587 performed depth-down-sampling for both the gene expression and chromatin accessibility

588 profiles of the multiome dataset. To down-sample the cell-by-gene count matrix for gene

589 expression, we used Scuttle::downsample [22] to reduce the sample depth to a percentage of

19

590    the original dataset. To down-sample the ATAC-seq depth, we performed down-sampling on the

591    fragment file and then regenerated the cell-by-peak count matrix. Specifically, we first counted

592    the number of fragments corresponding to the selected cells, then we calculated the target

593    depth by multiplying the original depth to the percentage factor. We randomly selected the

594    number of reads as calculated, without replacement, and saved this file as the new fragment file.

595    Then the down-sampled fragment file was sorted, recompressed, indexed with tabix and,

596    tabulated into peak counts with the original feature set with Signac:: FeatureMatrix [21] function.

597    This often resulted in less reduction in peak counts, as some of the fragments removed were

598    not previously assigned to the peaks.

599

600    Evaluated methods

601    We ran the unpaired integration methods with the multiome data appended to the single-

602    modality datasets as described above, the results were summarized under 'Unpaired (multiome-

603    split)'. We also ran the three multiome-guided methods.

604

605    Evaluations

606    The evaluation of annotation accuracy, cell type separation and batch mixing were calculated

607    using all cells present in simulated scRNA-seq, snATAC-seq, and the multiome datasets. Given

608    how the multiome data were split and appended to the single-modality datasets for the 'unpaired

609    (multiome-split)' category, it resulted in doubling the number of multiome cells. Thus, to ensure

610    a fair comparison between the two categories of methods, half of the multiome cells appended

611    to the RNA-seq were dropped while the other half of the multiome cells appended to the ATAC-

612    seq were dropped. As a result, the same number of cells was evaluated for the 'unpaired

613    (multiome-split)' and 'multiome-guided' methods.

614

615

616    **Scenario 3: evaluating the impact of batch effects on multi-omic data integration.**

617

618    Data simulation

619    The analysis of batch effects was only possible for the BMMC dataset. As mentioned before, the

620    BMMC dataset contains multiome cells generated at four different research sites and nine

621    donors. To create different types of batches, we used the multiome cells from donor 1 but

622    processed at three different sites (S1D1, S2D1, S4D1) as the data source to generate technical

623    batches. We used the multiome cells generated at research site 1 but from different donors

624 (S1D1, S1D2, S1D3) as the source of biological batches. To generate scenarios with mixed

625 technical and biological batch effects, we created more complex batch structures described as

626 'complex test' in Figure 4D using all samples that were either generated at research site 1 or

627 donor 1. After defining which sample each data type comes from and the number of cells, the

628 simulation is the same as described in 'Sceanrio 1', in which cells were randomly drawn from

629 the ground-truth multiome dataset to simulate scRNA-seq, snATAC-seq, and multiome samples.

630

631 Evaluated methods

632 The same seven methods, four from the 'unpaired (multiome-split)' and three from 'multiome-

633 guided' were ran. For situations were multiome were composed of two donors, an additional

634 variation of Seurat v4 was added, termed 'Seurat v4 integrate'. Specifically, the two multiome

635 datasets were first integrated across donors to generate one integrated reference before it was

636 used to integrate scRNA-seq and snATAC-seq datasets.

637

638 Evaluations

639 We calculated metrics measuring annotation accuracy, cell type separation, and batch mixing.

640 For batch mixing, we calculated both the mixing of data types, as well as the mixing of samples.

641 Similar to what was described in 'Scenario 2', to ensure that the same number of cells were

642 evaluated for the unpaired (multiome-split) methods and the multiom-guided methods, half of

643 multiome cells appended to the RNA-seq and the other half of the ATAC-seq dataset were

644 dropped.

645

21

# Declarations

**Availability of data and materials**

The source codes for simulation and evaluations are available online on GitHub at

https://github.com/myylee/benchmark_sc_multiomic_integration. For the multiome datasets

used to generate simulated data, the 10x PBMC dataset was downloaded from

https://www.genomics.com/resources/datasets/pbmc-from-a-healthy-donor-granulocytes-

removed-through-cell-sorting-10-k-1-standard-2-0-0 and the BMMC dataset was downloaded

from GEO accession GSE194122, and the fragment files were obtained from the authors [23].

**Competing interests**

M.L. receives research funding from Biogen Inc. The other authors declare no competing

interests.

**Authors' contributions**

M.Y.Y.L., M.L., K.H.K. conceived this project and designed the framework together. M.Y.Y.L.

performed the simulations and evaluations with guidance from M.L. All authors wrote and edited

the final manuscript. M.L. and K.H.K. supervised the study.

22

# References

1. Ma S, Zhang B, LaFave LM, Earl AS, Chiang Z, Hu Y, Ding J, Brack A, Kartha VK, Tay T, et al: **Chromatin Potential Identified by Shared Single-Cell Profiling of RNA and Chromatin.** *Cell* 2020, **183:**1103-1116 e1120.

2. **Chromium Single Cell Multiome ATAC + Gene Expression** [https://www.10xgenomics.com/products/single-cell-multiome-atac-plus-gene-expression#faqs]

3. Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, 3rd, Hao Y, Stoeckius M, Smibert P, Satija R: **Comprehensive Integration of Single-Cell Data.** *Cell* 2019, **177:**1888-1902 e1821.

4. Liu J, Gao C, Sodicoff J, Kozareva V, Macosko EZ, Welch JD: **Jointly defining cell types from multiple single-cell datasets using LIGER.** *Nat Protoc* 2020, **15:**3632-3662.

5. Kartha VK, Duarte FM, Hu Y, Ma S, Chew JG, Lareau CA, Earl A, Burkett ZD, Kohlway AS, Lebofsky R, Buenrostro JD: **Functional inference of gene regulation using single-cell multi-omics.** *Cell Genom* 2022, **2**.

6. Dou J, Liang S, Mohanty V, Miao Q, Huang Y, Liang Q, Cheng X, Kim S, Choi J, Li Y, et al: **Bi-order multimodal integration of single-cell data.** *Genome Biol* 2022, **23:**112.

7. Hao Y, Hao S, Andersen-Nissen E, Mauck WM, 3rd, Zheng S, Butler A, Lee MJ, Wilk AJ, Darby C, Zager M, et al: **Integrated analysis of multimodal single-cell data.** *Cell* 2021, **184:**3573-3587 e3529.

8. Tal Ashuach MIG, Michael I. Jordan, Nir Yosef: **MultiVI: deep generative model for the integration of multi-modal data.** *bioRxiv* 2021.

9. Gong B, Zhou Y, Purdom E: **Cobolt: integrative analysis of multimodal single-cell sequencing data.** *Genome Biol* 2021, **22:**351.

10. **PBMC from a Healthy Donor - Granulocytes Removed Through Cell Sorting (10k), single cell multiome atac + gene expression dataset by cell ranger arc 2.0.0.** [https://www.10xgenomics.com/resources/datasets/pbmc-from-a-healthy-donor-granulocytes-removed-through-cell-sorting-10-k-1-standard-2-0-0]

11. Lance CaL, Malte D. and Burkhardt, Daniel B. and Cannoodt, Robrecht and Rautenstrauch, Pia and Laddach, Anna and Ubingazhibov, Aidyn and Cao, Zhi-Jie and Deng, Kaiwen and Khan, Sumeer and Liu, Qiao and Russkikh, Nikolay and Ryazantsev, Gleb and Ohler, Uwe and , and Pisco, Angela Oliveira and Bloom, Jonathan and Krishnaswamy, Smita and Theis, Fabian J.: **Multimodal single cell data integration challenge: results and lessons learned.** *bioRxiv* 2022.

12. Lance C, Luecken MD, Burkhardt DB, Cannoodt R, Rautenstrauch P, Laddach A, Ubingazhibov A, Cao Z-J, Deng K, Khan S, et al: **Multimodal single cell data integration challenge: results and lessons learned.** *bioRxiv* 2022:2022.2004.2011.487796.

13. Wolf FA, Angerer P, Theis FJ: **SCANPY: large-scale single-cell gene expression data analysis.** *Genome Biol* 2018, **19:**15.

14. Hubert L, Arabie P: **Comparing Partitions.** *Journal of Classification* 1985, **2:**193-218.

15. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, et al: **Scikit-learn: Machine Learning in Python.** *Journal of Machine Learning Research* 2011, **12:**2825-2830.

722   16.     Buttner M, Miao Z, Wolf FA, Teichmann SA, Theis FJ: **A test metric for assessing single-**
723           **cell RNA-seq batch correction.** *Nat Methods* 2019, **16**:43-49.
724   17.     Korsunsky I, Millard N, Fan J, Slowikowski K, Zhang F, Wei K, Baglaenko Y, Brenner M,
725           Loh PR, Raychaudhuri S: **Fast, sensitive and accurate integration of single-cell data with**
726           **Harmony.** *Nat Methods* 2019, **16**:1289-1296.
727   18.     Luecken MD, Buttner M, Chaichoompu K, Danese A, Interlandi M, Mueller MF, Strobl DC,
728           Zappia L, Dugas M, Colome-Tatche M, Theis FJ: **Benchmarking atlas-level data**
729           **integration in single-cell genomics.** *Nat Methods* 2022, **19**:41-50.
730   19.     Wu H, Kirita Y, Donnelly EL, Humphreys BD: **Advantages of Single-Nucleus over Single-**
731           **Cell RNA Sequencing of Adult Kidney: Rare Cell Types and Novel Cell States Revealed**
732           **in Fibrosis.** *J Am Soc Nephrol* 2019, **30**:23-32.
733   20.     **Weighted Nearest Neighbor Analysis.**
734           [https://satijalab.org/seurat/articles/weighted_nearest_neighbor_analysis.html#wnn-
735           analysis-of-10x-multiome-rna-atac-1]
736   21.     Stuart T, Srivastava A, Madad S, Lareau CA, Satija R: **Single-cell chromatin state analysis**
737           **with Signac.** *Nat Methods* 2021, **18**:1333-1341.
738   22.     McCarthy DJ, Campbell KR, Lun AT, Wills QF: **Scater: pre-processing, quality control,**
739           **normalization and visualization of single-cell RNA-seq data in R.** *Bioinformatics* 2017,
740           **33**:1179-1186.
741   23.     Luecken M BD, Cannoodt R, Lance C, Agrawal A, Aliee H, Chen A,  Deconinck L, Detweiler
742           A, Granados A, Huynh S, Isacco, L, Kim Y, Klein D, De Kumar B, Kuppasani S, Lickert H,
743           McGeever A, Melgarejo J, Mekonen H, Morri M, and Muller M, Neff N, Paul S,  Rieck B,
744           Schneider K, Steelman S, Sterr M, Treacy D, Tong A, Villani A, Wang G, Yan J, Zhang C,
745           Pisco A, Krishnaswamy S, Theis F, Bloom JM: **A sandbox for prediction and integration**
746           **of DNA, RNA, and proteins in single cells.** In *Advances of Neural Information Processing*
747           *Systems*; 2021.
748

749

750  **Table 1:** Summary of the data used for simulation. Columns are number of cells (n_cells),
751  number of unique genes expressed per cell on average in the RNA profile (nGene_RNA), total
752  counts expressed per cell on average in RNA profile (nCount_RNA), number of unique
753  fragments per cell on average in the ATAC profile (nFrag_ATAC), number of peak counts per cell
754  on average in the ATAC profile (nPeakCount_ATAC).

755

| Source | n_cells | nGene_RNA | nCount_RNA | nFrag_ATAC | nPeakCount_ATAC |
|---|---|---|---|---|---|
| PBMC | 10085 | 2013 | 4463 | 15510 | 11305 |
| BMMC site 1 donor 2 (S1D2) | 6740 | 1365 | 2525 | 11064 | 7512 |
| BMMC site 1 or donor 1 | 29486 | 1205 | 2227 | 11798 | 7615 |

756
757

**Figure legends**

**Figure 1:** Outline of the benchmarking evaluations. (A) Scheme to evaluate if multiome data help the integration of single-modality data. (B) Scenarios simulated to evaluate multi-omic integration methods.

**Figure 2:** Comparison of integration performance without vs. with multiome cells. (A) The number of cells and cell types for each simulated dataset using the PBMC or BMMC multiome data as the ground truth. (B – C) Performance of cell type annotation and peak-gene association recovery in the PBMC-based simulations (B) and BMMC-based simulations (C). ARI and NMI measure agreement between predicted cell type and ground-truth labels. Peak-gene pair % recovered is the percentage of peak-gene pairs correctly identified comparing to the ground-truth list calculated using 10,412 paired PBMC cells (B) and 6,740 BMMC cells (C). F1 is the prediction accuracy normalized by the number of false positives and false negatives. Dashed line shows the percent recovery and F1 score calculated using 1,000 multiome cells. Error bar is mean ± standard deviation. (D) Runtime measured in seconds, for each method, in log2 scale. Error bar is mean ± standard deviation. (E) UMAP projection using integrated embedding for a select number of methods. UMAP projection for the other methods are shown in Supplementary Figures 3 (PBMC) and 4 (BMMC).

**Figure 3:** Evaluation of integration performance at varying sequencing depth for multiome cells. (A) Details of the simulation scheme. (B – C) Performance of cell type annotation and peak-gene association recovery in the PBMC-based simulations (B) and BMMC-based simulations (C: left panel, 2,000 multiome cells; right panel, 4,000 multiome cells). ARI and NMI measures agreement between predicted cell type and ground-truth labels. Peak-gene pair % recovered is the percentage of peak-gene pairs correctly identified comparing to the ground-truth list calculated using 10,412 paired PBMC cells (B) and 6,740 BMMC cells (C). F1 is the prediction accuracy normalized by the number of false positives and false negatives. (D) Performance of cell type annotation using Seurat v3 or Seurat v4 at increasing depth or increasing number of cells. (E) Performance of peak-gene association recovery using Seurat v3 or Seurat v4 at increasing depth or increasing number of cells. For all subplots, error bar is mean ± standard deviation.

**Figure 4:** Evaluation of integration performance in the presence of batch effects. (A) Simulation details for the constructed data with technical batches and biological batches. (B) Performance of cell type annotation and runtime in the presence of technical and biological batches shown in (A). ARI and NMI measure agreement between predicted cell type and ground-truth labels. Runtime is measured in seconds, for each method, in log2 scale. Error bar is mean ± standard deviation. (C) Simulation details for two datasets with more complex batch structures. (D) Performance of cell type annotation and runtime in the presence of technical and biological batches shown in (C). ARI and NMI measure agreement between predicted cell type and ground-truth labels. Runtime is measured in seconds, for each method, in log2 scale. Whisker is 1.5 times the inter-quartile range.

# A

**A** Does Multiome aid the analysis of single-modality datasets?

**B** What is the best integration method for scRNA, snATAC, and multiome?

## A

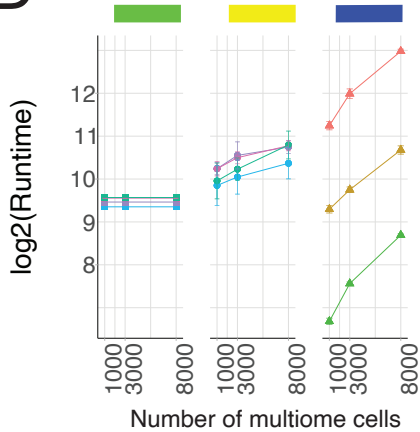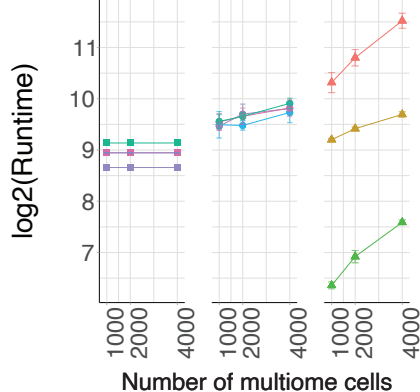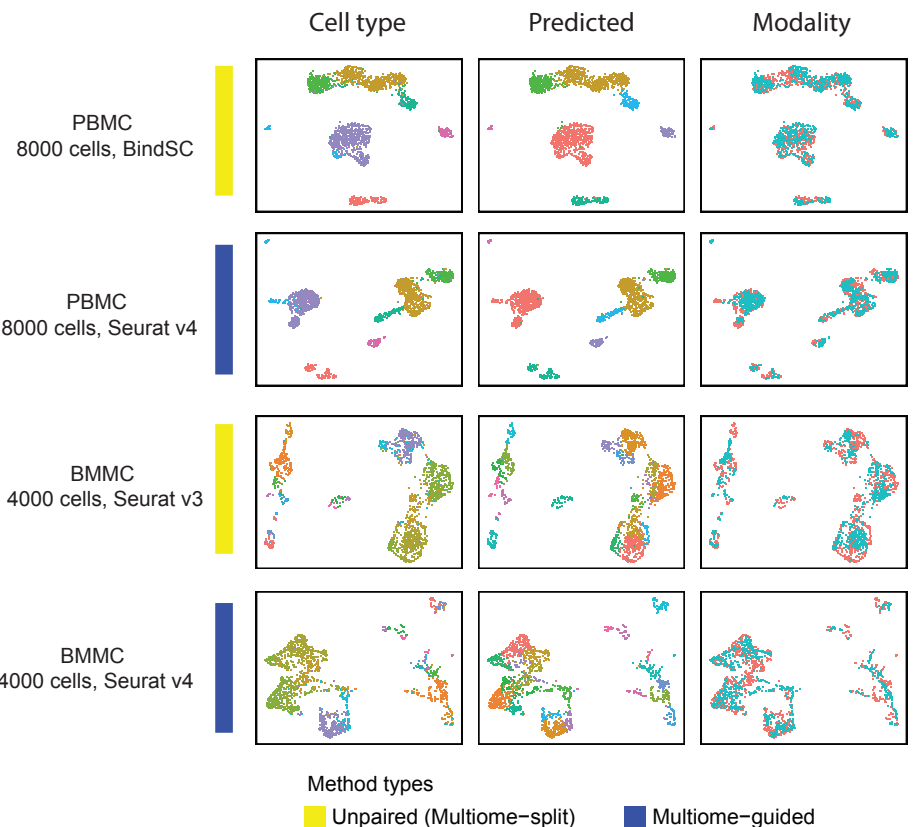| Source | scRNA | scATAC | Multiome | Number of cell types |
|--------|-------|--------|----------|---------------------|
| PBMC | 1000 | 1000 | 1000, 3000, 8000 | 7 |
| BMMC | 1000 | 1000 | 1000, 2000, 4000 | 21 |



**B** PBMC

**C** BMMC

**D**

**E**

Methods: MultiVI, Seurat v3, Seurat v4, BindSC, Cobolt, FigR, Liger

Method types: Unpaired, Unpaired (Multiome-split), Multiome-guided

Annotation: Observed accuracy