

1 Metrics Matter: Why We Need to Stop Using 2 Silhouette in Single-Cell Benchmarking

3 Pia Rautenstrauch (ORCID: 0000-0002-0070-4759)^{1,2} and Uwe Ohler (ORCID: 0000-0002-
4 0881-3116)^{1,2,3*}

5 ¹Humboldt-Universität zu Berlin, Department of Computer Science, 10099 Berlin, Germany.

6 ²Max-Delbrück-Center for Molecular Medicine in the Helmholtz Association (MDC), Berlin
7 Institute for Medical Systems Biology (BIMSB), Berlin, Germany.

8 ³Humboldt-Universität zu Berlin, Department of Biology, 10099 Berlin, Germany.

9 *Corresponding author(s). E-mail(s): uwe.ohler@hu-berlin.de; uwe.ohler@mdc-berlin.de;

10 Contributing authors: pia.rautenstrauch@gmail.com;

11 Abstract

12 Current-day single-cell studies comprise complex data sets affected by nested batch effects
13 caused by technical and biological factors, relying on advanced integration methods. Silhouette
14 is an established metric for assessing clustering results, comparing within-cluster cohesion to
15 between-cluster separation, and adaptations of it have emerged as the dominant choice to
16 evaluate the success of these integration methods. However, silhouette's assumptions are often
17 violated in single-cell data integration scenarios. We demonstrate that silhouette-based metrics
18 can neither reliably assess batch effect removal nor biological signal conservation and are thus
19 inherently unsuitable for data with (nested) batch effects. We propose alternative, robust
20 evaluation strategies that enable accurate integration method assessment and call to update
21 benchmarking practices.

22 Main text

23 Integrating single-cell data remains a key challenge of single-cell analysis due to the increasing
24 complexity and volume of data sets generated. These data sets often include intricate, nested
25 batch effects from both technical and biological factors, requiring rigorous evaluation of
26 integration methods to ensure accurate integration and interpretation. Silhouette-based
27 evaluation metrics have become widely adopted to address this challenge. As an integral part of
28 current data integration benchmarking, they are used for scoring both biological signal
29 conservation (bio-conservation) and batch removal. However, we demonstrate that these
30 metrics cannot reliably score data integration.

31
32 The metric "silhouette" scores clustering quality by comparing within-cluster cohesion to
33 between-cluster separation (Rousseeuw, 1987), and was originally developed for evaluating
34 unsupervised clustering results of unlabeled data (internal evaluation). In the single-cell field,
35 silhouette was thus quickly taken up for determining the optimal number of clusters in single-cell
36 data sets (Wagner et al., 2016; Scialdone et al. 2016). More recently, silhouette has been
37 adapted for evaluating horizontal data integration (Argelaguet et al., 2021), for instance, to score

38 bio-conservation by assessing how well cell type annotations (based on labeled data, i.e.,
39 external evaluation) from distinct batches co-cluster (Haghverdi et al., 2018; Tran et al., 2020;
40 Luecken et al., 2022). From 2017 onwards, silhouette-based metrics have also been employed
41 for scoring batch effect removal, another key challenge of horizontal data integration (Risso et
42 al., 2018; Büttner et al., 2019; Cole et al., 2019). Here, the silhouette concept is, however,
43 inverted for scoring how well cells from distinct batches (external labels) mix. Fueled by a large-
44 scale single-cell benchmark and accompanying toolbox, silhouette-based batch removal metrics
45 have become a predominant score to evaluate and claim the success of many new single-cell
46 integration methods (Luecken et al., 2021; Luecken et al., 2022).

47
48 Unfortunately, it appears to have gone unnoticed that silhouette-based batch removal metrics
49 completely fail when scoring data integration in even modestly challenging scenarios. To
50 illustrate this, consider a simplified, illustrative example: we simulate four single-cell RNA-seq
51 samples with three cell types. The samples are split into two groups, mimicking that they were
52 sequenced at two distinct sites (Figure 1(a)). This corresponds to data with batch effects nested
53 in groups with decreasing levels of between-group batch effects (or, conversely, increasing
54 levels of successful data integration), which we complement with an overcorrected scenario. To
55 evaluate the behavior of silhouette scores for evaluating batch removal, we chose the 'ASW
56 batch' metric, a commonly used cell-type dependent implementation of a silhouette-based batch
57 removal metric (scib package (Luecken et al., 2022)). We find that Batch ASW results in near
58 maximal, close to identical scores for every scenario - no matter whether data was actually
59 integrated or not. Silhouette scores only consider the nearest neighboring clusters - here,
60 assigned by sample - and when samples from the same group are highly similar, batch effects
61 between the groups cannot be captured (Figure 1 (b)). Given the increasing prevalence of
62 nested batch effects in single-cell studies, addressing this limitation is pivotal for ensuring
63 reliable data integration. As we will see, the problem results from the underlying definition of the
64 silhouette score, thus extending to **every** silhouette-based metric for batch removal.

65
66 The silhouette is defined as follows. For a cell i assigned to a cluster C_k . Given a_i : the mean
67 distance between a cell i and all other cells in the same cluster C_k . With: b_i : the mean distance
68 between a cell i and all other cells in the **nearest** (neighboring) other cluster C_l where $l \neq k$, the
69 silhouette coefficient of a single cell i , denoted as s_i is given by:

$$70 \quad s_i = \frac{b_i - a_i}{\max(a_i, b_i)} \quad (1).$$

71 Note that this is only defined for $2 \leq \# \text{clusters} \leq \# \text{cells} - 1$ and ranges between -1 and 1,
72 with 1 indicating good cluster separation ($a_i \ll b_i$), values near 0 indicating cluster overlapping
73 ($a_i = b_i$), and -1 wrong cluster assignment ($a_i \gg b_i$). In contrast to the use of silhouette for
74 internal clustering evaluation (unsupervised clustering), for scoring data integration in the single-
75 cell field, cells are not assigned to clusters in a data-driven manner, e.g., by the result of a
76 clustering algorithm, but by external information, such as cell type or batch labels.

77
78 To illustrate why silhouette is inadequate for evaluating batch removal, consider integrating
79 multiple data sets (samples) with a single cell type. In this context, the aim is to score cluster
80 overlap and not separation. Silhouette-based batch removal metrics first assign the cells of

81 distinct samples to corresponding clusters. The assumption is that silhouette values s_i around 0
82 indicate a high level of cluster overlap and, hence, batch effect removal. However, an important
83 detail goes unnoticed: silhouette (cf. equation (1)) considers the mean distance between a cell
84 i and all other cells in only the **nearest** (neighboring) other cluster C_i (b_i). A value for s_i around 0
85 is thus attainable if a given cluster overlaps with just a single other cluster and could still be very
86 distinct from all other remaining ones. This behavior is highly problematic in the presence of
87 nested batch effects, where samples within groups are a lot more similar to each other than
88 between groups. If samples within groups overlap, but differences remain between samples of
89 distinct groups, silhouette-based metrics can result in maximal scores despite remaining strong
90 batch effects, in the worst case, even favoring suboptimal methods. In practice, data sets
91 usually comprise a multitude of cell types. Silhouette-based batch removal scores are
92 commonly computed per cell type label and later aggregated to account for differences in cell
93 type composition between samples (Luecken et al., 2022). Additionally, they are transformed to
94 range between 0 and 1, with 1 indicating best performance. The same caveats apply - in the
95 presence of nested batch effects, maximal scores are reached even if data is insufficiently
96 integrated.

97
98 This behavior is not limited to toy examples but, in fact, painfully obvious on real data sets. We
99 empirically discovered this issue for 'Batch ASW' in the context of the NeurIPS 2021 challenge
100 (Lance et al., 2022). The benchmark data is rich in nested batch effects of samples sequenced
101 at different sites (intra-site differences smaller than inter-site) from bone marrow mononuclear
102 cells (Luecken et al., 2021). Choosing a scRNA-seq subset with four batches nested into two
103 groups (sites) for clarity, we compare metric performance on unintegrated, suboptimally
104 integrated, effectively integrated, and optimized integrated data (Figure 1(c)). Here, the
105 silhouette-based batch removal metric Batch ASW even favors worse solutions with stronger
106 batch effects (Figure 1(d)), with the same observations applying to the full data set
107 (Supplementary Figure 2(b)). While we demonstrate this behavior with scRNA-seq data, this
108 finding generalizes to any data with nested batch effects.

109
110 Single-cell integration benchmarking is an area of active research, which has seen large-scale
111 coordinated efforts (Tran et al., 2020; Luecken et al., 2021; Luecken et al., 2022; Hu et al.,
112 2024; Maan et al., 2024). When first introduced, silhouette-based batch removal metrics were
113 applied to small data sets without nested batch effects (Büttner et al., 2019), with the limitations
114 not becoming apparent. However, given the prevalence of nested batch effects in current-day
115 data sets, silhouette's inability to account for nested batch effects is a real concern. It is
116 especially problematic when they are not combined with metrics that could indicate insufficient
117 integration, but also when evaluation results are aggregated into a single summary score that
118 obscure possible discrepancies. Two classes of metrics should be considered to score
119 horizontal data integration: Batch removal and bio-conservation metrics (Tran et al., 2020;
120 Luecken et al., 2022; Maan et al., 2024). Among alternatives to silhouette, some batch removal
121 metrics score local batch mixing and are thus not prone to the same behavior, either without cell
122 type labels (iLISI (Korsunsky et al., 2019), kBET (Büttner et al., 2019)) or accounting for cell type
123 imbalance if cell type labels are available (CiLISI (Andreatta et al., 2024)). Concerning bio-
124 conservation, many clustering metrics have been applied to cell type labels (ARI, NMI, cell type

125 ASW). Evaluating performance on a high confidence subset, e.g., samples from the same donor
126 or technical replicates, can be a valuable option (Rautenstrauch & Ohler, 2024).

127

128 Combining local mixing batch removal with bio-conservation metrics on a cell type level has
129 proven to be a successful strategy for evaluating integration performance (Andreatta et al.,
130 2024). For example, applying CiLISI with ARI is robust to nested batch effects, leading to
131 accurate rankings in our simulated and real data scenarios while flagging overcorrection (high
132 batch removal but low bio-conservation scores) (Figure 1(b) and (d)). It is also possible to "fix"
133 the silhouette-based metric Batch ASW to be robust to nested batch effects by redefining b_i as
134 the mean distance between a cell i and all other cells in **any** other cluster C_l with $l \neq k$.

135 Changing euclidean to cosine distance results in higher discriminative power (cf. Methods for
136 further details). This adaptation, which we call batch removal adapted silhouette (BRAS), could
137 also be employed in other metric variants. Like CiLISI, the BRAS metric also accurately ranks
138 simulated and real data (cf. Figure 1(b) and (d) and Supplementary Figures 1(a) and (b) and
139 (2)).

140

141 Silhouette score problems are not limited to batch integration but also arise in scores adapted
142 for bio-conservation. As such, the Cell type ASW score shows significant limitations in
143 discriminating between scenarios (Figure 1(b) and (d); details concerning other bio-conservation
144 metrics can be found in Supplementary Note 1). This limitation also goes back to repurposing
145 the silhouette score - originally intended for internal - to external evaluation, which imposes
146 cluster labels on the data. Highly non-convex cluster shapes, particularly in the presence of
147 strong batch effects, cause unintended behavior as silhouette's comparison of within-cluster
148 cohesion to between-cluster separation becomes erratic, which can also affect batch removal
149 metrics. Arguably, such edge cases can and have been flagged by complementing Cell type
150 ASW with batch removal metrics (Haghverdi et al., 2018), similarly to the strategy that we show
151 to flag the overcorrection scenario (Figure 1(b)). However, current benchmarking practices often
152 aggregate scores across different metrics without identifying outliers. This practice can lead to
153 misleading evaluations, as high scores from unreliable metrics can disproportionately influence
154 the overall assessment of a method's performance.

155

156 Single-cell data integration remains a key computational challenge and an active area of
157 research. Our investigation reveals the inadequacy of currently prevalent silhouette-based
158 evaluation metrics for assessing data integration. In the presence of nested batch effects, these
159 metrics can produce near-maximal scores even when data integration fails, as they focus solely
160 on the nearest neighboring samples. We propose a robust evaluation strategy that combines
161 local batch mixing with bio-conservation metrics, along with modifications to the silhouette
162 metric to address its current issues. In any case, including a baseline model, such as
163 unintegrated data, is essential for meaningful evaluation of integration. The limitations of
164 silhouette metrics extend to bio-conservation assessments, where non-convex cluster shapes
165 resulting from batch effects lead to erratic behavior. In summary, silhouette-based integration
166 metrics are inadequate and should not be used to evaluate integration. Benchmarking practices
167 need to discontinue the use of silhouette-based metrics, especially when aggregating results.

168 This is required to ensure reliable assessments of integration methods, as method choice
169 impacts downstream analyses.

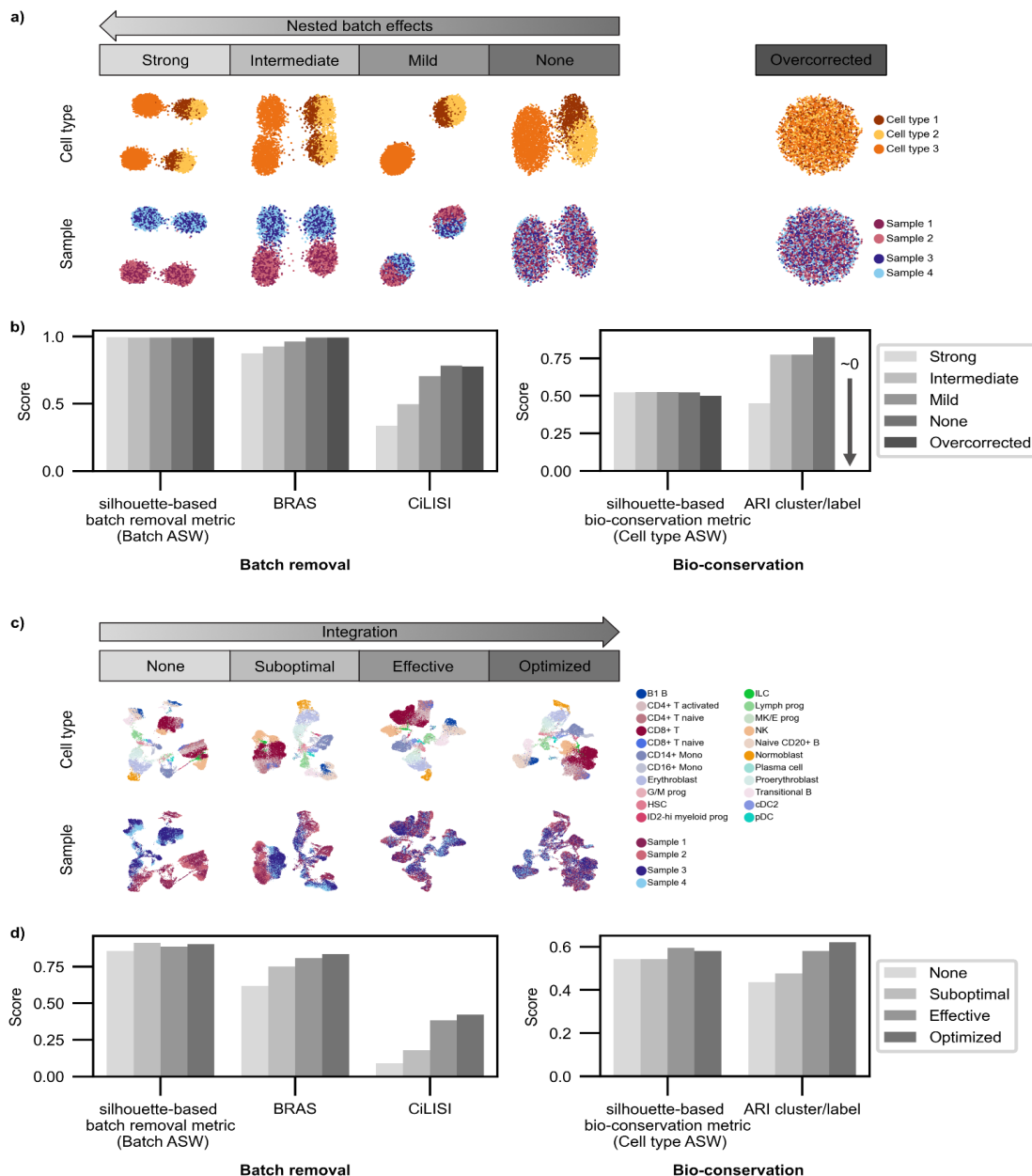
170 Acknowledgements

171 We wish to thank Michael I. Love from UNC-Chapel Hill for constructive feedback and
172 encouragement.

173 Funding

174 This project was funded in part by a grant from the Chan Zuckerberg Initiative ‘Single-Cell
175 Biology Data Insights’ program, and the DFG international research training group IRTG2403.

176 Figures



177
178 **Figure 1: Silhouette-based metrics (Batch ASW) are unreliable with nested batch effects,**
179 **failing single-cell data integration evaluation.**

180 (a) UMAPs of simulated data with nested batch effects between groups of samples with
181 decreasing levels of batch effects between groups. Colored by cell type and sample. (b) Batch
182 removal metrics: Unreliable metric (Batch ASW), reimplementing fixing erratic behavior called
183 batch removal adapted silhouette (BRAS), and an alternative cell type-dependent diversity
184 score: CiLISI. Bio-conservation metrics: Cell type ASW and ARI. (c) UMAPs of NeurIPS data
185 minimal example with nested batch effects integrated with increasing success, colored by cell
186 type and sample. (d) Metrics as in (b).

187 Online methods

188 **Data**

189 *Simulated data*

190 Drawing inspiration from Andreatta et al. (2024) and a recommendation of the Splatter
191 developer (<https://github.com/Oshlack/splatter/issues/99>), we simulate five scenarios with
192 decreasing levels of nested batch effects with the Splatter package (Zappia et al., 2017)
193 (version 1.26.0). Each scenario is composed of four samples with three cell types nested in two
194 groups, meaning that the samples within a group are more similar to each other than between
195 the groups. The scenarios are "Strong", "Intermediate", and "Mild", as well as "None" - with no
196 (nested) batch effects, and an "Overcorrected" scenario, with neither nested batch effects nor
197 biological cell type signal. We first simulate data with two samples of 2000 cells stemming from
198 three distinct cell types with varying proportions. We vary the nested batch effect for the
199 different scenarios via the `batch.facLoc` and `batch.facScale` parameters. We then select half of
200 the cells of the two samples, and add small noise factors to them, resulting in four samples
201 nested into two groups of 1000 cells each. The noise factor stems from another simulated data
202 matrix without batch and cell type structure where we use a small library size parameter
203 `lib.scale`. In the "Overcorrected" scenario, we choose no differential expression between cell
204 types and samples.

205

206 *Real data*

207 We employ a benchmarking data set from the NeurIPS 2021 Multimodal Single-Cell Data
208 Integration competition, specifically designed to contain nested batch effects for evaluating
209 integration. In particular, Luecken et al. (2021) profiled bone marrow mononuclear cells from
210 multiple donors across distinct sites, with inter-site batch effects being larger than intra-site
211 batch effects between donors. For demonstration purposes, we only use the scRNA-seq data of
212 the Multiome data accessible via GEO accession: GSE194122, in particular, a preprocessed
213 AnnData object provided as a supplementary file. We further used a minimal data subset
214 (minimal example) to illustrate the unreliable behavior of silhouette-based metrics with nested
215 batch effects with four samples from four donors from two distinct sites `s1d1`, `s1d3`, `s4d8`, and
216 `s4d9`, for our main figure panels, which we renamed to Sample 1, 2, 3, and 4, respectively. We
217 also consider the full data set, with results shown in Supplementary Figure 2.

218

219 **Data integration**

220 *Simulated data*

221 No integration was performed, as we have simulated differing levels of nested batch effects,
222 which can in turn be interpreted as varying success at batch effect removal.

223

224 *Real data*

225 To demonstrate the insensitivity of silhouette-based batch removal metrics to differing levels of
226 nested batch effects, we aimed to obtain integration results with varying success. The data was
227 first normalized to median total counts and logarithmized, and then dimensionality reduced with
228 PCA. No integration ("None") serves as a baseline. A naive, mild batch correction
229 ("Suboptimal") was achieved by batch-aware selection of highly variable genes (hvg), prioritizing

230 genes that are highly variable across batches, which is applied before PCA (carried out with
231 scanpy (Wolf et al., 2018)). To obtain different batch removal strengths, we used our tunable
232 model liam (Rautenstrauch & Ohler, 2024), which gives us control over distinct batch removal
233 strengths. In particular, we applied liam, to the raw scRNA-seq data of the BMCC Multiome data
234 set with default parameters ("Effective"), and increased batch removal by setting the adversarial
235 scaling parameter to 5 ("Optimized"). Of note, the findings related to the metrics are not specific
236 to the integration models used.

237

238 **Evaluation**

239 **Overview**

240 We assess horizontal data integration using a broad selection of metrics, in particular, Batch
241 ASW, iLISI, CiLISI, BRAS, and BRAS variants for batch removal, and cLISI, Cell type ASW,
242 NMI cluster/label and ARI cluster/label for bio-conservation.

243

244 For most metrics, we use the scib package, except for the implementations for the custom
245 CiLISI and newly proposed BRAS metrics (detailed below).

246

247 All metrics are scaled to range between 0 and 1, with 1 being optimal. For the silhouette-based
248 metric Cell type ASW this implies that original silhouette scores around 0 correspond to
249 transformed scores of approximately 0.5. We use low-dimensional embeddings as input: PCA
250 embeddings for simulated data, and PCA or liam embeddings for the NeurIPS data.

251

252 **Custom implementations of batch removal metrics robust to nested batch effects**

253 *CiLISI*: We implement a custom version of CiLISI (Andreatta et al., 2024), a cell-type aware
254 version of iLISI. First, we compute iLISI (range 0-1, scib implementation) per given cell type
255 label, which is summarized into a weighted mean (weighted by number of cells per cell type
256 label).

257

258 *Batch removal adapted silhouette (BRAS)*: To account for nested batch effects in single-cell
259 data, we modify the silhouette score s_i as described in equation 1. Specifically, we redefine b_i
260 as the mean distance between a cell i and all other cells in **any** other cluster (default in BRAS).
261 We also test a version with b_i as the distance between a cell i and all other cells in the **furthest**
262 other cluster.

263

264 The modified silhouette score is computed per cell i assigned to a cluster C_k . Following Luecken
265 et al.'s (2022) implementation:

266

267 $s_i = |s_i|$, with s_i computed as in equation 1.

268 Then, for each cell type label k corresponding to cluster C_k we define the BRAS score as:

$$269 \text{BRAS}_k = \frac{1}{|N_k|} \sum_{i \in N_k} 1 - s_i$$

270 where N_k denotes the set of cells assigned to cluster C_k and $|N_k|$ the number of cells in that set.

271 For the final *BRAS* score, we average over the set of unique cell labels M .

272

$$BRAS = \frac{1}{|M|} \sum_{k \in M} BRAS_k$$

273 We use cosine distance as the default for BRAS, finding it provides higher discriminative power
274 than euclidean distance (Supplementary Figure 1(a) and (b) and Supplementary Figure 2(b)).

275 We also compute Batch ASW and Cell type ASW with cosine distance.

276

277 *Details on ARI cluster/label and NMI cluster/label.*

278 Following Luecken et al. (2022), we optimized (Leiden) clustering with respect to the ARI and
279 NMI metric across a range of clustering resolutions (0-2, step 0.1) and show these results in
280 Figure 1 and Supplementary Figure 1 and 2 (Leiden is now the current default in scib, in the
281 original publication the Louvain algorithm was used). For a discussion on potential limitations of
282 this strategy, its impact on our results and alternative strategies see Supplementary Note 1 and
283 Supplementary Figures 3-5.

284

285 Code availability

286 The scripts and notebooks for data preprocessing, analyses, and figure generation are publicly
287 available at https://github.com/ohlerlab/metrics_matter_manuscript_reproducibility and will be
288 deposited in Zenodo upon acceptance.

289 References

- 290 Andreatta, M., Hérault, L., Gueguen, P., Gfeller, D., Berenstein, A. J., & Carmona, S. J. (2024). Semi-
291 supervised integration of single-cell transcriptomics data. *Nature Communications*, *15*(1), 872.
292 <https://doi.org/10.1038/s41467-024-45240-z>
293
- 294 Argelaguet, R., Cuomo, A. S. E., Stegle, O., & Marioni, J. C. (2021). Computational principles and
295 challenges in single-cell data integration. *Nature Biotechnology*, *39*(10), 1202–1215.
296 <https://doi.org/10.1038/s41587-021-00895-7>
297
- 298 Büttner, M., Miao, Z., Wolf, F. A., Teichmann, S. A., & Theis, F. J. (2019). A test metric for assessing
299 single-cell RNA-seq batch correction. *Nature Methods*, *16*(1), 43–49. [https://doi.org/10.1038/s41592-018-](https://doi.org/10.1038/s41592-018-0254-1)
300 [0254-1](https://doi.org/10.1038/s41592-018-0254-1)
301
- 302 Cole, M. B., Risso, D., Wagner, A., DeTomaso, D., Ngai, J., Purdom, E., Dudoit, S., & Yosef, N. (2019).
303 Performance Assessment and Selection of Normalization Procedures for Single-Cell RNA-Seq. *Cell*
304 *Systems*, *8*(4), 315–328.e8. <https://doi.org/10.1016/j.cels.2019.03.010>
305
- 306 Haghverdi, L., Lun, A. T. L., Morgan, M. D., & Marioni, J. C. (2018). Batch effects in single-cell RNA-
307 sequencing data are corrected by matching mutual nearest neighbors. *Nature Biotechnology*, *36*(5), 421–
308 427. <https://doi.org/10.1038/nbt.4091>
309
- 310 Hu, Y., Wan, S., Luo, Y., Li, Y., Wu, T., Deng, W., Jiang, C., Jiang, S., Zhang, Y., Liu, N., Yang, Z., Chen,
311 F., Li, B., & Qu, K. (2024). Benchmarking algorithms for single-cell multi-omics prediction and integration.
312 *Nature Methods*, *21*(11), 2182–2194. <https://doi.org/10.1038/s41592-024-02429-w>
313
- 314 Korsunsky, I., Millard, N., Fan, J., Slowikowski, K., Zhang, F., Wei, K., Baglaenko, Y., Brenner, M., Loh,
315 P., & Raychaudhuri, S. (2019). Fast, sensitive and accurate integration of single-cell data with Harmony.
316 *Nature Methods*, *16*(12), 1289–1296. <https://doi.org/10.1038/s41592-019-0619-0>
317
- 318 Lance, C., Luecken, M. D., Burkhardt, D. B., Cannoodt, R., Rautenstrauch, P., Laddach, A.,
319 Ubingazhibov, A., Cao, Z.-J., Deng, K., Khan, S., Liu, Q., Russkikh, N., Ryazantsev, G., Ohler, U.,
320 NeurIPS 2021 Multimodal data integration competition participants, Pisco, A. O., Bloom, J.,
321 Krishnaswamy, S., & Theis, F. J. (2022). Multimodal single cell data integration challenge: Results and
322 lessons learned. In D. Kiela, M. Ciccone, & B. Caputo (Eds.), *Proceedings of the NeurIPS 2021*
323 *Competitions and Demonstrations Track* (Vol. 176, pp. 162–176). PMLR.
324 <https://proceedings.mlr.press/v176/lance22a.html>
325
- 326 Luecken, M. D., Burkhardt, D. B., Cannoodt, R., Lance, C., Agrawal, A., Aliee, H., Chen, A. T., Deconinck,
327 L., Detweiler, A. M., Granados, A., Huynh, S., Isacco, L., Joon Kim, Y., Klein, D., de Kumar, B.,
328 Kuppasani, S., Lickert, H., McGeever, A., Mekonen, H., ... Bloom, J. M. (2021). A sandbox for prediction
329 and integration of DNA, RNA, and proteins in single cells. *Thirty-Fifth Conference on Neural Information*
330 *Processing Systems Datasets and Benchmarks Track (Round 2)*.
331 <https://openreview.net/forum?id=gN35BGa1Rt>
332
- 333 Luecken, M. D., Büttner, M., Chaichoompu, K., Danese, A., Interlandi, M., Mueller, M. F., Strobl, D. C.,
334 Zappia, L., Dugas, M., Colomé-Tatché, M., & Theis, F. J. (2022). Benchmarking atlas-level data

- 335 integration in single-cell genomics. *Nature Methods*, 19(1), 41–50. [https://doi.org/10.1038/s41592-021-](https://doi.org/10.1038/s41592-021-01336-8)
336 01336-8
337
338 Maan, H., Zhang, L., Yu, C., Geuenich, M. J., Campbell, K. R., & Wang, B. (2024). Characterizing the
339 impacts of dataset imbalance on single-cell data integration. *Nature Biotechnology*.
340 <https://doi.org/10.1038/s41587-023-02097-9>
341
342 Rautenstrauch, P., & Ohler, U. (2024). Liam tackles complex multimodal single-cell data integration
343 challenges. *Nucleic Acids Research*, 52(12), e52–e52. <https://doi.org/10.1093/nar/gkae409>
344
345 Risso, D., Perraudeau, F., Gribkova, S., Dudoit, S., & Vert, J.-P. (2018). A general and flexible method for
346 signal extraction from single-cell RNA-seq data. *Nature Communications*, 9(1), 284.
347 <https://doi.org/10.1038/s41467-017-02554-5>
348
349 Rousseeuw, P. J. (1987). Silhouettes: a graphical aid to the interpretation and validation of cluster
350 analysis. In *Journal of Computational and Applied Mathematics* (Vol. 20).
351
352 Scialdone, A., Tanaka, Y., Jawaid, W., Moignard, V., Wilson, N. K., Macaulay, I. C., Marioni, J. C., &
353 Göttgens, B. (2016). Resolving early mesoderm diversification through single-cell expression profiling.
354 *Nature*, 535(7611), 289–293. <https://doi.org/10.1038/nature18633>
355
356 Tran, H. T. N., Ang, K. S., Chevrier, M., Zhang, X., Lee, N. Y. S., Goh, M., & Chen, J. (2020). A
357 benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome Biology*,
358 21(1), 12. <https://doi.org/10.1186/s13059-019-1850-9>
359
360 Wagner, A., Regev, A., & Yosef, N. (2016). Revealing the vectors of cellular identity with single-cell
361 genomics. *Nature Biotechnology*, 34(11), 1145–1160. <https://doi.org/10.1038/nbt.3711>
362
363 Wolf, F. A., Angerer, P., & Theis, F. J. (2018). SCANPY: large-scale single-cell gene expression data
364 analysis. *Genome Biology*, 19(1), 15. <https://doi.org/10.1186/s13059-017-1382-0>
365
366 Zappia, L., Phipson, B., & Oshlack, A. (2017). Splatter: simulation of single-cell RNA sequencing data.
367 *Genome Biology*, 18(1), 174. <https://doi.org/10.1186/s13059-017-1305-0>