

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

4 Pia Rautenstrauch (ORCID: 0000-0002-0070-4759)<sup>1,2</sup> and Uwe Ohler (ORCID: 0000-0002-  
5 0881-3116)<sup>1,2,3\*</sup>

6 <sup>1</sup>Humboldt-Universität zu Berlin, Department of Computer Science, 10099 Berlin, Germany.

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

9 <sup>3</sup>Humboldt-Universität zu Berlin, Department of Biology, 10099 Berlin, Germany.

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

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

## 12 Supplementary Note 1

### 13 Impact of clustering strategy on bio-conservation metrics ARI and NMI

14 To compute ARI or NMI, used here to score bio-conservation, we need to compare a clustering  
15 for any given input to ground truth labels (cell type labels). The choice of clustering algorithm  
16 and hyperparameters affects results. Luecken et al. (2022) opted to optimize clustering for the  
17 Louvain algorithm with respect to the NMI and ARI metrics across a range of clustering  
18 resolutions. This strategy can lead to cluster numbers strongly deviating from the number of  
19 ground truth cell type labels and distinct number of clusters for any given scenario, complicating  
20 comparisons and potentially favoring unrealistic solutions. Recently, Maan et al. (2024) chose to  
21 optimize clustering based on the actual number of ground truth clusters (cell types).

22

23 A recent study proved that the NMI metric can exhibit biased behavior when the number of  
24 detected clusters exceeds the true number of clusters (Mahmoudi & Jemielniak, 2024). In light  
25 of this, and due to the potential limitations of optimizing with little constraints, we sought to  
26 assess the impact of different strategies for deriving a clustering to compare to ground truth  
27 labels with ARI and NMI.

28

29 We compare the results of choosing the maximum score in the full range of tested resolutions  
30 (0-2, step 0.1) of the Leiden clustering algorithm with choosing a maximum score only for results  
31 whose number of clusters is within  $\pm 20\%$  (bounded) of the ground truth (cell type labels).  
32 Supplementary Figures 3-5 a) show at which resolution and respective number of clusters  
33 maximal scores were reached in the full range and in the bounded region. Supplementary  
34 Figures 3-5 b) illustrate how this impacts the overall ranking of distinct scenarios for the different  
35 data sets.

36

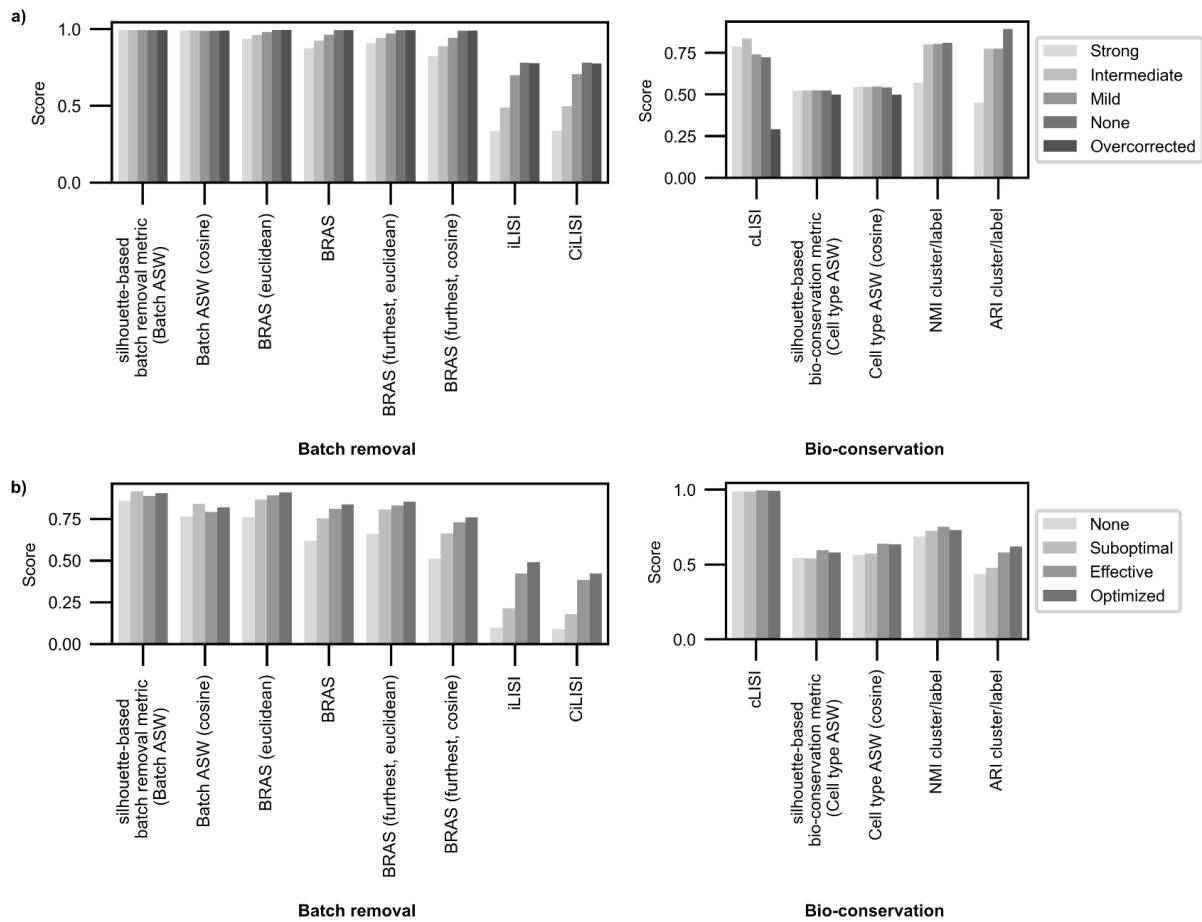
37 We find that this choice impacts results. For example, for the full NeurIPS data, scenario  
38 “Suboptimal”, the maximum scores for ARI and NMI in the full range corresponds to a clustering  
39 output (12 clusters) that strongly deviates from the number of ground truth clusters (22)  
40 (Supplementary Figure 5 a)). In several cases, the choice of strategy even led to different  
41 rankings (e.g., ARI for NeurIPS data minimal example, Supplementary Figure 4 b)).

42

43 These findings do not affect the main conclusions of our paper regarding silhouette-based  
44 metrics, but rather underscore that exploring optimization strategies based on the number of  
45 ground truth clusters needs further investigation.

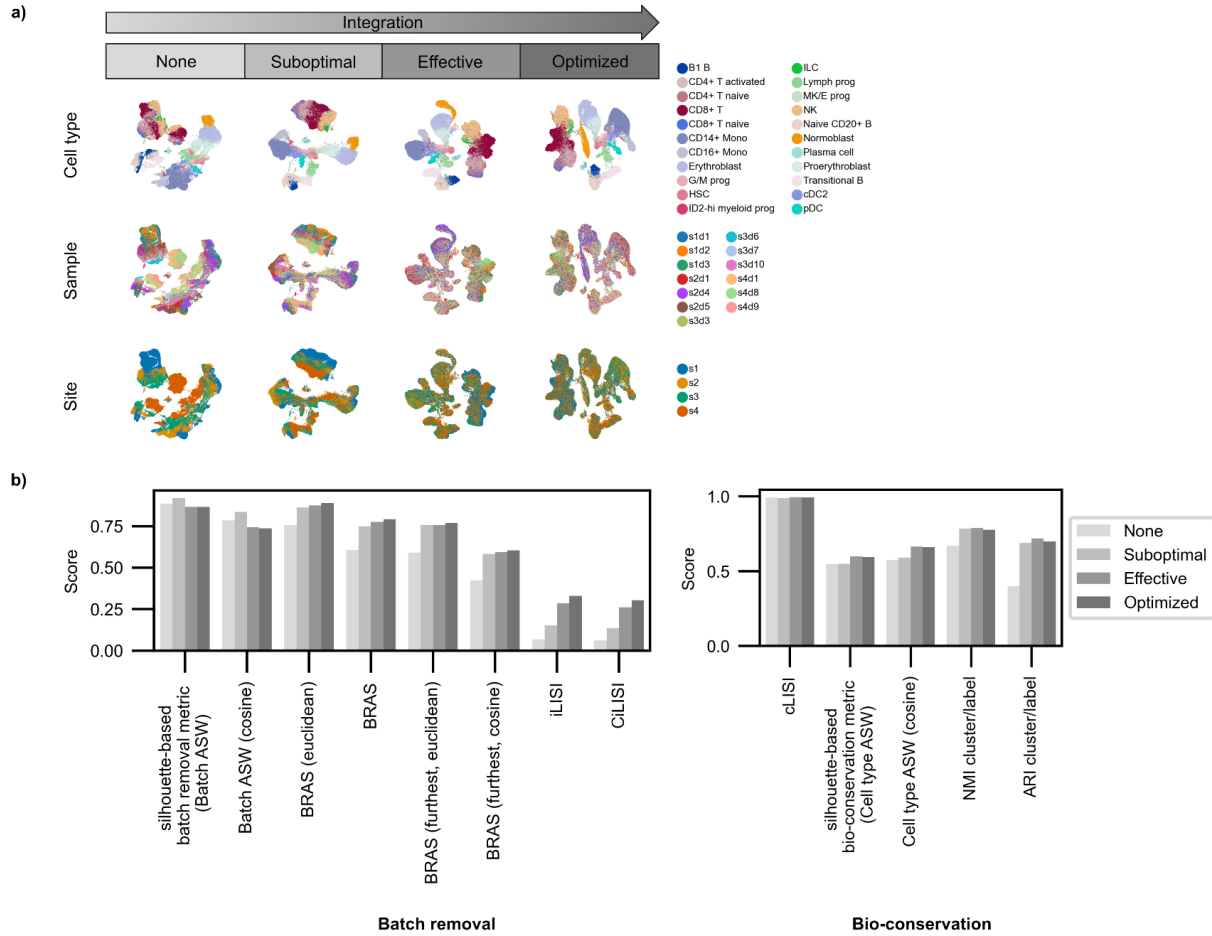
46 Supplementary Figures

47

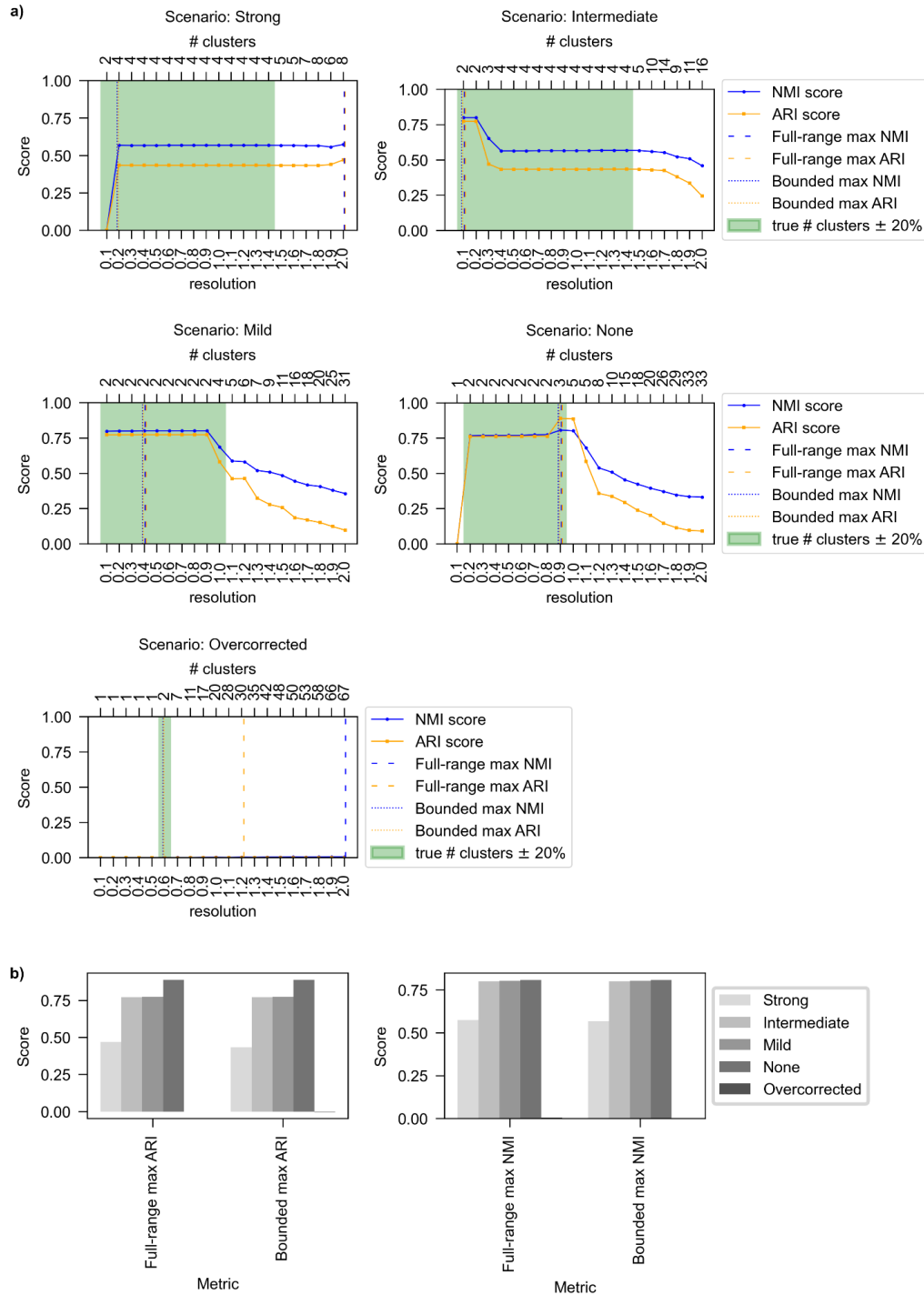


48

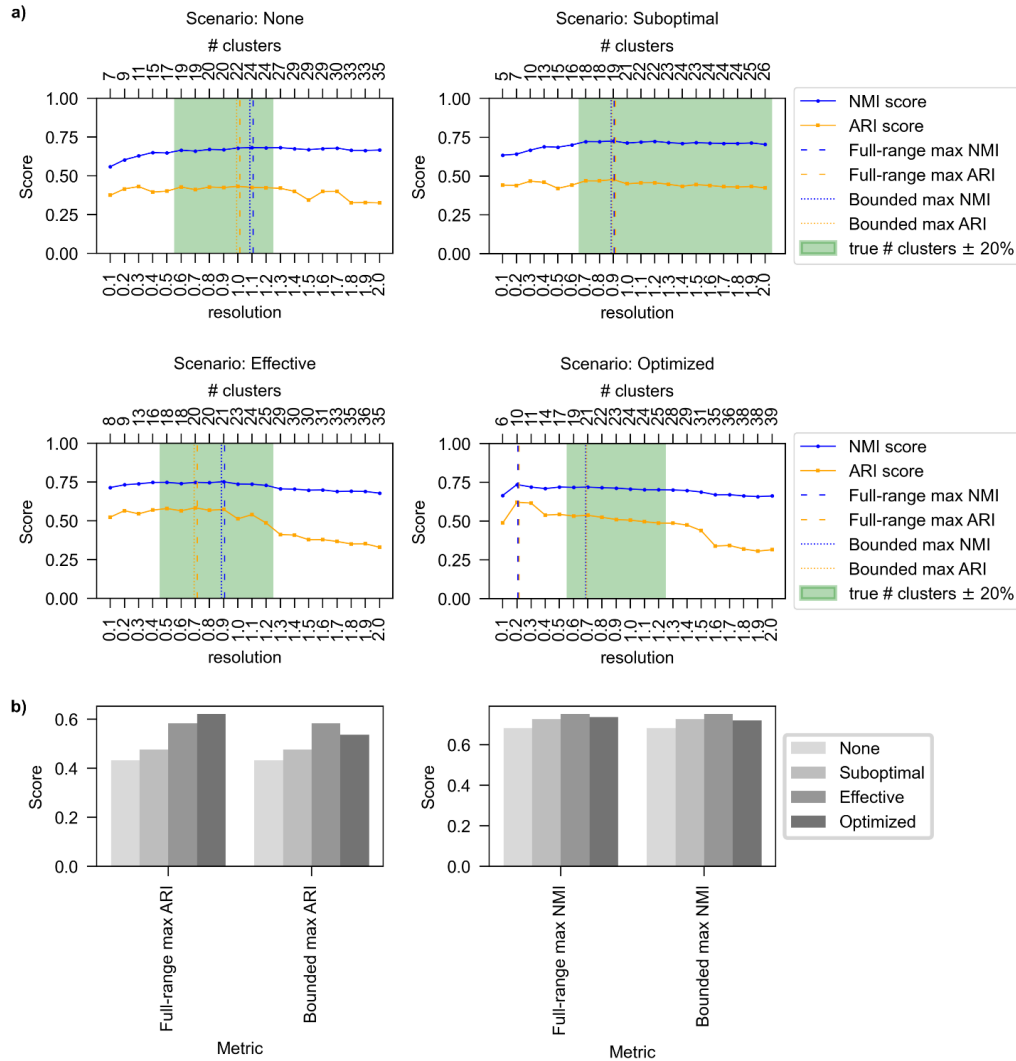
49 **Supplementary Figure 1:** Extended evaluation metrics. Batch removal and bio-conservation  
 50 metrics (a) for simulated and (b) for real data minimal example (cf. Figure 1).



51  
 52 **Supplementary Figure 2: Silhouette-based metrics (Batch ASW) are unreliable with**  
 53 **nested batch effects, failing single-cell data integration evaluation (2).**  
 54 **a)** UMAPs of full NeurIPS data set with nested batch effects integrated with increasing success,  
 55 colored by cell type, sample, and site. **(b)** Extended evaluation metrics.



56  
57 **Supplementary Figure 3: Impact of clustering strategy on ARI and NMI bio-conservation**  
58 **metrics for simulated data. a)** Relationship between Leiden clustering resolution (bottom x-  
59 axis), resulting cluster count (top x-axis), and corresponding ARI and NMI scores. Dashed lines  
60 indicate resolution and cluster count for maximum metric score across full range (0-2, step 0.1).  
61 Green area highlights results within  $\pm 20\%$  of true cluster count. Dotted lines show resolution and  
62 cluster count for maximum score within bounded range. True cluster count: 3. **b)** Comparison of  
63 max scores from different clustering strategies shown in a).



65

66 **Supplementary Figure 4: Impact of clustering strategy on ARI and NMI bio-conservation**

67 **metrics for NeurIPS data minimal example. a)** Relationship between Leiden clustering

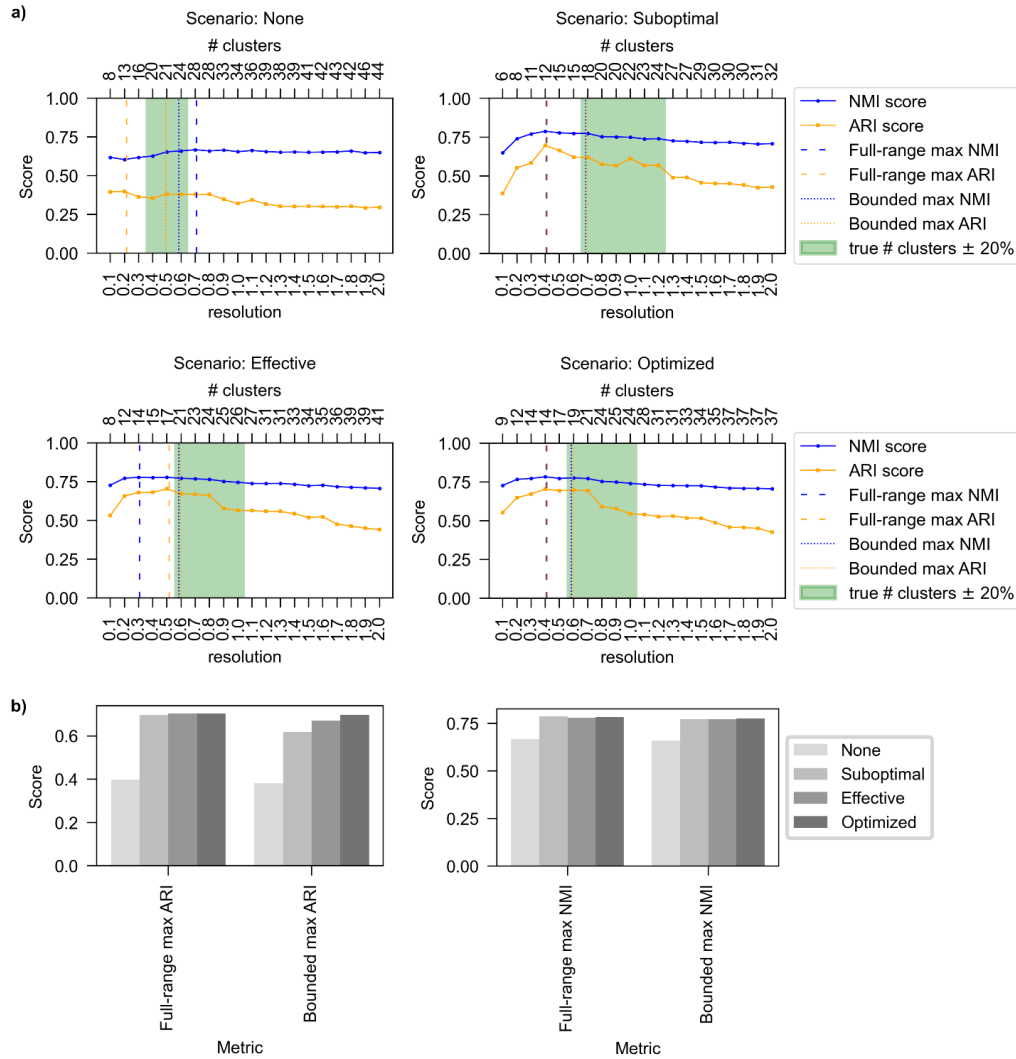
68 resolution (bottom x-axis), resulting cluster count (top x-axis), and corresponding ARI and NMI

69 scores. Dashed lines indicate resolution and cluster count for maximum metric score across full

70 range (0-2, step 0.1). Green area highlights results within  $\pm 20\%$  of true cluster count. Dotted

71 lines show resolution and cluster count for maximum score within bounded range. True cluster

72 count: 22. **b)** Comparison of max scores from different clustering strategies shown in a).



73  
 74 **Supplementary Figure 5: Impact of clustering strategy on ARI and NMI bio-conservation**  
 75 **metrics for full NeurIPS data. a)** Relationship between Leiden clustering resolution (bottom x-  
 76 axis), resulting cluster count (top x-axis), and corresponding ARI and NMI scores. Dashed lines  
 77 indicate resolution and cluster count for maximum metric score across full range (0-2, step 0.1).  
 78 Green area highlights results within  $\pm 20\%$  of true cluster count. Dotted lines show resolution and  
 79 cluster count for maximum score within bounded range. True cluster count: 22. **b)** Comparison  
 80 of max scores from different clustering strategies shown in a).

## 81 Supplementary References

- 82 Luecken, M. D., Büttner, M., Chaichoompu, K., Danese, A., Interlandi, M., Mueller, M. F., Strobl, D. C.,  
83 Zappia, L., Dugas, M., Colomé-Tatché, M., & Theis, F. J. (2022). Benchmarking atlas-level data  
84 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)  
85 [01336-8](https://doi.org/10.1038/s41592-021-01336-8)
- 86
- 87 Maan, H., Zhang, L., Yu, C., Geuenich, M. J., Campbell, K. R., & Wang, B. (2024). Characterizing the  
88 impacts of dataset imbalance on single-cell data integration. *Nature Biotechnology*.  
89 <https://doi.org/10.1038/s41587-023-02097-9>
- 90
- 91 Mahmoudi, A., & Jemielniak, D. (2024). Proof of biased behavior of Normalized Mutual Information.  
92 *Scientific Reports*, *14*(1), 9021. <https://doi.org/10.1038/s41598-024-59073-9>