

# Improving replicability in single-cell RNA-Seq cell type discovery with Dune

Hector Roux de Bézieux<sup>1,9</sup>, Kelly Street<sup>2,3</sup>, Stephan Fischer<sup>4</sup>, Koen Van den Berge<sup>5,6</sup>,  
Rebecca Chance<sup>7</sup>, Davide Risso<sup>8</sup>, Jesse Gillis<sup>4</sup>, John Ngai<sup>7</sup>, Elizabeth Purdom<sup>5</sup>,  
Sandrine Dudoit<sup>1,5,9,\*</sup>

<sup>1</sup> Division of Biostatistics, School of Public Health, University of California, Berkeley, CA, USA

<sup>2</sup> Department of Data Sciences, Dana-Farber Cancer Institute, Boston, MA, USA

<sup>3</sup> Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA

<sup>4</sup> Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA

<sup>5</sup> Department of Statistics, University of California, Berkeley, CA, USA

<sup>6</sup> Department of Applied Mathematics, Computer Science and Statistics, Ghent University, Ghent, Belgium

<sup>7</sup> Department of Molecular and Cell Biology, University of California, Berkeley, CA, USA

<sup>8</sup> Department of Statistical Sciences, University of Padova, Padova, Italy

<sup>9</sup> Center for Computational Biology, University of California, Berkeley, CA, USA

\* To whom correspondence should be addressed: sandrine@stat.berkeley.edu

March 30, 2020

## Abstract

Single-cell transcriptome sequencing (scRNA-Seq) has allowed many new types of investigations at unprecedented and unique levels of resolution. Among the primary goals of scRNA-Seq is the classification of cells into potentially novel cell types. Many approaches build on the existing clustering literature to develop tools specific to single-cell applications. However, almost all of these methods rely on heuristics or user-supplied parameters to control the number of clusters identified. This affects both the resolution of the clusters within the original dataset as well as their replicability across datasets. While many recommendations exist to select these tuning parameters, most of them are quite ad hoc. In general, there is little assurance that any given set of parameters will represent an optimal choice in the ever-present trade-off between cluster resolution and replicability. For instance, it may be the case that another set of parameters will result in more clusters that are also more replicable, or in fewer clusters that are also less replicable.

Here, we propose a new method called *Dune* for optimizing the trade-off between the resolution of the clusters and their replicability across datasets. Our method takes as input a set of clustering results on a single dataset, derived from any set of clustering algorithms and associated tuning parameters, and iteratively merges clusters within partitions in order to maximize their concordance between partitions. As demonstrated on a variety of scRNA-Seq datasets from different platforms, *Dune* outperforms existing techniques, that rely on hierarchical merging for reducing the number of clusters, in terms of replicability of the resultant merged clusters. It provides an objective approach for identifying replicable consensus clusters most likely to represent common biological features across multiple datasets.

Improvements in single-cell transcriptome sequencing (scRNA-Seq) over the last decade have allowed the characterization of gene expression in collections of thousands to hundreds of thousands of cells. While datasets have grown in size by several orders of magnitude, cell type identification remains a primary step in the analysis process [1]. We will focus here on unsupervised clustering, which can be broadly defined as partitioning observations into clusters based on a set of features, without using any prior knowledge on the groupings. In the scRNA-Seq context, clustering aims to identify groups of cells that are defined by a unique and consistent transcriptomic signature. Such groups of cells can represent both transient features, such as cellular states, or more permanent features, such as cellular types.

Many clustering algorithms have been proposed for scRNA-Seq, most of these being adaptations from the clustering literature at large. Popular methods include SC3 [2], Seurat [3], and Monocle [4]. However,

36 clustering remains a complex task. Kiselev et al. [5] outlined the various challenges – both biological and  
37 computational – of this step, including technical noise, biological heterogeneity, and the impact of tuning  
38 parameters for the clustering algorithms. In particular, obtaining replicable clusters can be difficult. In  
39 this work, we declare clusters as replicable if running the exact same clustering algorithm on a related  
40 dataset yields similar clusters. Duò et al. [6] offers a recent review and benchmark of some scRNA-Seq  
41 clustering algorithms, identifying SC3 and Seurat as the best methods overall. The selection of tuning  
42 parameters, however, remains an open question. While some methods, SC3 for example, provide a way  
43 to estimate the optimal value of its main tuning parameter, most do not, leaving the choice to the user.  
44 Consensus methods try to bypass this issue [2, 7], but they also rely on meta-parameters which can still  
45 have substantial impact on the results.

46 The aforementioned clustering algorithms identify a pre-specified number of clusters either directly,  
47 as in  $k$ -means, or indirectly, through another tuning parameter. They rely on the assumption that there  
48 is only one relevant level of clustering resolution, i.e., an optimal number of clusters, in the dataset.  
49 We argue that this is often not the case, since cell types usually have a hierarchy. For example, Tasic  
50 et al. [8] propose a tree structure for the mouse anterolateral motor (ALM) and primary visual (VISp)  
51 cortical areas. At the higher levels, cells can be clustered as neurons and non-neurons. Then, neurons  
52 can be further split into GABAergic and glutamatergic neurons and so on and so forth. This hierarchical  
53 structure means that the concept of an “optimal” number of clusters is not appropriate. Instead, many  
54 datasets can be better characterized with ever-finer levels of resolution. At the higher levels, cells are  
55 grouped into large clusters that are quite coarse, but are easily identifiable and very replicable across  
56 datasets. As the resolution increases, there are more and more clusters, but these are less and less certain,  
57 meaning that they are less likely to represent real biological cell types and more likely to be reflecting over-  
58 partitioning (cf. overfitting) of the data or the presence of transient states. This resolution-replicability  
59 trade-off is not obvious to quantify and is heavily dataset-dependent: it is not only influenced by the  
60 biological setting under study and its complexity, but also depends on technical properties of the data,  
61 such as sequencing depth and number of cells [1].

62 By far the most common method to establish a hierarchy for pre-defined clusters is agglomerative  
63 hierarchical clustering, a bottom-up method in which clusters are merged one-by-one until they are all  
64 merged into a single cluster. This procedure yields a tree structure linking clusters that are merged  
65 together. The tree can also be defined by merging clusters according to the fraction of differentially  
66 expressed (DE) genes between them [7, 8]. While several extensive benchmarks of clustering methods  
67 have been proposed [6, 9], these only focus on the resulting partitions rather than the full hierarchical  
68 structure. Zappia and Oshlack [10] proposed a representation of clustering trees to visually describe  
69 hierarchies but this type of analysis heavily relies on user-supervision.

70 Here, we present Dune, a method that aims to reconcile multiple clustering results and extract the  
71 common structure that they all identify. Dune takes as input a set of clustering results (i.e., results  
72 from a variety of clustering algorithms and associated tuning parameters applied to a given dataset) and  
73 produces hierarchies of clusters by merging clusters within each partition using information borrowed  
74 from the other partitions. While different clustering algorithms run with different tuning parameters will  
75 naturally provide discrepant clusters, all good clustering methods should be able to identify a common  
76 higher-level clustering that is robust to the choice of tuning parameters. Dune identifies this common  
77 higher level of resolution shared by all methods without requiring any tuning by the user. Examining  
78 this level can provide both useful biological insight and help to compare various clustering methods.

79 In this manuscript, we first introduce the Dune algorithm. Then, using a variety of scRNA-Seq and  
80 snRNA-Seq datasets from different sequencing platforms, we show that Dune outperforms agglomerative  
81 merging methods in navigating the trade-off between resolution and replicability and in identifying gold-  
82 standard high-level clusterings. Finally, we assess Dune’s robustness to poor inputs and to sample size.

## 83 Results

### 84 The Dune algorithm

85 The Dune algorithm is a general framework that increases the agreement between different clusterings  
86 of the same dataset through iterative merging. It takes as input  $R$  sets of clustering results, generally  
87 produced from running  $R$  clustering algorithms (or the same algorithm with different tuning parameter  
88 values) on the same dataset. An example can be seen in Figure 1a, where a small subset of the AIBS  
89 snRNA-Smart dataset [11] (see the “Methods, Case Studies” section) is used to demonstrate some of  
90 the main concepts underlying Dune. The first row displays three examples of clusterings (i.e., sets of

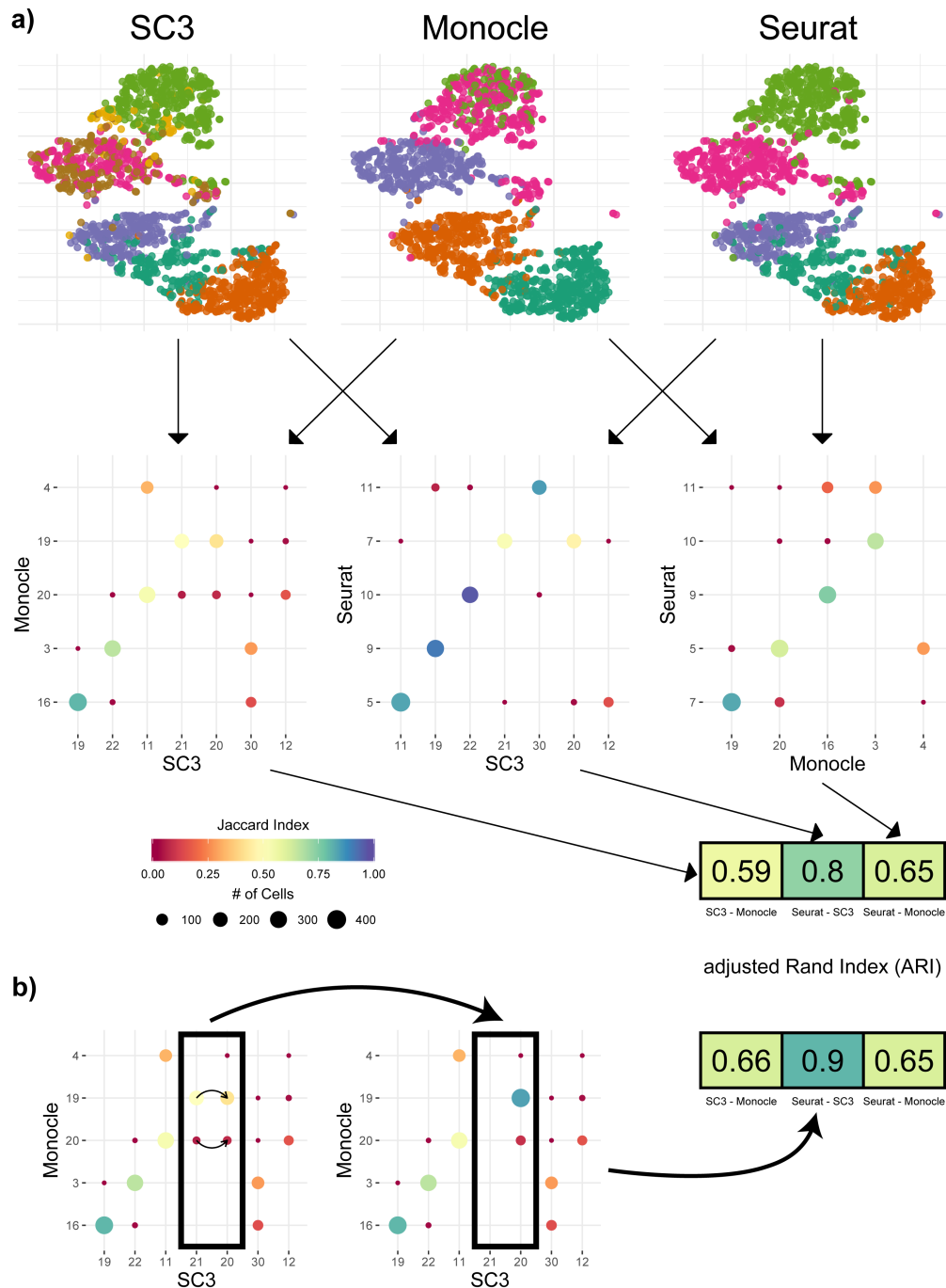


Figure 1: *Measuring and improving the concordance between clusterings.* We used a subset of the AIBS snRNA-Smart dataset as an example. Panel **a**. SC3, Monocle, and Seurat were run on the dataset and their results are displayed using scatterplots of the first two t-SNE components, where the color of the plotting symbol corresponds to the cluster label. Each pair of clusterings was then compared using a confusion matrix, resulting in three such matrices. For a pair of clusterings/partitions, a confusion matrix is a contingency table, where each entry corresponds to the number of observations in both a cluster from the first partition and a cluster from the second. The size of the dot represents the number of observations in both clusters and the color corresponds to the Jaccard index. Each confusion matrix produces one ARI value. Panel **b**. Merging clusters 20 and 21 from SC3 into one cluster changes the confusion matrix and increases the ARI.

91 cluster labels) produced by three different clustering algorithms applied to the same dataset, reduced

92 to two dimensions using t-SNE[12–14]. All three methods identify similar, but not identical clusters.  
93 Indeed, the algorithms output partitions with different levels of resolution. For example, Monocle splits  
94 the bottom region (in reduced dimension) into two clusters, while the other two methods find three  
95 clusters. Likewise, Monocle and SC3 find two clusters in the top region, while Seurat only finds one.  
96 These differences can be displayed using confusion matrices (second row of Figure 1a), where the overlap  
97 between two clusters from any pair of clusterings is displayed both in terms of the number of cells in  
98 the intersection and by the Jaccard index (i.e., the cardinality of the intersection of the two clusters  
99 over the cardinality of their union; [15]). Rows and columns are ordered so as to maximize, as much  
100 as possible, the sum of the diagonal entries. Confusion matrices can be further summarized using the  
101 adjusted Rand index (ARI). The ARI [16, 17] is a commonly used measure for the agreement between  
102 two sets of clustering labels, see the “Methods, ARI” section for more details. As can be seen in the  
103 confusion matrices, SC3 and Seurat have the highest level of agreement. Indeed, this is also reflected in  
104 the fact that they have the highest ARI of any pair.

105 Dune merges together the clusters within each of the  $R$  partitions so that the  $R$  clustering results  
106 more closely match each other. An example of the merging is displayed in Figure 1. Clusters 20 and 21  
107 from SC3 are merged together, resulting in one larger cluster named 20. Doing so increases the agreement  
108 between SC3 and Monocle in the confusion matrix, as reflected by an increase in ARI from 0.59 to 0.66.  
109 This merge also improves the ARI between SC3 and Seurat (from 0.8 to 0.9) and hence increases the  
110 overall agreement between the three clusterings. This is the main idea behind Dune. Specifically, Dune  
111 performs an iterative search, where, at each iteration, it identifies the partition and pair of clusters within  
112 this partition that, when merged, most improve the average of the adjusted Rand index over all pairs of  
113 clusterings (ARI). Thus, the Dune algorithm can be viewed as an iterative algorithm for maximizing the  
114 average pairwise ARI of a collection of clustering results. A more formal definition of the algorithm is  
115 provided in the “Methods, Dune” section.

116 We demonstrate how the Dune algorithm works in Figure 2, using the AIBS scRNA-Smart dataset, a  
117 scRNA-Seq dataset of 6,300 mouse brain cells further described in the “Methods, Case Studies” section.  
118 For this example, we ran SC3, Seurat, and Monocle to obtain our initial clustering results for input into  
119 Dune ( $R = 3$ ). Figure 2a displays the confusion matrix for a pair of clusterings (SC3 and Monocle)  
120 before any merging and Figure 2b displays a pseudocolor image of the matrix of all pairwise ARIs for the  
121 three clusterings before any merging. The overlap between the three methods is moderate. Indeed, the  
122 pairwise ARIs vary between 0.55 and 0.68 in Fig. 2b. However, as can be seen in the confusion matrix,  
123 the clusterings do capture a shared underlying structure, which will serve as grounding for the Dune  
124 merging. Figure 2d shows the confusion matrix for the same two partitions as in 2a, after merging with  
125 Dune. We can see that we have, by design, fewer clusters in both partitions, but also that the concordance  
126 between the two partitions is greatly improved (as indicated by the color of the plotting symbols, which  
127 represents the Jaccard Index). This is further evidenced in Figure 2e, where the pairwise ARIs between  
128 the three partitions are displayed. The average ARI after all merging steps increased from  $\sim 0.6$  to  
129  $\sim 0.89$ . Figures 2c and 2f demonstrate the evolution of the average ARI and of the number of clusters  
130 per partition through the Dune merging process. At each step, we merge the pair of clusters that leads  
131 to the greatest increase in average ARI. Hence, at each step, the average ARI increases (Fig. 2c) and  
132 the number of clusters in one of the partitions decreases by one (2f). The final partitions are achieved  
133 when the average ARI can no longer be improved.

134 In the following sections, we evaluate Dune and compare it to two hierarchical tree merging methods,  
135 using four datasets: two mouse brain datasets from the Allen Institute \*\*\* HRB: waiting for main paper  
136 and two human pancreas datasets [18, 19]. We then discuss the value of Dune’s stopping rule. Finally,  
137 we investigate the stability of the Dune algorithm to the clustering inputs and the sample size.

## 138 Dune outperforms other methods in recovering known biological subtypes

139 To evaluate Dune, we first considered how well the resulting merged clusters compare to known biological  
140 subtypes. We used the output of Dune on the  $R = 3$  clustering methods (namely, SC3, Seurat, and  
141 Monocle) applied to the AIBS scRNA-Smart dataset, as described above. For this dataset, we treated  
142 the labels from the original publication as the gold standard. At each merge (i.e., iteration), we computed  
143 the ARI between the the known subtypes and the Dune clusters. Figure 3a displays the ARI evolution  
144 for the clusters from SC3 as they are merged with Dune (blue curve). As merging occurs, the resolution  
145 (i.e., number of clusters) decreases and the ARI with the known cell subtypes increases. The entire ARI  
146 curve can be summarized by computing the the area under it, referred to herein as the area under the  
147 ARI curve (AUARIC), as depicted in Figure 3b.

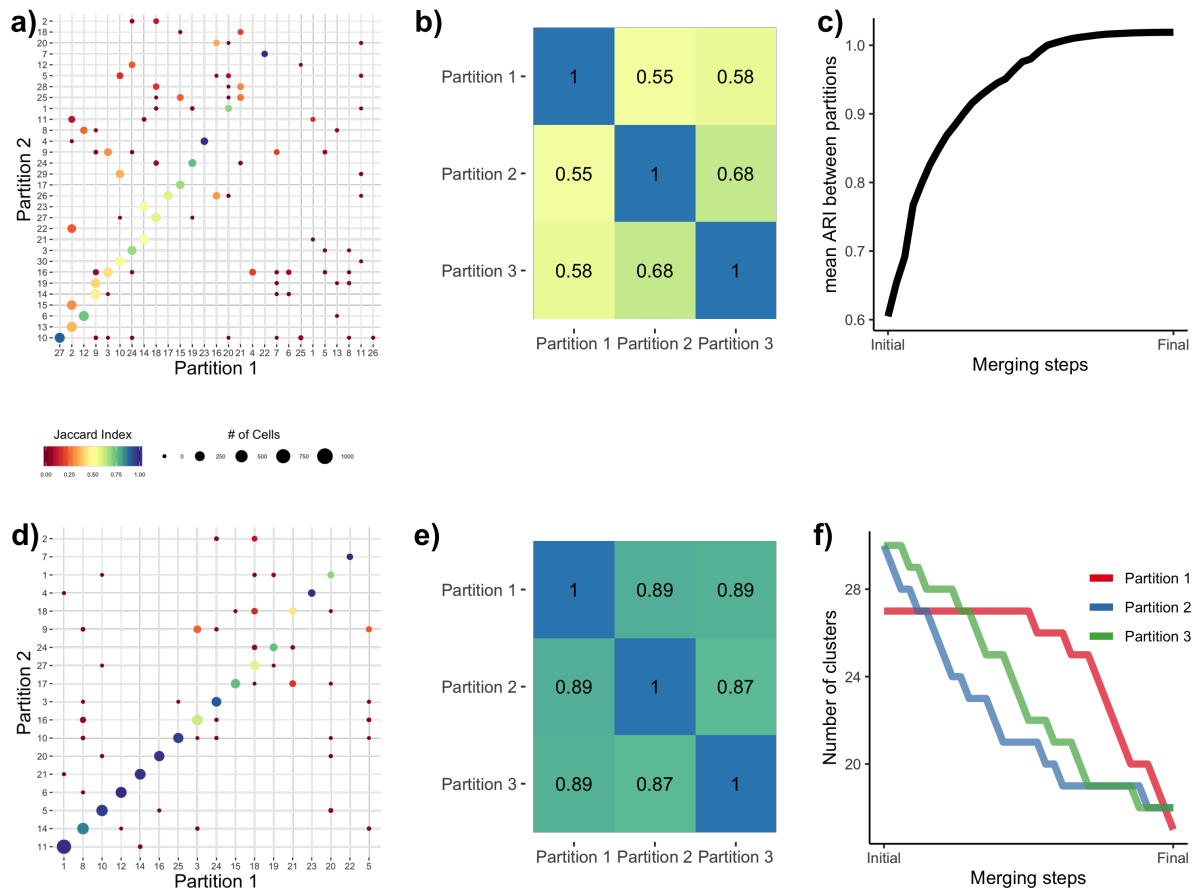


Figure 2: *Illustrating Dune on a dataset with three sets of clusters.* We used the AIBS scRNA-Smart dataset [11] as an example. Before any merging, the sets of cluster labels – or partitions – resulting from running SC3, Seurat, and Monocle have a moderate agreement. Panel a displays the confusion matrix between two of the partitions, where each entry corresponds to the number of observations in both a cluster from Partition 1 and a cluster from Partition 2. The confusion matrix shows that while many cells are clustered in similar clusters, i.e., along the main diagonal, many others are not. This can be summarized by the ARI between Partitions 1 and 2. Panel b displays a pseudocolor image of the matrix of all pairwise ARIs between the three partitions. Panel c illustrates that the average ARI between partitions increases as pairs of clusters are merged when applying Dune. After running Dune, the confusion matrix in Panel d and the pairwise ARI matrix in Panel e both show that the partitions are indeed more similar. Panel f shows that, at each merging step, the number of clusters in one of the partitions is decreased by one, in Dune’s greedy procedure to improve the average ARI by merging pairs of clusters.

148 We compared the performance of Dune to other methods of merging, referred to as Dist and DE (red  
 149 and green curves in Figure 3a, respectively). Both are hierarchical methods, that start by building a tree  
 150 between the clusters. The Dist method then merges clusters in a bottom-up manner, starting with the  
 151 two clusters that are closest in the tree and then iteratively until all clusters are merged. The second  
 152 approach, DE, follows the method implemented in RSEC and merges clusters bottom-up based on the  
 153 percentage of DE genes between clusters. It uses the limma package [20], where a gene is declared DE  
 154 if its nominal false discovery rate (FDR) adjusted  $p$ -value is below 0.05 [21]. Pairs of clusters with less  
 155 than a certain fraction of DE genes are merged. Increasing this threshold from 0 to 1 leads to an iterative  
 156 merging procedure. More details on these two procedures can be found in the [Method section](#).

157 In Figure 3a, we see that Dune consistently outperformed the other two integration methods in terms  
 158 of concordance with BICCN-curated clusters throughout the merging process and therefore also in term  
 159 of AUARIC. We note that while Dune stops merging when the average ARI can no longer be improved,  
 160 the hierarchical merging procedures have no meaningful stopping point and continue merging until only

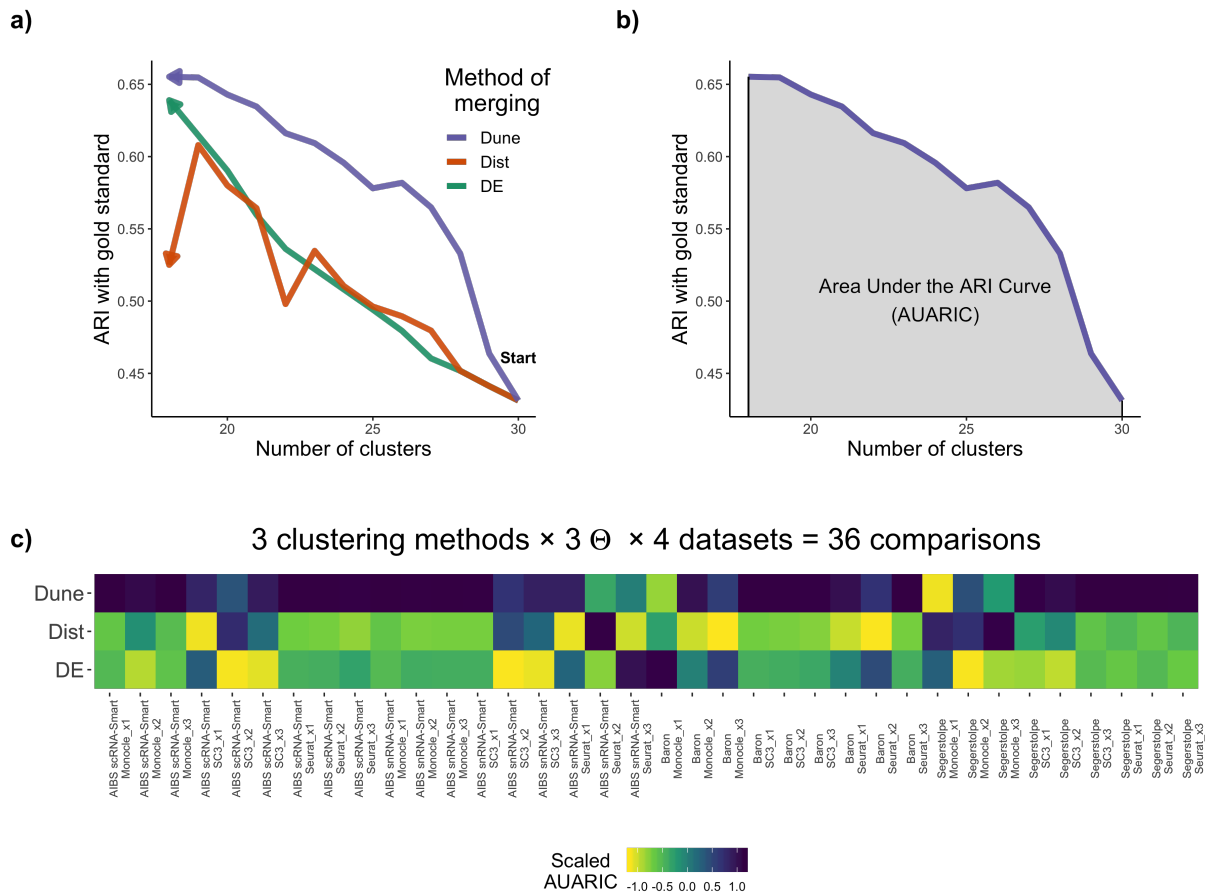


Figure 3: *Dune outperforms other methods in recovering known biological subtypes.* Panel **a**. SC3 was run on the AIBS scRNA-Smart dataset for  $\theta_{sc3} = 0$  and merged down with either DE, Dist, or Dune (with  $\theta_{Monocle} = 45$  and  $\theta_{Seurat} = 1.2$ , for Dune). The ARI with the labels from the original publication, treated as gold standard, was computed at each step of all three merging procedures. Panel **b**. For each merging method from **a**, the area under the ARI curve (AUARIC) was computed. This was repeated for three clustering methods, each with three different values of their respective tuning parameter  $\theta$ , and four datasets. The resulting 36 AUARIC are displayed in the pseudocolor image of Panel **c**. The AUARIC values are scaled to have a column mean of zero and column variance of 1. This was done to make AUARIC values comparable across datasets, clustering methods, and parameter values, since the AUARIC can have different scales across scenarios.

161 one cluster is left. To provide a reasonable stopping point, we stopped the other methods when merging  
 162 no longer improves the ARI, similar to the requirement of Dune, which means we did not penalize the  
 163 other methods for not providing a natural stopping point. For each merging method, we computed an  
 164 area under its ARI curve (AUARIC), as depicted in Figure 3b for the merging of the SC3 clusters of the  
 165 AIBS scRNA-Smart dataset using Dune.

166 Figure 3c show the results when repeating this process over a multiplicity of scenarios. Dune and  
 167 the other merging methods rely on one or multiple clustering results – in this work, clusterings from  
 168 SC3, Seurat, and Monocle. Because each of these methods have tuning parameters than can affect  
 169 their performance, we ran each of the three clustering methods on a grid of tuning parameter values  
 170 for all 4 datasets, as described in the “Methods, Data analysis” section. The AUARIC for the three  
 171 merging methods across these 36 scenarios are displayed in Figure 3c and Table S2. Overall, Dune clearly  
 172 outperformed the other two merging methods. Table S2 recapitulates all rankings. In particular, in 29  
 173 out of the 36 evaluations, Dune resulted in the highest ARI increase and was the lowest performer only  
 174 twice.

## 175 **Dune outperforms other methods in terms of the resolution-replicability trade-** 176 **off**

177 We then considered the replicability of the clusters found by **Dune** compared to the other two merging  
178 strategies. We measured replicability by evaluating whether the method finds similar clusters for multiple  
179 independent datasets – for example, datasets on the same biological system but from different labs or  
180 technologies. We considered two pairs of datasets: The two mouse brain AIBS Smart datasets from the  
181 Allen Institute and the two human pancreas datasets **Baron** and **Segerstople**. To measure replicability,  
182 we relied on the **MetaNeighbor** algorithm from Crow et al. [22], which identifies replicable clusters between  
183 pairs of datasets (see “[Methods, metaneighbour](#)” for description). The replicability of a set of clusters  
184 was then defined as the fraction of cells in replicable clusters. We used this measure to compare **Dune** to  
185 other merging procedures.

## 186 **Illustration of the trade-off between resolution and replicability**

187 Figure 4a displays replicability vs. resolution for a wide range of clustering results, where three clustering  
188 methods (SC3, Seurat, and Monocle) were run with a large grid of tuning parameter values, on the pair  
189 of mouse brain datasets. This clearly demonstrates the trade-off between replicability and resolution:  
190 As the number of clusters increased, the fraction of cells in replicable clusters decreased, regardless of  
191 the clustering method used. While the actual trade-off is specific to the biological context and the pair  
192 of datasets that are being considered, it should be noted that a similar trade-off is clearly visible when  
193 applying the same type of analysis to the human pancreas datasets (Figure S2). Note that although  
194 it might be tempting to use this figure to contrast and benchmark clustering methods, this would not  
195 appropriate. Indeed, pre-processing steps were not identical between the three methods – as described  
196 in “[Methods, Data analysis](#)” – and, as such, no direct comparison is possible.

## 197 **Comparison of merging methods**

198 As pairs of clusters are merged, the resolution decreases, so a well-performing merging method is one  
199 that improves the replicability of the clusters. Therefore, a natural way to benchmark merging methods  
200 is to measure how and if replicability improves as the number of clusters is reduced. For example, in  
201 Figure 4b, Seurat was run with  $\theta_{Seurat} = 1.7$  on each of the two AIBS Smart datasets. The two sets  
202 of clusters were then merged using the three different merging methods, independently on each dataset.  
203 **Dune** also used the clusterings from SC3 ( $\theta_{SC3} = 15$ ) and Monocle ( $\theta_{Monocle} = 15$ ). At each step of  
204 the merging, we then tracked how replicability evolves. All three merging methods outputted sets of  
205 clusters with increasing replicability as resolution decreases, but **Dune** produced clusters that have higher  
206 replicability compared to the other two. The area under the replicability curve (AURC) was computed  
207 for each merging method. This was repeated for the three clustering methods, each with three values of  
208 their respective tuning parameter  $\theta$ , and two pairs of datasets, which lead to 18 comparisons, depicted  
209 in the pseudocolor image of Figure 4c. **Dune** outperformed the other two merging methods in all 18  
210 comparisons. Note that, as in the previous section, merging for the other methods was stopped at the  
211 resolution level where **Dune** stopped, which provided these methods with more information than they  
212 would otherwise have had.

## 213 **Dune has a natural stopping point**

214 Unlike other merging methods, **Dune** provides a meaningful stopping point, i.e., it keeps merging clusters  
215 until no improvement in average ARI occurs. By contrast, the two hierarchical merging methods continue  
216 to merge until there is only one cluster, which is not biologically meaningful or interesting.

217 Each clustering method has some strengths and drawbacks: **Dune**’s stopping point identifies the level  
218 of resolution where all clustering algorithms are close to full agreement. Furthermore, at the stopping  
219 point, the clusters overlap very well with gold-standard clusters. In Figure S3a, the outputs from SC3,  
220 Seurat, and Monocle were used as inputs to **Dune** on the AIBS snRNA-Smart dataset. After merging  
221 with **Dune**, the clusters from SC3 overlap well with the Allen Institute subclass labels. Indeed, the ARI  
222 between the SC3 clusters and the subclasses increases from  $\sim .63$  before merging to  $\sim .83$  after merging.

## 223 **Dune robustness analysis**

224 **Robustness to poor clustering inputs** Since **Dune** takes as input the results from clustering algo-  
225 rithms, it is sensitive to the quality of the clusterings produced by these algorithms. In general, **Dune** will

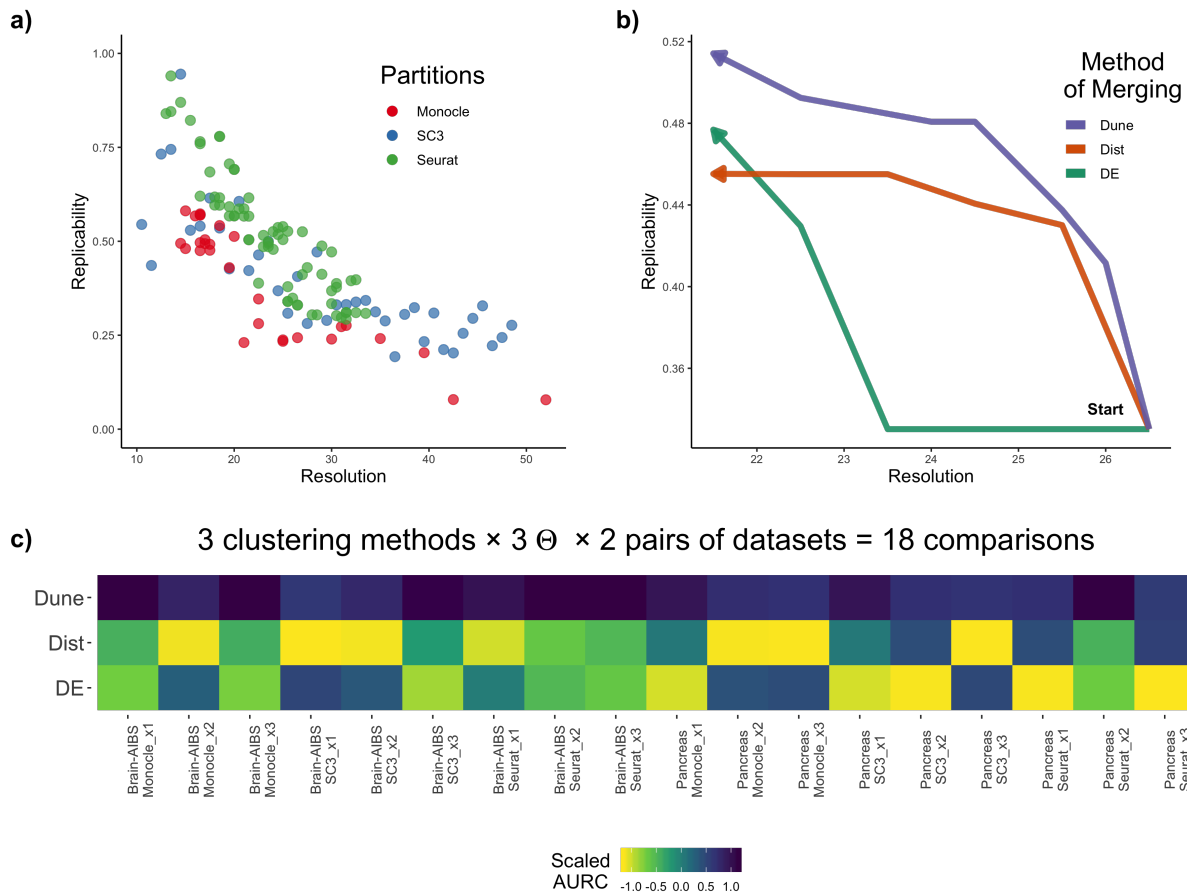


Figure 4: *Dune* correctly navigates the resolution-replicability trade-off. Panel **a**. SC3, Seurat, and Monocle were run on the two AIBS snRNA-Smart datasets, as described in Methods, for a wide range of tuning parameter values. Then, the MetaNeighbor method was used to find the clusters that replicate between these two datasets. Replicability was then computed as the fraction of cells in replicable clusters. There is an apparent trade-off between resolution and replicability. Panel **b**. For a given point from **a**, we merged the clusters and tracked how replicability evolved as we decreased resolution. Panel **c**. For each of the curves in **b**, we computed an area under the replicability curve (AURC). This was repeated over the three clustering methods, each with three different values of their respective tuning parameter  $\theta$ , and for the two pairs of datasets. AURC were scaled column-wise for display in the pseudocolor image.

226 not be able to produce good clusters when merging only clusters that capture no underlying biological  
 227 signal. However, we showed that *Dune* is robust to a mix of “good” clustering inputs and “bad” clus-  
 228 tering inputs. We used as “good” inputs the results of SC3, Seurat, and Monocle and as “bad” inputs  
 229 fully random clusters (see the “Methods, Data analysis” section). Then, the replicability of the “good”  
 230 clusterings was measured as merging happened and the AURC was computed and compared to the AURC  
 231 when there was no “bad” inputs. As more and more “bad” clusters were added (Figure S3b), *Dune* still  
 232 improved the replicability of the “good” clusters as it merged them, even when half of the clusters used  
 233 as inputs were random. Hence, *Dune* can recover from very poor clustering inputs.

234 **Robustness to sample size** We investigated how *Dune* handles datasets with an ever-smaller number  
 235 of cells. To simulate such datasets, we downsampled the two pancreas datasets. Downsampling could  
 236 affect both the quality of input clusters and the merging procedure of *Dune*. To disentangle these two  
 237 effects, we downsampled the two human pancreas datasets after running SC3, Seurat, and Monocle, but  
 238 before running *Dune*. We then measured how and whether merging still improved the cluster replicability  
 239 by computing the AURC and contrasting it to its value without downsampling (see the “Methods, Data  
 240 analysis” section for more details).

241 When the datasets were downsampled to between 90% and as low as 10% of the original number of



242 cells, Dune still correctly navigated the trade-off between resolution and replicability (Fig. S3c). Only  
243 when fewer than 10% of the cells were used (which amounts to datasets of fewer than 200 cells) did Dune’s  
244 capacity to improve cluster replicability worsen noticeably. This demonstrates that the method is very  
245 stable to the number of cells.

## 246 Discussion

247 We have introduced Dune, a new method for navigating the resolution-replicability trade-off in cluster  
248 analysis and for aggregating clustering results from multiple algorithms. We stress that Dune is not a new  
249 clustering algorithm; instead, it relies on different clustering methods to identify the highest resolution  
250 at which cluster quality (i.e., replicability across datasets) remains high. In doing so, Dune identifies the  
251 commonalities of the input clusterings and uses this to improve each of these clusterings. The method is  
252 stable with respect to the quality of the input clusterings as well as to the number of cells/observations  
253 to be clustered. Furthermore, as a result of merging clusters, Dune provides a sensible hierarchy on the  
254 clusters based on their commonality across different methods. As we go up in this hierarchy, the number  
255 of clusters is reduced, but their replicability improves. In this regard, Dune outperforms more commonly  
256 used hierarchical merging methods.

257 Dune automatically stops at a meaningful resolution level, where all clustering algorithms are in  
258 agreement, while the other methods either keep merging until all clusters are merged into one or require  
259 user supervision to stop early. This feature helps users in identifying reliable structure in their scRNA  
260 and snRNA datasets. The manual choice of a stopping point is difficult since, in practice, it is often  
261 impossible to measure replicability given the lack of a second appropriate dataset.

262 Dune relies on the adjusted Rand index (ARI) to decide which clusters to merge. Because of this,  
263 it currently cannot be used with clustering methods that do not cluster all cells unambiguously, e.g.,  
264 with soft or fuzzy clustering methods which could assign some cells to multiple clusters based on weights.  
265 Other approaches, such as RSEC, leave some cells unclustered. For now, using such methods as input to  
266 Dune would require forcing a hard assignments of the cells (possibly to their nearest cluster) or excluding  
267 ambiguous/unclustered cells. Extensions of the ARI to fuzzy clustering have been proposed [23, 24] and  
268 would need to be evaluated.

269 This manuscript focuses on the question of unsupervised clustering. Recent work in supervised cluster-  
270 ing [25–28] has proposed labeling cells in a new dataset by relying on information contained in other  
271 datasets or even cell atlases. In practice, these methods define marker genes for known cell types and  
272 build classifiers to assign new cells to these cell types. In particular, Garnett [29] allows a hierarchical  
273 clustering structure, but one that needs to be predefined, and scClassify [30] uses the HOPACH [31] al-  
274 gorithm to establish a hierarchy in the training dataset. Most of these algorithms can also identify new  
275 cell types not present in the reference. It is therefore possible to use a supervised clustering method to  
276 identify the cells of a dataset that have a known cell types. If these cells do not provide information to  
277 help cluster the rest of the cells, we can remove them, and then use unsupervised clustering methods and  
278 Dune on the remaining cells.

279 While the method we propose has only been benchmarked on scRNA-Seq and snRNA-Seq datasets,  
280 it is a general framework that can be applied to any clustering setting.

## 281 Acknowledgments

282 We thank Quentin Dumont and Frank Herbert [32] for inspiration for the Dune name.

283  
284 This work used the Extreme Science and Engineering Discovery Environment or XSEDE (which  
285 is supported by National Science Foundation grant number ACI-1548562) PSC, Bridges Regular and  
286 Large Memory at the Pittsburgh Supercomputing Center through allocations TG-IBN180019 and TG-  
287 IBN190010. This work was supported by NIH grant U19MH114830 (JN), U19MH114821 (JG) and  
288 R01MH113005(JG). KVdB is a postdoctoral fellow of the Belgian American Educational Foundation  
289 (BAEF) and is supported by the Research Foundation Flanders (FWO), grant 1246220N.

## 290 Authors’ contributions

291 HRB, KS, JN, EP, and SD conceived and designed the study. HRB and KS developed and implemented  
292 the method. HRB, KS, SF, and DR analyzed the data. RC provided resources. HRB and SF wrote the  
293 initial draft of the manuscript, and KS, KVdB, RC, DR, JG, JN, EP, and SD contributed to revisions.

## 294 Data availability

295 The Pancreas datasets were downloaded from the Hemberg group website: <https://hemberg-lab.github.io/scRNA.seq.datasets/human/pancreas/> on October 1<sup>st</sup>, 2018. The AIBS datasets can be obtained from Neuroscience Multi-omics Archive (*RRID* : *SCR\_002001*; [nemoarchive.org](https://nemoarchive.org)), (*Zeng sn Ssv4* <https://assets.nemoarchive.org/dat-k7p82j4> and *Zeng sc Ssv4* <https://assets.nemoarchive.org/dat-55mowp9>).

## 300 Code availability

301 The results from this paper can be reproduced using code from the following GitHub repository: [https://github.com/HectorRDB/Dune\\_Paper](https://github.com/HectorRDB/Dune_Paper). The Dune method is implemented in an open-source R package available on Github: <https://github.com/HectorRDB/Dune> (*RRID* : *SCR\_018218*) and to be released through the Bioconductor Project (<http://www.bioconductor.org>).

## 305 Competing interests

306 The authors declare that they have no competing interests.

## 307 Methods

308 Consider a – possibly high-dimensional – dataset of  $n$  observations,  $\mathbf{X} = \{x_1, \dots, x_n\}$ , where  $x_i \in \mathbb{R}^J$ ,  
309  $i = 1, \dots, n$ . For instance, in scRNA-Seq,  $x_i$  corresponds to the  $J$  gene expression measures (i.e.,  
310 normalized read counts) of cell  $i$ . Represent the results of any (non-fuzzy) clustering method as a  
311 partition,  $\mathbf{P}$ , which splits the set of  $n$  observations into  $k$  disjoint subsets or clusters,  $\{\mathcal{C}_1, \dots, \mathcal{C}_k\}$ , where:  
312 **1)**  $\mathcal{C}_i \cap \mathcal{C}_j = \emptyset, \forall i, j \in \{1, \dots, k\}$ , and **2)**  $\cup_{i \in \{1, \dots, k\}} \mathcal{C}_i = \mathbf{X}$ . Accordingly, a collection of  $R$  clustering  
313 results may be represented as multiple partitions,  $\mathbf{P}_1, \dots, \mathbf{P}_R$ , with partition  $\mathbf{P}_r$  containing  $k_r$  clusters,  
314  $r = 1, \dots, R$ . For each observation  $x_i$ , denote by  $c_{i,r} \in \{\mathcal{C}_1^r, \dots, \mathcal{C}_{k_r}^r\}$  the cluster to which it belongs in  
315 partition  $\mathbf{P}_r$ .

316 The focus of the present manuscript is to develop a general approach to combine clusters within  
317 the different partitions,  $\mathbf{P}_1, \dots, \mathbf{P}_R$ , in order to balance the trade-off between cluster resolution and  
318 replicability. In the remainder of this section, we first present the Rand index, a well-known measure of  
319 concordance between two partitions, and its adjusted version. We also review popular clustering methods  
320 in the scRNA-Seq literature and alternative approaches to merge clusters. Finally, we formalize the two  
321 key notions of cluster resolution and cluster replicability.

## 322 Adjusted Rand index

323 The Rand index [16] measures the concordance between two partitions  $\mathbf{P}_1$  and  $\mathbf{P}_2$ . Denote by  $a =$   
324  $|\{(x_i, x_j) \in \mathbf{X}^2 | (c_{i,1} = c_{j,1}) \& (c_{i,2} = c_{j,2})\}|$  the number of pairs of observations that are in the same  
325 cluster for both partitions  $\mathbf{P}_1$  and  $\mathbf{P}_2$  and by  $b = |\{(x_i, x_j) \in \mathbf{X}^2 | (c_{i,1} \neq c_{j,1}) \& (c_{i,2} \neq c_{j,2})\}|$  the number  
326 of pairs of observations that are in different clusters for both partitions  $\mathbf{P}_1$  and  $\mathbf{P}_2$ . The Rand index is  
327 then the ratio of  $a + b$  over the total number of pairs of observations

$$\text{RI}(\mathbf{P}_1, \mathbf{P}_2) = \frac{a + b}{\binom{n}{2}} \in [0, 1]. \quad (1)$$

328 Thus, intuitively, the Rand index is the proportion of pairs of observations for which the two partitions  
329 are in agreement.

330 However, the Rand index does not account for the fact that a pair of observations might be in the  
331 same (different) cluster(s) in the two partitions purely by chance. The adjusted Rand index (ARI) [17]  
332 adjusts for the level of concordance expected by chance, yielding a value between  $-1$  and  $+1$ . Specifically,  
333 considering  $\mathbf{P}$  a fixed partition and  $\mathbf{R}$  a random permutation of  $\mathbf{P}$ , then  $\mathbb{E}[\text{ARI}(\mathbf{P}, \mathbf{R})] = 0$ , where the  
334 expected value is over all cluster permutations (i.e., permutations of the cluster assignments of the  
335 observations, while keeping the number of clusters and the sizes of the clusters fixed). Negative values  
336 indicate less than the expected level of concordance and positive values indicate more than the expected  
337 level of concordance. The ARI relies on the contingency table of two partitions  $\mathbf{P}_1$  and  $\mathbf{P}_2$ , with the  
338  $(i, j)^{\text{th}}$  entry  $n_{i,j}$  defined as the number of observations both in cluster  $i$  of partition  $\mathbf{P}_1$  and cluster

Table 1: *Adjusted Rand index*. Contingency table for two partitions  $\mathbf{P}_1$  and  $\mathbf{P}_2$ .

	$\mathcal{C}_1^2$	$\mathcal{C}_2^2$	$\dots$	$\mathcal{C}_{k_2}^2$	Sums
$\mathcal{C}_1^1$	$n_{1,1}$	$n_{1,2}$	$\dots$	$n_{1,k_2}$	$a_1$
$\mathcal{C}_2^1$	$n_{2,1}$	$n_{2,2}$	$\dots$	$n_{2,k_2}$	$a_2$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$\mathcal{C}_{k_1}^1$	$n_{k_1,1}$	$n_{k_1,2}$	$\dots$	$n_{k_1,k_2}$	$a_{k_1}$
Sums	$b_1$	$b_2$	$\dots$	$b_{k_2}$	

339  $j$  of partition  $\mathbf{P}_2$  (Table 1). Examples of contingency tables between two partitions can be found in  
 340 Figures 1a, 1b, 2a, and 2d.

341 Given the contingency table notation, the adjusted Rand index is defined as

$$\text{ARI}(\mathbf{P}_1, \mathbf{P}_2) = \frac{\sum_{i,j} \binom{n_{i,j}}{2} - \frac{1}{\binom{n}{2}} \sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}}{\frac{1}{2} (\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2}) - \frac{1}{\binom{n}{2}} \sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}}. \quad (2)$$

342 For  $R$  partitions, the level of concordance can be quantified by the average ARI for all possible pairs  
 343 of partitions

$$\overline{\text{ARI}}(\mathbf{P}_1, \dots, \mathbf{P}_R) = \frac{1}{\binom{R}{2}} \sum_{\{(r,s) \in \{1, \dots, R\} | r < s\}} \text{ARI}(\mathbf{P}_r, \mathbf{P}_s). \quad (3)$$

344 Note that, in the case of  $R = 2$  partitions, this is simply the ARI between the two partitions. If one  
 345 considers the matrix of pairwise ARIs between partitions, such as displayed in Figures 2b and e, then the  
 346 average ARI is defined as the mean of the upper(or lower)-triangular matrix.

### 347 ARI merging with Dune

348 Given  $R$  partitions (possibly the result of different clustering algorithms or different tuning parameter  
 349 values for the same clustering algorithm or both),  $\mathbf{P}_1, \dots, \mathbf{P}_R$ , with  $\mathbf{P}_r$  containing  $k_r$  clusters,  $r =$   
 350  $1, \dots, R$ , Dune seeks to improve the overall agreement among these, as measured by the average ARI,  
 351 through an iterative process of merging clusters within partitions.

352 Specifically, Dune searches over each partition  $\mathbf{P}_r$  and over each of  $\binom{k_r}{2}$  pairs of clusters in  $\mathbf{P}_r$  for the  
 353 pair which produces the largest improvement in ARI when merged, i.e.,

$$(r^*, i^*, j^*) := \arg \max_{\substack{r \in \{1, \dots, R\} \\ i, j \in \{1, \dots, k_r\}}} \sum_{\{s \in \{1, \dots, R\} | s \neq r\}} \text{ARI}(\mathbf{P}_r^{i \cup j}, \mathbf{P}_s) - \text{ARI}(\mathbf{P}_r, \mathbf{P}_s), \quad (4)$$

where  $\mathbf{P}_r^{i \cup j}$  is the partition created by merging clusters  $\mathcal{C}_i^r$  and  $\mathcal{C}_j^r$  in partition  $\mathbf{P}_r$

$$\begin{aligned} \mathbf{P}_r^{i \cup j} &:= \mathbf{P}_r \setminus \{\mathcal{C}_i^r, \mathcal{C}_j^r\} \cup \{\mathcal{C}_i^r \cup \mathcal{C}_j^r\} \\ &= \{\mathcal{C}_1^r, \dots, \mathcal{C}_{i-1}^r, \mathcal{C}_{i+1}^r, \dots, \mathcal{C}_{j-1}^r, \mathcal{C}_{j+1}^r, \dots, \mathcal{C}_{k_r}^r, \mathcal{C}_i^r \cup \mathcal{C}_j^r\}. \end{aligned}$$

Dune amounts to a greedy algorithm for maximizing the average ARI,  $\overline{\text{ARI}}$ . At each step, we find the  
 pair of clusters that, when merged, lead to the greatest improvement in ARI. Once we have identified  
 this pair of clusters, we update the collection of partitions:  $\{\mathbf{P}_1, \dots, \mathbf{P}_R\} \rightarrow \{\mathbf{P}_1, \dots, \mathbf{P}_{r^*}^{i^* \cup j^*}, \dots, \mathbf{P}_R\}$ .  
 We continue iterating until no beneficial merge can be identified, that is, we stop updating when

$$\max_{r, i, j} \sum_{s \neq r} \text{ARI}(\mathbf{P}_r^{i \cup j}, \mathbf{P}_s) - \text{ARI}(\mathbf{P}_r, \mathbf{P}_s) < 0.$$

354 This greedy approach means that each update step is constrained to merging a single pair of clusters  
 355 from a single partition. As such, we never merge three clusters together in one iteration or two pairs of  
 356 clusters in the same or in separate partitions. This ensures that, in our applications, we do not converge  
 357 to the naive optimal solution of merging all clusters, which does represent a full agreement between the  
 358 partitions but is of no practical interest.

359 While Dune provides a natural stopping point for merging, it is also possible to stop earlier in the  
 360 merging process, by tuning the merging parameter  $m_{\text{Dune}}$ , which is defined as the fraction of ARI improve-  
 361 ment over the total ARI improvement. For example,  $m_{\text{Dune}} = .5$  means that Dune returns the merged  
 362 partitions that have a mean ARI halfway between the mean ARI of the original partitions and the mean  
 363 ARI of the final ones.

## 364 Computational implementation and run time

365 The Dune algorithm has been implemented in an open-source R package available on Github: <https://github.com/HectorRDB/Dune>. It is implemented in a fully-parallel and efficient manner. Run time for  
366 a large dataset of  $\sim 100,000$  cells with 3 partitions is under 15 minutes with 10 CPUs. The package also  
367 contains plotting functions used to create many panels of the paper, as well as options to create GIFs  
368 and track the evolution of mean ARI or confusion matrices across the merging steps.

## 370 Clustering algorithms for scRNA-Seq data

371 Any combination of clustering algorithms and associated tuning parameters, applied to an appropriate  
372 dataset, can produce a set of partitions that can be used as input to Dune. However, as our work was  
373 motivated by the classification of cells based on transcriptomic signatures, we will focus on this particular  
374 setting to benchmark Dune.

375 In the descriptions below, we use the notation from the original papers to describe the tuning param-  
376 eters of each method; the same notation may therefore correspond to different parameters depending on  
377 the algorithm.

378 SC3 [2] is a consensus clustering method that involves performing  $k$ -means clustering on different  
379 dimensionality reductions of the input dataset. A hierarchical clustering method is then applied to the  
380 resulting consensus matrix. The main parameter is the number of clusters  $k$ , which is used both in  $k$ -  
381 means and to cut the hierarchical clustering tree. The method provides an estimate of the optimal value  
382 of this parameter,  $k_0$ , based on the number of eigenvalues of the centered and scaled distance matrix that  
383 are significantly different from 0 (see Kiselev et al. [2] for more details). For large datasets, there exists  
384 a hybrid version of the algorithm, where the full SC3 clustering method is run on only a fraction of the  
385 cells to identify the clusters and the rest of the cells are assigned to the clusters using a support vector  
386 machine (SVM) algorithm.

387 Seurat's clustering algorithm (*SEURAT*, *RRID* : *SCR\_007322*) has evolved over the different versions  
388 of the software; here, we focus on version 3 [3] (we specifically use version 3.1.1). The algorithm first  
389 reduces the dimension of the data by selecting the first  $p$  principal components (PCs) and then computes  
390 a  $k$ -nearest neighbor ( $k$ -NN) graph. After refining the graph, it groups cells together using, as default,  
391 the Louvain algorithm [33]. The two main tuning parameters are the number of neighbors  $k$  used to build  
392 the  $k$ -NN graph and the resolution parameter for the Louvain algorithm.

393 Monocle's clustering algorithm has also changed and we focus on version 3 [4] (implemented in the  
394 Monocle3 package, although we keep the name Monocle for simplicity; we specifically use version 0.1.3).  
395 Monocle's clustering algorithm is similar to the one implemented in Seurat, with a few differences. After  
396 initial dimensionality reduction based on principal component analysis (PCA), Monocle performs another  
397 dimensionality reduction step using uniform manifold approximation and projection (UMAP) [34, 35] and  
398 relies on that representation to build the  $k$ -NN graph. It then clusters cells using, by default, the Leiden  
399 algorithm [36].

400 Resampling-based sequential ensemble clustering (RSEC [7]) is a consensus method over user-supplied  
401 clustering algorithms and their associated tuning parameters. In order to improve the stability and tight-  
402 ness of the clusters, it also provides the option to perform clustering on subsamples of the observations, as  
403 well as sequential clustering. However, in this paper, we mainly use RSEC for its final step of hierarchical  
404 merging, see section [Existing methods to merge clusters](#).

## 405 Method parameters

406 For each method, we only tune the main parameter. For Seurat, however, there are two main tuning  
407 parameters. The  $k$  parameter controls the number of neighbors used to build the  $k$ -NN graph, while the  
408 resolution parameter defines the neighborhood in the Louvain clustering algorithm. In practice, the  $k$   
409 parameter has much less impact than the resolution parameter (see Figure S1). Moreover, depending on  
410 the value of the resolution, increasing  $k$  either increases or decreases the final number of clusters. As a  
411 result, we only consider changing the resolution parameter.

412 For ease and generality of notation, we will denote each method's main tuning parameter by  $\theta$  and  
413 define  $\theta$  such that increasing  $\theta$  increases the number of clusters. Thus, for the methods described above,  
414  $\theta_{SC3} = k$ ,  $\theta_{Seurat} = \text{Resolution}$ , and  $\theta_{Monocle} = -k$ . Each combination  $\Theta = \{\theta_{SC3}, \theta_{Seurat}, \theta_{Monocle}\}$  of  
415 the three parameters defines a set of partitions that serves as input for Dune.

## 416 Existing methods to merge clusters

417 Once a set of clusters has been identified, one can build a hierarchical tree for these clusters and then  
418 merge clusters that are similar. This involves specifying a measure of distance or similarity between  
419 individual observations (i.e., cells) as well as between clusters. It should be noted that the distance used  
420 to build the tree of clusters need not be the same as the distance used to merge clusters.

421 For scRNA-Seq datasets, commonly used between-cell distance measures include the Euclidean distance  
422 and one minus the Spearman correlation coefficient. Between-cluster distances include classical  
423 linkage measures used in hierarchical clustering, e.g., maximum/minimum/average of all pairwise dis-  
424 tances between observations in two clusters or distance between the cluster averages or medoids. For  
425 scRNA-Seq, another sensible between-cluster distance measure is the proportion of differentially expressed  
426 (DE) genes between clusters [7, 8]. A detailed discussion of such measures is out of the scope of this  
427 manuscript[37].

428 Here, we consider two possible ways of merging. In both cases, we compute the cluster medoids  
429 (median of the cluster) based on the log-transformed count matrix (adding 1 to avoid taking the log  
430 of zero). We then build a hierarchical tree of clusters using the Euclidean distance between the cluster  
431 medoids. The first merging approach directly uses this tree to decide how to merge clusters. Specifically,  
432 clusters are merged bottom-up, starting with the two clusters that are closest in the tree and then  
433 iteratively until all clusters are merged. The parameter  $m_{Dist} = n_{merges}$ , the number of merges (between  
434 0 and the initial number of clusters minus one), controls the amount of merging. The second approach  
435 follows the method implemented in RSEC. It computes the percentage of DE genes between clusters, using  
436 the limma package [20] (*LIMMA*, *RRID* : *SCR.010943*), where a gene is declared DE if its nominal FDR  
437 adjusted  $p$ -value is below 0.05 [21]. The main tunable parameter is  $m_{DE} = \alpha \in [0, 1]$ , the threshold for the  
438 percentage of DE genes below which we merge. We name these two methods Dist and DE, respectively.

## 439 Cluster replicability using MetaNeighbor

440 We quantify the replicability of clusters across datasets by applying a modified version of unsuper-  
441 vised MetaNeighbor [22] (*MetaNeighbor*, *RRID* : *SCR.016727*). MetaNeighbor requires as input a set  
442 of unnormalized datasets, a set of cluster labels, and a set of highly variable genes. It uses a cross-  
443 dataset validation scheme to quantify how well clusters match across datasets. Given any two datasets,  
444 MetaNeighbor builds a cell-cell similarity network based on the Spearman correlation over the set of  
445 highly variable genes. One of the datasets is treated as a test dataset, where all cluster labels are hidden,  
446 the other dataset is treated as a training dataset, whose labels are propagated to the test dataset through  
447 the cell-cell similarity network. Each pair of clusters (one in the training dataset, the other in the test  
448 dataset) receives a score based on how well the training cluster predicts the labels from the test cluster.  
449 This score is the area under the receiver operator characteristic curve (one-vs-one AUROC). We define  
450 the best matching cluster as the test cluster which dominates all other test clusters (one-vs-one AUROC  
451  $> 0.5$ ). Finally, we reduce the test set to the two best matching clusters, recompute an AUROC, which  
452 we call one-vs-best AUROC, and record this as the pair’s final score. Then the role of the test and  
453 training datasets are reversed. A cluster is considered replicable if there is a cluster in the other dataset  
454 such that the clusters are reciprocal best hits with a high AUROC score (one-vs-best AUROC  $> 0.6$  both  
455 ways). See Crow et al. [22] for details.

456 The **replicability score of a cluster** is defined as the fraction of cells contained in replicable clusters.  
457 More specifically, for a comparison of two datasets, we enumerate replicable clusters in each dataset, then  
458 deduce the number of cells that are in replicable clusters, sum this number across datasets, and divide  
459 by the total number of cells.

460 We used MetaNeighbor’s **variableGenes** procedure to select genes that were detected as highly  
461 variable across all datasets. For performance reasons, the **variableGenes** procedure was applied to a  
462 random subset of 50,000 cells for datasets exceeding that size. However, the full datasets were use for  
463 the rest of the analysis. In the end, we obtained a set of 541 highly variable genes for the Allen brain  
464 datasets and 2, 147 genes for the pancreas datasets.

## 465 Case studies

### 466 AIBS Smart mouse brain datasets

467 We used the two AIBS Smart datasets produced as part of the Brain Initiative Cell Census Net-  
468 work (BICCN: *RRID* : *SCR.015820*) and described in Yao et al. [11], one is single-cell (*Zeng sn*

469 *SSv4* (<https://assets.nemoarchive.org/dat-k7p82j4>) and the other is single-nuclei (*Zeng sc SSv4*  
470 <https://assets.nemoarchive.org/dat-55mowp9>). We use the subclass labels as gold-standard cluster  
471 labels for these datasets. Those datasets can be downloaded from the Neuroscience Multi-omics  
472 Archive (*RRID* : *SCR\_002001*; [nemoarchive.org](https://nemoarchive.org)). More details on the parent data set (<https://assets.nemoarchive.org/dat-ch1nqb7>)  
473 and data access can be found in Yao et al. [11].

#### 474 Human pancreas datasets

475 We focus on two datasets from [18] (8,568 cells) and [19] (3,514 cells) which we name **Baron** and  
476 **Segerstople**, respectively. Both datasets were downloaded from the <https://hemberg-lab.github.io/scRNA.seq.datasets/>  
477 on October 1<sup>st</sup>, 2018. We use the clusters from the original publications as gold-standard cluster labels.

#### 478 Data analysis

479 Except when otherwise specified, all methods and algorithms were run with default parameters or, if no  
480 available default, with the parameters recommended in the vignette or tutorial.

481 **Pre-processing:** Count matrices were filtered to remove lowly-expressed genes with fewer than  $i$  reads  
482 in  $j$  cells. See Table S1 for values of  $i$  and  $j$  for each dataset.

483 As indicated below, we follow different normalization strategies before running Seurat and Monocle  
484 in order to obtain more diverse clustering results. This is appropriate, as the goal of the manuscript is  
485 not to compare different clustering methods, but rather different merging methods for given clustering  
486 results. The merging methods that Dune is compared to rely on only one clustering input; we therefore  
487 seek to benchmark merging methods using a variety of clustering inputs.

488 **Seurat:** Following the tutorial, we run `FindVariableFeatures` and `ScaleData` to normalize the data.  
489 Counts are log-transformed (adding 1 to avoid taking the log of zero) and normalized by sequencing depth.  
490 For the two pancreas datasets, batches are also normalized using the `scaleData` function. Following  
491 principal component analysis, `FindNeighbors` and `FindClusters` are run for a number of neighbors  $k$  in  
492  $\{30, 50, 70\}$  and resolution  $\theta$  from 0.3 to 2.5 in increments of 0.1.

493 **SC3:** The algorithm is run on a dataset normalized as above with the Seurat pipeline. The optimal  
494 value of  $k$ ,  $k_0$ , is computed using the `sc3.estimate_k` function. The parameter  $\theta$  is transformed to be  
495  $\theta_{SC3} = k - k_0$ . SC3 is then run for values of  $\theta$  ranging from  $-15$  to  $+15$ .

496 **Monocle:** `zinbwave` [7] is first used for normalization and dimensionality reduction on the filtered count  
497 data. For the two pancreas datasets, batches are included as model covariates. We select  $K$ , the number  
498 of reduced dimensions, based on a visual representation for each dataset, see Table S1. This first step  
499 of dimensionality reduction is followed by another using UMAP [35] with two dimensions. The resulting  
500 two-dimensional representation is then used to build the  $k$ -NN graph, with  $k$  ranging from 10 to 150 in  
501 increments of 10.

502 **Dune:** For a given set of values for  $\Theta = \{\theta_{SC3}, \theta_{Seurat}, \theta_{Monocle}\}$ , we get three sets of cluster labels that  
503 we can use as input to Dune.

504 **Building the hierarchical tree:** The output of each clustering method is used as input to RSEC's  
505 `makeDendrogram` function. Then, we either cut the tree using R's `cutree` function or RSEC's `mergeClusters`  
506 function.

507 **Producing "bad" clusters:** For each value of the tuning parameters  $\Theta$ , on the pancreas datasets, we  
508 add fully random inputs to Dune. That is, we create "bad" clusterings by randomly assigning each cell a  
509 number (or cluster label) between 1 and  $(k_{SC3} + k_{Monocle} + k_{Seurat})/3$ , where  $k$  denotes the number of  
510 clusters for a particular clustering algorithm. Since cells are assigned randomly, the size of the clusters  
511 will vary, but all clusters have the same expected size. To account for the stochastic nature of this  
512 procedure, we repeat this 10 times.

513 **Downsampling:** Downsampling the number of cells at the beginning of the analysis pipeline would  
514 affect both the quality of the clustering results and the quality of the merging with Dune. As such, to test  
515 only the stability of Dune to the number of cells, we downsample the cells just before running Dune, that  
516 is, the clustering algorithms are run on the full dataset but only a subset of the dataset is used to decide  
517 which clusters to merge and in which order. Afterwards, cells that are not in the subsample are assigned  
518 to the merged clusters based on their original cluster labels. That is, if Cluster 1 and 2 are merged, all  
519 cells that were originally in Cluster 1 and 2, even those not selected in the downsampling and used as  
520 input to Dune, are assigned to the merged cluster.  
521 Most of the code was run using xsede [38].

## 522 References

- 523 [1] Valentine Svensson and Eduardo da Veiga Beltrame. A curated database reveals trends in single cell  
524 transcriptomics. *bioRxiv*, page 742304, 2019. doi: 10.1101/742304.
- 525 [2] Vladimir Yu Kiselev, Kristina Kirschner, Michael T Schaub, Tallulah Andrews, Andrew Yiu, Tamir  
526 Chandra, Kedar N Natarajan, Wolf Reik, Mauricio Barahona, Anthony R Green, and Martin Hem-  
527 berg. SC3: consensus clustering of single-cell RNA-seq data. *Nature Methods*, 14(5):483–486, may  
528 2017. ISSN 1548-7091. doi: 10.1038/nmeth.4236. URL [http://www.nature.com/articles/nmeth.](http://www.nature.com/articles/nmeth.4236)  
529 [4236](http://www.nature.com/articles/nmeth.4236).
- 530 [3] Tim Stuart, Andrew Butler, Paul Hoffman, Christoph Hafemeister, Efthymia Papalexi, William M  
531 Mauck, Yuhan Hao, Marlon Stoekius, Peter Smibert, and Rahul Satija. Comprehensive Integration  
532 of Single-Cell Data. *Cell*, 177(7):1888–1902.e21, jun 2019. ISSN 10974172. doi: 10.1016/j.cell.2019.  
533 05.031. URL <http://www.ncbi.nlm.nih.gov/pubmed/31178118><http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC6687398>.
- 534 [4] Junyue Cao, Malte Spielmann, Xiaojie Qiu, Xingfan Huang, Daniel M. Ibrahim, Andrew J. Hill, Fan  
535 Zhang, Stefan Mundlos, Lena Christiansen, Frank J. Steemers, Cole Trapnell, and Jay Shendure. The  
536 single-cell transcriptional landscape of mammalian organogenesis. *Nature*, 566(7745):496–502, feb  
537 2019. ISSN 0028-0836. doi: 10.1038/s41586-019-0969-x. URL [http://www.nature.com/articles/](http://www.nature.com/articles/s41586-019-0969-x)  
538 [s41586-019-0969-x](http://www.nature.com/articles/s41586-019-0969-x).
- 539 [5] Vladimir Yu Kiselev, Tallulah S. Andrews, and Martin Hemberg. Challenges in unsupervised clus-  
540 tering of single-cell RNA-seq data, may 2019. ISSN 14710064. URL [http://www.nature.com/](http://www.nature.com/articles/s41576-018-0088-9)  
541 [articles/s41576-018-0088-9](http://www.nature.com/articles/s41576-018-0088-9).
- 542 [6] Angelo Duò, Mark D. Robinson, and Charlotte Soneson. A systematic  
543 performance evaluation of clustering methods for single-cell RNA-seq data.  
544 *F1000Research*, 7:377–382, 2018. ISSN 1759796X. doi: 10.5256/f1000research.17093.  
545 r36544. URL [https://f1000researchdata.s3.amazonaws.com/manuscripts/17687/](https://f1000researchdata.s3.amazonaws.com/manuscripts/17687/49bacf17-03b0-4f80-bde4-e892c8c3e22f_{_}15666_{_}-_{_}charlotte_{_}soneson_{_}v2.pdf?doi=10.12688/f1000research.15666.2{&}numberOfBrowsableCollections=17{&}numberOfBrowsableInstitutionalCollections=4{&}numberOfBrows)  
546 [49bacf17-03b0-4f80-bde4-e892c8c3e22f\\_{\\_}15666\\_{\\_}-\\_{\\_}charlotte\\_{\\_}soneson\\_{\\_}v2.](https://f1000researchdata.s3.amazonaws.com/manuscripts/17687/49bacf17-03b0-4f80-bde4-e892c8c3e22f_{_}15666_{_}-_{_}charlotte_{_}soneson_{_}v2.pdf?doi=10.12688/f1000research.15666.2{&}numberOfBrowsableCollections=17{&}numberOfBrowsableInstitutionalCollections=4{&}numberOfBrows)  
547 [pdf?doi=10.12688/f1000research.15666.2{&}numberOfBrowsableCollections=](https://f1000researchdata.s3.amazonaws.com/manuscripts/17687/49bacf17-03b0-4f80-bde4-e892c8c3e22f_{_}15666_{_}-_{_}charlotte_{_}soneson_{_}v2.pdf?doi=10.12688/f1000research.15666.2{&}numberOfBrowsableCollections=17{&}numberOfBrowsableInstitutionalCollections=4{&}numberOfBrows)  
548 [17{&}numberOfBrowsableInstitutionalCollections=4{&}numberOfBrows.](https://f1000researchdata.s3.amazonaws.com/manuscripts/17687/49bacf17-03b0-4f80-bde4-e892c8c3e22f_{_}15666_{_}-_{_}charlotte_{_}soneson_{_}v2.pdf?doi=10.12688/f1000research.15666.2{&}numberOfBrowsableCollections=17{&}numberOfBrowsableInstitutionalCollections=4{&}numberOfBrows)
- 549 [7] Davide Risso, Liam Purvis, Russell B. Fletcher, Diya Das, John Ngai, Sandrine Dudoit, and Elizabeth  
550 Purdom. clusterExperiment and RSEC: A Bioconductor package and framework for clustering of  
551 single-cell and other large gene expression datasets. *PLoS Computational Biology*, 14(9):e1006378,  
552 sep 2018. ISSN 15537358. doi: 10.1371/journal.pcbi.1006378. URL [http://dx.plos.org/10.1371/](http://dx.plos.org/10.1371/journal.pcbi.1006378)  
553 [journal.pcbi.1006378](http://dx.plos.org/10.1371/journal.pcbi.1006378).
- 554 [8] Bosiljka Tasic, Zizhen Yao, Lucas T. Graybuck, Kimberly A. Smith, Thuc Nghi Nguyen, Darren  
555 Bertagnolli, Jeff Goldy, Emma Garren, Michael N. Economo, Sarada Viswanathan, Osnat Penn,  
556 Trygve Bakken, Vilas Menon, Jeremy Miller, Olivia Fong, Karla E. Hirokawa, Kanan Lathia,  
557 Christine Rimorin, Michael Tieu, Rachael Larsen, Tamara Casper, Eliza Barkan, Matthew Kroll,  
558 Sheana Parry, Nadiya V. Shapovalova, Daniel Hirschstein, Julie Pendergraft, Heather A. Sullivan,  
559 Tae Kyung Kim, Aaron Szafer, Nick Dee, Peter Groblewski, Ian Wickersham, Ali Cetin, Julie A.  
560 Harris, Boaz P. Levi, Susan M. Sunkin, Linda Madisen, Tanya L. Daigle, Loren Looger, Amy  
561 Bernard, John Phillips, Ed Lein, Michael Hawrylycz, Karel Svoboda, Allan R. Jones, Christof

- 563 Koch, and Hongkui Zeng. Shared and distinct transcriptomic cell types across neocortical areas. *Nature*, 563(7729):72–78, nov 2018. ISSN 14764687. doi: 10.1038/s41586-018-0654-5. URL <http://www.nature.com/articles/s41586-018-0654-5>.
- 564
- 565
- 566 [9] Saskia Freytag, Luyi Tian, Ingrid Lönnstedt, Milica Ng, and Melanie Bahlo. Comparison of clustering tools in R for medium-sized 10x genomics single-cell RNA-sequencing data. *F1000Research*, 7, 2018. ISSN 1759796X. doi: 10.12688/f1000research.15809.1.
- 567
- 568
- 569 [10] Luke Zappia and Alicia Oshlack. Clustering trees: a visualization for evaluating clusterings at multiple resolutions. *GigaScience*, 7(7):1–9, 2018. ISSN 2047217X. doi: 10.1093/gigascience/giy083. URL <http://orcid.org/0000-0001-9788-5690Address:>.
- 570
- 571
- 572 [11] Zizhen Yao, Hanqing Liu, Fangming Xie, Stephan Fischer, Ricky S Adkins, Andrew I Aldrige, Seth A Ament, Ann Bartlett, M. Margarita Behrens, Koen Van den Berge, Darren Bertagnolli, Tommaso Biancalani, A. Sina Boeshaghi, Hector Corrada Bravo, Tamara Casper, Carlo Colantuoni, Jonathan Crabtree, Heather Creasy, Kirsten Crichton, Megan Crow, Nick Dee, Elizabeth L Dougherty, Wayne I Doyle, Sandrine Dudoit, Rongxin Fang, Victor Felix, Olivia Fong, Michelle Giglio, Jeff Goldy, Michael Hawrylycz, Hector Roux de Bezieux, Brian R. Herb, Ronna Hertzano, Xiaomeng Hou, Qiwen Hu, Z. Josh Huang, Jayaram Kancherla, Matthew Kroll, Kanan Lathia, Yang Eric Li, Jacinta D. Lucero, Chongyuan Luo, Anup Mahurkar, Delissa McMillen, Naeem M. Nadaf, Joseph R. Nery, Thuc Nghi Nguyen, Sheng-Yong Niu, Vasilis Ntranos, Joshua Orvis, Julia K. Osteen, Thanh Pham, Antonio Pinto-Duarte, Olivier Poirion, Sebastian Preissl, Elizabeth Purdom, Christine Rimorin, Davide Risso, Angeline C. Rivkin, Kimberly Smith, Kelly Street, Josef Sulc, Valentine Svensson, Michael Tieu, Amy Torkelson, Herman Tung, Eeshit Dhaval Vaishnav, Charles R. Vanderburg, Cindy van Velthoven, Xinxin Wang, Owen White, Jesse Gillis, Peter V. Kharchenko, John Ngai, Lior Pachter, Aviv Regev, Bosiljka Tasic, Joshua D Welch, Joseph R. Ecker, Evan Macosko, Bing Ren, BRAIN Initiative Cell Census Network (BICCN), Hongkui Zeng, and Eran A Mukamel. An integrated transcriptomic and epigenomic atlas of mouse primary motor cortex cell types. *bioRxiv*, page 2020.02.29.970558, mar 2020. doi: 10.1101/2020.02.29.970558.
- 573
- 574
- 575
- 576
- 577
- 578
- 579
- 580
- 581
- 582
- 583
- 584
- 585
- 586
- 587
- 588
- 589 [12] L.J.P. van der Maaten and G.E. Hinton. Visualizing high-dimensional data using t-sne. *Journal of Machine Learning Research*, 9:2579–2605, 2008.
- 590
- 591 [13] L.J.P. van der Maaten. Accelerating t-sne using tree-based algorithms. *Journal of Machine Learning Research*, 15:3221–3245, 2014.
- 592
- 593 [14] Jesse H. Krijthe. *Rtsne: T-Distributed Stochastic Neighbor Embedding using Barnes-Hut Implementation*, 2015. URL <https://github.com/jkrijthe/Rtsne>. R package version 0.15.
- 594
- 595 [15] Paul Jaccard. Distribution de la flore alpine dans le bassin des dranses et dans quelques rgions voisines. *Bulletin de la Societe Vaudoise des Sciences Naturelles*, 37:241–72, 01 1901. doi: 10.5169/seals-266440.
- 596
- 597
- 598 [16] William M. Rand. Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical Association*, 66(336):846–850, 1971. doi: 10.1080/01621459.1971.10482356. URL <https://www.tandfonline.com/doi/abs/10.1080/01621459.1971.10482356>.
- 599
- 600
- 601 [17] Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of Classification*, 2(1):193–218, Dec 1985. ISSN 1432-1343. doi: 10.1007/BF01908075. URL <https://doi.org/10.1007/BF01908075>.
- 602
- 603
- 604 [18] Maayan Baron, Adrian Veres, Samuel L. Wolock, Aubrey L. Faust, Renaud Gaujoux, Amedeo Vetere, Jennifer Hyoje Ryu, Bridget K. Wagner, Shai S. Shen-Orr, Allon M. Klein, Douglas A. Melton, and Itai Yanai. A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas Reveals Inter- and Intra-cell Population Structure. *Cell Systems*, 3(4):346–360.e4, oct 2016. ISSN 24054720. doi: 10.1016/j.cels.2016.08.011. URL <https://www.sciencedirect.com/science/article/pii/S2405471216302666?via%3Dihub>.
- 605
- 606
- 607
- 608
- 609
- 610 [19] Åsa Segerstolpe, Athanasia Palasantza, Pernilla Eliasson, Eva Marie Andersson, Anne Christine Andréasson, Xiaoyan Sun, Simone Picelli, Alan Sabirsh, Maryam Clausen, Magnus K. Bjursell, David M. Smith, Maria Kasper, Carina Ämmälä, and Rickard Sandberg. Single-Cell Transcriptome Profiling of Human Pancreatic Islets in Health and Type 2 Diabetes. *Cell Metabolism*, 24(4):593–607,
- 611
- 612
- 613



- 614 oct 2016. ISSN 19327420. doi: 10.1016/j.cmet.2016.08.020. URL <https://www.sciencedirect.com/science/article/pii/S1550413116304363?via%3Dihub>.  
615
- 616 [20] Matthew E. Ritchie, Belinda Phipson, Di Wu, Yifang Hu, Charity W. Law, Wei Shi, and  
617 Gordon K. Smyth. limma powers differential expression analyses for RNA-sequencing and  
618 microarray studies. *Nucleic Acids Research*, 43(7):e47–e47, apr 2015. ISSN 1362-4962.  
619 doi: 10.1093/nar/gkv007. URL [http://academic.oup.com/nar/article/43/7/e47/2414268/  
620 limma-powers-differential-expression-analyses-for](http://academic.oup.com/nar/article/43/7/e47/2414268/limma-powers-differential-expression-analyses-for).
- 621 [21] Yosef Benjamini, Yoav ; Hochberg. Controlling the False Discovery Rate - a Practical and Power-  
622 ful Approach to Multiple Testing. Journal of the Royal Statistical Society Series B-Methodological  
623 1995.pdf. *Journal of the Royal Statistical Society Series B (Methodological)*, 57(1):289–300, 1995.  
624 doi: 10.2307/2346101. URL [https://www.jstor.org/stable/2346101http://www.jstor.org/  
625 stable/2346101](https://www.jstor.org/stable/2346101http://www.jstor.org/stable/2346101).
- 626 [22] Megan Crow, Anirban Paul, Sara Ballouz, Z. Josh Huang, and Jesse Gillis. Characterizing the  
627 replicability of cell types defined by single cell RNA-sequencing data using MetaNeighbor. *Nature  
628 Communications*, 9(1):884, dec 2018. ISSN 20411723. doi: 10.1038/s41467-018-03282-0. URL  
629 <http://www.nature.com/articles/s41467-018-03282-0>.
- 630 [23] R. J.G.B. Campello. A fuzzy extension of the Rand index and other related indexes for clustering  
631 and classification assessment. *Pattern Recognition Letters*, 28(7):833–841, may 2007. ISSN 01678655.  
632 doi: 10.1016/j.patrec.2006.11.010.
- 633 [24] Roelof K. Brouwer. Extending the rand, adjusted rand and jaccard indices to fuzzy partitions.  
634 *Journal of Intelligent Information Systems*, 32(3):213–235, jun 2009. ISSN 09259902. doi: 10.1007/  
635 s10844-008-0054-7.
- 636 [25] Allen W. Zhang, Ciara O’Flanagan, Elizabeth A. Chavez, Jamie L.P. Lim, Nicholas Ceglia, Andrew  
637 McPherson, Matt Wiens, Pascale Walters, Tim Chan, Brittany Hewitson, Daniel Lai, Anja Mottok,  
638 Clementine Sarkozy, Lauren Chong, Tomohiro Aoki, Xuehai Wang, Andrew P. Weng, Jessica N.  
639 McAlpine, Samuel Aparicio, Christian Steidl, Kieran R. Campbell, and Sohrab P. Shah. Probabilistic  
640 cell-type assignment of single-cell RNA-seq for tumor microenvironment profiling. *Nature Methods*,  
641 16(10):1007–1015, oct 2019. ISSN 15487105. doi: 10.1038/s41592-019-0529-1.
- 642 [26] Ze Zhang, Danni Luo, Xue Zhong, Jin Huk Choi, Yuanqing Ma, Stacy Wang, Elena Mahrt, Wei  
643 Guo, Eric W. Stawiski, Zora Modrusan, Somasekar Seshagiri, Payal Kapur, Gary C. Hon, James  
644 Brugarolas, and Tao Wang. Scina: Semi-supervised analysis of single cells in silico. *Genes*, 10(7):531,  
645 jul 2019. ISSN 20734425. doi: 10.3390/genes10070531. URL [https://www.mdpi.com/2073-4425/  
646 10/7/531](https://www.mdpi.com/2073-4425/10/7/531).
- 647 [27] Sergii Domanskyi, Anthony Szedlak, Nathaniel T Hawkins, Jiayin Wang, Giovanni Paternostro, and  
648 Carlo Piermarocchi. Polled Digital Cell Sorter (p-DCS): Automatic identification of hematological  
649 cell types from single cell RNA-sequencing clusters. *BMC Bioinformatics*, 20(1):369, dec 2019. ISSN  
650 14712105. doi: 10.1186/s12859-019-2951-x. URL [https://bmcbioinformatics.biomedcentral.  
651 com/articles/10.1186/s12859-019-2951-x](https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-019-2951-x).
- 652 [28] Florian Wagner and Itai Yanai. Moana: A robust and scalable cell type classification framework  
653 for single-cell RNA-Seq data. *bioRxiv*, page 456129, 2018. doi: 10.1101/456129. URL [https:  
654 //www.biorxiv.org/content/10.1101/456129v1](https://www.biorxiv.org/content/10.1101/456129v1).
- 655 [29] Hannah A. Pliner, Jay Shendure, and Cole Trapnell. Supervised classification enables rapid  
656 annotation of cell atlases. *Nature Methods*, 16(10):983–986, oct 2019. ISSN 15487105. doi:  
657 10.1038/s41592-019-0535-3.
- 658 [30] Yingxin Lin, Yao Cao, Hani J Kim, Agus Salim, Terence P. Speed, Dave Lin, Pengyi Yang, and  
659 Jean Yee Hwa Yang. scClassify: hierarchical classification of cells. *bioRxiv*, page 776948, 2019. doi:  
660 10.1101/776948.
- 661 [31] Mark van der Laan and Katherine Pollard. Hybrid clustering of gene expression data with visual-  
662 ization and the bootstrap. *Mark J. van der Laan*, 117, 01 2001.
- 663 [32] Frank Herbert. *Dune*. Chilton Books, Philadelphia, NY, 1965.

- 664 [33] Vincent D Blondel, Jean Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast  
665 unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and*  
666 *Experiment*, 2008(10):P10008, oct 2008. ISSN 17425468. doi: 10.1088/1742-5468/2008/  
667 10/P10008. URL [http://stacks.iop.org/1742-5468/2008/i=10/a=P10008?key=crossref.](http://stacks.iop.org/1742-5468/2008/i=10/a=P10008?key=crossref.46968f6ec61eb8f907a760be1c5ace52)  
668 [46968f6ec61eb8f907a760be1c5ace52](http://stacks.iop.org/1742-5468/2008/i=10/a=P10008?key=crossref.46968f6ec61eb8f907a760be1c5ace52).
- 669 [34] Etienne Becht, Leland McInnes, John Healy, Charles-Antoine Dutertre, Immanuel W H Kwok,  
670 Lai Guan Ng, Florent Ginhoux, and Evan W Newell. Dimensionality reduction for visualizing  
671 single-cell data using UMAP. *Nature Biotechnology*, 37(1):38–44, jan 2019. ISSN 1087-0156. doi:  
672 10.1038/nbt.4314. URL <http://www.nature.com/articles/nbt.4314>.
- 673 [35] Leland McInnes, John Healy, and James Melville. UMAP: Uniform Manifold Approximation and  
674 Projection for Dimension Reduction. *arxiv*, feb 2018. URL <http://arxiv.org/abs/1802.03426>.
- 675 [36] V. A. Traag, L. Waltman, and N. J. van Eck. From Louvain to Leiden: guaranteeing well-  
676 connected communities. *Scientific Reports*, 9(1):5233, dec 2019. ISSN 2045-2322. doi: 10.1038/  
677 s41598-019-41695-z. URL <http://www.nature.com/articles/s41598-019-41695-z>.
- 678 [37] Taiyun Kim, Irene Rui Chen, Yingxin Lin, Andy Yi-Yang Wang, Jean Yee Hwa Yang, and  
679 Pengyi Yang. Impact of similarity metrics on single-cell RNA-seq data clustering. *Briefings in*  
680 *Bioinformatics*, 20(6):2316–2326, nov 2019. ISSN 1467-5463. doi: 10.1093/bib/bby076. URL  
681 <https://academic.oup.com/bib/article/20/6/2316/5077112>.
- 682 [38] J. Towns, T. Cockerill, M. Dahan, I. Foster, K. Gaither, A. Grimshaw, V. Hazlewood, S. Lathrop,  
683 D. Lifka, G. D. Peterson, R. Roskies, J. R. Scott, and N. Wilkins-Diehr. Xsede: Accelerating scientific  
684 discovery. *Computing in Science & Engineering*, 16(5):62–74, Sept.-Oct. 2014. ISSN 1521-9615. doi:  
685 10.1109/MCSE.2014.80. URL [doi.ieeecomputersociety.org/10.1109/MCSE.2014.80](https://doi.ieeecomputersociety.org/10.1109/MCSE.2014.80).

## 686 Supplementary Material

### 687 Supplementary methods

Table S1: *Parameters for processing the datasets.* Each dataset is filtered such that we keep all genes with a least  $i$  reads in  $j$  samples. Then, zinbwave is run with  $K$  dimensions.

Dataset	$i$	$j$	$K$
AIBS scRNA-Smart	50	50	30
AIBS snRNA-Smart	50	50	14
Baron	5	5	10
Segerstople	5	5	20

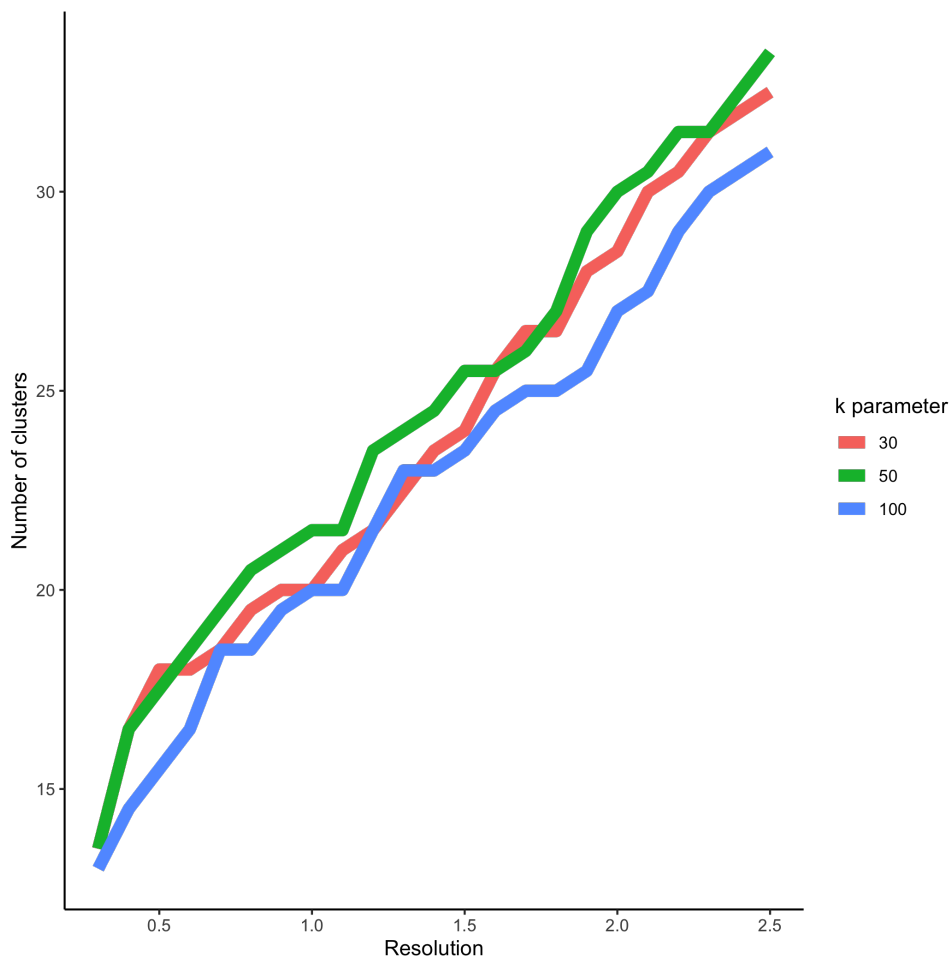


Figure S1: *Impact of Seurat's two main tuning parameters on the number of clusters.* The Seurat algorithm is run on the two AIBS snRNA-Smart datasets, for a grid of tuning parameter values. The average number of clusters found in both datasets is then computed. For increasing values of the resolution parameter and fixed values of the  $k$  parameter, the number of clusters is always increasing. On the other hand, for increasing values of the  $k$  parameter and fixed values of the resolution parameter, the number of clusters can either increase or decrease. This can be seen in the fact that the curves are all increasing but intersect multiple times.

### 688 Supplementary results

Table S2: Ranking of merging methods over all 36 comparisons for improving ARI with gold standard. See Figure 3

	1	2	3
DE	2	21	13
Dist	5	10	21
Dune	29	5	2

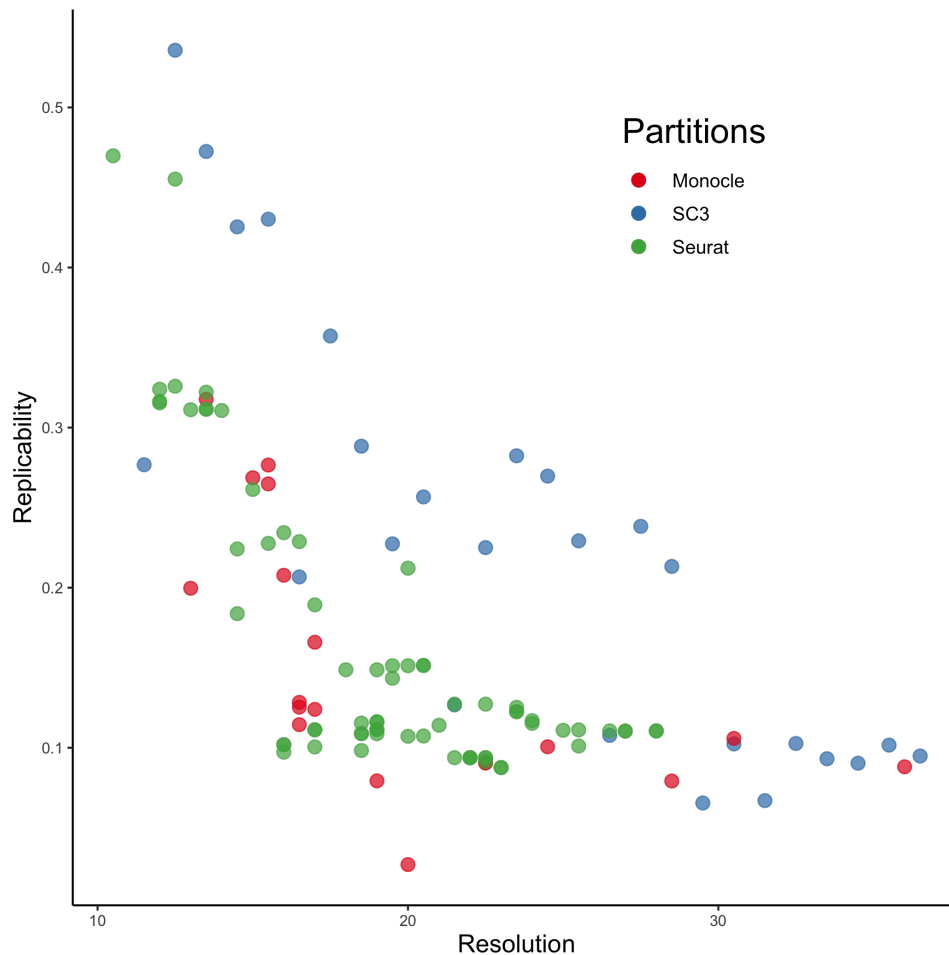


Figure S2: Resolution-replicability trade-off on the Pancreas datasets. Seurat, SC3, and Monocle are run on the two Pancreas datasets, as described in Methods, for a wide range of tuning parameter values. Then, the MetaNeighbor method is used to compute replicability scores for the resulting clusters between these two datasets. An apparent trade-off between replicability and resolution is visible.

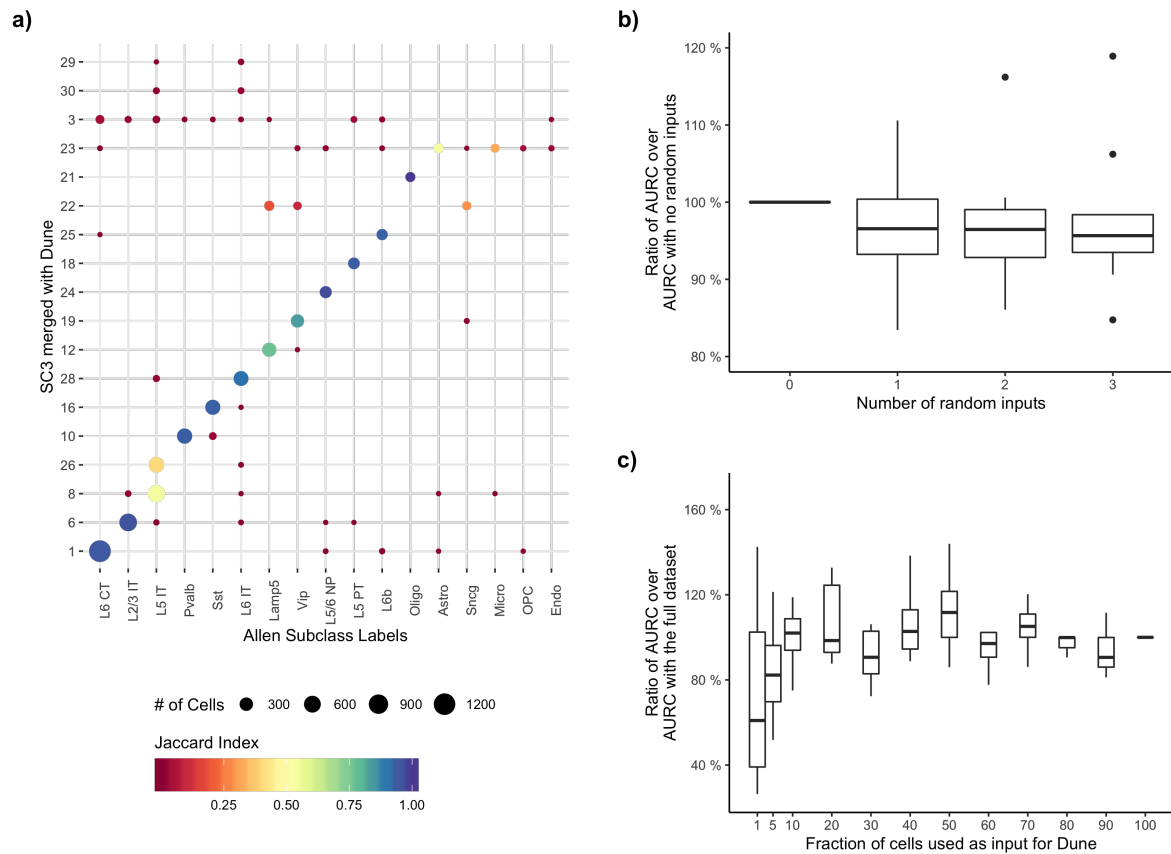


Figure S3: *Dune* robustness analysis. Panel **a**. Fully merging SC3 with *Dune* produces meaningful high-level biological clusters, as can be seen by the overlap between the clustering and the Allen subclass labels. Panel **b**. Adding an increasing number of random clustering inputs to *Dune* impacts only slightly the resolution-replicability area under the curve when merging the other correct clusters. Panel **c**. Likewise, *Dune* is stable to decreasing the number of input cells, as low as 10% of the original sample size.