



A joint deep learning model enables simultaneous batch effect correction, denoising and clustering in single-cell transcriptomics

Justin Lakkis, David Wang, Yuanchao Zhang, et al.

Genome Res. published online May 25, 2021

Access the most recent version at doi:[10.1101/gr.271874.120](https://doi.org/10.1101/gr.271874.120)

P<P	Published online May 25, 2021 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **A joint deep learning model enables simultaneous batch effect correction,** 2 **denoising and clustering in single-cell transcriptomics**

3 Justin Lakkis^{1,*}, David Wang², Yuanchao Zhang¹, Gang Hu³, Kui Wang⁴, Huize Pan⁵, Lyle Ungar⁶, Muredach
4 P. Reilly⁵, Xiangjie Li^{7,*}, Mingyao Li^{1,*}

5 1. Department of Biostatistics, Epidemiology and Informatics, Perelman School of Medicine, University
6 of Pennsylvania, Philadelphia, PA 19104, USA.

7 2. Graduate Group in Genomics and Computational Biology, Perelman School of Medicine, University of
8 Pennsylvania, Philadelphia, PA 19104, USA.

9 3. School of Statistics and Data Science, Key Laboratory for Medical Data Analysis and Statistical
10 Research of Tianjin, Nankai University, Tianjin 300071, China.

11 4. Department of Information Theory and Data Science, School of Mathematical Sciences and LPMC,
12 Nankai University, Tianjin 300071, China.

13 5. Division of Cardiology, Department of Medicine, Columbia University Irving Medical Center, New
14 York, NY 10032, USA.

15 6. Department of Computer and Information Science, School of Engineering and Applied Sciences,
16 University of Pennsylvania, Philadelphia, PA 19104, USA.

17 7. State Key Laboratory of Cardiovascular Disease, Fuwai Hospital, National Center for Cardiovascular
18 Diseases, Chinese Academy of Medical Sciences and Peking Union Medical College, Beijing 100037,
19 China.

20 ***Correspondence:**

21 Justin Lakkis (jlakkis@pennmedicine.upenn.edu)

22 Xiangjie Li (xiangjie631@outlook.com)

23 Mingyao Li, Ph.D. (mingyao@pennmedicine.upenn.edu)

24 **Abstract**

25 Recent developments of single-cell RNA-seq (scRNA-seq) technologies have led to enormous biological
26 discoveries. As the scale of scRNA-seq studies increases, a major challenge in analysis is batch effects,
27 which are inevitable in studies involving human tissues. Most existing methods remove batch effects in a
28 low-dimensional embedding space. Although useful for clustering, batch effects are still present in the
29 gene expression space, leaving downstream gene-level analysis susceptible to batch effects. Recent
30 studies have shown that batch effect correction in the gene expression space is much harder than in the
31 embedding space. Methods such as Seurat 3.0 rely on the mutual nearest neighbor (MNN) approach to
32 remove batch effects in gene expression, but MNN can only analyze two batches at a time, and it
33 becomes computationally infeasible when the number of batches is large. Here we present CarDEC, a
34 joint deep learning model that simultaneously clusters and denoises scRNA-seq data while correcting
35 batch effects both in the embedding and the gene expression space. Comprehensive evaluations
36 spanning different species and tissues showed that CarDEC outperforms Scanorama, DCA + Combat,
37 scVI, and MNN. With CarDEC denoising, non-highly variable genes offer as much signal for clustering as
38 the highly variable genes (HVGs), suggesting that CarDEC substantially boosted information content in
39 scRNA-seq. We also showed that trajectory analysis using CarDEC's denoised and batch corrected
40 expression as input revealed marker genes and transcription factors that are otherwise obscured in the
41 presence of batch effects. CarDEC is computationally fast, making it a desirable tool for large-scale
42 scRNA-seq studies.

43 **Introduction**

44 Single-cell RNA sequencing (scRNA-seq) analysis has substantially advanced our understanding of
45 cellular heterogeneity and transformed biomedical research. However, the analysis of scRNA-seq data
46 remains confounded by batch effects, which are inevitable in analyses of human tissue and are
47 prevalent in many scRNA-seq studies in general (Hicks et al. 2018; Lahnemann et al. 2020). Several
48 methods have been developed to remove batch effects in scRNA-seq data analysis (Haghverdi et al.
49 2018; Lopez et al. 2018; Barkas et al. 2019; Korsunsky et al. 2019; Stuart et al. 2019b; Welch et al. 2019;
50 Li et al. 2020; Polanski et al. 2020). These methods can be divided into two categories: 1) batch
51 correction in the low-dimensional embedding space, and 2) batch correction in the original gene
52 expression space. Most published papers belong to the first category (Barkas et al. 2019; Korsunsky et
53 al. 2019; Welch et al. 2019; Li et al. 2020; Polanski et al. 2020). Although useful for profiling the overall
54 characteristics of cells such as clustering and trajectory reconstruction, these methods cannot be used
55 for downstream gene-level analysis like differential expression and co-expression analysis.

56
57 A recent benchmarking study has shown that correcting batch effects in the gene expression space is
58 much more challenging than in the embedding space (Lucken et al. 2020). Popular methods such as
59 Seurat 3.0 (Stuart et al. 2019b) rely on the mutual nearest neighbor (MNN) approach (Haghverdi et al.
60 2018) to remove batch effects in the gene expression space, but MNN can only analyze two batches at a
61 time. Its performance is affected by the order in which batches are corrected and it quickly becomes
62 computationally infeasible when the number of batches gets large. Moreover, our evaluations indicate
63 that MNN performs poorly for removing batch effects for genes that are not highly variable. Non-highly
64 variable genes represent the majority of genes in the genome, where batch effects constitute a larger
65 fraction of variance in the transcriptome and are much harder to correct.

66

67 To address this gap in the literature, we present CarDEC (Count adapted regularized Deep Embedded
68 Clustering), a joint deep learning framework for simultaneous batch effect correction, denoising, and
69 clustering of scRNA-seq data. Rather than explicitly modeling batch effect, CarDEC jointly optimizes its
70 reconstruction loss with a self-supervised clustering loss. By minimizing a clustering loss iteratively, the
71 batch effects in the embedding are reduced and cell type signal is improved (Li et al. 2020). The
72 denoised gene expression values, computed from this embedding using a decoder, are then corrected
73 for batch effects as well. To address the difficulty of batch correcting genes that are not highly variable,
74 which suffer from a lower cell type signal-to-noise ratio, we designed CarDEC using a branching
75 architecture that treats highly variable genes (HVGs) and the remaining genes, which we designate as
76 lowly variable genes (LVGs), as distinct feature blocks.

77
78 CarDEC is unique among batch effect correction methods in that it implicitly corrects for batch effects
79 through joint optimization of its dual objective function, rather than explicitly modeling batch effects
80 using batch indicators as in methods such as MNN (Haghverdi et al. 2018) and scVI (Lopez et al. 2018).
81 Moreover, it corrects batch effects both in the low-dimensional embedding space and the original gene
82 expression space. CarDEC’s architecture is founded on the idea of treating HVGs and LVGs as different
83 “feature blocks,” which enables CarDEC to use the HVGs to drive the clustering loss, while still allowing
84 the LVG reconstructions to depend on the rich, batch corrected embedding learned from the HVGs, to
85 help remove batch effects in the LVGs.

86

87 **Results**

88 **Overview of CarDEC and evaluation**

89 An outline of the CarDEC workflow is shown in Fig. 1. CarDEC starts with data preprocessing and
90 pretraining of an autoencoder using HVGs with a mean squared error reconstruction loss function. After

91 pretraining, the weights learned from the pretrained autoencoder are transferred over to the main
92 CarDEC model, which treats HVGs and LVGs as different feature blocks. The main CarDEC loss function is
93 a weighted combination of the reconstruction losses for the HVGs and the LVGs, and a self-supervised
94 clustering loss function driven by the HVGs. This combined loss function allows CarDEC to preserve local
95 structure of the data during clustering (Guo et al. 2017). By minimizing this self-supervised combined
96 loss function, CarDEC not only improves the low-dimensional embedding for clustering but the
97 reconstructed genewise features, which are computed as a function of the low-dimensional embedding,
98 are also denoised and batch effect corrected, leading to improved gene expression quality.

99

100 We evaluated CarDEC on a diverse set of challenging real datasets that range from human to mouse and
101 have different flavors of batch effects. In our evaluations, we wish to assess two properties of CarDEC: 1)
102 its ability to recover biological signals in the data, and 2) its ability to remove spurious technical signals
103 driven by batch effects. An ideal method should strive to remove batch effects while maintaining true
104 biological variations. We compared CarDEC with several state-of-the-art scRNA-seq methods for
105 denoising, batch effect correction, and clustering. scVI (Lopez et al. 2018) and DCA (Eraslan et al. 2019)
106 are multi-use methods that provide denoised counts in the gene expression feature space, and also a
107 low-dimensional embedding that can be used for tasks like clustering and visualization. scVI also
108 attempts to correct for batch effects by conditioning on batch annotation when modeling the denoised
109 counts with a zero-inflated negative binomial distribution. Since DCA is not designed for batch effect
110 correction, to make a fair comparison, we applied Combat (Johnson et al. 2007) to DCA denoised gene
111 expression. MNN (Haghverdi et al. 2018) is a batch correction method that merges batches in a pairwise
112 manner and generates batch corrected gene expression on a cosine scale. We used the implementation
113 of the MNN method provided in the R programming language (R Core Team 2020). Scanorama (Hie et al.
114 2019) is also based on the mutual nearest neighbor idea, but it finds matching elements among all

115 batches at once, thus speeding it up computationally and making the method invariant to batch order.
116 scDeepCluster (Tian et al. 2019) is a clustering method that also draws inspiration from the self-
117 supervised clustering loss (Guo et al. 2017). We provide the exact software implementations of these
118 packages that we used in Supplemental Table 3.

119

120 To measure the degree of batch mixing, we examined the batchwise centroids after denoising and/or
121 batch correction with each method by calculating a coefficient of variation (CV) metric. For each gene, a
122 CV is calculated for a given cell type using the centroid of each batch in that cell type. Then, to obtain a
123 single CV score, we take the weighted average of the cell-type-specific CVs, where the weight of a cell
124 type is the fraction of the dataset's cells that belong to that cell type. The reason for computing CV
125 scores within cell types is that we expect minimal biological heterogeneity within cell types, so in the
126 absence of batch effects we expect batch centroids to be similar to one another within, but not between
127 cell types. A higher value of CV corresponds to greater variation of gene expression among batches and
128 less batch mixing, whereas a good batch effect removal method should drive the CV value close to zero.

129

130 **Application to human pancreatic islet data from four protocols**

131 A unique feature of CarDEC is the branching architecture for both the HVGs and the LVGs. To
132 demonstrate that this architecture is key in removing batch effects, we combined four datasets
133 consisting of scRNA-seq expression data in human pancreas generated using Fluidigm C1 (Lawlor et al.
134 2017), Smart-seq2 (Segerstolpe et al. 2016), CEL-Seq (Grun et al. 2016), and CEL-Seq2 (Muraro et al.
135 2016). This is a challenging task as there are strong batch differences among these different scRNA-seq
136 protocols, and analysis using the raw data as input yielded low Adjusted Rand Indexes (ARIs)
137 (Supplemental Fig. S1). The branching architecture was designed with two objectives in mind. First, we
138 wish to show that when correcting batch effects and denoising both the HVGs and the LVGs, using a

139 branching model that treats these feature blocks differently improves the quality of denoised expression
140 values relative to a naïve architecture that treats these feature blocks the same. Second, we hope to
141 design a model architecture such that including the LVGs in the model does not worsen denoising and
142 batch effect correction quality of the HVGs, relative to a naïve model that only denoises the HVGs and
143 does not attempt to denoise LVGs.

144

145 As shown in Fig. 2, the branching architecture posts significant performance boosts over the naïve
146 architecture that treated all genes as the same feature block in the input. The branching architecture
147 performed better for denoising both the HVGs (ARI of 0.93 over 0.72) and the LVGs (ARI of 0.83 over
148 0.67) relative to the naïve architecture (Fig. 2A,B), underscoring the necessity of using the branching
149 architecture to denoise all genes as efficiently as possible. We also observed that for the purpose of
150 denoising and batch correcting only the HVGs, the branching architecture performed just as well as a
151 naïve model that only included the HVGs and completely discarded the LVGs (ARI of 0.93 vs 0.95)
152 (Supplemental Fig. S2). This verifies that the branching architecture does not trade off denoising and
153 batch correction effectiveness on the HVGs at all to denoise the LVGs. Additionally, the denoised counts
154 from CarDEC showed less batch effects compared to denoised expression from Scanorama, scVI, and
155 batch corrected expression from MNN (Fig. 2A,B). Fig. 2C shows that genewise CV scores obtained from
156 CarDEC are the closest to zero among all compared methods. We also noticed that the clustering
157 accuracy obtained when clustering is done using denoised values in the gene expression space is similar
158 to the clustering accuracy obtained when using the embedding of CarDEC to do clustering (Fig. 2D).

159

160 **Application to macaque retina data with multi-level batch effect**

161 After finalizing the CarDEC architecture, we next evaluated the performance of CarDEC on a macaque
162 retina dataset (Peng et al. 2019). This dataset poses a great challenge for batch effect correction and

163 denoising because it features strong, multi-level batch effects, with cells sequenced from two different
164 regions, four different macaques, and thirty different samples (Supplemental Fig. S3).

165

166 For the task of denoising and batch effect correction in the gene expression space, CarDEC was again the
167 best performing method (Fig. 3, Supplemental Fig. S4, Supplemental Fig. S5), followed by scVI whose
168 ARIs are close to CarDEC. CarDEC and scVI not only removed the multi-level batch effects but also
169 preserved inter-cell type variation. Notably, the ARI for clustering using the LVG denoised and batch
170 effect corrected counts from CarDEC is 0.98 and is 0.97 for scVI (Fig. 3B), which is as high as that when
171 using the HVGs to do clustering (Fig. 3A). As a comparison, the ARI is only 0.15 using the LVG raw counts
172 as input for clustering (Supplemental Fig. S3). This suggests that the denoising and batch correction in
173 CarDEC and scVI substantially boosted the signal-to-noise ratio in the LVGs. Moreover, CarDEC's and
174 scVI's genewise CVs are the closest to zero, providing evidence that cells were mixed well by batch by
175 these two methods (Fig. 3C, Supplemental Fig. S6).

176

177 The other methods all struggled with batch effects in the denoised counts. Scanorama largely failed to
178 correct batch effects: when using Scanorama batch corrected gene expression as input, the cells were
179 separated primarily by batch rather than by cell type. For the LVGs, its ARI is as low as that when using
180 the LVG raw counts as input for clustering (Scanorama 0.21 vs raw 0.15) (Fig. 3B, Supplemental Fig. S3).
181 DCA had slightly higher ARIs than Scanorama for both the HVGs and the LVGs, although both are
182 significantly lower than CarDEC and scVI (Fig. 3A,B). MNN performed much better than Scanorama and
183 DCA for batch correcting the HVG counts, achieving an ARI of 0.91 (Fig. 3A). However, it still fell short of
184 CarDEC and scVI for this evaluation (CarDEC ARI 0.98 and scVI ARI 0.96). Looking more closely at the
185 HVG UMAP plots, the batches were mixed less thoroughly with MNN than they were for CarDEC and
186 scVI, and the cells were separated less by cell type indicating that MNN failed to completely recover cell

187 type variation. This is further confirmed by the genewise CV density plot in which the MNN density curve
188 is further away from zero than CarDEC and scVI (Fig. 3C). For removing batch effects in the LVG counts,
189 MNN was the worst performing method because it removed nearly all biological variations, leaving only
190 batch effects.

191
192 For the simpler task of clustering using the low-dimensional embedding representation, existing
193 methods did considerably better than they did at batch effect correction in the gene expression space,
194 but still fell short of CarDEC and scVI (Fig. 3D, Supplemental Fig. S7). CarDEC and scVI achieved ARIs of
195 nearly 1 for clustering using the embedding. DCA and scDeepCluster performed better than Louvain's
196 algorithm using raw HVGs, but still fell short of achieving as high of an ARI as CarDEC and scVI.
197 Scanorama struggled on this dataset with an ARI of only 0.57.

198

199 **Application to mouse cortex and PBMC data from four protocols**

200 We next compared different methods using a mouse cortex dataset (Ding et al. 2020). This dataset
201 poses a great challenge for batch correction and denoising on two fronts (Supplemental Fig. S9). First, it
202 exhibits very serious batch effects because cells were generated using four different scRNA-seq
203 protocols. Furthermore, this dataset is heavily dominated by excitatory and inhibitory neurons, and the
204 other cell types are rare, so preserving biological variation is especially imperative for detecting and
205 analyzing these rarer subpopulations.

206

207 For the task of denoising and batch correcting the gene expression space, CarDEC and scVI performed
208 considerably better than the other methods (Fig. 4). CarDEC performed the best, followed closely by
209 scVI, at balancing between removing batch effects while preserving as much cell type variability as
210 possible. For both CarDEC and scVI, the ARIs are similar when using the HVG denoised counts and the

211 LVG denoised counts as input for clustering (Fig. 4A,B). As a comparison, the ARIs are only 0.28 and 0.26
212 when using the HVG and the LVG raw counts as input for clustering, respectively (Supplemental Fig. S9).
213 The relatively low ARI when using the HVG raw counts as input for clustering demonstrates the strong
214 batch effects in this dataset. However, for this challenging dataset, using CarDEC denoised and batch
215 corrected LVG counts, the ARI increased to 0.78, suggesting that CarDEC substantially boosted the
216 signal-to-noise ratio in the LVGs by simultaneous denoising and batch effect removal. We observed a
217 similar signal boost by scVI in which the ARI for the LVGs is 0.73. The genewise CVs for CarDEC are also
218 the closest to zero among all methods (Fig. 4C). Although scVI achieved high ARIs for both the HVG and
219 LVG denoised gene expression, the cells appeared to be less well mixed than CarDEC (Fig. 4A,B), which is
220 in agreement with its larger CVs than CarDEC.

221

222 DCA largely failed for this dataset (Fig. 4A,B), even though its denoised gene expression was also batch
223 corrected by Combat. For both the HVGs and the LVGs, DCA separated the cells purely by scRNA-seq
224 protocol with no mixing of cells from different batches. After denoising using DCA and batch correction
225 with Combat, cell variation was driven entirely by batch, rendering the denoised counts ineffective for
226 downstream analyses. For batch correcting the HVGs, Scanorama was the third-best performer (Fig. 4A),
227 followed by MNN. MNN did not merge batches to the extent that CarDEC did and failed to preserve as
228 much cell type variability, causing cell types to mix more. For removing batch effects in the LVGs, MNN
229 did considerably worse than CarDEC and scVI (Fig. 4B). It suffered from the same problems as DCA for
230 the LVGs in that cell type variation was lost and all variability was driven by batch. We also noticed that
231 the CVs for DCA are the closest to zero, however, the small CVs are mainly due to the overcorrection of
232 Combat as the ARIs are low for both the HVGs and the LVGs.

233

234 Due to the strong batch effects in this dataset, even the simpler task of clustering using embedding was
235 very difficult (Fig. 4D, Supplemental Fig. S10). In particular, scDeepCluster performed poorly at this task,
236 scoring a lower ARI than a straightforward application of Louvain's algorithm to the raw data. The ARI
237 for Scanorama is only slightly better than that obtained from the raw data. For this task, scVI was the
238 leader, achieving an ARI of 0.74, followed closely by CarDEC with an ARI of 0.73.

239

240 We also analyzed a dataset of human PBMCs from the same paper (Ding et al. 2020) as the mouse
241 cortex data. This dataset was similar to the cortex dataset: featuring eight batches spanning five scRNA-
242 seq protocols and the results were largely the same: CarDEC and scVI were the best performers for
243 denoising/batch correcting the HVGs and also the best for denoising/batch correcting the LVGs
244 (Supplemental Fig. S12-S14).

245

246 **Application to human monocyte data with pseudotemporal structure**

247 We next show the utility of CarDEC for improving trajectory analysis for cells with pseudotemporal
248 structure. We analyzed an scRNA-seq dataset generated from monocytes derived from human
249 peripheral blood mononuclear cells by Ficoll separation followed by CD14⁻ and CD16⁺ cell
250 selection (Li et al. 2020). This dataset includes 10,878 monocytes from one healthy subject. The cells
251 were processed in three batches from blood drawn on three different days. Although monocytes can be
252 classified as classical (CD14⁺⁺/CD16⁻), intermediate (CD14⁺⁺/CD16⁺), and nonclassical patrolling (CD14⁻
253 /CD16⁺⁺) subpopulations based on surface markers, our previous analysis based on scRNA-seq data
254 indicates that these cells show continuous transcriptional characteristics and trajectory analysis is an
255 appropriate approach to characterize them (Li et al. 2020). This dataset has strong batch effects
256 (Supplemental Fig. S16). To reconstruct the trajectories of these cells, for each method, we first

257 denoised and/or batch corrected the gene expression matrix, which was then fed into Monocle 3 (Cao et
258 al. 2019) to estimate the pseudotime of each cell.

259

260 Fig. 5A shows that CarDEC yields the best pseudotime analysis results with cells from the three batches
261 well mixed, and a clear pseudotemporal path emerged. The batchwise density plots show that the three
262 batches have similar pseudotime distributions, suggesting that CarDEC successfully removed batch
263 effects. The plots of *FCGR3A* (known marker gene for nonclassical monocytes) and *S100A8* (known
264 marker gene for classical monocytes) gene expression also showed expected patterns (Supplemental
265 Fig. S17). There are two key points of evidence from these marker gene plots suggesting that CarDEC
266 recovered biological signal. First, for each marker gene, the expression levels are virtually identical
267 across batches for all pseudotime points, which indicates that batch effects were removed for each gene
268 expression and pseudotime relationship. Also, *FCGR3A* gene expression decreases monotonically with
269 pseudotime, while *S100A8* expression increases monotonically with pseudotime. This is exactly the kind
270 of behavior we expect from these marker genes. Since *FCGR3A* and *S100A8* are markers for the
271 nonclassical and classical monocytes, respectively, we expect a good pseudotime analysis to segment
272 the monocytes from nonclassical to classical (or vice versa) and for *FCGR3A* and *S100A8* expressions to
273 be monotonic functions of pseudotime with opposite trends. By denoising and batch correcting gene
274 counts, CarDEC successfully mixed batches and recovered biological signal down to the individual
275 marker gene level.

276

277 By contrast, no other methods were able to achieve CarDEC's success in improving pseudotime analysis.
278 Scanorama (Fig. 5B), DCA (Fig. 5C), and MNN (Fig. 5E) all failed to mix batches in the UMAP embedding
279 from Monocle 3. Although scVI (Fig. 5D) mixed batches well in the UMAP embedding, the pseudotime
280 distributions showed substantial variation across batches. For Scanorama, DCA, and MNN, neither of the

281 marker genes show strong monotonic trends as a function of the pseudotime, and for each marker
282 gene, the relationship between expression and pseudotime varied by batch. These issues suggest that
283 Scanorama, DCA, and MNN confounded biological signal, and obscured signals from canonical markers
284 of established subpopulations, *FCGR3A* and *S100A8*, as marker genes.

285 There are other approaches to using these denoising and batch correction methods for pseudotime
286 analysis. For example, one can subset the denoised and batch corrected matrix to include only the HVGs
287 and then feed this into Monocle 3 (Supplemental Fig. S18, Supplemental Fig. S20). Alternatively, one can
288 use the embedding from CarDEC, Scanorama, DCA, or scVI as the reduced dimension space to build the
289 Monocle 3 pseudotime graph (Supplemental Fig. S19, Supplemental Fig. S21). In both of these other
290 cases, the conclusions are largely the same, CarDEC is the best method for improving pseudotime
291 analysis.

292

293 Next, we examined whether the denoised and batch corrected gene expression values can help improve
294 gene expression quality for biological discovery. We focused our analyses on 61 transcription factors
295 (TFs) that were expressed in the monocyte data and also found to be differentially expressed among
296 classical, intermediate, and nonclassical monocytes by Wong *et al.* (Wong et al. 2011). Among these 61
297 TFs, 23 were selected as HVGs and the remaining 38 were designated LVGs. Fig. 6A shows that the
298 CarDEC denoised gene expression revealed a gradually decreasing trend from nonclassical to classical for
299 TFs that are known to be highly expressed in nonclassical monocytes, e.g., *TCF7L2*, *POU2F2*, *CEBPA*, and
300 *HSBP1*. We also observed expected gene expression increase from nonclassical to classical for TFs that
301 are known to be highly expressed in classical monocytes, e.g., *NFE2*, *CEBPD*, *GAS7*, and *MBD2*. Notably,
302 some of the TFs with these expected expression patterns were not selected as HVGs, suggesting that
303 denoising and batch correction in CarDEC helped recover the true biological variations in both HVGs and

304 LVGs. By contrast, when using raw UMI counts as input, the heatmap did not reveal any meaningful
305 biological patterns even for those TFs that were selected as HVGs (Fig. 6B).

306

307 An important task in trajectory analysis is to identify genes whose expression values change over
308 pseudotime and whether the expression patterns are different between conditions (e.g., healthy vs
309 diseased) over pseudotime. Avoiding false positive results is critical as failure of doing so may lead to
310 follow-up of a wrong signal. Since the three batches were obtained from the same subject, we do not
311 expect to detect significant gene expression differences over pseudotime among them. To this end, we
312 performed differential expression analysis and compared the distributions of gene expression changes
313 over pseudotime across the three batches. We performed hypothesis tests using the '*gam*' function in R
314 package *mgcv* and tested whether gene expression patterns for the three batches are significantly
315 different over pseudotime. Fig. 6C shows the p-values from this differential expression analysis for each
316 method. CarDEC is much more effective in removing batch effects than Scanorama, DCA, scVI, and MNN.
317 For the 23 HVG TFs, the median $-\log_{10}$ p-value for CarDEC is 0.30, whereas the median $-\log_{10}$ p-values for
318 Scanorama, DCA, scVI, and MNN are 18.37, 5.29, 25.90, and 28.52, respectively. For the 38 LVG TFs, the
319 median $-\log_{10}$ p-value for CarDEC is increased to 3.70, but still much lower than the other methods (9.81
320 for Scanorama, 6.36 for DCA, 39.19 for scVI, and 9.07 for MNN). These results indicate that failure to
321 correct for batch effects could lead to a severe inflation of false positive results. Fig. 6D shows four
322 selected TFs, where the denoised and batch corrected gene expression for CarDEC agreed well among
323 the three batches, further confirming the effectiveness of CarDEC in removing batch effects in the gene
324 expression space.

325

326 **CarDEC is scalable to large dataset**

327 As the scale of scRNA-seq continues to grow, it becomes increasingly important for a method to be
328 scalable to large datasets. To evaluate the scalability of CarDEC, we leveraged a dataset of 104,694
329 human fetal liver cells (Popescu et al. 2019). Since we are principally interested in the problem of
330 denoising and batch correcting in the full gene expression space, we retained all 21,521 genes after
331 initial filtering for this analysis. For CarDEC we benchmarked two variations: a version that provides only
332 denoised/batch corrected expression in the Z-score space (CarDEC Z-score) and a version that provides
333 denoised/batch corrected expression in the count space (CarDEC Count).

334 We evaluated the runtime needed to process 10%, 20%, 40%, 60%, 80%, and 100% of cells in the human
335 fetal liver dataset for CarDEC, Scanorama, DCA, scVI, and MNN. All evaluations were done on a 2019
336 edition MacBook Pro with 2.4 GHz 8-Core Intel Core i9 CPU and 32 GB of memory. CarDEC, DCA, and
337 scVI were all trained with early stopping to halt training upon convergence. The results are shown in
338 Supplemental Fig. S22. Both versions of CarDEC as well as DCA scaled approximately linearly with the
339 number of cells and all three of these methods finished the analysis in less than 3.5 hours. We were
340 unable to train scVI on more than 60% of the cells, as the jupyter kernel crashed mid-training for 80% of
341 the cells. scVI also took considerably longer to run. At 60% of the data, CarDEC Z-score, DCA, and CarDEC
342 Count took about 40 minutes, 1 hour, and 1.5 hours, respectively. By contrast, scVI took about 3 hours
343 and 45 minutes to analyze this data. Scanorama had serious scalability issues. Like scVI, we could not run
344 Scanorama for more than 60% of the data. Among the data points that we do have, Scanorama's run
345 time is clearly not $O(n_{cells})$, while CarDEC, DCA, and scVI scale roughly linearly due to the minibatch
346 gradient descent algorithm that trains all of them (any non-monotonicity is due to early stopping
347 variation). Scanorama's run time appears to be parabolic as a function of sample size. MNN, also has
348 serious scalability issues, which is consistent with a recent benchmarking study (Haghverdi et al. 2018).
349 It took over 12 hours to analyze 20% of the data, and over 47 hours to analyze 40% of the data. We
350 could not run MNN in under 48 hours using more than 40% of the data.

351

352 CarDEC is robust to hyperparameters

353 CarDEC involves several parameters that need to be tuned. To evaluate if CarDEC is robust to different
354 choices of these parameters, we conducted additional analyses. First, we evaluated the performance of
355 CarDEC for the α weight parameter, which is used when combining the self-supervised clustering loss
356 and the reconstruction losses. The α parameter can be set to any value in the range of [0, 2]. Since a
357 weight of 1 is the midpoint of this range, we chose it as a natural default. To evaluate the robustness of
358 CarDEC to the choice of α , we varied its value across the set [0.1, 0.55, 1.0, 1.45, 1.9] and measured how
359 the ARI changes for both the HVGs and the LVGs for each chosen value of α with all other parameters
360 held equal. As shown in Supplemental Fig. S23, CarDEC is robust to this choice of parameter.

361

362 Next, we evaluated the performance of CarDEC when varying the number of clusters for clustering. To
363 show that CarDEC is robust to this parameter, we cannot rely on ARI, since the maximum achievable ARI
364 for any clustering method decreases the more the number of clusters is misspecified relative to the
365 number of labels/cell types in the set of gold standard labels. Instead, we chose to show that increasing
366 the number of clusters specified for CarDEC just splits existing clusters, without substantially changing
367 cell type signals/separation. To do this, we fit CarDEC on the pancreas dataset repeatedly, varying the
368 number of clusters with each fit. The numbers of clusters selected were 5, 9, 11, 14, 18, respectively. As
369 shown in Supplemental Fig. 24A, even when the number of clusters is increased from 5 to 18, the cells
370 were still separated mainly by their gold standard cell type labels. Closer examination of the UMAPs
371 revealed that increasing the number of clusters usually just split cell types into two or more cluster label
372 bins to accommodate surplus cluster labels, without changing the underlying structure of the UMAP
373 plot. The CV score distributions shown in Supplemental Fig. 24B further supports that CarDEC is robust
374 to the choice of the number of clusters since the CV score distributions were largely unchanged when

375 the number of clusters varied. We also visualized how cells in a cluster were split when CarDEC was refit
376 with an increasing number of clusters using a Sankey Plot. As shown in Supplemental Fig. 24C, in nearly
377 all cases, increasing the number of clusters just split existing clusters.

378

379 **Discussion**

380 We developed CarDEC, a joint deep learning model, that removes batch effects not only in the low-
381 dimensional embedding space but also across the entire gene expression space. As demonstrated in our
382 evaluations and a recent benchmarking study (Lucken et al. 2020), it is considerably harder to correct for
383 batch effects in the gene expression space than in the embedding space, and especially hard to correct
384 for batch effects in LVGs, which constitute the majority of the transcriptome. CarDEC was built to tackle
385 these challenges. To remove batch effects in the gene expression space, we minimize a loss function
386 that combines clustering and reconstruction losses. The self-supervised clustering loss, driven by HVGs,
387 regularizes the embedding and removes batch effects in the embedding. The rich, batch corrected
388 embedding is then used to compute an effectively batch corrected representation in the original gene
389 expression space. To address the difficulty associated with batch correcting LVGs, we implemented a
390 branching architecture, where embeddings are computed separately for HVGs and LVGs and where only
391 the HVG embedding is used to compute the clustering loss. Using the pancreatic islet datasets generated
392 from four scRNA-seq protocols, we demonstrated that this branching architecture substantially
393 improved batch effect removal on both the HVG and LVG gene expression spaces, as compared to the
394 naïve architecture.

395

396 Across a variety of datasets, with batch effects spanning multiple complexities in level and strength, we
397 demonstrated that CarDEC consistently led in its ability to remove batch effects. CarDEC was
398 consistently the best for removing batch effects in all capacities: in the embedding space, the HVG

399 expression space, and the LVG expression space. We showed that with appropriate denoising and batch
400 correction, the LVGs offer as much signal for clustering as the HVGs, suggesting that CarDEC has
401 substantially boosted the amount of information content in scRNA-seq. We also demonstrated that by
402 batch correcting gene expression counts, CarDEC improved pseudotemporal analysis of human
403 monocytes, an example of how batch correction can be used to improve downstream analyses.

404

405 CarDEC removes batch effects in the LVG denoised gene expression based on the embedding layer that
406 concatenates the HVG embedding. A potential concern of this concatenation is that the denoised LVG
407 expression values might contain artificially introduced signals from the HVGs. To examine this, we
408 focused on the bottom 16215 genes in the Pancreas dataset by variance and evaluated the performance
409 of CarDEC when varying the number of HVGs. Then, we fit CarDEC for various numbers of HVGs and
410 evaluated the ARI obtained from clustering the denoised expression for these 16215 genes. If the LVGs
411 contain artificial signals from the HVGs, then we would expect the ARI computed from the bottom
412 16125 genes to vary by how many HVGs were used to train CarDEC. As demonstrated in Supplemental
413 Fig. S25, this is not the case. We observed that ARI on the LVG set increased modestly when the number
414 of HVGs was increased from 500 to 1000, likely because 500 HVGs were too few highly variable features.
415 From 1000 up to 3000 inclusive, the ARI on the LVG set was very stable. After 3000 HVGs, there was a
416 modest drop in the LVG set's ARI, but this is likely because too many noisy features were introduced into
417 the HVG set, which hurt the quality of the KL-Divergence gradients. Since the 16125 LVG set's ARI is
418 fairly invariant to the number of HVGs used, especially in the range of 1000 to 3000 HVGs, we can infer
419 that the model architecture is not adding significant artificial signals from the HVGs to the LVGs.

420

421 In this paper, we focused on the analyses when using all genes as input. Since Scanorama, DCA, scVI, and
422 MNN are intended to be used with HVGs only, we reanalyzed every dataset with only HVGs selected as

423 the input using every method. As expected, many methods performed better when only HVGs were
424 modeled (Supplemental Fig. S2, Supplemental Fig. S8, Supplemental Fig. S11, Supplemental Fig. S15).
425 However, for the macaque retina data, Scanorama still struggled to remove its batch effects even when
426 only HVGs were considered. For this challenging dataset, the ARI from Scanorama is only 0.18, which is
427 much lower than all of the other methods (Supplemental Fig. S8).

428

429 Current scRNA-seq studies often include a large number of cells generated from many samples, across
430 multiple conditions, and possibly using different protocols. Removing batch effects is critical for data
431 integration. Since CarDEC provides efficient batch correction in the full gene expression space, it can be
432 used for a wide array of analyses to facilitate biological discovery. Harmonized counts in the gene
433 expression space can be used to estimate unbiased, batch corrected log fold changes, which can be used
434 to identify marker genes for different cell types. These counts can also be used to reconstruct
435 trajectories and identify genes showing pseudotemporal patterns. Lastly, CarDEC is computationally fast
436 and memory efficient, making it a desirable tool for analyses of complex data in large-scale single-cell
437 transcriptomics studies.

438 **Figure legends**

439 **Figure 1. The workflow of CarDEC.** The CarDEC workflow can be summarized in four steps that are
440 depicted here: preprocessing, pretraining, denoising, and optionally, denoising on the count scale.

441
442 **Figure 2. Justification for the branching architecture in CarDEC.** The CarDEC API splits the input matrix
443 into HVGs and LVGs and treats them separately with a “Branching” architecture as in Fig. 1.
444 Alternatively, we can use a “Naïve” model that treats all features the same regardless of gene
445 expression variance. The naïve model consists of an autoencoder with a clustering loss in addition to the
446 reconstruction loss. Here we demonstrate the utility of the Branching architecture. The HVGs and LVGs
447 were clustered separately, and the ARI of assignments is provided along with a UMAP plot. First row was
448 colored by cell type, second by scRNA-seq protocol. (A) UMAP embedding computed from the denoised
449 HVG counts for each method: CarDEC with Branching architecture, CarDEC with Naïve Architecture,
450 Scanorama, scVI, and MNN. (B) UMAP embedding computed from the denoised LVG counts for each
451 method. Figure legends are the same as those in (A). (C) Density plot of genewise coefficient of variation
452 (CV) among batch centroids. Density plots are provided for HVGs and LVGs separately. (D) Clustering
453 accuracy metrics were obtained using the embedding based methods to cluster the data, rather than
454 running Louvain on the full gene expression space. Results for “Raw” were obtained by using Louvain’s
455 algorithm on the original HVG counts and are provided as a baseline with which to compare embedding
456 based clustering results to.

457
458 **Figure 3. Comparison of different methods on the macaque retina dataset.** (A) UMAP embedding
459 computed from the denoised HVG counts for each method. Top row was colored by cell type; bottom
460 was colored by Macaque ID. Cells were also clustered with Louvain’s algorithm. (B) UMAP embedding
461 computed from the denoised LVG counts for each method. Figure legends are the same as those in (A).

462 (C) Density plot of genewise coefficient of variation (CV) among batch centroids. Density plots are
463 provided for HVGs and LVGs separately. Centroids computed with sample id as batch definition. (D)
464 Clustering accuracy metrics obtained using the embedding based methods to cluster the data, rather
465 than running Louvain on the full gene expression space. Results for “Raw” were obtained by using
466 Louvain’s algorithm on the original HVG counts and are provided as a baseline with which to compare
467 embedding based clustering results to.

468

469 **Figure 4. Comparison of different methods on the mouse cortex dataset.** (A) UMAP embedding
470 computed from the denoised HVG counts for each method. Top row was colored by cell type; bottom
471 was colored by batch. Cells were also clustered with Louvain’s algorithm, and resultant ARI is provided.
472 (B) UMAP embedding computed from the denoised LVG counts for each method. Figure legends are the
473 same as those in (A). (C) Density plot of genewise coefficient of variation (CV) among batch centroids.
474 Density plots are provided for HVGs and LVGs separately. (D) Clustering accuracy metrics obtained using
475 the embedding-based methods to cluster the data, rather than running Louvain on the full gene
476 expression space. Results for “Raw” were obtained by using Louvain’s algorithm on the original HVG
477 counts and are provided as a baseline with which to compare embedding based clustering results to.

478

479 **Figure 5. Comparison of different methods for pseudotime analysis in the human monocyte data.** The
480 analysis is done on monocytes derived from three technical replicates from the same subject. For each
481 method, the full dataset was denoised/batch corrected and then fed to Monocle 3 for pseudotime
482 analysis. We show the UMAP embedding colored by batch (column 1) and estimated pseudotime
483 (column 2). We also visualize the kernel density distribution of pseudotime by batch (column 3) and plot
484 the distributions of marker genes *FCGR3A* and *S100A8* against pseudotime (columns 4 and 5,
485 respectively). (A) Pseudotime analysis when using denoised/batch corrected gene expression matrix

486 from CarDEC as input. (B) Pseudotime analysis when using batch corrected gene expression matrix from
487 Scanorama as input. (C) Pseudotime analysis when using denoised gene expression matrix from DCA but
488 with Combat post-hoc batch correction as input. (D) Pseudotime analysis when using denoised/batch
489 corrected gene expression matrix from scVI as input. (E) Pseudotime analysis when using batch
490 corrected gene expression matrix from MNN as input.

491

492 **Figure 6. Comparison of different methods for differential expression analysis of transcription factors**

493 **in the human monocyte data.** (A) Heatmap of scaled gene expression for CarDEC. Pseudotime was

494 inferred based on embedding obtained from CarDEC using Monocle 3. (B) Heatmap of scaled raw UMI

495 counts. Pseudotime was inferred based on embedding obtained from the scaled raw UMI counts using

496 Monocle 3. (C) p-values obtained from differential expression analysis among the three batches over

497 pseudotime. For each method, the pseudotime was inferred based on embedding obtained from the

498 corresponding method. The red dotted line corresponds to p-value = 0.01. The top panel is for the 23

499 HVG TFs and the bottom panel is for the 38 LVG TFs. (D) Denoised and batch corrected gene expression

500 for CarDEC, batch corrected gene expression for Scanorama, denoised gene expression for DCA with

501 Combat post-hoc batch correction, denoised and batch corrected gene expression for scVI, and batch

502 corrected gene expression for MNN over pseudotime for four selected TFs, *HIF1A*, *HES4*, *HSBP1*, and

503 *GAS7*. For each method, the pseudotime was inferred based on embedding from the corresponding

504 method.

505 **Methods**

506 The CarDEC workflow (Fig. 1) involves four steps: preprocessing, pretraining, gene expression denoising
507 in Z-score space, and (optionally) denoising in count space. Below we briefly describe each of these
508 steps. Details of the implementation is described in Supplemental Note 1, and the hyperparameters of
509 CarDEC are shown in Supplemental Table 1.

510

511 **Step 1: preprocessing**

512 We first remove any cells expressing less than 200 genes, and then remove any genes expressed in less
513 than 30 of the remaining cells. Let \mathbf{X} be an $n \times p$ gene count matrix with n cells and p genes after
514 filtering. The gene expression values are normalized. In the first step, cell level normalization is
515 performed in which gene expression for a given gene in each cell is divided by the total gene expression
516 across all genes in the cell, multiplied by the median total expression across all cells, and then
517 transformed to a natural log scale. In the second step, gene level normalization is performed in which
518 the cell level normalized values for each gene are standardized by subtracting the mean and dividing by
519 the standard deviation across all cells within the same batch for the given gene. Highly variable genes
520 (HVGs) are selected based on the log-normalized counts using the approach introduced by Stuart and
521 Butler (Stuart et al. 2019a) and implemented in the “pp.highly_variable_genes” function with
522 “batch_key” parameter in the SCANPY package (version ≥ 1.4) (Wolf et al. 2018). The remaining genes
523 that are not selected as HVGs are considered lowly variable genes (LVGs). We note that many of the
524 LVGs still show cell-to-cell variability and are useful for clustering analysis after appropriate denoising
525 and batch effect correction. We select 2,000 HVGs for all analyses in this paper.

526

527 **Step 2: pretraining using the HVGs**

528 The pretraining step is a straightforward implementation of an autoencoder. Let p_{HVG} be the number of
 529 HVGs selected in Step 1, and \mathbf{Y}_{HVG} be the corresponding $n \times p_{HVG}$ matrix of normalized expression,
 530 subsetting to include only the HVGs. Define a standard autoencoder for \mathbf{Y}_{HVG} with encoder and decoder
 531 represented by $f_{E,HVG}(\cdot; W_{E,HVG})$ and $f_{D,HVG}(\cdot; W_{D,HVG})$, respectively. The weights $W_{E,HVG}$ and
 532 $W_{D,HVG}$ are randomly initialized using the glorot uniform approach and are tuned during pretraining. We
 533 use the tanh activation for the output of the encoder, and the linear activation function for the output
 534 of the decoder. For all intermediate hidden layers in the encoder and decoder, we use the ReLU
 535 activation function. The autoencoder is pretrained with mean squared error loss using minibatch
 536 gradient descent with the Adam optimizer (Kingma and Ba 2015).

537

538 **Step 3: denoising Z-scores**

539 In this step, we use an expanded, branching architecture to accommodate LVGs, and introduce a
 540 clustering loss that regularizes the embedding and improves batch mixing and denoising especially in the
 541 gene space. Let p_{LVG} be the number of LVGs selected in Step 1, and \mathbf{Y}_{LVG} be the corresponding $n \times p_{LVG}$
 542 matrix of normalized expression, subsetting to include only the LVGs, and $\mathbf{y}_{i,HVG}$ and $\mathbf{y}_{i,LVG}$ be the
 543 vectors of HVGs and LVGs, respectively in cell i . We retain the encoder and decoder mappings for HVGs,
 544 $f_{E,HVG}(\cdot; W_{E,HVG})$ and $f_{D,HVG}(\cdot; W_{D,HVG})$ from Step 2, including the learnt weights $W_{E,HVG}$ and
 545 $W_{D,HVG}$. We introduce a clustering layer that takes the HVG embedding
 546 $\mathbf{z}_{i,HVG} = f_{E,HVG}(\mathbf{y}_{i,HVG}; W_{E,HVG})$ as input and returns for each cell a vector of cluster membership
 547 probabilities for h clusters, where h is a user-specified number. For this clustering layer, we introduce an
 548 $h \times d$ matrix of trainable weights/cluster centroids \mathbf{M} , where the j^{th} row of \mathbf{M} is a cluster centroid $\boldsymbol{\mu}_j$,
 549 and d is dimension of the embedding.

550

551 To initialize \mathbf{M} , we run Louvain’s algorithm on the embeddings $\{\mathbf{z}_{i,HVG}: i \in \{1, 2, \dots, n\}\}$ learnt from the
 552 pretrained autoencoder and find the cluster centroid for each cluster. The clustering layer computes a
 553 vector of cluster membership probabilities for cell i , denoted by \mathbf{q}_i . Let q_{ij} , the j^{th} element of \mathbf{q}_i ,
 554 denote the probability that cell i belongs to cluster j . Then the membership probabilities are computed
 555 using a t -distribution kernel as follows,

556

$$q_{ij} = \frac{\left(1 + \|\mathbf{z}_{i,HVG} - \mu_j\|^2\right)^{-1}}{\sum_{j'} \left(1 + \|\mathbf{z}_{i,HVG} - \mu_{j'}\|^2\right)^{-1}}$$

557

558 Since we do not have cell type labels in an unsupervised analysis, we create “pseudo-labels” that can be
 559 used in place of real labels for optimizing clustering weights. Inspired by Xie *et al.* (Xie et al. 2016), these
 560 pseudo-labels are computed from the membership probabilities q_{ij} as follows,

561

$$p_{ij} = \frac{q_{ij}^2 / \sum_i q_{ij}}{\sum_{j'} q_{ij'}^2 / \sum_i q_{ij'}}$$

562

563 Let \mathbf{p}_i be an h -dimensional vector whose j^{th} element is p_{ij} . Then the clustering loss for cell i is defined
 564 as the following Kullback–Leibler divergence (KLD),

565

$$l_{i,c} = KLD(\mathbf{p}_i || \mathbf{q}_i) = \sum_j p_{ij} \log\left(\frac{p_{ij}}{q_{ij}}\right)$$

566

567 This loss is a component of the total loss defined later. Since it takes the embedding vectors $\mathbf{z}_{i,HVG}$ as
 568 input, minimizing this objective function can refine the embedding and help to remove batch effects

569 from denoised counts computed using this embedding as input. The use of this KLD loss function was
 570 inspired by DESC(Li et al. 2020), which has shown that batch effects can be gradually removed over
 571 iterations. The intuition is to use “easy-to-cluster” cells, i.e., cells with “pseudo-labels”, to guide the
 572 neural network to learn cluster-specific gene expression features while ignoring other unwanted noises
 573 such as batch effects. Over iterations, the algorithm ignores information unrelated to cell clustering by
 574 adjusting the network weights using the gradient descent algorithm and learns cluster-specific
 575 information. During the iterative update procedure, cells that are initially assigned to the same cluster
 576 are moved closer and closer to the cluster centroid, hence removing batch effects in the embedding
 577 space. Since the algorithm learns information on cell clusters from those “easy-to-cluster cells” while
 578 ignoring other irrelevant information by constructing the auxiliary distribution \mathbf{P} and optimizing the KL-
 579 divergence between \mathbf{P} and \mathbf{Q} , then as long as technical differences (e.g., between batches) are smaller
 580 than biological differences (e.g., between cell types), it can remove batch effect successfully.

581
 582 We also introduce encoder and decoder mappings $f_{E,LVG}(\cdot; W_{E,LVG})$ and $f_{D,LVG}(\cdot; W_{D,LVG})$ to address
 583 the problem of denoising and batch correction for the LVGs. Unlike the HVG decoder $f_{D,HVG}$, the LVG
 584 decoder $f_{D,LVG}(\cdot; W_{D,LVG})$ does not map the low-dimension embedding $\mathbf{z}_{i,LVG}$ alone to reconstruct
 585 $\hat{\mathbf{y}}_{i,LVG}$ in the original p_{LVG} -dimension space. Rather, we concatenate the HVG and LVG embeddings
 586 together, and feed the combined vector $[\mathbf{z}_{i,HVG} \ \mathbf{z}_{i,LVG}]$ into the decoder to denoise and batch correct
 587 LVG expression in the original p_{LVG} -dimension space. That is,

$$588 \quad \hat{\mathbf{y}}_{i,LVG} = f_{D,LVG}([\mathbf{z}_{i,HVG} \ \mathbf{z}_{i,LVG}]; W_{D,LVG}).$$

589
 590 This concatenated embedding is critical because it allows CarDEC to only use the high signal-to-noise
 591 ratio HVGs to drive the clustering loss, while still using the rich, batch corrected embedding that is

592 refined using this clustering loss to denoise and batch correct LVGs. The activation functions for the
 593 encoder and decoder of the LVGs are similarly defined as the autoencoder in Step 2.

594

595 To train this branching model, we first introduce two reconstruction losses, one for the HVGs and one
 596 for the LVGs computed as follows for cell i ,

597

$$l_{i,HVG} = \frac{1}{p_{HVG}} \|\hat{\mathbf{y}}_{i,HVG} - \mathbf{y}_{i,HVG}\|^2$$

$$l_{i,LVG} = \frac{1}{p_{LVG}} \|\hat{\mathbf{y}}_{i,LVG} - \mathbf{y}_{i,LVG}\|^2$$

598

599 Then the total loss is calculated as a multi-component loss function as follows,

600

$$l_i = \alpha l_{i,c} + (2 - \alpha) \frac{l_{i,HVG} + l_{i,LVG}}{2},$$

601

602 where α is a hyperparameter ranging from 0 to 2 that balances reconstruction loss with clustering loss.

603 We set α at 1 as default value. The total loss is minimized in an iterative fashion until certain

604 convergence criteria are satisfied.

605

606 **Step 4: denoising gene expression counts**

607 In Step 3, the denoised expression values obtained from the decoder are on a Z-score scale and are not

608 naturally comparable to raw UMI counts. To remedy this, we offer an optional downstream modeling

609 step that provides denoised expression values on the original count scale. This strategy involves finding

610 mean and dispersion parameters that maximize a negative binomial likelihood. We choose the negative

611 binomial distribution because previous studies have shown that UMI counts are not zero-inflated, and
 612 the negative binomial distribution fits the data well(Chen et al. 2018; Wang et al. 2018; Svensson 2020).

613

614 After the training in Step 3, we have obtained batch corrected low-dimension embeddings, $\mathbf{z}_{i,HVG}$ and
 615 $\mathbf{z}_{i,LVG}$ for each cell i from the fine-tuned HVG and LVG encoders. We will use two separate neural
 616 networks to maximize the negative binomial losses: one for the HVGs and one for the LVGs. These
 617 models are completely separate from one another but are trained almost identically with only minor
 618 differences. The goal is to map the embeddings into the full gene space to obtain mean and dispersion
 619 parameters for each gene. Without loss of generality, we use the HVGs as an example to illustrate how
 620 the neural network is built.

621

622 The vector of genewise means $\boldsymbol{\mu}_{i,HVG}$ and vector genewise dispersions $\boldsymbol{\theta}_{i,HVG}$ are given below,

623

$$\begin{aligned}\boldsymbol{\mu}_{i,HVG} &= s_i \times \exp(\mathbf{W}_{\mu,HVG} \times \tilde{\mathbf{z}}_{i,HVG}), \\ \boldsymbol{\theta}_{i,HVG} &= \text{softplus}(\mathbf{W}_{\theta,HVG} \times \tilde{\mathbf{z}}_{i,HVG}),\end{aligned}$$

624

625 where s_i is the size factor for cell i , $\mathbf{W}_{\mu,HVG}$ and $\mathbf{W}_{\theta,HVG}$ are trainable weight matrices, and \exp and
 626 softplus are activation functions that are applied elementwise. For each gene j in cell i , we compute
 627 the negative log likelihood of the negative binomial distribution as

628

$$l_{ij} = -\log\left(\frac{\Gamma(x_{ij} + \theta_{ij})}{\Gamma(\theta_{ij})} \left(\frac{\theta_{ij}}{\theta_{ij} + \mu_{ij}}\right)^{\theta_{ij}} \left(\frac{\mu_{ij}}{\theta_{ij} + \mu_{ij}}\right)^{x_{ij}}\right),$$

629

630 where x_{ij} is the original count in HVG gene j for cell i , and where μ_{ij} and θ_{ij} are the j^{th} elements of
631 $\boldsymbol{\mu}_{i,HVG}$ and $\boldsymbol{\theta}_{i,HVG}$, respectively. For the HVG count model, the full loss for cell i is then $l_i =$
632 $\frac{1}{p_{HVG}} \sum_{j=1}^{p_{HVG}} l_{ij}$. The loss for the LVG count model can be similarly defined. Both the HVG and LVG count
633 models are trained using their own early stopping and learning rate decay convergence monitoring.

634

635 **Evaluation of batch effect removal in the gene expression space and the embedding space**

636 Here, we briefly describe the workflow to evaluate batch effect removal and comparison between
637 different methods (details see Supplemental Note 2). First, we evaluated the performance of different
638 methods in removing batch effect in the gene expression space. For this evaluation, we considered
639 CarDEC, Scanorama, DCA with post-hoc batch effect correction by Combat, scVI, and MNN. We ran all
640 denoising/batch correction methods on the full data matrix and then present clustering results for
641 denoised HVGs and denoised LVGs separately by subsetting the HVGs and LVGs from the full
642 denoised/batch corrected expression matrix. The subsetting matrix that includes only the HVGs (or LVGs)
643 is then passed down to the Louvain's clustering algorithm. All steps in this workflow are identical for
644 both the HVGs and the LVGs, and all methods used the same HVGs and LVGs as input for clustering.
645 Furthermore, on a given dataset, we benchmarked all methods with the same number of clusters.
646 Second, we evaluated batch effect removal for the embedded representations of scRNA-seq. For this
647 evaluation we considered CarDEC, Scanorama, DCA with post-hoc batch effect correction by Combat,
648 scVI, and scDeepCluster. We excluded MNN since it has no embedding functionality. We included "raw"
649 as a control method for comparison, which is just subsetting the raw data to include only the HVGs, and
650 then running the clustering workflow.

651

652 **Coefficient of variation (CV) analysis**

653 To measure batch mixing, we examined the batchwise centroids before and after denoising. Let \mathbf{X}' be a
 654 matrix of gene expression counts (including both HVGs and LVGs). \mathbf{X}' can be the matrix of raw counts, or
 655 the denoised/batch corrected counts from any of CarDEC, Scanorama, DCA with post-hoc batch effect
 656 correction by Combat, scVI, or MNN. For CarDEC, we only considered denoised counts, not denoised
 657 expression in the Z-score space. If \mathbf{X}' consists of MNN corrected expression, then we did not preprocess
 658 the data since MNN denoised expression is on a cosine scale. In the case of MNN, Scanorama and DCA
 659 after Combat correction, a fraction of expression counts can be negative, which poses difficulties when
 660 computing coefficients of variation. To circumvent this issue, any expression values after correction by
 661 these methods that are negative were truncated to zero for the CV analysis. For all other methods we
 662 have denoised expression in the non-negative count space, so we performed cell normalization and log
 663 normalization on \mathbf{X}' , exactly in the same way described in the Step 1 (preprocessing) of CarDEC.

664
 665 Below we describe how the CV scores are calculated. Let x_{ij}' be the expression value in gene j of cell i in
 666 \mathbf{X}' . Let S_{ab} be a set of integers defined such that $i \in S_{ab}$ if and only if cell i belongs to cell type a and is
 667 sequenced from batch b . Furthermore, let $c_{ajb} = \sum_{i \in S_{ab}} x_{ij}' / \sum_{i \in S_{ab}} 1$ be the centroid (mean
 668 expression) of batch b for gene j of cell type a . Let $C_{aj} = \{c_{ajb}\}$ be the collection of batch centroids for
 669 gene j in cell type a . We can now define a **cell-type-specific CV** as follows:

670

$$CV_{aj} = \frac{\sqrt{\text{Var}(C_{aj})}}{\text{Mean}(C_{aj}) + \gamma}$$

671

672 where $\gamma = 10^{-12}$ is a small constant to mitigate computational instability. A higher value of CV_{aj}
 673 corresponds to greater variation among batches and less batch mixing for cell type a , so a good batch
 674 effect removal method should drive CV_{aj} closer to zero. Normalizing by mean expression adjusts the

675 metric for how highly expressed the gene is, so that CVs from more highly expressed genes are
 676 comparable to CVs from less highly expressed genes.

677

678 To consolidate our results across cell types, so that we can present a single distribution of CV scores
 679 across all cell types, we combine the cell-type-specific CV scores. Then we summarize the cell-type-
 680 specific CV scores by taking a weighted average of these scores across cell types, where the weight of a
 681 cell type is proportional to that cell type's frequency in the dataset. Specifically, the weight for cell type
 682 a is computed as $c(a) = \sum_i I(a_i = a) / \sum_i 1$, where a_i is the cell type of cell i . The combined CV score is
 683 calculated as follows:

684

$$CV_j = \sum_a c(a) CV_{aj}$$

685

686 **Evaluation metrics for clustering**

687 For all of our benchmark datasets, we used the cell type labels reported in the original papers as the
 688 gold standard. The clustering performance of each method was mainly evaluated using the adjusted
 689 rand index (ARI), calculated as below,

690

$$ARI = \frac{\sum_{ij} \binom{n_{ij}}{2} - [\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}] / \binom{n}{2}}{\frac{1}{2} [\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2}] - [\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}] / \binom{n}{2}},$$

691

692 where n_{ij} is the number of cells in both cluster i from the cluster assignments obtained when
 693 benchmarking and in cell type j according to the gold standard cell type labels. a_i is the total number of
 694 cells in cluster i from the cluster assignments obtained when benchmarking, b_j is the total number of

695 cells in cell type j according to the gold standard cell type labels from the original study, and n is the
 696 total number of cells. Additionally, we also computed normalized mutual information (NMI) and purity,
 697 calculated as below,

698

$$NMI = 2 \times \frac{\sum_{ij} \frac{n_{ij}}{n} \log \left(\frac{n \times n_{ij}}{a_i \times b_j} \right)}{\sum_i \frac{a_i}{n} \log \left(\frac{n}{a_i} \right) + \sum_j \frac{b_j}{n} \log \left(\frac{n}{b_j} \right)},$$

699

$$Purity = \frac{1}{n} \sum_i \max_j n_{ij}.$$

700

701 **Data sets**

702 We analyzed multiple published scRNA-seq datasets, which are available through the accession numbers
 703 reported in the original papers. 1) Human pancreatic islet data: CelSeq (Gene Expression Omnibus
 704 GSE81076), CelSeq2 (Gene Expression Omnibus GSE85241), Fluidigm C1 (Gene Expression Omnibus
 705 GSE86469), and Smart-seq2 (Array Express E-MTAB-5061); 2) Bipolar cells from mouse retina (Gene
 706 Expression Omnibus GSE81904); 3) Bipolar cells from macaque retina (Gene Expression Omnibus
 707 GSE118480); 4) mouse cortex data (Single Cell Portal SCP425); 5) human PBMC data (Single Cell Portal
 708 SCP424); 6) human monocyte data GEO (GSE146974); 7) human fetal liver data (Array Express E-MTAB-
 709 7407). Details of these datasets were described in Supplemental Table 2.

710

711 **Software availability**

712 An open-source implementation of the CarDEC algorithm can be downloaded from GitHub
 713 (<https://github.com/jlakkis/CarDEC>) and is available as Supplemental Code. Code to reproduce all of our

714 analyses can be found on GitHub (https://github.com/jlakkis/CarDEC_Codes) and is available as
715 Supplemental Code as well.

716 **Competing interest statement**

717 The authors declare no competing interests.

718

719 **Acknowledgments**

720 This work was supported by the following grants: R01GM125301 (to M.L.), R01EY030192 (to M.L.),
721 R01EY031209 (to M.L.), R01HL113147 (to M.L. and M.P.R.), and R01HL150359 (to M.L. and M.P.R.). We
722 thank Sean Simmons and Jiarui Ding for their help on the mouse cortex and human PBMC data analysis.
723 We also thank Romain Lopez and Adam Gayoso for letting us know the scVI update, which has greatly
724 improved its performance. *Author contributions:* This study was conceived of and led by M.L.. J.L.
725 designed the model and algorithm, implemented the CarDEC software, and led data analysis with input
726 from M.L. and X.L.. X.L. led data analysis for the human monocyte data and designed the workflow
727 figure. D.W., G.H., and K.W. participated the early stage of algorithm design and testing. Y. Z. provided
728 input on memory management and analysis for the human monocyte data. L. U. provided input on the
729 model and algorithm design. H.Z. and M.P.R. provided input on the human monocyte data analysis. J.L.
730 and M.L. wrote the paper with feedback from all coauthors.

731

732 **References**

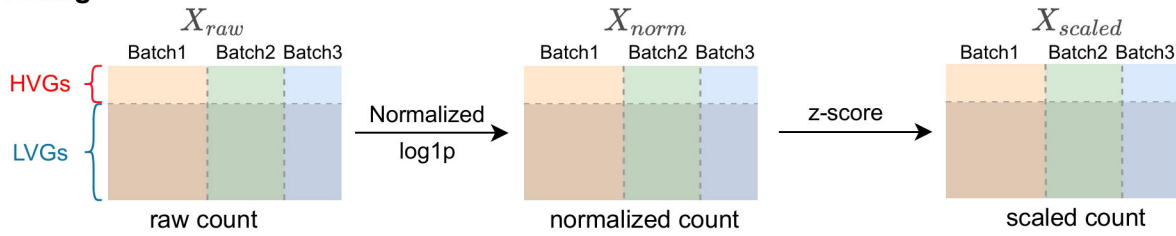
- 733 Barkas N, Petukhov V, Nikolaeva D, Lozinsky Y, Demharter S, Khodosevich K, Kharchenko PV. 2019. Joint
734 analysis of heterogeneous single-cell RNA-seq dataset collections. *Nat Methods* **16**: 695-698.
- 735 Cao J, Spielmann M, Qiu X, Huang X, Ibrahim DM, Hill AJ, Zhang F, Mundlos S, Christiansen L, Steemers FJ
736 et al. 2019. The single-cell transcriptional landscape of mammalian organogenesis. *Nature* **566**:
737 496-502.
- 738 Chen W, Li Y, Easton J, Finkelstein D, Wu G, Chen X. 2018. UMI-count modeling and differential
739 expression analysis for single-cell RNA sequencing. *Genome Biol* **19**: 70.

- 740 Ding J, Adiconis X, Simmons SK, Kowalczyk MS, Hession CC, Marjanovic ND, Hughes TK, Wadsworth MH,
741 Burks T, Nguyen LT et al. 2020. Systematic comparison of single-cell and single-nucleus RNA-
742 sequencing methods. *Nat Biotechnol* **38**: 737-746.
- 743 Eraslan G, Simon LM, Mircea M, Mueller NS, Theis FJ. 2019. Single-cell RNA-seq denoising using a deep
744 count autoencoder. *Nat Commun* **10**: 390.
- 745 Grun D, Muraro MJ, Boisset JC, Wiebrands K, Lyubimova A, Dharmadhikari G, van den Born M, van Es J,
746 Jansen E, Clevers H et al. 2016. De Novo Prediction of Stem Cell Identity using Single-Cell
747 Transcriptome Data. *Cell Stem Cell* **19**: 266-277.
- 748 Guo X, Gao L, Liu X, Yin J. 2017. Improved deep embedded clustering with local structure preservation.
749 *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*: 1753-
750 1759.
- 751 Haghverdi L, Lun ATL, Morgan MD, Marioni JC. 2018. Batch effects in single-cell RNA-sequencing data
752 are corrected by matching mutual nearest neighbors. *Nat Biotechnol* **36**: 421-427.
- 753 Hicks SC, Townes FW, Teng M, Irizarry RA. 2018. Missing data and technical variability in single-cell RNA-
754 sequencing experiments. *Biostatistics* **19**: 562-578.
- 755 Hie B, Bryson B, Berger B. 2019. Efficient integration of heterogeneous single-cell transcriptomes using
756 Scanorama. *Nature biotechnology* **37**: 685-691.
- 757 Johnson WE, Li C, Rabinovic A. 2007. Adjusting batch effects in microarray expression data using
758 empirical Bayes methods. *Biostatistics* **8**: 118-127.
- 759 Kingma DP, Ba JL. 2015. ADAM: a method for stochastic optimization. *International Conference on*
760 *Learning Representation*.
- 761 Korsunsky I, Millard N, Fan J, Slowikowski K, Zhang F, Wei K, Baglaenko Y, Brenner M, Loh PR,
762 Raychaudhuri S. 2019. Fast, sensitive and accurate integration of single-cell data with Harmony.
763 *Nat Methods* **16**: 1289-1296.
- 764 Lahnemann D, Koster J, Szczurek E, McCarthy DJ, Hicks SC, Robinson MD, Vallejos CA, Campbell KR,
765 Beerenwinkel N, Mahfouz A et al. 2020. Eleven grand challenges in single-cell data science.
766 *Genome Biol* **21**: 31.
- 767 Lawlor N, George J, Bolisetty M, Kursawe R, Sun L, Sivakamasundari V, Kycia I, Robson P, Stitzel ML.
768 2017. Single-cell transcriptomes identify human islet cell signatures and reveal cell-type-specific
769 expression changes in type 2 diabetes. *Genome Res* **27**: 208-222.
- 770 Li X, Wang K, Lyu Y, Pan H, Zhang J, Stambolian D, Susztak K, Reilly MP, Hu G, Li M. 2020. Deep learning
771 enables accurate clustering with batch effect removal in single-cell RNA-seq analysis. *Nat*
772 *Commun* **11**: 2338.
- 773 Lopez R, Regier J, Cole MB, Jordan MI, Yosef N. 2018. Deep generative modeling for single-cell
774 transcriptomics. *Nat Methods* **15**: 1053-1058.
- 775 Lucken MD, Buttner M, Chaichoompu K, Danese A, Interlandi M, Mueller MF, Strobl DC, Zappia L, Dugas
776 M, Colome-Tatche M et al. 2020. Benchmarking atlas-level data integration in single-cell
777 genomics. *bioRxiv* doi:<https://doi.org/10.1101/2020.05.22.111161>.
- 778 Muraro MJ, Dharmadhikari G, Grun D, Groen N, Dielen T, Jansen E, van Gurp L, Engelse MA, Carlotti F, de
779 Koning EJ et al. 2016. A Single-Cell Transcriptome Atlas of the Human Pancreas. *Cell Syst* **3**: 385-
780 394 e383.
- 781 Peng YR, Shekhar K, Yan W, Herrmann D, Sappington A, Bryman GS, van Zyl T, Do MTH, Regev A, Sanes
782 JR. 2019. Molecular Classification and Comparative Taxonomics of Foveal and Peripheral Cells in
783 Primate Retina. *Cell* **176**: 1222-1237 e1222.
- 784 Polanski K, Young MD, Miao Z, Meyer KB, Teichmann SA, Park JE. 2020. BBKNN: fast batch alignment of
785 single cell transcriptomes. *Bioinformatics* **36**: 964-965.

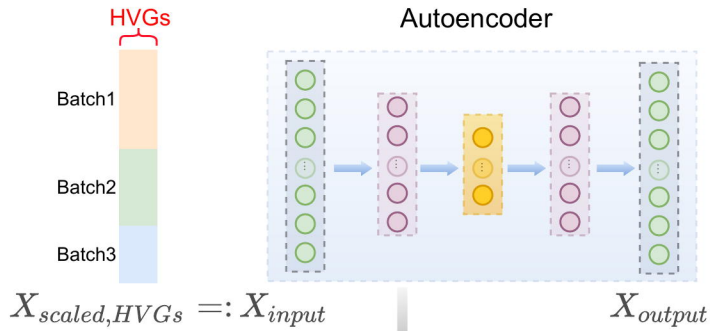
- 786 Popescu DM, Botting RA, Stephenson E, Green K, Webb S, Jardine L, Calderbank EF, Polanski K, Goh I,
787 Efremova M et al. 2019. Decoding human fetal liver haematopoiesis. *Nature* **574**: 365-371.
- 788 R Core Team (2020). R: A language and environment for statistical computing. R Foundation for
789 Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- 790 Segerstolpe A, Palasantza A, Eliasson P, Andersson EM, Andreasson AC, Sun X, Picelli S, Sabirsh A,
791 Clausen M, Bjursell MK et al. 2016. Single-Cell Transcriptome Profiling of Human Pancreatic
792 Islets in Health and Type 2 Diabetes. *Cell Metab* **24**: 593-607.
- 793 Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck III WM, Hao Y, Stoeckius M, Smibert P,
794 Satija R. 2019a. Comprehensive integration of single-cell data. *Cell* **177**: 1888-1902. e1821.
- 795 Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, 3rd, Hao Y, Stoeckius M, Smibert P,
796 Satija R. 2019b. Comprehensive Integration of Single-Cell Data. *Cell* **177**: 1888-1902 e1821.
- 797 Svensson V. 2020. Droplet scRNA-seq is not zero-inflated. *Nat Biotechnol* **38**: 147-150.
- 798 Tian T, Ji W, Qi S, Wei Z. 2019. Clustering single-cell RNA-seq data with a model-based deep learning
799 approach. *Nature Machine Intelligence* **1**: 191-198.
- 800 Wang J, Huang M, Torre E, Dueck H, Shaffer S, Murray J, Raj A, Li M, Zhang NR. 2018. Gene expression
801 distribution deconvolution in single-cell RNA sequencing. *Proc Natl Acad Sci U S A* **115**: E6437-
802 E6446.
- 803 Welch JD, Kozareva V, Ferreira A, Vanderburg C, Martin C, Macosko EZ. 2019. Single-Cell Multi-omic
804 Integration Compares and Contrasts Features of Brain Cell Identity. *Cell* **177**: 1873-1887 e1817.
- 805 Wolf FA, Angerer P, Theis FJ. 2018. SCANPY: large-scale single-cell gene expression data analysis.
806 *Genome biology* **19**: 15.
- 807 Wong KL, Tai JJ, Wong WC, Han H, Sem X, Yeap WH, Kourilsky P, Wong SC. 2011. Gene expression
808 profiling reveals the defining features of the classical, intermediate, and nonclassical human
809 monocyte subsets. *Blood* **118**: e16-31.
- 810 Xie J, Girshick R, Farhadi A. 2016. Unsupervised deep embedding for clustering analysis. *Proceedings of*
811 *International Conference on Machine Learning*: 478-487.
- 812

Workflow of CarDEC

Step 1: Preprocessing



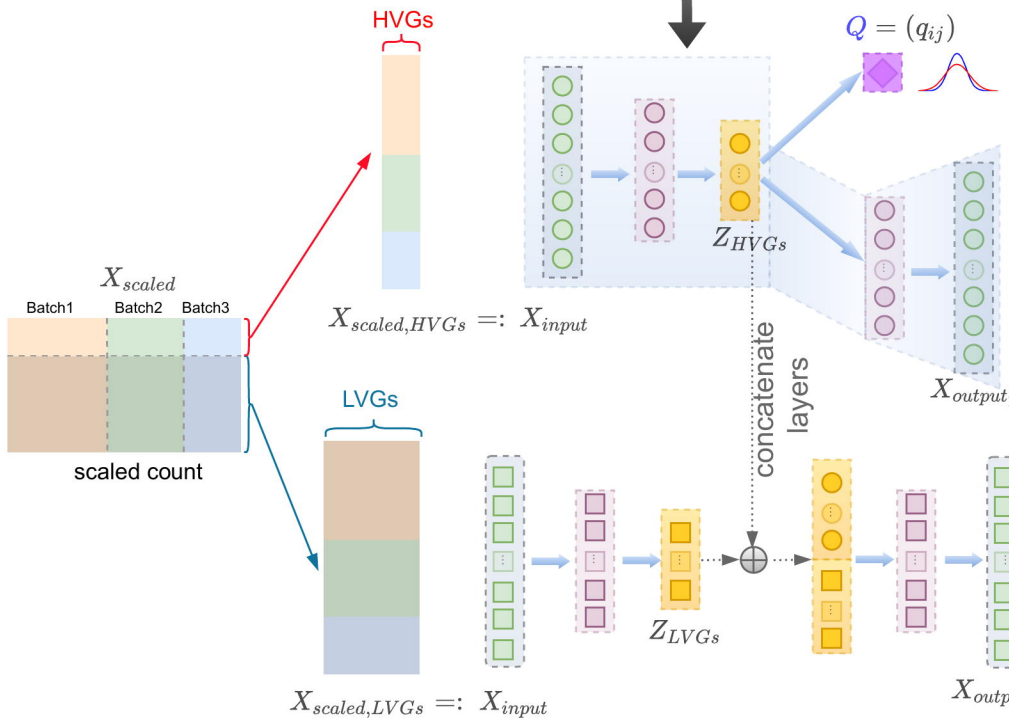
Step 2: Pretraining



Reconstruction loss with only HVGs:

$$\ell = \frac{1}{p_{HVGs}} \|X_{scaled,HVGs} - X_{output}\|^2$$

Step 3: Denoised z-score



Clustering loss:

$$\ell_c = KL(P||Q) = \sum_{i=1}^n \sum_{j=1}^K p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

Reconstruction loss with only HVGs:

$$\ell_{HVGs} = \frac{1}{p_{HVGs}} \|X_{scaled,HVGs} - X_{output_1}\|^2$$

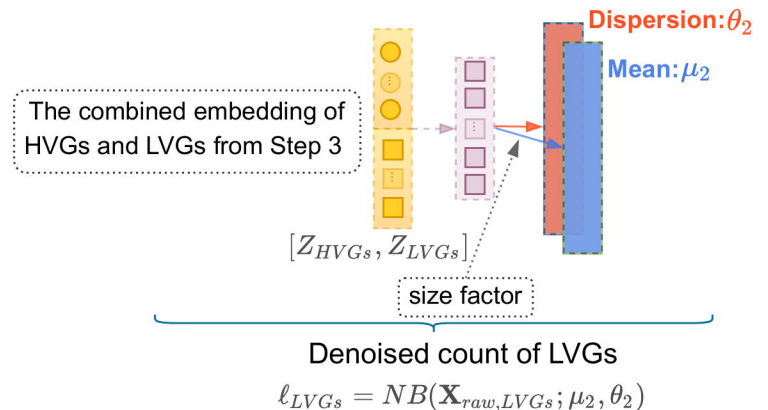
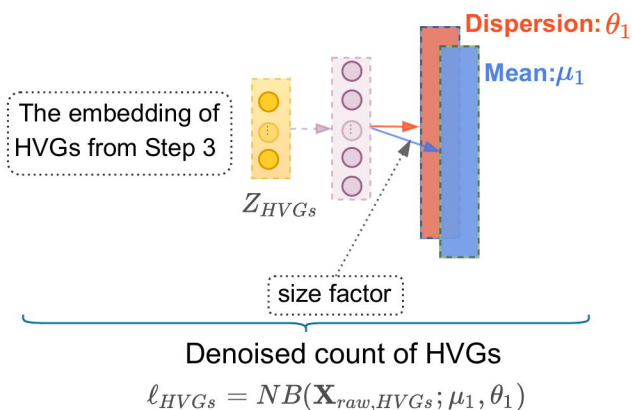
Reconstruction loss with only LVGs:

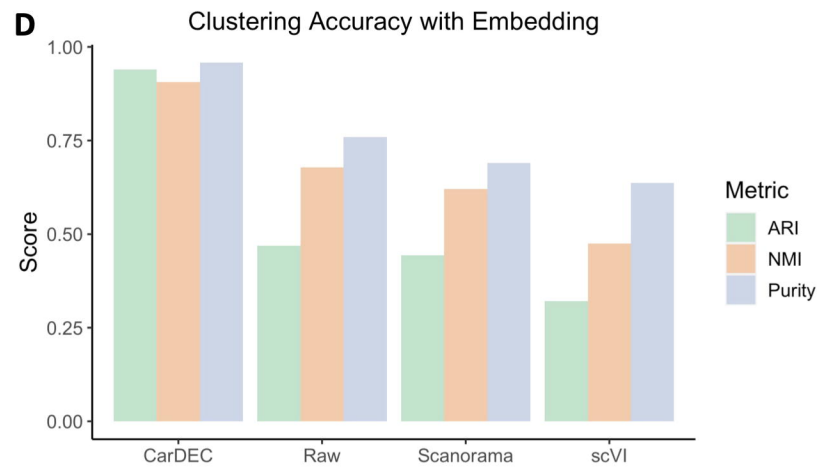
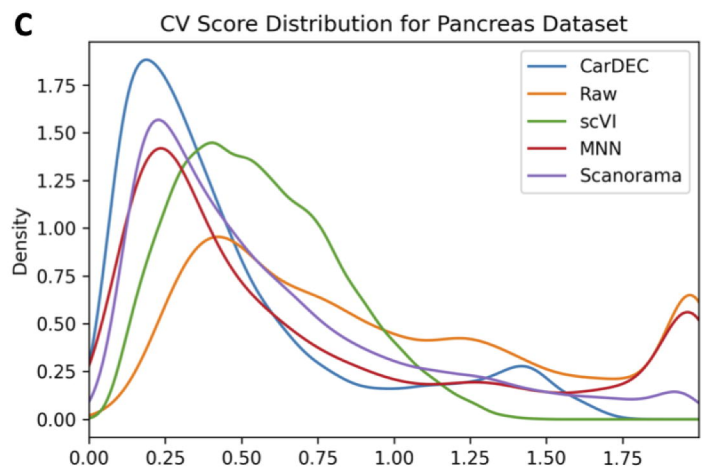
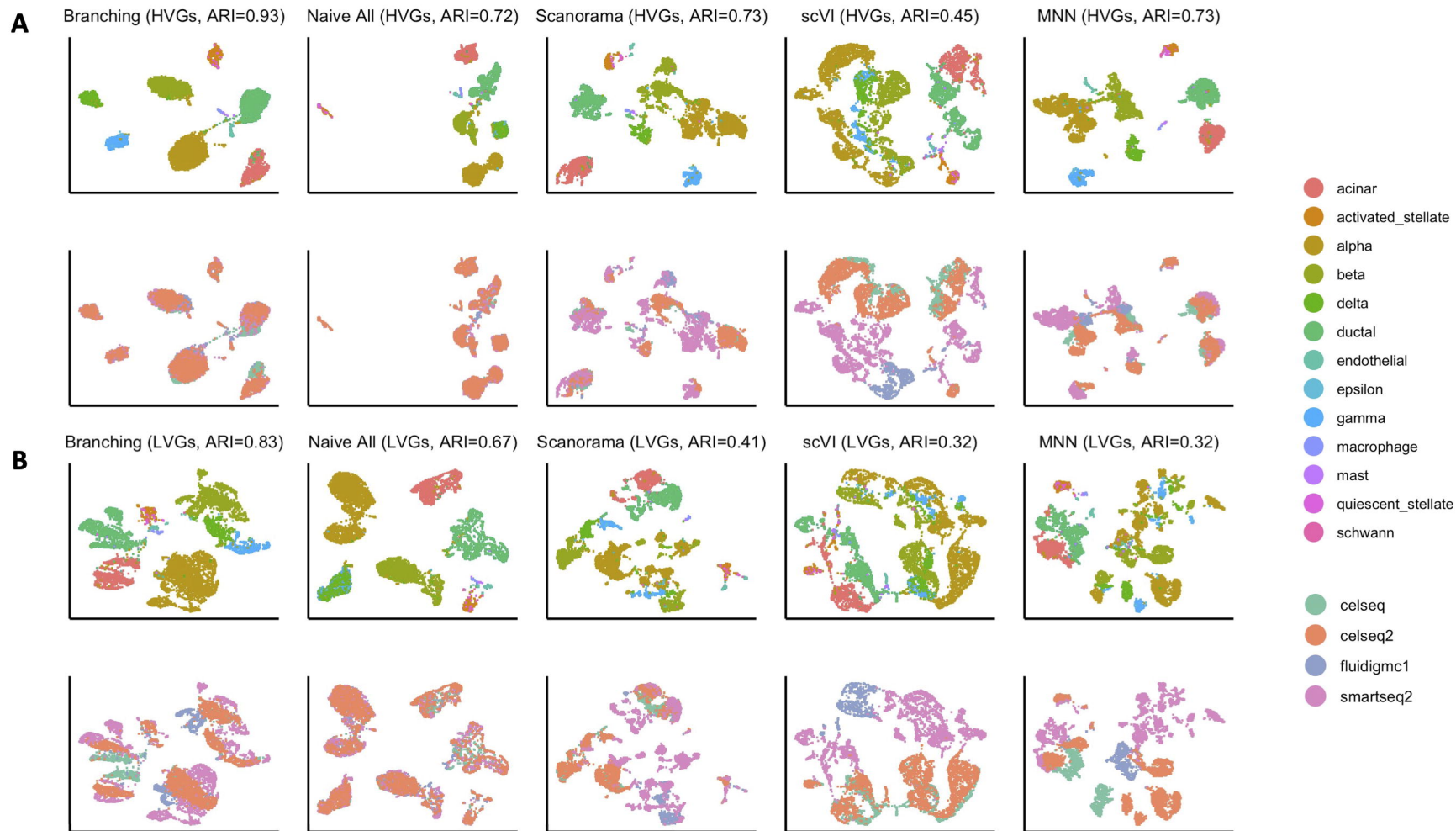
$$\ell_{LVGs} = \frac{1}{p_{LVGs}} \|X_{scaled,LVGs} - X_{output_2}\|^2$$

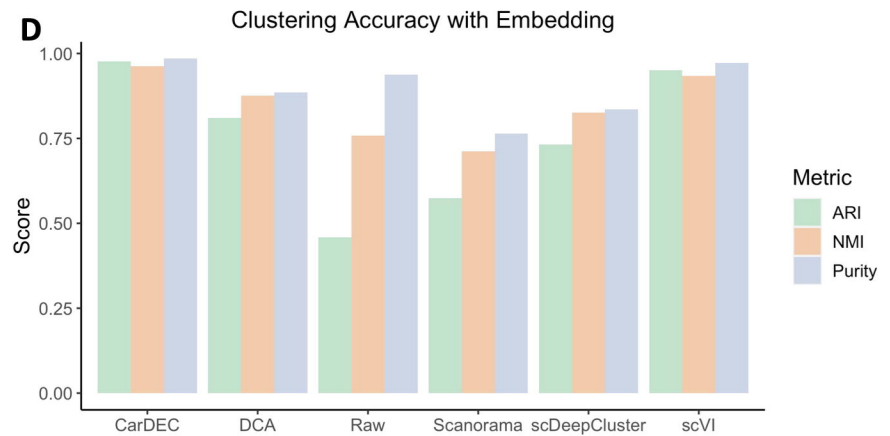
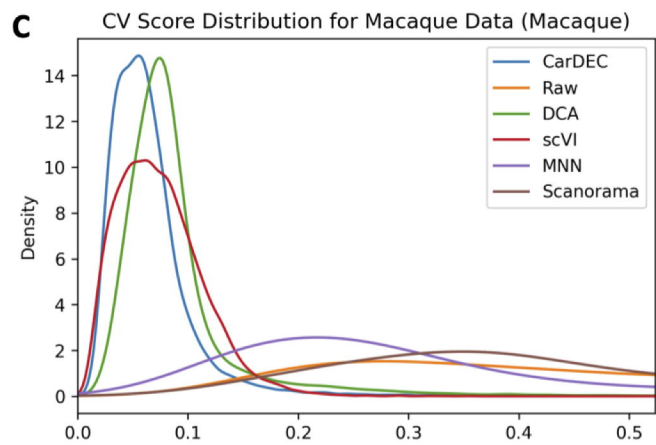
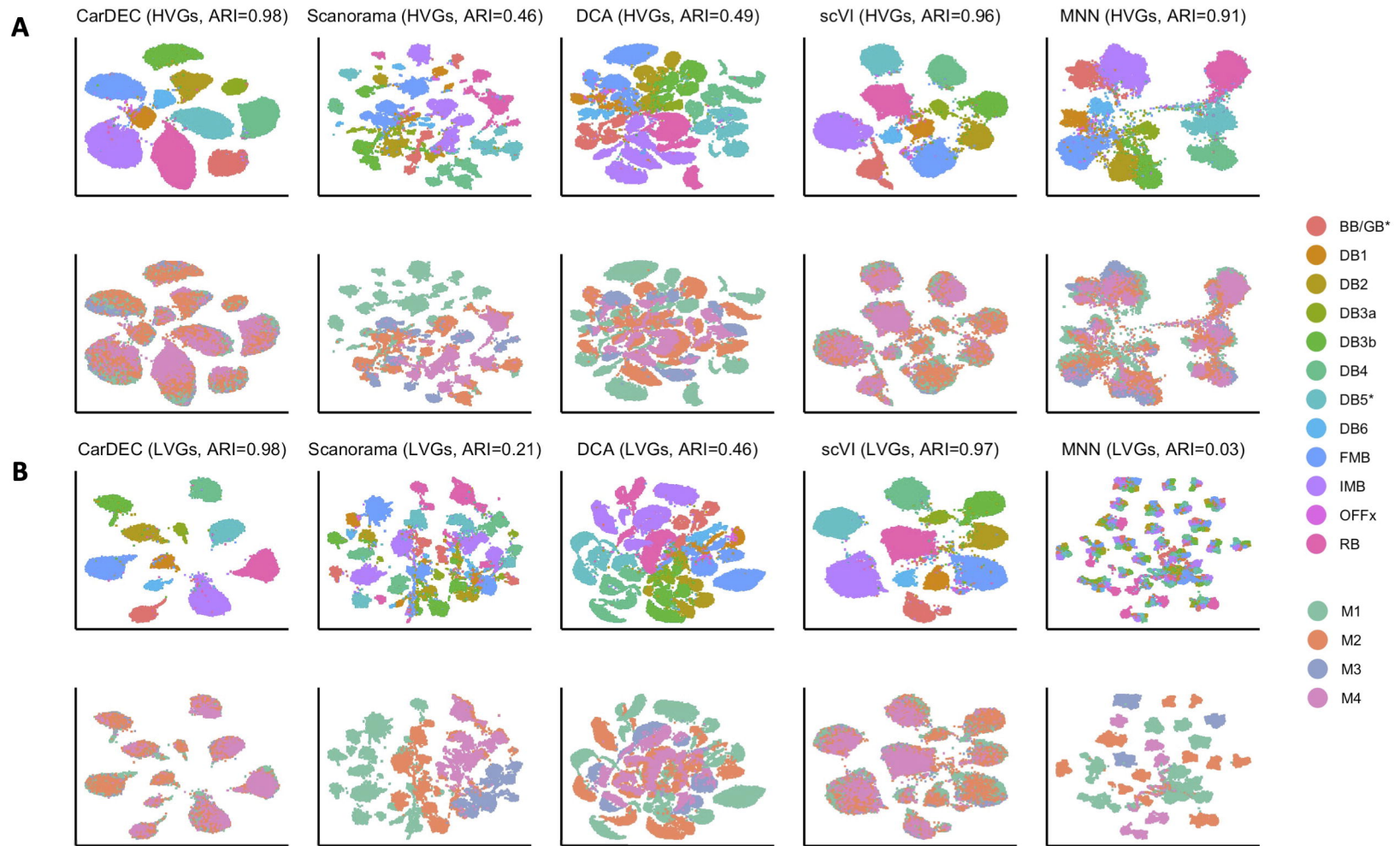
Total loss:

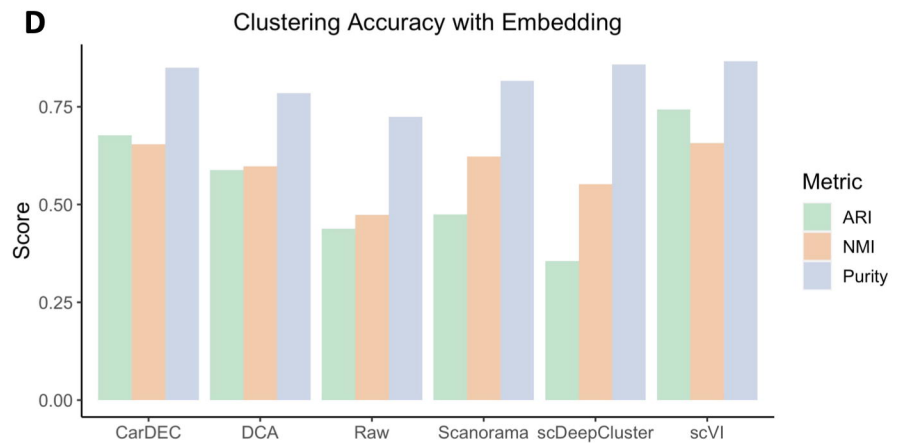
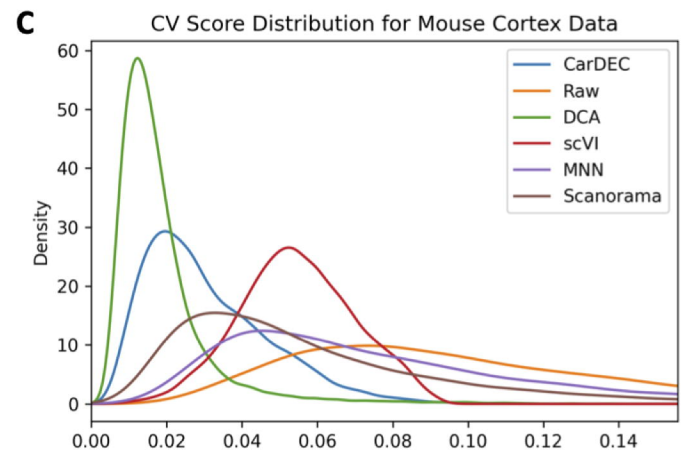
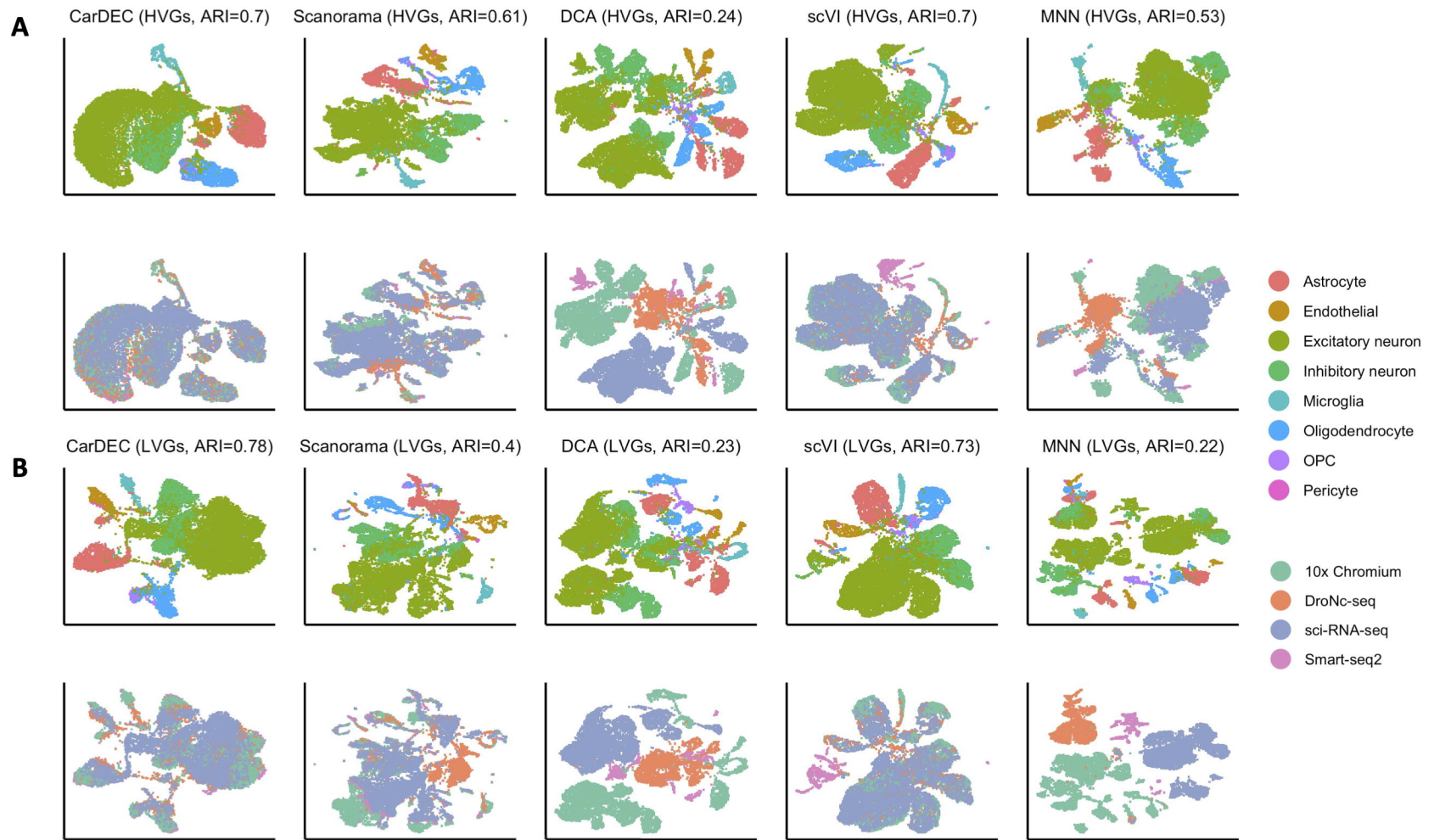
$$\ell_{total} = (2 - \lambda)\ell_c + \lambda(\ell_{HVGs} + \ell_{LVGs})/2$$

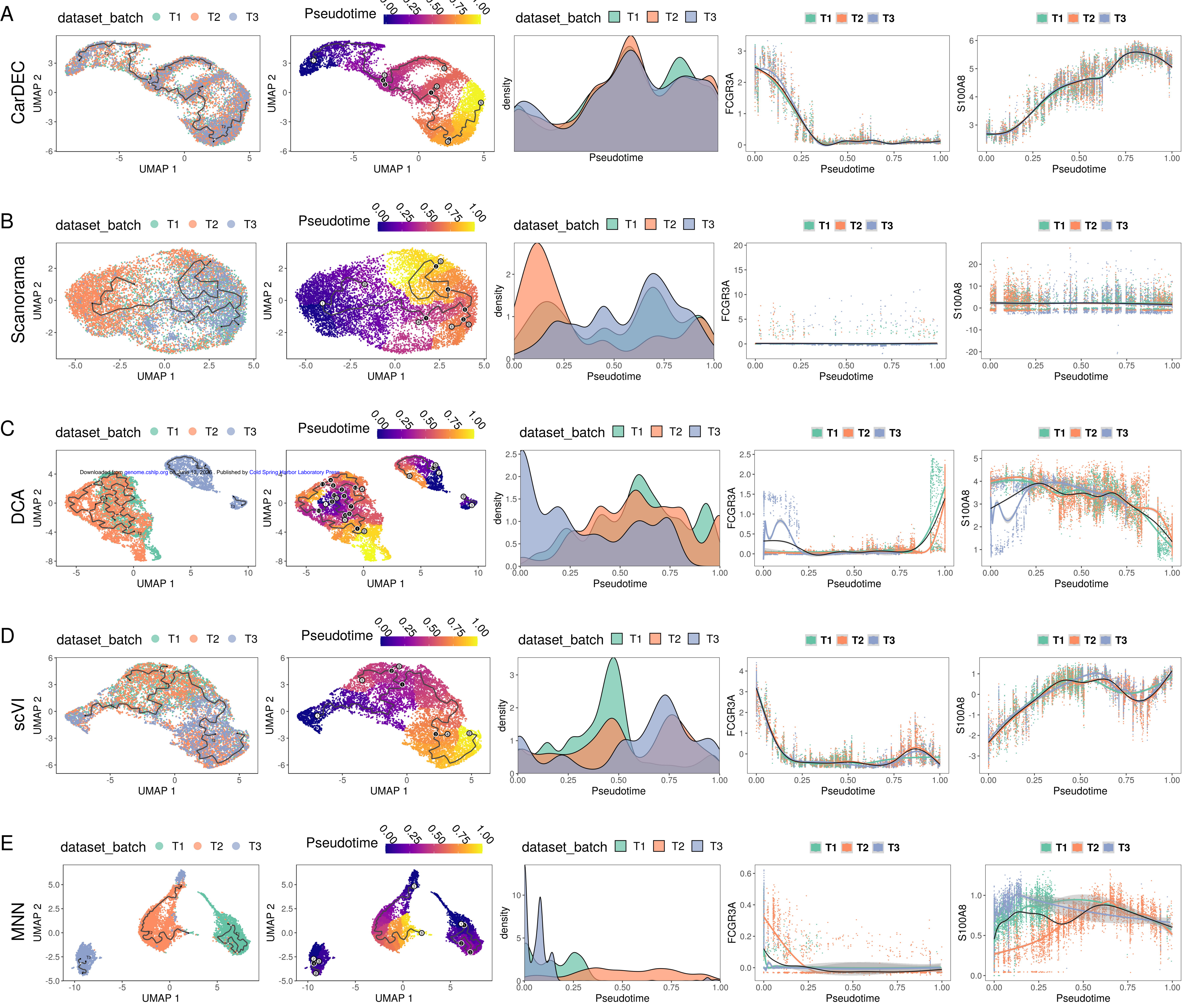
Step 4: Denoised raw count (optional)



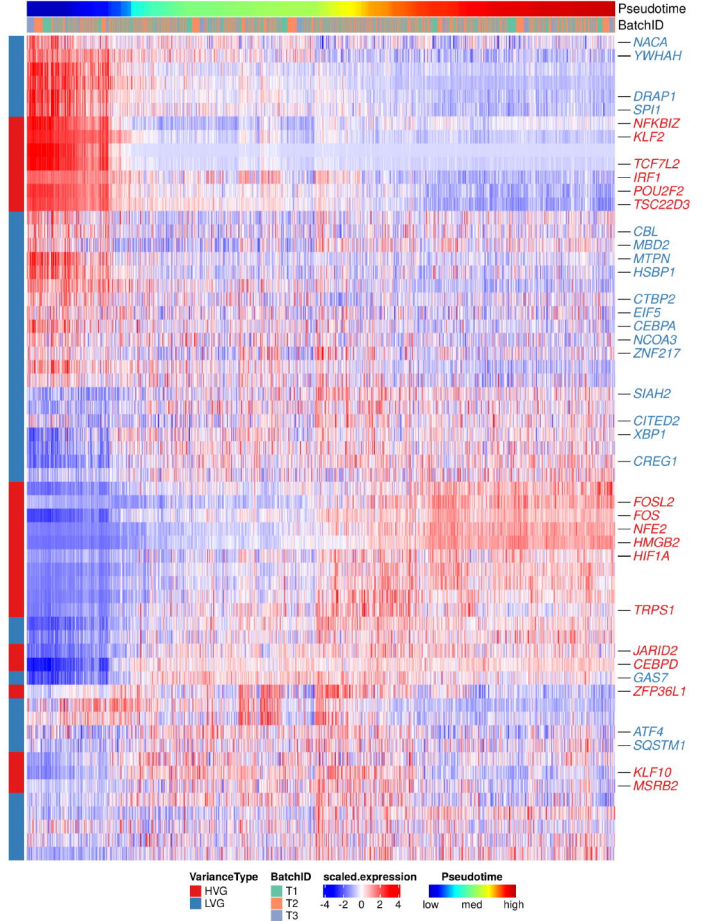




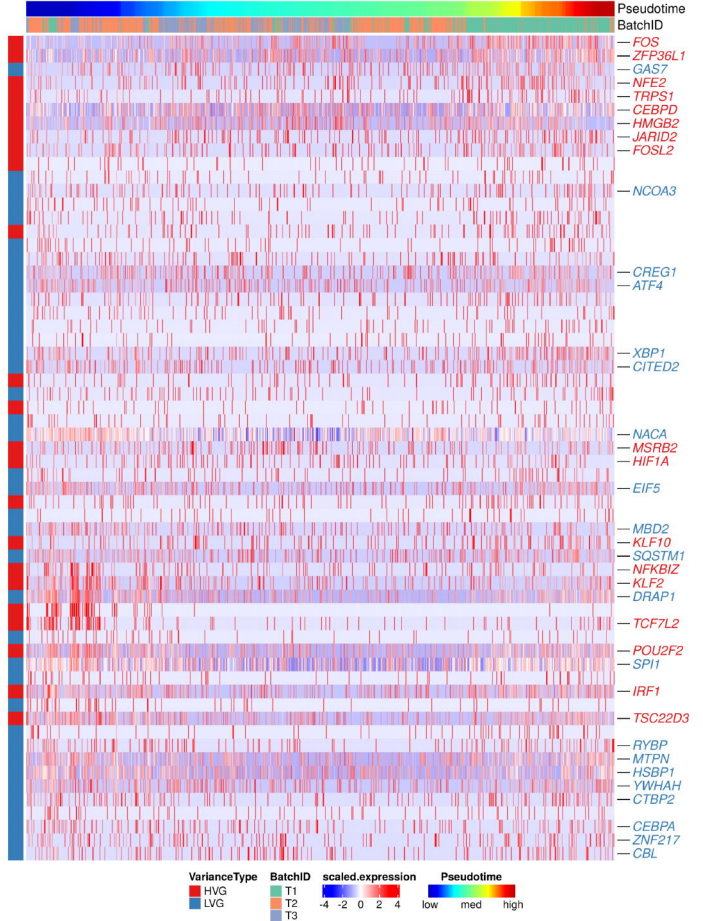




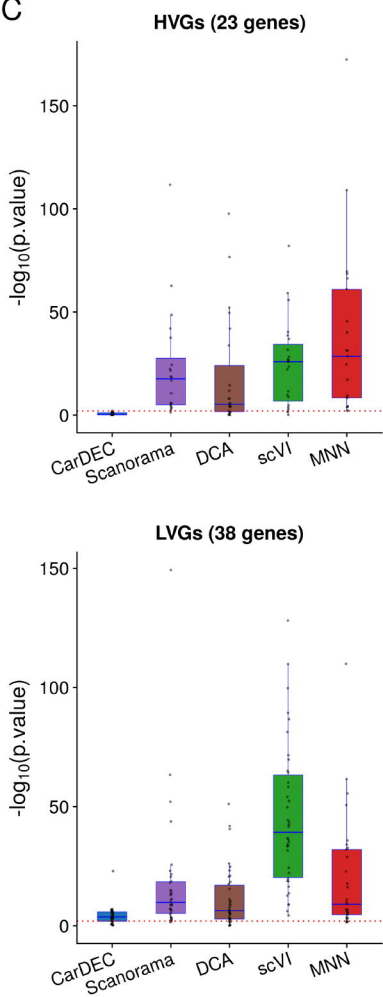
A Denoised gene expression (pseudotime by CarDEC)



B Raw gene expression (pseudotime by Raw)



C



D

