

DiSCERN - Deep Single Cell Expression ReconstructioN for improved cell clustering and cell subtype and state detection.

Fabian Hausmann^{a,b,1}, Can Ergen-Behr^{a,1}, Robin Khatri^{a,b}, Mohamed Marouf^a, Sonja Hänzelmann^{a,b}, Nicola Gagliani^{c,d,e}, Samuel Huber^{c,d}, Pierre Machart^{a,b,*}, Stefan Bonn^{a,b,*}

^a*Institute of Medical Systems Biology, University Medical Center Hamburg-Eppendorf, Hamburg, Germany.*

^b*Center for Biomedical AI, University Medical Center Hamburg-Eppendorf, Hamburg, Germany.*

^c*Section of Molecular Immunology and Gastroenterology, I. Department of Medicine, University Medical Center Hamburg-Eppendorf, 20246 Hamburg, Germany*

^d*Hamburg Center for Translational Immunology (HCTI), University Medical Center Hamburg-Eppendorf, 20246 Hamburg, Germany*

^e*Department of General, Visceral and Thoracic Surgery, University Medical Center Hamburg-Eppendorf, 20246 Hamburg, Germany*

Abstract

Single cell sequencing provides detailed insights into biological processes including cell differentiation and identity. While providing deep cell-specific information, the method suffers from technical constraints, most notably a limited number of expressed genes per cell, which leads to suboptimal clustering and cell type identification. Here we present DISCERN, a novel deep generative network that reconstructs missing single cell gene expression using a reference dataset. DISCERN outperforms competing algorithms in expression inference resulting in greatly improved cell clustering, cell type and activity detection, and insights into the cellular regulation of disease. We used DISCERN to detect two unseen COVID-19-associated T cell types, cytotoxic CD4⁺ and CD8⁺ Tc2 T helper cells, with a potential role in adverse disease outcome. We utilized T cell fraction information of patient blood to classify mild or severe COVID-19 with an AUROC of 81 % that can serve as a biomarker of disease stage. DISCERN can be easily integrated into existing single cell sequencing workflows and readily adapted to enhance various other biomedical data types.

Keywords: Single cell RNA-seq, RNA sequencing, imputation, cell clustering, cell type identification, expression reconstruction, Deep Learning, Machine

*Corresponding authors

Email addresses: pierre.machart@neclab.eu (Pierre Machart), sbonn@uke.de (Stefan Bonn)

¹Authors contributed equally

Learning, auto encoder, batch effect correction, transfer learning, probabilistic modeling, reference atlas mapping, COVID-19, T helper cell, transcription factor analysis, single nuclear RNA-seq
2010 MSC: 00-01, 99-00

1 Introduction

Single-cell RNA sequencing (scRNA-seq) technologies allow the dissection of gene expression at single-cell resolution, which improves the detection of known and novel cell types and the understanding of cell-specific molecular processes [1, 2]. The extension of the basic scRNA-seq technology with epitope sequencing of cell-surface protein levels (CITE-seq), allows for the simultaneous surveillance of the gene and protein surface expression of a cell [3]. Another recent technological innovation was TCR-seq, which enables the simultaneous sequencing of essential immune cell features and the variable segments of T cell antigen receptors (TCRs) that confer antigen specificity [4, 5].

While several commercial platforms have enabled researchers to use single cell sequencing methods with relative ease and at reasonable cost, the analysis of the high-dimensional scRNA-seq data still remains challenging [6, 7]. The main technical downside of single cell sequencing that impedes downstream analysis is the sparsity of gene expression information and high technical noise. Depending on the platform used, single cell sequencing detects around three thousand genes per cell, giving almost an order of magnitude less genes detected than bulk RNA-sequencing [8]. The term ‘dropout’ refers to genes that are expressed by a cell but cannot be observed in the corresponding scRNA-seq data, a technical artifact that afflicts predominantly lowly to medium expressed genes, as their transcript number is insufficient to reliably capture and amplify them. This missing expression information limits the resolution of downstream analyses, such as cell clustering, differential expression, marker gene and cell type identification [9].

To improve the lack and stochasticity of gene expression information in single cell experiments, several *in silico* gene imputation methods have been designed based on different principles. Gene imputation infers gene expression in a given cell type or state, based on the information from other biologically similar cells of the same dataset. Several methods utilizing this principle have been developed [10], amongst them DCA, MAGIC, scImpute, DeepImpute and CarDEC [11, 12, 13, 14, 15]. DCA is an autoencoder-based method for denoising and imputation of scRNA-seq data using a zero-inflated negative binomial model of the gene expression. MAGIC uses a nearest neighbor diffusion graph to impute gene expression and scImpute estimates gene expression and drop-out probabilities using linear regression. DeepImpute is an ensemble method, splitting the expression data into multiple pieces and trying to learn imputation of highly correlated genes using deep learning. CarDEC uses a two step procedure of imputation and batch correction using a neural network. All of these algorithms use information from similar cells with measured expression of the same dataset

40 for imputation. Another class of imputation algorithms use bulk RNA-seq data
41 to constrain scRNA-seq expression imputation. Bfimpute [16] uses Bayesian fac-
42 torization, SCRABBLE [17] matrix regularization, and SIMPLEs [18] a prior
43 distribution on the bulk data to impute scRNA-seq expression. Unfortunately,
44 SCRABBLE and Bfimpute do not scale beyond small single cell datasets and
45 few genes (3000 cells and genes in our hands), and SIMPLEs requires matching
46 single cell and bulk RNA-seq samples, severely constraining their usability.

47 Similarly, methods (e.g. multigrate[19]) were developed, which use scRNA-
48 seq in combination with complementary, matching data (e.g. CITE-seq, ATAC-
49 seq) to improve imputation. While complementary CITE-seq information is
50 available for many scRNA-seq datasets, other information such as ATAC-seq
51 data of the same sample is usually missing.

52 While current imputation methods provide improved gene expression infor-
53 mation, they still rely on the comparison of similar cells with largely absent gene
54 expression information, for example by using clustering approaches. Genes that
55 are not expressed in neighboring cells cannot be imputed, limiting the value of
56 classical imputation. In an ideal case, it would be possible to obtain information
57 of the expected true gene expression per cell, or at least expression information
58 with less technical noise, to reconstruct the true expression at single cell level.
59 Additionally, recent studies question the number of technical dropouts in UMI-
60 based sequencing technologies [20, 21] and thus challenge classical imputation
61 based methods. However, there are still batch specific changes, e.g. capture
62 rate of specific genes and differences in sample processing, which affect the sin-
63 gle cell data, beyond dropout. These changes can be wanted (enforced by the
64 experimental setup) or unwanted (stochastic changes in the experimental setup,
65 material).

66 Recent work has shown the effectiveness of deep generative models (e.g. Au-
67 toencoders and Generative Adversarial Networks) to infer realistic scRNA-seq
68 data and augment scarce cell populations using Generative Adversarial Net-
69 works [22] or the prediction of perturbation response using Autoencoders [23].
70 We hypothesized that a deep generative model could allow for the reconstruc-
71 tion of missing single cell gene expression information (low quality - lq) by
72 using related data with more genes expressed (high-quality - hq) as a reference,
73 a reference-based approach to gene expression inference (Figure 1A). In other
74 words, lq data with many missing gene expression values and bad clustering
75 could be transformed into data with few missing genes and improved clustering
76 if the “style” of a related hq dataset could be transferred to it. In the best case,
77 it would be possible to infer gene expression information for single cell data (lq)
78 by using purified bulk RNA-seq data (hq), obtaining over ten thousand genes
79 expressed per cell. We envision that this approach, when properly calibrated,
80 gains deep mechanistic insights into data beyond what is currently measurable.
81 It is important to note that the concept of using hq data to reconstruct gene
82 expression in lq data is different from classical imputation algorithms that infer
83 gene expression based on nearby cells from the same dataset, as outlined above.

84 Based on the above considerations, we developed DISCERN, a novel deep
85 generative neural network for directed single cell expression reconstruction. DIS-

86 CERN allows for the realistic reconstruction of gene expression information by
87 transferring the style of hq data onto lq data, in latent and gene space. Our ex-
88 periments on real and simulated data show that DISCERN outperforms several
89 existing algorithms in gene expression inference across a wide array of single
90 cell datasets and technologies, improving cell clustering, cell type and activity
91 detection, and pathway and gene regulation identification. To obtain deep in-
92 sights into the cellular changes underlying COVID-19, we reconstructed single
93 cell expression data of patient blood and lung immune data. While in our ini-
94 tial analysis [24] of blood data we detected few immune cell types, expression
95 reconstruction with DISCERN resulted in the detection of 28 cell types and
96 states in blood, including two unseen disease-associated T cell types, cytotoxic
97 CD4⁺ and CD8⁺ Tc2 T helper cells. Reconstructing a second COVID-19 blood
98 dataset with disease severity information, we were able to classify mild and se-
99 vere COVID-19 with an AUROC of 81 %, obtaining a potential biomarker of
100 disease stage. DISCERN can be easily integrated into existing workflows, as an
101 additional step after count mapping. Given that DISCERN is not limited by
102 a predefined distribution of data, we believe that it can be readily adapted to
103 enhance various other biomedical data types, especially other omics data such
104 as proteomics and spatial transcriptomics.

105 2. Results

106 2.1. The DISCERN algorithm for directed expression reconstruction

107 We aim to realistically reconstruct gene expression in scRNA-seq data by
108 using a related hq dataset. Ideally, this expression reconstruction algorithm
109 should meet several requirements [7]. First, it needs to be **precise** and model
110 gene expression values realistically. It shouldn't remove information of cellular
111 identity to form 'average cells' or collapse different cell types or states into one.
112 Second, the network should be **robust** to the presence of different cell types
113 in hq and lq data, or an imbalance in their relative ratios. It shouldn't, for
114 instance, 'hallucinate' hq-specific cells into the lq data. Lastly, the network
115 should be directional, as the user should be able to choose the target (reference)
116 dataset.

117 With these prerequisites in mind, we designed a deep neural network for
118 directed single cell expression reconstruction (DISCERN) (Figure S1B) that is
119 based on a modified Wasserstein Autoencoder [25]. A unique feature of DIS-
120 CERN is that it transfers the "style" of hq onto lq data to reconstruct missing
121 gene expression, which sets it apart from other batch correction methods such
122 as [26], which operate in a lower dimensional representation of the data (e.g.
123 PCA, CCA). To allow DISCERN to accurately reconstruct single cell RNA-
124 seq expression based on reference data, the structure of the network had to be
125 adapted in several ways. First, we implemented Conditional Layer Normaliza-
126 tion (CLN) [27, 28, 22] to allow for directed expression reconstruction of lq data
127 based on reference hq data (Figure S1B & S2). Second, we added two decoder
128 heads to the network to enable it to model dataset-specific dropout rates and

129 gene expression separately. Lastly, we extended DISCERN’s loss function with
130 a binary cross-entropy term for learning the probability of dropouts to increase
131 general inference fidelity. Further algorithmic details of DISCERN can be found
132 in the methods and Figure S1.

133 We first demonstrate DISCERN’s capabilities to faithfully reconstruct gene
134 expression using five pancreas single cell expression datasets from 5 different
135 studies [29, 30, 31, 32, 33], with varying quality (Tables S1 and S2). The pan-
136 creas data is widely used for benchmarking and it is ideal to evaluate expression
137 reconstruction for many cell types and sequencing technologies. We consider a
138 dataset as hq when the average number of genes detected per cell (GDC) (e.g.
139 smartseq2, GDC 6214) is much higher than in a comparable lq dataset (Ta-
140 ble S2). Conversely, a dataset is lq when the average cell has lower counts and
141 fewer genes expressed than a comparable hq dataset (e.g. indrop, GDC 1887).
142 Throughout this text, we will name sequencing technologies with capital (e.g.
143 Smart-Seq2, InDrop) and datasets with lower case first letters (smartseq2, in-
144 drop). We trained DISCERN on these five pancreatic single cell datasets and
145 assessed the integration of data in gene space and the expression reconstruction
146 per cell type. While uncorrected data cluster by batch and not by cell type,
147 DISCERN-integrated data show good batch mixing and clustering of cells by
148 cell type across all five datasets (Figure 1B & Figure S2). To get a clearer
149 picture of DISCERN’s expression reconstruction capabilities we next calculated
150 correlation coefficients of measured expression between the lowest quality in-
151 Drop and highest quality Smart-Seq2 data, before and after expression recon-
152 struction using DISCERN. The mean expression reconstruction of indrop-lq to
153 smartseq2-hq and smartseq2-hq to indrop-lq data is very accurate, showing a
154 Pearson correlation of $r = 0.95$, while mean expression correlation between un-
155 corrected indrop-lq and smartseq2-hq data is only $r = 0.77$ due to strong batch
156 effects (Figure 1C & D, Figures S3 and S4). The improved quality of indrop-
157 lq data reconstructed to smartseq2-hq level is validated by the strong increase
158 of genes expressed per cell, ranging from ≈ 2000 genes per cell in the uncor-
159 rected indrop-lq data to ≈ 6000 genes in the indrop-lq data after reconstruction
160 (Figure S5).

161 We next investigated the effect of reconstruction of three cell type-specific
162 genes, before and after correction across the five pancreas datasets (Figure S6).
163 Insulin expression in the pancreas should be largely restricted to beta cells [34],
164 which can be observed in the uncorrected smartseq2-hq and celseq2 datasets,
165 while the indrop-lq batch shows a diffuse pattern of insulin expression across
166 cell types (Figure S6A left panel). This diffuse insulin expression is corrected
167 by reconstructing the smartseq2-hq expression pattern from the indrop-lq data
168 (Figure S6A middle panel). In general, the expected specificity of insulin ex-
169 pression in beta cells can be recovered for all datasets when using DISCERN’s
170 reconstruction using the smartseq2-hq reference. Conversely, the reconstruction
171 from hq to the indrop-lq reference results in diffuse insulin expression across all
172 reconstructed datasets (Figure S6A right panel). We obtained similar results for
173 the pancreatic acinar cell-specific gene REG1A and the delta cell-specific gene
174 SST, both of which show diffuse expression across cell types in the uncorrected

175 inDrop data and cell-specific expression after reconstruction using smartseq2-hq
176 reference (Figure S6B & C). Interestingly, DISCERN can not only recover bio-
177 logical expression information, but it is also able to apply sequencing method-
178 specific effects after reconstruction. The smartseq2-hq dataset, for instance,
179 displays nearly no ribosomal protein coding gene expression after sequencing as
180 previously reported by [8], while data sequenced using InDrop, Cel-Seq, or Cel-
181 Seq2 shows prominent ribosomal protein coding gene expression (Figure S6D,
182 left panel). When reconstructing smartseq2-hq data to indrop-lq data, riboso-
183 mal protein coding gene expression is re-instantiated (Figure S6D, right panel).

184 We further corroborated DISCERN's capability to integrate and reconstruct
185 gene expression in the more complex difftec dataset (Tables S1 and S2), consist-
186 ing of 14 single cell peripheral blood mononuclear cell (PBMC) datasets across a
187 wide range of technologies. Similar to pancreas, the difftec dataset is widely used
188 for benchmarking and it is ideal to evaluate expression reconstruction for even
189 more cell types and sequencing technologies. The different single cell technolo-
190 gies show large variation in quality, with an GDC ranging from 422 in Seq-Well
191 to 2795 in Smart-seq2. We trained DISCERN on these 14 PBMC single cell
192 datasets and observed very good integration in gene space (Figure S7). We
193 then reconstructed chromium-v2-lq (GDC 795) using a chromium-v3-hq refer-
194 ence (GDC 1514) and observed high mean gene expression correlation between
195 the reconstructed and reference datasets (Figures S8 and S9). These results
196 across 19 single cell datasets provide first evidence for the high-quality data in-
197 tegration and expression reconstruction that can be obtained with DISCERN.

198 *2.2. Specific and robust gene expression inference*

199 We next investigated the precision and robustness of DISCERN's expression
200 reconstruction in more detail and compared DISCERN's performance to several
201 state-of-the-art algorithms for expression imputation and data integration.

202 We explored the robustness of DISCERN to the choice of its hyperparameter
203 by testing various non-default combinations of the four hyperparameters influ-
204 encing the model training. In all combinations DISCERN was able to achieve a
205 pearson correlation of > 0.94 and a correlation of 0.95 with the default param-
206 eter when reconstructing the indrop-lq batch to the smartseq-hq batch of the
207 pancreas dataset (Figure S10). This provides strong evidence that DISCERN's
208 performance is robust to the choice of hyperparameters.

209 Since expression reconstruction can be seen as a generalization of expression
210 imputation, we compared DISCERN to DCA, MAGIC, and scImpute, CarDEC,
211 and DeepImpute, five state-of-the-art imputation algorithms [11, 12, 13, 14, 15].
212 Expression reconstruction can also be viewed as a batch correction task in gene
213 space, which is why we additionally compared DISCERN to scGEN, Seurat, tr-
214 VAE and scVI [23, 26, 35, 36]. It is important to note, however, that these batch
215 correction methods were not designed for the expression reconstruction task and
216 use a lower dimensional representation to align different batches. Seurat uses
217 canonical correlation analysis and scGEN uses the bottleneck layer representa-
218 tion of an autoencoder to calculate and apply linear transformations. trVAE

219 and scVI explicitly encode the conditional information in the autoencoder ar-
220 chitecture.

221 We compared the ability of these models to adjust expression information
222 on the pancreas dataset by reconstructing the indrop-lq expression based on
223 the smartseq2-hq expression. Generally deep learning methods, which allow for
224 projection (scGEN, scVI, trVAE, DISCERN), show the best performance, with
225 DISCERN showing the lowest deviation between cell types (Figure S11). We
226 also investigated the gene expression standard deviation on the same data, show-
227 ing that DISCERN reconstructs the variation in the indrop-lq best, with scVI
228 showing only slightly worse performance (Figure S12). A factor which has a high
229 impact on the variation is the number of dropouts found in each gene. While
230 most imputation methods try to remove them, we think they contain useful in-
231 formation as well [37]. DISCERN is able to capture the batch-specific dropout
232 rate much better compared to other batch correction or imputation methods
233 (Figure S13). Interestingly deep learning methods, scVI, scGEN, DeepImpute
234 and DCA for example, achieve a similar correlation of the dropout rate than
235 classical methods, for example Seurat and MAGIC, even if deep learning meth-
236 ods seem to be better in reconstruction of mean expression (Figure S11). It is
237 important to highlight that the proper estimation of expression variation and the
238 dropout rate is pivotal for the reliable computation of differentially expressed
239 genes. Since DISCERN displays the best variance estimation, it also achieves
240 the best median correlation of the differentially expressed genes (Figure S14).

241 To investigate the precision of gene expression reconstruction, we created an
242 artificial dataset by dividing the smartseq2-hq pancreas data into two batches,
243 smartseq-lq and smartseq2-hq. In the smartseq-lq batch, the top one KEGG
244 pathways per cell type were removed by setting the expression of genes con-
245 tained in these pathways to zero, while the smartseq2-hq remained unaltered.
246 Therefore, a reconstruction of smartseq-lq data using smartseq2-hq reference
247 (reconstructed-hq) should ideally recover the smartseq-lq expression to its orig-
248 inal state, prior to the removal of the genes. DISCERN is able to reconstruct
249 the mean expression for all cell types, achieving a correlation $r = 0.99$ (Fig-
250 ure 2A). DCA ($r = 0.66$), MAGIC ($r = 0.34$), scImpute ($r = 0.80$), Deep-
251 Impute $r = 0.89$ and Seurat ($r = 0.76$) have significantly lower correlation
252 between the smartseq2-hq and reconstructed-hq gene expression (Figure 2A).
253 scGen ($r = 0.98$), scVI ($r = 0.99$) and trVAE ($r = 0.99$) show similar perfor-
254 mance compared to DISCERN. Moreover, scGEN and trVAE however perform
255 worse in reconstruction of highly expressed genes, while scVI slightly overes-
256 timates the expression in general (Figure 2A). We obtained similar results on
257 the difftec dataset, with DISCERN ($r = 0.98$) outperforming DCA ($r = 0.47$),
258 MAGIC ($r = 0.21$), scImpute ($r = 0.04$), Seurat ($r = 0.92$), scVI ($r = 0.96$),
259 trVAE ($r = 0.95$), DeepImpute ($r = 0.58$), and scGEN ($r = 0.94$) (Figure S15).
260 To further investigate gene expression reconstruction specificity, we compared
261 the correlation of reconstructed-hq to smartseq2-hq data after performing dif-
262 ferential gene expression (DEG) for each cell type against all other cell types
263 (Figure 2B, upper panel). DISCERN is able to recover the correct DEG t-
264 statistics with a median correlation of 0.92, improving over state-of-the-art tools

265 by more than 6 percentage points. In the corresponding experiment using the
266 difftec dataset, DISCERN achieves a median correlation of 0.86, which is a 21
267 percentage point improvement over competing methods (Figure S16).

268 Since the genes were initially selected using KEGG gene set enrichment
269 analysis, the reconstruction of the corresponding pathways was investigated by
270 performing KEGG gene set enrichment analysis on the DEG results. DISCERN
271 is able to recover the pathway expression enrichment scores with a median cor-
272 relation of 0.88, exceeding the performance of scVI by more than 3 percentage
273 points on median (Figure 2B, lower panel). In the corresponding experiment
274 using the difftec dataset, DISCERN achieves a median correlation of 0.77, out-
275 performing Seurat and scGen by more than 16 percentage points (Figure S17).

276 While DISCERN outperforms competing algorithms in expression and path-
277 way reconstruction correlation, it achieves the fourth-best correlation for the
278 DEG fold-change (FC) of reconstructed-hq to smartseq2-hq data for the pan-
279 creas (Figure S18) and reconstructed-hq to chromium-v3-hq difftec datasets
280 (Figure S19). In both cases Seurat, scVI and CarDEC achieve better correla-
281 tion, which is due to the fact that DISCERN slightly underestimates FC in
282 favor of superior DEG variance estimation.

283 Next, we show DISCERN's expression reconstruction robustness with re-
284 spect to varying sizes of lq to hq data. It is conceivable to assume that a large
285 amount of hq data would benefit the expression reconstruction of the lq data,
286 which makes it important to understand at what ratio good results can be ex-
287 pected. Interestingly, DISCERN seems to be very robust across a wide range
288 of smartseq2-lq to smartseq2-hq ratios, with correlations of 0.98 (ratio of lq/hq
289 0.14) to 0.93 (ratio of lq/hq 18.4), while the second-best performing algorithm
290 scGen showed a 11 percentage point decrease in performance (0.82 for ratio of
291 lq/hq 18.4) (Figure 2C, Figure S20). We observed similar results for the correla-
292 tion of t-statistics, showing a slight dependence of DISCERN's performance
293 on the lq/hq ratio (Figure S21). In general, all methods show better perfor-
294 mance with a small ratio of lq/hq data, while DISCERN and scVI shows least
295 dependence and outperform other algorithms in the correlation of expression
296 and t-statistics, especially in the case of high lq/hq ratio.

297 Another aspect of expression reconstruction robustness is the dependence of
298 the algorithm on the cell type or cell state similarity of the lq and hq datasets.
299 In the optimal case, DISCERN would not require that the lq and hq datasets
300 have overlapping cell types to perform an accurate expression reconstruction,
301 which is theoretically possible if the network learns the general gene-regulatory
302 expression logic of the hq data (see discussion). To understand the dependence
303 on dataset similarity, we removed a complete cell type, pancreas alpha cells,
304 from the smartseq2-hq data and left the alpha cells in the smartseq2-lq data.
305 We then additionally varied the number of common cells in the lq and hq data,
306 starting with no overlapping cells (only alpha cells in the lq and all cells except
307 alpha in the hq data) and ending with almost complete overlap (all cells overlap
308 between the smartseq2-hq and -lq data, except for the alpha cells only present
309 in lq data) (Figure 2D). When evaluating DEG correlation, DISCERN was
310 the only method consistently achieving better performance than uncorrected

311 data, outperforming scVI by 2 to 17 percentage points (Figure 2D). Similarly,
312 DISCERN was consistently achieving better performance than uncorrected data
313 in the FC correlation task (Figure S22).

314 We next took a closer look at the integration and expression reconstruction
315 performance when no cell types overlap between the lq (alpha cells only) and
316 hq (all other cells) data. Notably, Seurat seems to over-integrate cell types,
317 mixing smartseq2-hq beta and gamma cells with reconstructed-hq alpha cells
318 from other batches (Figure S23), while all other methods keep the smartseq2-hq
319 and reconstructed-hq exclusive cell types separate (Figure 2E & Figure S23).
320 This over-integration seems to be causal for Seurat's poor DEG correlation per-
321 formance ($r = 0.19$), while DISCERN ($r = 0.55$) is the only method achieving
322 better performance than uncorrected cells ($r = 0.52$) (Figure 2F). Thus, DIS-
323 CERN is able to keep existing expression correlations and improves the detec-
324 tion of cell type specific genes by reconstruction using an hq batch as reference.
325 In conclusion, DISCERN is both a precise and robust method for expression
326 reconstruction that outperforms existing methods by a significant margin.

327 *2.3. Improving cell cluster, type, and trajectory identification*

328 The comparison to competing methods provided evidence for DISCERN's
329 superior expression reconstruction. Now, we will delineate how DISCERN's
330 expression reconstruction improves downstream cell clustering, cell type and
331 activity state identification, marker gene determination, and gene regulatory
332 network and cell trajectory analysis.

333 Batch correction algorithms are usually evaluated by comparing their ability
334 to integrate cells coming from the same cell type but different batches, using
335 the silhouette score, the adjusted rand index (ARI), and adjusted mutual infor-
336 mation (AMI). DISCERN often outperforms all competing methods across all
337 metrics, achieving state-of-the-art performance in batch mixing and cell type
338 clustering (Figures S24 to S26).

339 To understand if cell-determining gene expression and pathways could be
340 recovered with expression reconstruction, we used a single nuclear sequencing
341 (sn-lq) and scRNA-seq (sc-hq) data pair that was prepared from the same liver
342 metastasis biopsy [38]. We reconstructed sn-lq data using the sc-hq reference,
343 obtaining reconstructed-hq data. While single nuclear sequencing provides re-
344 duced expression information in the average counts per cell as compared to
345 scRNA-seq (Table S2) [38], it is still the method of choice to obtain cell-specific
346 expression information when intact single cells cannot be recovered from a tis-
347 sue (e.g. after tissue fixation or freezing). It is important to note that nuclear
348 transcripts reflect current gene activity, which in part might not correlate with
349 transcripts that have lifetimes of up to days. Before integration, the sn-lq and
350 sc-hq datasets cluster by batch and not by cell type, while after expression
351 reconstruction with DISCERN cells cluster by type and not by batch (Fig-
352 ure S27). This is reflected in an expression correlation of 0.49 (sc-hq vs. sn-lq)
353 before and 0.97 after reconstruction (sc-hq vs. reconstructed-hq) (Figure S28).
354 DISCERN reconstruction resulted in the expression of T cell receptor signaling
355 genes in reconstructed T cells (Figure S29) and antigen presentation genes in

356 macrophages (Figure S30), providing evidence that DISCERN faithfully recreates
357 cell-determining genes and pathways based on the hq data. Seurat, CarDEC
358 and scImpute are not able to reconstruct the expression information and show
359 a similar expression pattern as the uncorrected sn-lq dataset. In their recon-
360 structions (seurat-hq, CarDEC-hq, scImpute-hq and sn-lq) the expression of
361 important T cell marker genes such as *CD3E*, *CD3D* and *CD8A* is largely ab-
362 sent, while in sc-hq and DISCERN-hq the expression is easily detectable (Fig-
363 ure S29). DCA, scVI,scGEN, MAGIC and trVAE show a strongly disturbed
364 expression pattern, where many genes show a much larger expression than in
365 the sc-hq or the sn-lq datasets (Figures S29 and S30).

366 To further corroborate the advantage of single nuclear expression recon-
367 struction, we next aimed to increase the T cell subtype resolution of human
368 single nucleus acute kidney injury data (kidney-lq) by using matching single
369 cell data (kidney-hq). Only 1% of kidney-lq nuclei show *CD3D*, *CD3E* or
370 *CD3G* expression, compared to 7% of the cells in the kidney-hq dataset. Seu-
371 rat and DISCERN were able to detect T cells in the reconstructed kidney-lq
372 (reconstructed-hq) and the kidney-hq data with notable *CD3D* expression in
373 this cluster (Figure S31). The reconstructed-hq and the kidney-hq T cells were
374 further classified into T cell subtypes and activation states (Figure S31C). While
375 a large proportion of T cells detected in Seurat reconstructed data could not
376 be annotated due to missing *CD3D*, *CD4*, and *CD8A* expression, DISCERN
377 reconstructed data does not present these limitations.

378 It is intriguing to observe that many marker genes are still hard to detect in
379 kidney single cell RNA-seq data but also in the antigen presentation pathway
380 in macrophages (Figure S30). This is most probably due to dropout. Thus, we
381 rationalized that bulk RNA sequencing (RNA-seq) data of purified cell types
382 (e.g. FACS sorted immune cells) is a suitable hq proxy for the expected gene
383 expression per cell. RNA-seq data of purified cells is readily available from
384 public repositories, making it possible to obtain thousands of purified immune
385 cell RNA-seq samples (see methods). We therefore set out to increase cluster,
386 cell type, gene regulatory network, and trajectory identification of scRNA-seq
387 data by reconstructing gene expression using a related RNA-seq reference (Fig-
388 ure S32). For the scRNA-seq data we chose a cord blood mononuclear cite-
389 seq dataset (cite-lq) that was labeled with 15 antibodies (Table S3) to allow
390 for surface protein-based cell type discovery [39]. The CITE-seq information
391 allowed us to confirm expression reconstruction by DISCERN in cases where
392 gene expression is absent but protein expression and cell identity are validated
393 via antibody labeling. For the RNA-seq data, we selected 9852 purified immune
394 samples (bulk-hq) and proceeded to reconstruct cite-lq (GDC 798) using a bulk-
395 hq (GDC 13 104) reference to obtain reconstructed-hq data with DISCERN. We
396 first investigated the correspondence of gene expression prior (cite-lq) and post
397 reconstruction (bulk-hq) with antibody-based surface protein labeling of *CD3D*,
398 *CD4*, *CD8A*, *CD2*, *B3GAT1*, *FCGR3A*, *CD14*, *ITGAX* and *CD19* (Figure 3A,
399 Figure S33). For several proteins (CD8A, B3GAT1, CD4), the corresponding
400 cite-lq gene expression was absent and cell type-specifically re-instantiated in
401 the reconstructed-hq expression data with DISCERN (Figure 3A, Figure S33).

402 In cases where cell type-specific gene and protein expression matched cite-lq
403 data (*CD3D*, *CD14*) the expression in reconstructed-hq data was left unaltered
404 (Figure S33). In some instances, we observed low cell type-specific expression
405 in the cite-lq data (*CD8A*, *CD2*, *FCGR3A*, *CD19*) that matched protein ex-
406 pression (Figure S33). In these cases, gene expression was increased in the cor-
407 rect cell types in the reconstructed-hq data. In general, we observed increased
408 agreement between cell type-specific surface protein and gene expression af-
409 ter reconstruction, showing that DISCERN doesn't invent or 'hallucinate' cell
410 types but reconstructs the expected expression specific for each cell type. We
411 further corroborated these results by selecting eight known cell type-specific
412 cytosolic proteins and investigated their expression before and after expression
413 reconstruction. *MS4A1* (B cells), *IL7R* (CD4⁺ T cells), *MS4A7* (Monocytes),
414 *GNLY* and *NKG7* (NK cells) showed consistent expression before and after
415 reconstruction (Figure S34). The chemokine receptors *CCR2* (Monocytes, ac-
416 tivated T cells), *CXCR1* (NK cells), and *CXCR6* (CD8⁺ T cells) showed the
417 correct cell type-specific expression only after expression reconstruction (Fig-
418 ure S34) [40]. It is notoriously hard to obtain cell subtype-specific information
419 from blood mononuclear scRNA-seq data, especially for CD4⁺ T helper cells due
420 to their limited activation status in healthy individuals. This doesn't mean that
421 polarized CD4⁺ T helper cells do not exist in healthy blood, as they are com-
422 monly detected after stimulation using FACS (Table S3) [41]. This lack of reso-
423 lution in scRNA-seq impedes clustering, marker gene, and trajectory analyses, a
424 drawback that could be overcome using DISCERN's expression reconstruction.
425 We therefore compared CD4⁺ T cell (gene expression of *CD4* > 1 and *CD3E*
426 > 2.5) clustering and subtype identification using cite-lq and reconstructed-hq
427 data. While clustering with the leiden algorithm [42] using highly variable genes
428 of cite-lq data resulted in an unstructured distribution of CD4⁺ T cell subtypes
429 (Figure 3B), clustering of reconstructed-hq data yields detailed insights into
430 T helper cell subtypes of blood mononuclear data (Figure 3C). After recon-
431 struction, we were able to characterize TH17, TH2, TH1, HLA-DR express-
432 ing TREG (Active.TREG), naive CD4⁺ T cells (CD4_naive), effector-memory
433 CD4⁺ T cells (CD4_EM), central-memory CD4⁺ T cells (CD4_CM), and effector
434 cells expressing IFN-regulated genes (IFN_regulated) (Figure 3C). We selected
435 published cell-determining marker genes and observed that many of them were
436 dropped out in the uncorrected data but were present after reconstruction (Fig-
437 ure S35). The absence of marker genes in uncorrected data results in poor
438 clustering and cell type identification, while single positive cells are detectable
439 in the respective neighborhood identified by reconstructed counts (Figure S35).
440 Importantly, we observed that in all cases the DISCERN-estimated proportions
441 of T helper subsets fall within the range of expected proportions as assessed by
442 previous FACS studies (Table S3, Figure S36). These findings are important,
443 as they prove once more that DISCERN discovers the correct cell subtypes and
444 cell proportions, in this case substantially outperforming the available CITE-seq
445 information in cell subtype resolution.

446 To further verify the cell type annotations, we extracted the top cluster-
447 determining genes from the reconstructed-hq data. Members of the TNF-

448 receptor superfamily are known to be expressed in T helper cell subtypes [43],
449 which can be observed after reconstruction in TH17 cells and partially in TH1,
450 TH2, Active_TREG and IFN_regulated cells (Figure S37). Similarly, recon-
451 structed TH1 cells show the expected high expression of granzymes *GZMK* and
452 *GZMA* [44], while *MIAT* and *HLA* expression are found in activated TREG
453 cells after reconstruction (Active_TREG cluster, Figure S37) [45, 46]. *NOG* ex-
454 pression is detected in reconstructed CD4_naive cells, as previously described
455 [47]. In addition, reconstructed CD4_naive, CD4_EM and CD4_CM show low
456 expression of the genes important for the T helper subtypes TH1, TH2, TH17,
457 Active_TREG and IFN_regulated. We further corroborated our cell type anno-
458 tation of reconstructed-hq data by observing the expected expression of several
459 established T cell subtype markers (Figure S38). We compared these newly
460 found clusters to representations found with Seurat, multigrade, and in uncor-
461 rected cite-lq data. The uncorrected cite-lq data manifests cluster separation
462 for some cell types, most notably IFN_regulated and Active_TREG cells (Fig-
463 ure S39A). Seurat reconstruction and multigrade imputation with CITE-seq
464 information results in the mixing of cell types and clusters (Figure S39B & C).
465 A further comparison to Bfimpute and SCRABBLE was impossible due to the
466 dataset size, as outlined in the introduction.

467 Similar to improved clustering and cell subtype detection, DISCERN reconstructed-
468 hq data resulted in improved gene regulatory network inference with SCENIC
469 [48]. SCENIC infers transcription factor-regulated gene expression modules
470 of single cell data. While cite-lq data resulted in a scattered distribution of
471 transcription factor networks across several T helper cell subtypes, SCENIC
472 with reconstructed-hq data showed transcription factor regulation in the cor-
473 rect subtypes (Figure 3D). After expression reconstruction the IKZF2 regulon
474 is detected in activated TREG cells [49] and the MAF regulon is found in differ-
475 entiated CD4⁺ T cells but not in naive CD4⁺ T cells [50]. A weak signal of the
476 MAF regulon is already detectable in the cite-lq data, yet strongly increased in
477 reconstructed-lq, while maintaining differentiated T helper cell specificity (Fig-
478 ure 3D). Furthermore, after reconstruction with DISCERN we could identify
479 the TH17 associated master transcriptional regulators RORC(+) and RORA(+)
480 [51], which were scattered over all TH17 cells before reconstruction (Figure S40).
481 Seurat is able to partially reconstruct the expression of the RORC(+) regulon
482 but fails to detect the more specific RORA(+) expression (Figure S40).

483 Finally, we wanted to investigate if DISCERN could also enhance cell trajec-
484 tory analyses with Slingshot of the citeseq data [52]. We focused on the differen-
485 tiation of effector and other T helper cell subtypes and found five lineages that
486 either pass through or terminate in the effector cell cluster in reconstructed-hq
487 data (Figure 3C). Two trajectories were of special interest to us: Lineage1 from
488 CD4_naive to TH1 cells (Figure S41) and Lineage2 from CD4_naive to TH17
489 cells (Figure S42). While the expression change along the trajectory in uncor-
490 rected data (Figure S41A, Figure S42A) is hardly visible, cell type-specific clus-
491 ters can be easily observed after DISCERN reconstruction (for lineage details
492 see Figure S41B, Figure S42B). The detailed insights into cell differentiation
493 that we obtained with reconstructed data are in stark contrast to the Slingshot

494 results obtained with cite-lq data. While terminal effector molecules can be de-
495 tected with cite-lq data and seurat-hq data, intermediate stages remain hidden,
496 which prohibits the detection of trajectories and results in a shuffling of marker
497 gene expression (Figures S41 and S42). Taken together these results highlight
498 how expression reconstruction using DISCERN improves downstream analyses
499 and yields deeper biological insights into cell type and state identification, gene
500 regulation, and developmental trajectories of cells.

501 *2.4. Discovering COVID-19 disease-relevant cells in lung and blood*

502 The previous sections have demonstrated DISCERN's utility to reconstruct
503 single cell expression data based on an hq reference, vastly improving the detec-
504 tion of cell (sub-) types and their signaling. Given these advantages, we won-
505 dered if DISCERN's expression reconstruction could deepen our understanding
506 of cell type-composition and signaling changes of immune cells in COVID-19
507 disease (Figure S43), using two published datasets [53, 24]. To obtain best re-
508 construction results, we again resorted to using bulk-hq immune reference data
509 (Table S1) [54], as outlined in the previous section.

510 First, we used a COVID-19 blood dataset (covid-blood-lq) with limited cell
511 type resolution, which was originally analyzed by our group using Seurat (Ta-
512 ble S1) [24]. While CD4⁺, CD8⁺, and NK cells formed separate clusters we
513 were unable to visibly distinguish subpopulations of these cells in covid-blood-
514 lq data [24]. Reconstruction of gene expression using bulk-hq data led to the
515 identification of 24 subtypes of CD4⁺ and CD8⁺ T cells in covid-blood-hq data
516 (Figure S44). Several cell clusters identified in covid-blood-hq data showed the
517 correct cell type-specific marker gene expression in covid-blood-lq data, albeit
518 in fewer cells, reduced in magnitude, and in some cases less specific (Figures S45
519 and S46). Reconstruction also led to the identification of CD4⁺ TH17 helper
520 cells that express *RORC* (Figure 4A & B, Figure S47). Based on the molecular
521 footprint of these TH17 cells they were further subdivided into TH17_cluster1
522 that exhibits a memory T cell phenotype with elevated *IL7R* expression and
523 TH17_cluster2 that exhibits an activated T cell phenotype with elevated *MHC-*
524 *II*, *CCR4* and *RBPJ* expression (Figure 4B, Figure S47). The expression of
525 *RBPJ* is of particular interest, as it is linked to TH17 cell pathogenicity, sug-
526 gesting a role of pathogenic TH17 cells in COVID-19 [55]. It is common practice
527 to stimulate memory T cells in vitro to trigger IL-17A production and a shift
528 towards a TH17 phenotype was previously described in COVID-19 [56]. With
529 DISCERN we are able to distinguish these cells in COVID-19 patient blood
530 without stimulation, identifying cytokine producing memory cells with a TH17-
531 like phenotype (Figure S47).

532 To further validate the existence of activated TH17 cells in COVID-19 pa-
533 tient blood, we next analyzed the corresponding lung data (covid-lung) of the
534 patients for shared T cell receptor clones (Figure S48). The underlying assump-
535 tion is that cells with the same T cell receptor in lung and blood originate
536 from the same progenitor and therefore have a high probability of belonging
537 to the same cell type. For this comparison we used the cell type annotation
538 and representation of our original analysis of the covid-lung data, in which

539 memory T and TH17 cells were readily observed without reconstruction [24].
540 TH17_cluster1 cells showed strong clonal overlap with covid-lung CD4⁺ memory
541 T cells (Figure S48) and expressed comparable levels of *RORC* to covid-lung
542 effector memory TH17 cells (Figure S49), indicating that these CD4⁺ central
543 memory T cells could be TH17 (-like) cells. TH17_cluster2 in blood exhibited
544 strong clonal overlap with effector memory and resident memory TH17 cells
545 in covid-lung data (Figure S48) that express *RORC* and *IL-17A* (Figure S49).
546 Using the clonotype information of resident memory cells producing *IL-17A* in
547 inflamed lung (TRM17), we further corroborated the existence of the newly
548 identified population of IL-17A-producing TH17 cells in reconstructed COVID-
549 19 blood data (Figure S48). In general, the T cell receptor clonal information in
550 blood and lung therefore corroborated our cell type annotation in covid-blood-
551 hq data.

552 To understand the role of T cell subtypes in COVID-19 disease progression
553 we analyzed a second blood single cell dataset (covid-blood-severity-lq) contain-
554 ing disease-severity information for 130 COVID-19 patients [53]. To obtain opti-
555 mal cell type resolution, we combined the covid-blood-severity-lq T cell data[53]
556 with CD3⁺ covid-blood-lq cells [24] and reconstructed gene expression for the
557 combined dataset using bulk T cell sequencing reference data[54], resulting in
558 covid-blood-severity-hq data. Many of the 15 CD4⁺ T cell clusters identified in
559 covid-blood-severity-hq data (Figure S50) were also present in the covid-blood-
560 hq data, further validating the consistency of our cell type identification. This is
561 also corroborated by the available surface protein data for covid-blood-severity
562 data, substantiating that naive cells are CD45RA, memory cells are CD45RO,
563 and effector cell types are CD45RO positive (further details in Figure S51). We
564 compared the clusters that we identified in the covid-blood-hq with clusters iden-
565 tified in the covid-blood-severity-hq data and found confined and overlapping
566 regions of TFH, TH17_cluster1, and TH17_cluster2 cells (Figure S52). We also
567 compared the identified clusters to clusters defined in the original publication
568 (Figure S53). Cells identified as TFH in the original publication show signif-
569 icant overlap with naive CD4⁺ T cells (defined on transcriptome and protein
570 level) and CD4⁺ IL22⁺ cells (CD4.IL22) show marked overlap with TREG cells.
571 These results confirm once more the precise and robust cell type identification
572 that can be achieved with DISCERN.

573 Interestingly, we also identified two rather unexpected cell types after re-
574 construction. One cluster is positive for *CD4* and negative for *CD8A* while
575 otherwise expressing a signature of CD8⁺ effector memory cells with high ex-
576 pression of *GZMB*, *GZMH* and *PRF1* (Figure 4D & 4E). This signature points
577 to a CD4⁺ cytotoxic phenotype and indeed virus-reactive CD4⁺ cytotoxic cells
578 were described to be increased in blood during COVID-19 [57]. The other cell
579 type expresses *CD8*, *IL6R*, and *GATA3*, while being negative for *SLAMF7* (Fig-
580 ure 4D & 4E). These cells were described in the literature to be CD8⁺ T helper
581 cells [58], exert T helper function, and have been shown to lack cytotoxicity.
582 They lack expression of a significant number of cytokines and key transcription
583 factors pointing to a TH17 or TH22 phenotype. On a protein level these cells
584 express *CCR4*, while being negative for *CCR6*, making them cytolytic CD8⁺ T

585 helper type 2 cells (Tc2) cells. Part of this cluster overlaps with CD4 single-
586 positive cells and might explain why T helper type 2 cells are missing in the
587 CD4 cell clustering.

588 Overall, the highly specific and sensitive cell type identification in covid-
589 blood-severity-hq data enabled us to correlate the five COVID-19 disease sever-
590 ity categories to shifts in cell type and activity information. We first validated
591 the decrease in TFH cells with increasing disease severity, as described in the
592 original work (Figure S54) [53]. TH17 cells have been extensively studied using
593 flow cytometry and in accordance with our results MHC-II positive as well as
594 *CCR4* positive cells were described in COVID-19 patients (Figure 4B) [56]. We
595 observed a strong decrease in naive T helper cells in severe disease, most pro-
596 nounced for naive TREGs, while the fraction of TH17 cells showed little correla-
597 tion with disease severity (Figure S54). Of the two mixed cell types we detected
598 in COVID-19 data, cytotoxic CD4⁺ cells were increased in moderate and severe
599 disease (Figure S55). A similar increase is visible in patients with severe respi-
600 ratory disease without COVID-19 (Figure S56) and these cells might therefore
601 be a general marker of severe respiratory illness. Cytolytic CD8⁺ Tc2 cells are
602 increased in patients with severe symptoms and in those who died from COVID-
603 19 (Figure S55) and are described to be reduced after recovery from COVID-19
604 [59]. This positive correlation and the known role of Tc2 cells in fibroblast
605 proliferation induction and tissue remodeling could pinpoint a mechanistic role
606 of these cells in lung fibrosis as witnessed in severe COVID-19 patients. The
607 possibility to observe these cells in reconstructed single cell data may pave the
608 way to study the functional role of these cells in adverse COVID-19 outcome.

609 The relatively strong correlation of some cell types with COVID-19 out-
610 come suggests that blood cell fraction information might be used for patient
611 severity prediction. We trained a Gradient Boosting Machine (GBM) using
612 leave-one-out-cross-validation (LOOCV) on the fractions of all T cell types and
613 performed a forward feature elimination, to obtain a sparse, optimal model for
614 patient blood-based severity prediction. We first classified patients into three
615 groups, mild (union of asymptomatic and mild, $n = 26$), moderate ($n = 26$),
616 and severe (union of severe and critical, $n = 19$), reaching an AUROC of 0.63
617 (Table S4). We noticed that the mild and moderate groups were indistinguish-
618 able for the classifier (Figure S57). Training a GBM classifier on mild and severe
619 cases substantially increased classification performance, reaching an AUROC of
620 0.81 and accuracy, and F1 score of 0.82 (Table S4, Figure 4F &G). Compared
621 to the original T cell types and fractions reported (accuracy 0.61) [53], DIS-
622 CERN reconstructed T cell fractions are 33 % more accurate in the prediction
623 of COVID-19 disease severity (Figure 4G, Table S4). This classification improve-
624 ment is remarkable, given that DISCERN has no notion of disease severity when
625 it reconstructs gene expression. These results further demonstrate DISCERN's
626 precise and robust expression reconstruction that enabled the discovery of a
627 potential new blood-based biomarker for COVID-19 severity prediction.

628 3. Discussion

629 The sparsity of gene expression information and high technical noise in sin-
630 gle cell sequencing technologies limits the resolution of cell clustering, cell type
631 identification, and many other analyses. Several algorithms such as scImpute,
632 MAGIC, DeepImpute, and DCA have addressed this problem by imputing miss-
633 ing gene expression in single cell data by borrowing expression information from
634 similar cells within the same dataset. While gene imputation clearly improves
635 gene expression by inferring values for dropped out genes, it comes with several
636 shortcomings. Andrews and Hemberg (2018) showed that several state-of-the-
637 art imputation tools increase the number of false positives [60] by imputing
638 biological absent genes. Additionally the data generated by imputation meth-
639 ods often violate the statistical assumptions made by downstream algorithms,
640 e.g. negative binomial distribution. Furthermore, imputation relies on the com-
641 parison of similar cells with largely absent gene expression information in the
642 same dataset. With DISCERN we approach to gene expression inference of single
643 cell data, by realistic reconstruction of missing gene expression in scRNA-seq
644 data using a related dataset (single cell or bulk RNAseq) with more complete
645 gene expression information. We thus propose to call this procedure ‘expression
646 reconstruction’ to highlight the fundamental difference to classical imputation
647 and refer to the dataset with missing gene expression information as low qual-
648 ity (lq) and the reference dataset as high-quality (hq). We considered a dataset
649 high quality, if it showed a good tradeoff between the mean number of expressed
650 genes and the cell number. For example in the pancreas dataset the smartseq2
651 (6214.0 genes) and the fluidigm1 (8127.4 genes) show a the highest number of
652 expressed genes, but the fluidigm1 batch only consists of 638 cells compared
653 to the smartseq2 batch with 2394 cells, thus we selected the smartseq2 batch as
654 the high-quality batch. However, DISCERN does not require the definition of
655 a high quality batch a priori and it can depend on the scientific question, e.g. a
656 specific batch shows enriched expression of specific genes. In this case the eval-
657 uation of multiple reconstructions with different “high quality” batches can be
658 useful. Furthermore, the use of the dropout estimation procedure in the decoder
659 allows to achieve a single-cell data-like distribution of the reconstructed data
660 and thus is not as strongly violating statistical assumption of downstream anal-
661 ysis. Thus, we consider DISCERN as an approach for expression reconstruction
662 including batch correction, where the reference does not need to be defined *a*
663 *priori* and can come from single cell as well as bulk RNAseq experiments, which
664 enables DISCERN to improve over current state-of-the-art batch correction and
665 imputation methods.

666 We provide compelling evidence that our reference-based reconstruction out-
667 performs contemporary expression imputation algorithms as well as batch cor-
668 rection algorithms such as Seurat, scGen, scVI, and CarDEC when they are
669 repurposed for expression reconstruction. To obtain an objective and thorough
670 performance evaluation for expression inference, we used seven performance
671 metrics on 19 datasets, including 12 single cell sequencing technologies. These
672 datasets cover a range of differences, both technical and biological. While we do

673 not distinguish them in this work, DISCERN could be conditioned on technical
674 as well as biological differences to, for instance, generate ‘diseased’ expression
675 programs from healthy data. We focused our performance evaluation on three
676 scenarios with available ground-truth information, i) the in silico creation of
677 defined gene and pathway drop out events in scRNA-seq data, ii) published
678 hq and lq data pairs from the same tissue (pancreas, difftec, sn/scRNA-seq
679 datasets), and iii) CITE-seq protein expression as ground-truth for cell types
680 (citeseq dataset). In total, DISCERN achieved best performance in 21 out of 27
681 experiments. While DISCERN yields first place to other methods in FC expres-
682 sion correlation comparisons, it always obtains best results across all datasets
683 in gene expression, gene regulatory network analysis, pathway reconstruction,
684 and cell type and activity identification and is the most stable algorithm for
685 different lq to hq size ratios and cell type overlaps. Furthermore it reaches best
686 performance in several batch correction evaluation metrics.

687 It is important to note that DISCERN is a **precise** network that models
688 gene expression values realistically while retaining prior and vital biological in-
689 formation of the lq dataset after reconstruction. The network is also **robust**
690 to the presence of different cell types in hq and lq data, or an imbalance in
691 their relative ratios, and is robust to ‘hallucinating’ hq-specific cells into the lq
692 data. Thus, DISCERN evidently shows less increase in the number of false posi-
693 tives compared to other data smoothing and imputation algorithms. Several
694 algorithmic choices are the foundation of DISCERN’s precision and robustness.
695 The network was designed to model the sequencing-technology-specific and the
696 underlying biological signals in separate components of its architecture. Dis-
697 entanglement of those two components is necessary to accurately reconstruct
698 expression information in the case where lq and hq datasets have different con-
699 tent, i.e. cell type compositions. If the component designed to model the effect
700 of sequencing technology also captures the difference in the biological signal,
701 the reconstruction will lead to a lack of integration across the two datasets
702 where some cell types are still clustered by dataset (similar to scGen in Fig-
703 ure S27). On the contrary, if the component modeling the biological signal
704 captures sequencing-technology-specific features, the reconstruction will lead
705 to an over-integration of the datasets where cells of different types are mixed
706 together (similar to Seurat in Figure S23). The demonstrated ability of DIS-
707 CERN to avoid those shortcomings, even in scenarios where there is very little
708 to no overlap between cell types across datasets, lies in the carefully crafted
709 balance between the expressivity of its components. The representational capa-
710 bilities of DISCERN, achieved via batch normalization, five loss terms, and a
711 dual head decoder, would reduce DISCERN’s usability, if they would require fre-
712 quent dataset-specific tuning. The stability and usability was therefore a central
713 concern in the design and evaluation phase of DISCERN, which resulted in an
714 algorithm that gave very good results with a single set of default (hyper-) param-
715 eters. All comparisons to other algorithms, for instance, were performed with
716 default settings. Only the expression reconstruction of the exceptionally large
717 COVID-19 datasets required the fine-tuning of the learning rate, cross entropy
718 term, sigma, and the MMD penalty term. Another important technical feature

719 of DISCERN is that it can easily be integrated into existing workflows. It takes
720 a normalized count matrix, as created by nearly all existing single cell analysis
721 workflows, as input and produces a reconstructed expression matrix. This can
722 be used for most downstream applications (i.e. cell clustering, cell type identifi-
723 cation, cell trajectory analysis, and differential gene expression). DISCERN can
724 be trained on standard processors (CPU) for small and medium-sized datasets
725 and requires graphical processing units (GPU) for the expression reconstruction
726 of large datasets. Altogether, the usability and robustness of DISCERN should
727 enable even non-expert users to perform gene expression reconstruction.

728 A unique feature of DISCERN is the use of an hq reference to infer biolog-
729 ically meaningful gene expression. While we consider this a main strength of
730 DISCERN, the dependence on a suitable reference dataset might also limit its
731 application. We took great care in this manuscript to mitigate this concern by
732 showing how DISCERN is able to reconstruct gene expression for many differ-
733 ent types of lq and hq pairs, ranging from indrop - smartseq2 to single nucleus
734 - single cell data pairs. Remarkable in this context is DISCERN's robustness
735 to differences between the cell type compositions of lq and hq data pairs, with
736 DISCERN being the only algorithm obtaining robust expression reconstruction
737 when few or no cell types overlap. We have also shown that purified bulk RNA-
738 seq samples can be used as hq reference, as successfully applied to PBMC and
739 COVID-19 datasets in this study. We used 9852 FACS purified immune cell
740 bulk sequencing samples [54], comprising 27 cell types, to successfully recon-
741 struct single cell expression data. This implies that most single cell studies
742 involving immune cells (with or without other cell types present) can be re-
743 constructed with DISCERN using a single published bulk RNA-seq dataset.
744 Furthermore, public RNA-seq repositories such as NCBI GEO contain tens of
745 thousands of samples of immune and non-immune cells that could serve as refer-
746 ence for most expression reconstruction experiments. Conversely, pure cell type
747 or subtype bulk RNA-seq data could be hard to obtain as the sorting of cells
748 might have limited resolution or might be partially impure. In consequence,
749 the usage of bulk RNA-seq data as reference for expression reconstruction could
750 lead to a grouping or averaging of cell subtypes. While these potential caveats
751 might adversely affect expression reconstruction, we have not observed merging
752 or averaging effects of single cell subtypes when corresponding bulk RNA-seq
753 cell type information was not present or present at different proportions (Fig-
754 ure 3B & 3C, Figure S36). Importantly, cells do not necessarily cluster into
755 distinct classes but can build cell continua, as shown in the trajectory analy-
756 sis in Figure 3B & 3C, where T cells seem to differentiate into each other and
757 do not form clearly separable clusters. In general, handling continua of cell
758 types is challenging for imputation and batch correction algorithms, as many of
759 them, including for instance scGEN, Bfimpute, SIMPLEs, and cscGAN, require
760 or recommend cluster or cell type annotation. This might lead to under- or
761 over-integration of cell continua. DISCERN does not rely on cluster (or cell
762 type) information and seamlessly integrates and reconstructs cell clusters and
763 continua (Figure 3C, Figure S44). In conclusion, we provide strong evidence
764 that DISCERN is widely and easily applicable to many single cell experiments.

765 While DISCERN gave good reconstruction results using default parameters
766 for most datasets we analyzed, we would like to highlight that the immense
767 representational power of generative neural networks can remove or hallucinate
768 biological information if not properly handled [6]. This is true for data integra-
769 tion [61, 62] as well as for expression reconstruction algorithms and we would
770 highlight two guiding principles for optimal results. For non-expert users, we
771 would recommend the use of default settings and a careful selection of a re-
772 lated hq dataset. When datasets are large and complex, with many cell types
773 in the lq and several non-overlapping cell types in the hq data, one should al-
774 ways ensure that training does not merge or mix non-overlapping cell types with
775 other cells, by investigating that these cells keep their cell type-specific marker
776 gene expression. Keeping these ‘checks and balances’ will usually result in good
777 reconstruction results even for complex datasets such as covid-blood-severity.

778 To obtain novel insights into COVID-19 disease mechanisms and a new blood-
779 based biomarker for disease severity we reconstructed two published datasets
780 with DISCERN, Hamburg COVID-19 patients (covid-lung, -blood) and the
781 COVID-19 cell atlas (covid-blood-severity). The application of DISCERN to
782 the covid-blood dataset (COVID-19 patient blood) enabled us to detect 24 dif-
783 ferent immune cell types and activity states, which is quite remarkable given
784 that we find these cells in blood. Two TH17 subtypes caught our attention, as
785 they share the TCR clonality with the lung data from the same patients (covid-
786 lung), suggesting bloodstream re-entry of lung TH17 cells. We linked these two
787 subclusters to their functional role by separating them into a memory-like and
788 activated-like phenotype. The clonal overlap of activated TH17 cells in blood
789 with previously discovered lung-resident cells suggests that activated TH17 cells
790 in blood are resident T cells from the lung reentering circulation. These cells
791 might in part explain the multi-organ pathology observed in COVID-19, as
792 activated T cells might travel via the blood to secondary organs and cause in-
793 flammation and tissue damage. Future work might demonstrate the effect of
794 these activated T cells on tissue inflammation.

795 Given the detailed cell type and activity information we reached with gene
796 expression reconstruction, we wondered if changes in blood immune cell popu-
797 lations might be useful as a biomarker for disease severity prediction. We used
798 DISCERN to reconstruct the covid-blood and the covid-blood-severity datasets
799 and again identified a plethora of different T cell subtypes in the blood of pa-
800 tients with COVID-19. Using these cell proportions, we were able to classify
801 mild and severe disease using a GBM machine learning algorithm with 82 %
802 accuracy, outperforming classification with the originally published T cell types
803 by 21 percent points. This improvement is absolutely striking, as DISCERN
804 has no notion of the classification groups. It simply reconstructs gene expres-
805 sion and thereby improves cell type detection. These results are a convincing
806 implicit proof not only of the usefulness of DISCERN but more importantly of
807 its precision and robustness. While the use of this scRNA-seq-based biomarker
808 would be too expensive and time-consuming for clinical care, it strongly suggests
809 that FACS-based T cell fraction or count information from blood could be used
810 to trace and predict the severity state and potentially the disease trajectory of

811 COVID-19 patients.

812 Interestingly, we also discovered two atypical T cell types in reconstructed
813 COVID-19 patient blood single cell data. While cytotoxic CD4⁺ T cells have
814 been observed in COVID-19, we can show that this increase is not COVID-19
815 specific and is also observed in other types of pneumonia. Interestingly, we also
816 detected cytolytic CD8⁺ Tc2 cells that express *CD8A*, *GATA3*, *IL6R* and are
817 negative for *SLAMF6*. This cell type is linked to tissue fibrosis and steroid
818 refractory disease in asthma [63]. The increase in CD8⁺ Tc2 cells that we ob-
819 serve specifically in COVID-related death could be associated with COVID-19
820 patients that do not respond to steroids. Demonstration of increase of this cell
821 type in patients dying of COVID-19 points to a potential therapeutic inter-
822 vention with the drug Fevipiprant, which blocks CD8⁺ Tc2 cell activation and
823 its pro-fibrotic effects by inhibiting prostaglandin D2 signaling [64]. Functional
824 analysis of these cells has to demonstrate whether these cells are an early marker
825 of later death or whether it is a marker of already escalated treatment.

826 The basic concept of utilizing a high-quality reference to improve lower qual-
827 ity data might be applied to many other research areas where technological
828 limitations restrict biological insights. The usage of deep generative networks
829 and other artificial intelligence methodology to infer information beyond what
830 is technically measurable could be transformative in future biomedical research.

831 **Acknowledgements**

832 We thank Immo Prinz, Manuel Friese, Johannes Soeding, Robert Zinzen, Yu
833 Zhao and Stefan Kurtz for their helpful comments and suggestions. We thank
834 Rajasree Menon and Matthias Kretzler for providing kidney single cell and sin-
835 gular nuclear RNA-seq data and their support for corresponding analysis. FH,
836 RK, MM, PM, & SB were supported by the LFF-FV 78, EU ERare-3 Maxo-
837 mod, SFB 1286 Z2, FOR 296, and FOR 5068 research grants. Further support
838 was obtained from the UKE R3 reduction of animal testing grant. CE was sup-
839 ported by DFG ER 981/1-1 and the clinician scientist programme of university
840 Hamburg, NG was supported by ERC StG-715271, and SHu was supported by
841 ERC CoG-865466 and has an endowed Heisenberg-Professorship awarded by
842 the Deutsche Forschungsgemeinschaft. SHa was funded by the BMBF STOP-
843 FSGS-01GM1518C and SFB 1192 B08 research grants.

844 **Competing interests**

845 The authors declare no competing interests.

846 **Author contributions**

847 SB initiated and SB, PM, FH, and CE conceptualized the study with help
848 from MM. FH and CE implemented DISCERN, MM refactored the code, and
849 PM reviewed the DISCERN implementation. FH, CE, and RK performed the

850 analyses. SB, PM, NG, and SHu supervised the study. SB, FH, and CE wrote
 851 the manuscript. SHu, PM, NG, RK and SHa provided ideas, contributed to the
 852 manuscript text and critically reviewed the manuscript. All authors read and
 853 approved the final manuscript.

854 **Main figures**

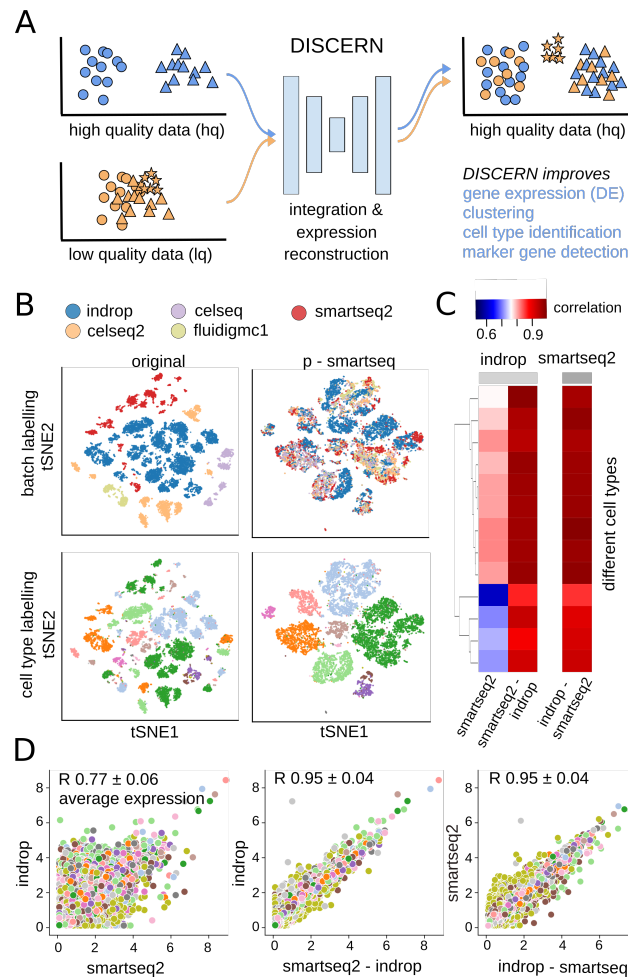


Figure 1: *Integration and expression reconstruction of single cell sequencing data.* **A:** DISCERN transfers the style of a high-quality (hq) dataset to a related low quality (lq) dataset, enabling gene expression reconstruction that results in improved clustering, cell type identification, marker gene detection, and mechanistic insights into cell function. The hq and lq datasets have to be related but not identical, containing for example several overlapping cell types but also exclusive cell types of cell activity states for one or the other dataset. **B:** t-SNE visualization of the pancreas dataset before reconstruction (original) and after transferring

the style of the smartseq2 dataset using DISCERN (p-smartseq2). The upper row shows the dataset of origin before and after projection colored by batch and the lower row colored by cell type annotation (details of 13 cell types in supplements). **C** and **D**: Average gene expression (over all the cells of a given type) of the pancreas indrop and smartseq2 datasets before (first column and panel) and after smartseq2 to indrop (second column and panel), and after indrop to smartseq2 projection (third column and panel). **C**: Gene correlation by cell type shown in colored heatmap. **D**: Each colored point represents a single gene colored by the cell type. The mean Pearson correlation with one standard deviation over all cell types is shown in the figure title.

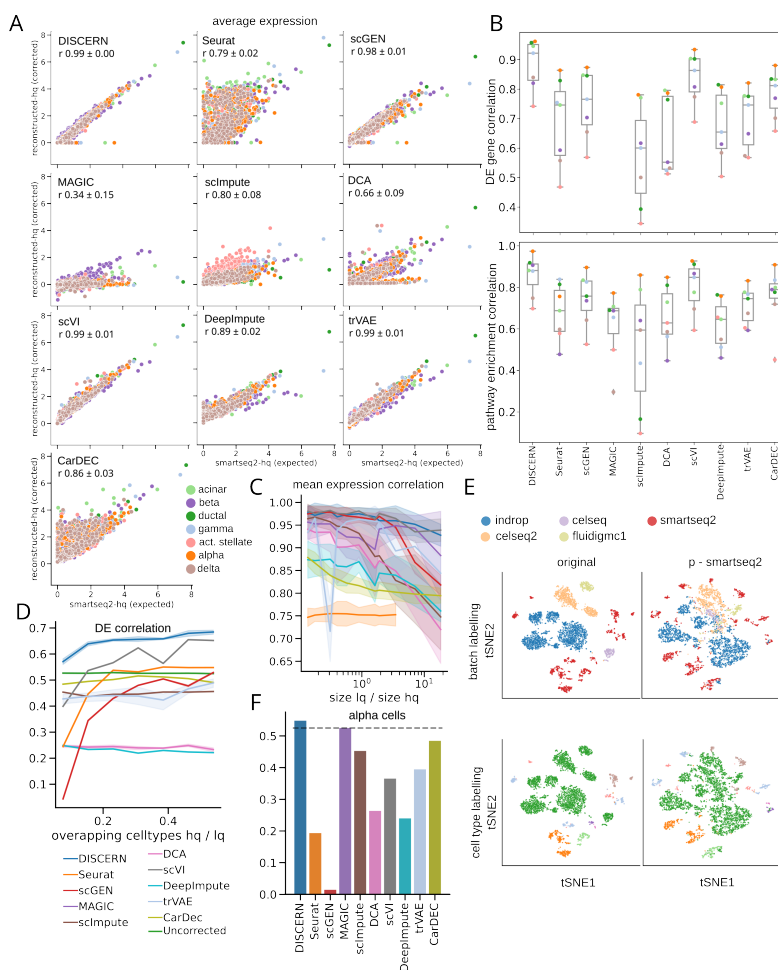


Figure 2: *Expression reconstruction benchmark of DISCERN and five state-of-the-art batch correction and imputation algorithms.* **A**: Comparison of the expression reconstruction performance of Seurat, scGEN, MAGIC, scImpute, DCA, scVI, trVAE, DeepImpute, and DISCERN using smartseq2 data. The smartseq2 data was split into a smartseq2-lq and a smartseq2-hq batch. The smartseq2-lq batch was modified such that the expression of all genes of a cell type determining pathway (top ranked by GSEA) was set to zero. The expression of the in silico altered pathway genes was then compared between reconstructed-hq data and the unaltered

smartseq2-hq data. **B**: Differential gene expression and pathway enrichment correlation of the reconstructed-hq to the expected values before removal. The smartseq2-lq data was the same as in **A**. The DEG analysis was restricted to genes which were removed in the smartseq2-lq batch. Correlation of the DEG analysis was based on the t-statistic and for the pathway enrichment analysis on the normalized enrichment scores. **C**: Mean expression correlation of reconstructed-hq with the expected expression in smartseq2-hq data for different ratios of lq to hq data. The standard deviation indicates the deviation in correlation of the cell types. The datasets were created as described in **A**. **D**: Alpha cells were removed from the smartseq2-hq batch and left in the low quality batches. The number of overlapping cell types between the hq and lq data was then altered by removing cell types, which overlap between lq and hq data, from the lq data before preprocessing and expression reconstruction. The ratio of the intersection size to the total number of cell types is shown on the x-axis. The y-axis shows the correlation of the t-statistics of alpha cells from lq-batches vs other cells from the smartseq2 batch with ground truth alpha cells from the smartseq2 batch vs other cells from the uncorrected smartseq2 batch. We used Spearman rank correlation for the comparison, since no gene subset was used. **E**: t-SNE visualization of the cell type removal experiment where alpha cells are removed from the smartseq2 batch and all non-alpha cells are removed from the lq-batches, such that there is no overlap between lq and hq. **F**: Spearman correlation of the t-statistics of alpha cells from lq-batches vs other cells from the smartseq2 batch with ground truth alpha cells from the smartseq2 batch vs other cells from the uncorrected smartseq2 batch. The dataset was the same as in **E** (no cell type overlap between hq and lq data). The dotted line indicates the correlation achieved without reconstruction.

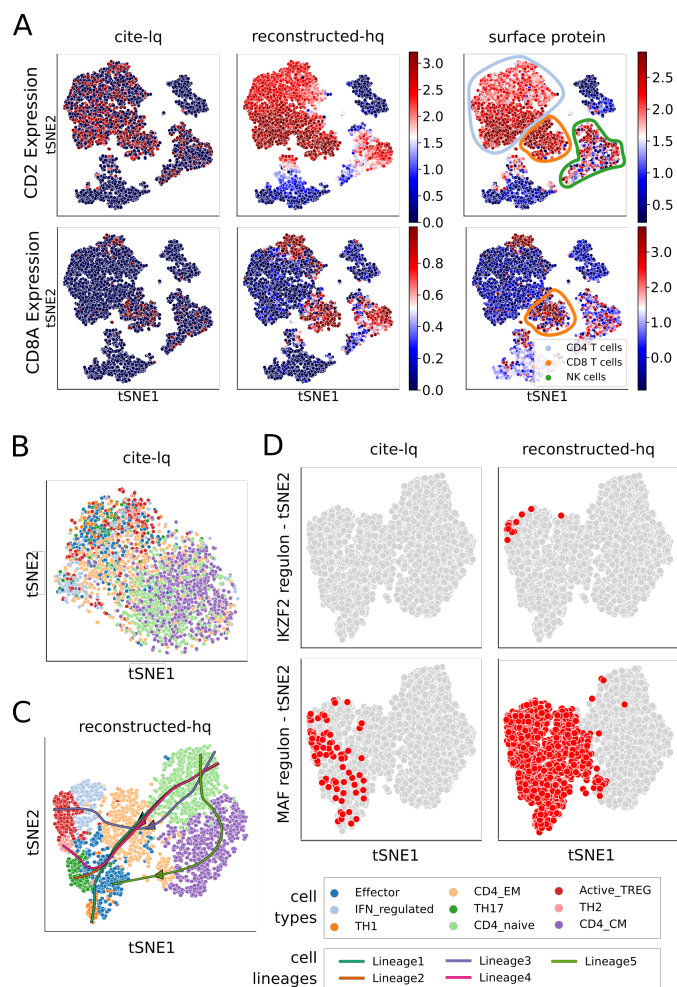


Figure 3: *Expression reconstruction improves downstream analyses including cell identification, gene regulation, and trajectory inference.* The cite-lq dataset was reconstructed using bulk-hq data and compared to ground truth CITE-seq (surface protein) information. The CITE-seq information was not used during training of DISCERN. **A**: t-SNE visualization of *CD2* (first row) and *CD8A* (second row) gene (first two columns) and protein (last column) expression. The first column depicts gene expression for uncorrected cite-lq, the second for reconstructed-hq, and the third protein surface expression ground truth information. Cell types commonly known to express these genes are highlighted with colored circles in the last column. **B**: t-SNE visualization of $CD4^+$ T cells in the cite-lq dataset. Cell types were assigned using louvain clustering on the reconstructed-hq data (see C) and show no clear clustering. **C**: t-SNE and trajectory information of $CD4^+$ T cell subtypes found by Slingshot analysis on reconstructed-hq data. While uncorrected data shows no clear cell type clustering (see B), reconstructed data shows a clear grouping of cell types. Trajectories were calculated using *CD4_naive* as starting point and *TH2*, *TH17*, *TH1*, *Active_TREG*, *CD4_CM* as endpoints. Lineage1 indicates *TH1*, Lineage2 *TH17*, Lineage3 *Active_TREG*, Lineage4 *TH2*, and Lineage5 *Effector* cell differentiation. **D**: Detection of regulons that are specific for $CD4^+$ T cell subtypes using pySCENIC. The first column shows regulons found in the uncorrected

cite-lq and the second column in reconstructed-hq data.

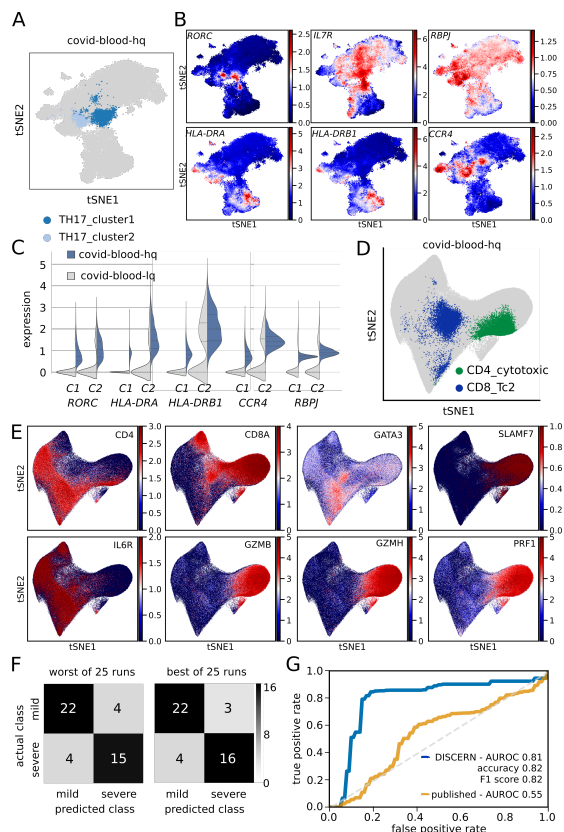


Figure 4: *Expression reconstruction improves COVID-19 cell type identification and allows for efficient disease severity prediction.* Two COVID-19 blood datasets were reconstructed and analyzed. Hamburg covid-blood-lq and covid-lung-lq data was reconstructed using bulk-hq data, resulting in the respective -hq datasets. Similarly, Cambridge covid-blood-severity-lq data, which contains disease severity information, was reconstructed using bulk-hq data. **A**: t-SNE representation of TH17 subclusters using reconstructed covid-blood-hq data. Clusters were defined using the leiden clustering algorithm on CD4⁺ T cells. **B**: t-SNE representation colored by expression of reconstructed genes distinguishing TH17_cluster1 and TH17_cluster2 cells. TH17_cluster1 displays a central memory and TH17_cluster2 a more activated phenotype. **C**: Violin plots of expression levels for genes distinguishing TH17_cluster1 (C1) and TH17_cluster2 (C2) cells before (covid-blood-lq) and after (covid-blood-hq) reconstruction with DISCERN. **D**: Rare and unexpected cell types found in the reconstructed covid-blood-hq data with covid-blood-severity and bulk data. Cytotoxic CD4⁺ T cells (CD4_cytotoxic) are displayed in green, CD8⁺ Tc2 helper cells (CD8_Tc2) in blue, and all other cells in gray color. **E**: t-SNE representation of key marker genes in covid-blood-hq data for CD4_cytotoxic and CD8_Tc2 cells displayed in **D**. **F**: Best and worst confusion matrix for disease severity prediction using GBM classifiers trained on fractions of five T cell types (CD4_CM, CD4_cytotoxic, CD4_naive, CD8_EM, CD8_effector) using reconstructed covid-blood-severity-hq data. Category “critical” was combined with “severe” and “asymptomatic” with “mild”. **G**: ROC curve of the GBM predictions outlined in **F** using reconstructed (blue color) covid-blood-

severity-hq (CD4_CM, CD4_cytotoxic, CD4_naive, CD8_EM, CD8_effector) and published T cell information from uncorrected (yellow color) data (CD4_CM, CD4_Tfh, CD8_EM, NKT, Treg). Confidence intervals (color shades) indicate one standard deviation.

855

856

857 4. Methods

858 4.1. Data availability

859 In this manuscript many different scRNA-seq and RNA-seq datasets were
860 used. A comprehensive overview of dataset, method, cell type, origin, size, and
861 naming convention can be found in Tables S1 to S3. All datasets are publicly
862 available as listed in Table S1.

863 4.2. Dataset description

864 *Pancreas.* The pancreas dataset is a collection of different scRNA-seq datasets,
865 profiling pancreas cells in the context of diabetes [65]. The pancreas dataset is
866 a widely used dataset for batch correction benchmark experiments and due to
867 its high number of cell types and sequencing technologies it allows to evaluate
868 differences between cells and sequencing technologies at the same time. The ex-
869 pression table, including the annotation, is available from SeuratData ([https://](https://github.com/satijalab/seurat-data)
870 github.com/satijalab/seurat-data) as `panc8.SeuratData` (v3.0.2) [65]. The
871 dataset was sequenced using five sequencing technologies (Smart-Seq2, Flu-
872 idigm C1, CelSeq, CEL-Seq2, inDrop) and consists of 13 cell types (alpha, beta
873 ,ductal, acinar, delta, gamma, activated_stellate, endothelial, quiescent_stellate,
874 macrophage, mast, epsilon, schwann). In total, before preprocessing, the dataset
875 contains 14 890 cells.

876 *difftec.* The difftec dataset was created for a systematic comparative analysis
877 of scRNA-seq methods [66]. Similar to pancreas, the difftec dataset is ideal
878 for the evaluation of expression reconstruction across many cell types and se-
879 quencing technologies. Seven sequencing technologies (10x Chromium v2, 10x
880 Chromium v3, Smart-Seq2, Seq-Well, inDrop, Drop-seq, CEL-Seq2) were used
881 with at least two replicates each. In this dataset 10 different cell types (Cy-
882 tototoxic T cell, CD4⁺ T cell, CD14⁺ monocyte, B cell, Natural killer cell,
883 Megakaryocyte, CD16⁺ monocyte, Dendritic cell, Plasmacytoid dendritic cell,
884 Unassigned) were annotated, and make up for 31 021 cells in total before filter-
885 ing. The expression table including the annotation is available from SeuratData
886 as `pbmcsca.SeuratData` (v3.0.0).

887 *snRNA & scRNA.* The dataset was created for the validation of a single cell
888 and single nuclei analysis toolbox [38]. Since snRNA-seq and scRNA-seq data
889 varies in the amount of counts per cell and the genes detected, we tested if
890 DISCERN could reconstruct snRNA-seq expression so that it would closely
891 resemble scRNA-seq expression, providing a biological ground-truth. While we
892 label snRNA-seq data as `lq` and scRNA-seq as `hq`, this distinction is incorrect

893 from a biological perspective, as gene expression should be in part different
894 between the nucleus and the cytosol. The dataset consists of a liver biopsy
895 sample (HTAPP-963) of metastatic breast cancer with single cell sequencing
896 and single nuclei sequencing. Eight cell types (Epithelial cells, Macrophages,
897 Hepatocytes, T cells, Endothelial cells, Fibroblasts, B cells, NK cells) were found
898 in the original publication in a total of 12 423 cells. The data was sequenced
899 using the Chromium V3 technology on a Illumina HiSeq X sequencer.

900 *covid-lung & covid-blood.* The COVID-19 dataset we have previously published
901 consists of blood and bronchoalveolar lavage (BAL) samples from four patients
902 with bacterial pneumonia and eight patients with SARS-CoV-2 infection[24].
903 In total 155 706 cells were sequenced using TCR-seq technology, which allows
904 for the comparison of clonal expansion in both tissues. While we investigated
905 the lung data in detail in the original publication, the analysis of the blood was
906 largely limited to cell type identification. Using DISCERN, we use the blood
907 data to find previously unobserved cell types, link them to cell clones found in
908 the lung, and derive a biomarker based on cell fractions (see also covid-blood-
909 severity data). Cell type annotations for the BAL samples were used as in the
910 original publication.

911 *citeseq.* This dataset contains CITE-seq information of healthy human PBMCs
912 for 6 cell types (B cells, CD4 T cells, NK cells, CD14⁺ Monocytes, FCGR3A⁺
913 Monocytes, CD8 T cells) [39]. In our analyses we used the cell type information
914 provided in the original publication [67]. The CITE-seq data is ideal to bench-
915 mark DISCERN, as the information of 13 surface proteins offers ground-truth
916 information on the cell types and a good proxy for the expression of the 13
917 corresponding genes.

918 *bulk.* We used this large dataset of 28 FACS sorted and bulk sequenced immune
919 cell types as ‘ultimate’ hq reference data for lq immune single cell sequencing
920 data. Each of the 9852 samples provides an average expression information for
921 13 104 genes for a specific immune cell type, providing a hq reference for e.g. lq
922 single cell PBMC CITE-seq data with only 798 expressed genes per cell. We
923 further assume that this dataset is large enough to provide enough per cell type
924 variability for our deep neural network to faithfully learn and represent its gene
925 expression. In more detail, the dataset consists of 28 sorted immune cell types
926 (Naive CD4, Memory CD4, TH1, TH2, TH17, Tfh, Fr. I nTreg, Fr. II eTreg,
927 Fr. III T, Naive CD8, Memory CD8, CM CD8, EM CD8, TEMRA CD8, NK,
928 Naive B, USM B, SM B, Plasmablast, DN B, CL Monocytes, Int Monocytes,
929 NC Monocytes, mDC, pDC, Neutrophils, LDG) with \geq 99% purity [54]. Total
930 RNA was extracted using RNeasy Micro Kits (QIAGEN). Libraries for RNA-seq
931 were prepared using SMART-seq v4 Ultra Low Input RNA Kit (Takara Bio).
932 In total, the dataset contains 9852 samples collected in two phases from 416
933 donors, out of which 79 are healthy. For training DISCERN, bulk TPM counts
934 and all cell types were used if not stated otherwise.

935 *covid-blood-severity*. This dataset is an aggregation of three COVID-19 sequenc-
936 ing studies using the 10X Genomics Chromium Single Cell 5' v1.1 technology.
937 It contains a large number of cell types with fine-grained cell type annotations
938 that are complemented with information on COVID-19 disease severity for each
939 patient sequenced. We used this dataset to obtain a blood-based biomarker of
940 COVID-19 disease severity, based on T cell fractions observed with DISCERN.
941 The data consists of PBMCs from 29 healthy, 89 COVID-19 and 12 LPS-treated
942 patients. The authors detected 51 cell types in their original work (see Ta-
943 ble S1) [53] and COVID-19 patients were classified by their disease severity
944 (worst clinical outcome) into 'asymptomatic', 'mild', 'moderate', 'severe', 'crit-
945 ical', and 'death'. Count data together with CITE-seq information was used
946 as provided in the original publication ([https://covid19.cog.sanger.ac.uk/
947 submissions/release1/haniffa21.processed.h5ad](https://covid19.cog.sanger.ac.uk/submissions/release1/haniffa21.processed.h5ad)).

948 *kidney-lq (snRNA-seq) & kidney-hq (scRNA-seq)*. The kidney dataset consists
949 of single cell RNA-seq and single nuclei RNA-seq data of 9 patients with acute
950 kidney injury sequenced using 10X Genomics Chromium technology. It contains
951 in total 82 701 cells with 52 934 cells sequenced using snRNA-seq and 29 767 cells
952 sequenced using scRNA-seq. The dataset does not contain cell type annotation,
953 but in initial analysis using a different subset [68] suggested that identification
954 of T cells in the snRNA-seq data is challenging. For this reason, the analysis
955 was focused on the detection of T cells and their subtypes.

956 4.3. Code availability

957 All original code has been deposited at github.com ([https://github.com/
958 imsb-uke/discern](https://github.com/imsb-uke/discern)) and is publicly available as of the date of publication. Any
959 additional information required to reanalyze the data reported in this paper is
960 available from the lead contact upon request.

961 4.4. Preprocessing

962 Raw expression data (Counts) preprocessing was performed as previously
963 described [69] using the scanpy (v1.6.1, [70]) implementation. In particular,
964 the intersection of genes between batches was used. The cells were filtered
965 to a minimum of 10 genes per cell and a minimum of 3 cells per gene. Li-
966 brary size normalization was performed to a value of 20 000 with subsequent
967 log-transformation. As model input for DISCERN the genes were scaled to
968 zero mean and unit variance. However, for all further evaluation the genes
969 were scaled to their uncorrected mean and variance not considering the batch
970 information.

971 4.5. Description of DISCERN

972 DISCERN is based on a Wasserstein Autoencoder with several added and
973 modified features. We will describe the details of DISCERN's architecture in
974 the next paragraphs and a compact representation can be found in Figure S1B.

975 *Wasserstein Autoencoder.* While neural network-based autoencoders have been
976 widely used for decades for dimensionality reduction [71, 72], recent advances
977 have also allowed their use to build a generative model of the distribution of
978 the data at hand[73]. More recently, leveraging results from optimal transport
979 [74], Wasserstein Generative Adversarial Networks (WGAN) [75] and Wasser-
980 stein Autoencoders (WAE) [25] have been designed to explicitly minimize the
981 (Wasserstein, or earth-mover) distance between the distribution of the input
982 data and their reconstruction. WGANs only implicitly encode their input into
983 a latent representation (called latent code), while WAE has the useful property
984 of using an explicit encoder, which makes it possible for the model to directly
985 manipulate the different representations of single-cell data. Finally, the WAE
986 framework, established in [25], allows the use of a wide range of architecture and
987 losses, which we are going to detail now. First of all, in order to effectively use a
988 number of latent dimensions that adaptively matches the intrinsic dimension of
989 the scRNA-seq data at hand, DISCERN uses a random encoder as prescribed
990 in [76].

991 *Architecture.* Autoencoders widely used for transcriptomics applications are
992 shown to perform well on several tasks, like drug perturbation prediction [23]
993 or dropout imputation [12]. Since the ordering of the genes in scRNA-seq count
994 matrices is mostly arbitrary, fully-connected layers are usually used in this task.
995 In our case, DISCERN consists of three fully connected layers in the encoder
996 and the decoder. The bottleneck of the autoencoder (or latent space) contains
997 48 neurons, which is sufficient to accurately model all the datasets we used in
998 our experiments. Additionally, we exploit a finding from [76] to let the net-
999 work learn the appropriate amount of latent dimensions. While the encoder
1000 will be tasked to transform the distribution of the input data into a fixed,
1001 low-dimensional prior distribution (i.e. a standard Gaussian), the decoder will
1002 perform the opposite, i.e. transforming the fixed, low-dimensional prior distri-
1003 bution into gene space. scRNA-seq data is known to display a high level of zero
1004 measurements, called dropout, which is essential to accurately model the count
1005 distribution. To describe scRNA-seq data in a parametric way, it is common to
1006 model the expression level of a gene with zero-inflated negative binomial dis-
1007 tribution [77]. Despite the several non-linearities in the decoder architecture,
1008 it is, however, difficult to learn an encoding function that maps a simple prior
1009 to the distribution leading to low quality modeling of low expressed genes. To
1010 address this issue, we scale the gene expression and attach a second head to the
1011 decoder (i.e. a second decoder sharing all weights with the first, except for the
1012 last layer). The task of the second decoder head is to predict, for each gene
1013 of a cell, the probability of its expression to be dropped out, giving rise to a
1014 random decoder. Thus, this second decoder head predicts dropout probabili-
1015 ties and models the dropout probabilities for different batches. This additional
1016 head allows modeling the dropout and the expression independently, to capture
1017 the specific distribution of single cell data without the need for further explicit
1018 assumption about the distribution. During inference the predicted expressions
1019 are randomly set to zero based on these predicted dropout probabilities. This

1020 sampling procedure does not have any trainable parameter, is therefore not part
1021 of the model training and only performed during inference.

Loss function. The loss optimized during the training of DISCERN is composed of four terms: a data-fitting (or reconstruction) loss, a dropout fitting (cross entropy) loss, a prior-fitting term (ensuring that DISCERN approximately minimizes the Wasserstein distance) and a variance penalty term (that controls the randomness of the encoder). Thus, DISCERN can be considered as a Wasserstein Autoencoder as introduced in [25]. For the reconstruction term, the framework introduced in [25] allows the use of any positive cost function. We elected to use the Huber loss [78] as it is well suited for modeling scaled scRNA-seq expression data, because it allows to select a threshold value to give lower weight to high differences in highly expressed genes and thus allows the model to learn a more robust expression estimate without focusing too much on outlier values. This reconstruction term is defined as

$$L_{\delta}(x, \hat{x}^{count}) = \frac{1}{d_x} \sum_{i=1}^{d_x} \begin{cases} \frac{1}{2}(x_i - \hat{x}_i^{count})^2 & \text{for } \|x_i - \hat{x}_i^{count}\| \leq \delta, \\ \delta (\|x_i - \hat{x}_i^{count}\| - \frac{1}{2}\delta), & \text{otherwise.} \end{cases}$$

1022 where x is the input expression matrix, \hat{x}^{count} the predicted expression matrix
1023 from one decoder head, d_x the number of genes, and δ a threshold deciding
1024 between the two conditions of the Huber loss.

1025 For the prior-fitting term, following [25], DISCERN uses the Maximum Mean
1026 Discrepancy (MMD) [79] between the aggregate posterior (i.e. the distribution
1027 of the input single-cells after encoding) and a standard Gaussian. We use the
1028 sum over an inverse multiquadratic kernel with different sizes for this task.

Similar to [79], we define the MMD as

$$\mathcal{D}_Z(P_Z, Q_Z) = \left\| \int_{\mathcal{Z}} k(z, \cdot) dP_Z(z) - \int_{\mathcal{Z}} k(z, \cdot) dQ_Z(z) \right\|_{\mathcal{H}_k}$$

1029 where P_Z is the gaussian prior distribution and Q_Z the aggregated posterior
1030 in the latent space for a positive-definite reproducing kernel $k: \mathcal{Z} \times \mathcal{Z} \rightarrow \mathcal{R}$
1031 and a corresponding real valued reproducing kernel hilbert space \mathcal{H}_k . For the
1032 implementation details please refer to [25] or the provided implementation.

Then, to prevent the random encoder (with diagonal covariance) from collapsing to a deterministic one, a penalty term that enforces that some components of the variance are close to 1. Intuitively, that means that the superfluous latent dimensions will only contain random noise (see [76] for more details). We define this penalty term as

$$S_{\sigma}(x) = \sum_{i=1}^{d_z} \left| \log(\sigma_i^2(x)) \right|$$

1033 where d_z is the number of latent dimensions and σ the function generating the
1034 components of the variance in the latent space, in our case, the encoder network.

Another loss term, namely the binary cross-entropy loss, on the second decoder head is used to enable the model to learn a dropout probability for each gene and sample. The loss on the dropout layer enables the model to capture the bimodal distribution of single cell data. We define the binary cross-entropy loss as

$$H(x^{dropout}, \hat{x}^{dropout}) = -\frac{1}{d_x} \sum_{i=1}^{d_x} x_i^{dropout} \log(\hat{x}_i^{dropout}) + (1 - x_i^{dropout}) \log(1 - \hat{x}_i^{dropout})$$

where $x^{dropout}$ is the binarized expression information, $\hat{x}^{dropout}$ is the predicted binarized expression (probability of dropout) and d_x the number of genes. Additionally, activity regularization is applied on the Conditional Layer Normalization (CLN), such that the weights of the conditional layers are only regularized in a batch-specific manner and the regularization is not applied for batches, which are not present in the current mini-batch. This has the advantage that the batch dependent weights are not influenced too much by different batch sizes. The four loss terms are added (and weighed using λ s) together to form the loss that DISCERN minimizes during training:

$$L = L_\delta(x, \hat{x}^{count}(z)) + \lambda_{prior} \cdot \mathcal{D}_Z(q_z, p_z) + \lambda_{sigma} \cdot S_\sigma(x) + \lambda_{dropout} \cdot H(\mathcal{I}_{>0}(x), \hat{x}^{dropout}(z))$$

1035 See also Figure S1B for a graphical depiction of the loss terms.

1036 *Conditional Layer Normalization.* The weights of those fully-connected layers are
1037 shared for all the batches that DISCERN is trained on. However, to model the batch-
1038 specific differences, we use a Conditional Layer Normalization (CLN) that applies the
1039 idea proposed in [27] to Layer Normalization [28] after each fully connected layer (see
1040 Figure S1B). In essence, for each batch, different sets of shifting factors are learned.
1041 Note that in DISCERN, no scaling factors are used to limit the expressivity of the
1042 conditioning and therefore reduce the chance of over integration. This allows not only
1043 to accurately model the batch-specific differences between batches, but also to trans-
1044 fer the batch effect from one dataset onto another, in the spirit of the style-transfer
1045 approach developed in [27]. To make things clear, DISCERN does not explicitly train
1046 to integrate datasets. Instead, it trains to accurately model the input data, capturing
1047 the batch-specific differences with the weights of the CLN layers (i.e. conditioning),
1048 and the biological signal (which is mostly shared across the batches to integrate) with
1049 the weights of the fully-connected layers. After training, we encode all the cells we
1050 want to reconstruct, conditioning the process on their batch of origin. Then, we take
1051 the batch chosen by the user and proceed to decode all the cells conditioning on that
1052 specific batch, effectively transferring the batch effect of one specific batch onto all
1053 of the batches we want to integrate and reconstruct. The training loss is computed
1054 over the complete minibatch, thus it is not different per batch (dataset). The weights
1055 of the conditional layer normalization are learned together with the weights of the
1056 feed-forward network using the same loss function.

1057 *Activations & dropout.* With the exception of the output layer, every other fully-
1058 connected layer of the encoder and the decoder was followed by a CLN, a Mish ([80])
1059 activation function, and dropout during model training to reduce overfitting.

1060 *Optimization.* To optimize the weights of our model, DISCERN uses Rectified Adam
1061 ([81], which addresses some of the shortcomings of the widely used Adam [82] and gen-
1062 erally yields more stable training. To prevent overfitting, the optimization is stopped
1063 early. It is implemented as a modification of the Keras EarlyStopping (with parameter
1064 minDelta set to 0.01 and the patience to 30) where the callback is delayed by a fixed
1065 number of 5 epochs. The delay was implemented to prevent too early stopping due to
1066 the optimization procedure.

1067 *Reconstruction.* The reconstruction (or projection) to a reference batch is not per-
1068 formed during training and thus the network is not optimized to it. However, during
1069 inference, the reconstruction can be performed by providing the correct batch label
1070 in the encoder part of the network, while only providing the reference batch label
1071 for the decoder part. Therefore, The network will encode the dataset to a batch-
1072 independent latent representation and decode it using only the reference label and
1073 therefore project the complete dataset to the reference batch. This can be done for
1074 any number of batches without re-training of DISCERN.

1075 *Running time and memory usage.* DISCERNs running time for training is linear in
1076 the number of cells and the number of training epochs. However, the use of the early
1077 stopping mechanism greatly reduces the running time and improves model perfor-
1078 mance. Additionally the running time, for training and inference, is dependent on the
1079 size of the mini-batches. The memory requirements are also linear in the number of
1080 cells and genes for training and inference. Since DISCERN is trained on mini-batches
1081 the memory requirements can also be slightly adjusted by changing the mini-batch
1082 size during training or inference.

1083 4.6. Hyperparameters

1084 As outlined in the architecture section of the methods and depicted in Figure S1,
1085 DISCERN features several learnable hyperparameters. The complexity of the hyper-
1086 parameter search space is a potential downside of DISCERN, if these hyperparameters
1087 would be unstable across different datasets or in other words, would require constant
1088 tuning. Fortunately, DISCERN’s hyperparameters are very stable across the multi-
1089 titude of datasets tested in this manuscript, which we will outline in this paragraph.
1090 Naturally, there is no rule without an exception, which in this manuscript are the
1091 COVID-19 datasets that required optimization for several hyperparameters.

1092 *Constant hyperparameters.* DISCERN features a number of hyper-parameters that
1093 can be tuned through hyperparameter optimization (see below for details). Most
1094 of them have default values that yield reasonable performance across the different
1095 datasets we used and are being kept constant across experiments, including the COVID-
1096 19 dataset. Those constant hyperparameters are: the choice of the reconstruction loss
1097 (Huber loss), activation functions (Mish), CLN for the conditioning, number of fully-
1098 connected layers (3) and their size (1024, 512, 256 and 256, 512, 1024 neurons for the
1099 encoder and the decoder respectively), number of latent dimensions (48), learning rate
1100 (1×10^{-3}), decay rates β_1 and β_2 of Rectified Adam (0.85 and 0.95 respectively), batch
1101 size (192), label smoothing for our custom cross entropy loss (0.1), dropout rates (0.4
1102 in the encoder and 0 in the decoder), delta parameter of the Huber loss (9.0), weight
1103 on the penalty on the randomness of the encoder λ_{sigma} (1×10^{-8}), weight on the cross
1104 entropy loss term $\lambda_{dropout}$ (1×10^5), weight on the MMD penalty term λ_{prior} (1500).

1105 *Dataset-specific hyperparameters.* The optimal value of the L2 regularization applied
1106 on the weights of our custom CLN highly depends on the dataset at hand and thus
1107 requires dataset-specific tuning. For datasets with a very small variance in cell compo-
1108 sitions the L2 CLN regularization can be turned off (weight set to 0). When datasets
1109 have different compositions the L2 CLN regularization requires higher values (typically
1110 between 1×10^{-3} and 0.2).

1111 *COVID-19 hyperparameters.* For the experiments with COVID-19 datasets slightly
1112 adjusted hyperparameters were used: learning rate of 6e-3, label smoothing for our
1113 custom crossentropy loss of 0.05, weight on the penalty on the randomness of the
1114 encoder λ_{sigma} (1e-4), weight on the cross entropy loss term λ_{dropout} (2e3), weight on
1115 the MMD penalty term λ_{prior} (2000).

1116 *Hyperparameter optimization.* DISCERN implements different techniques for hyper-
1117 parameter optimization by using the ray[tune] library [83]. For most use cases the
1118 model does not require hyperparameter tuning and the default parameter should be
1119 sufficient. However, DISCERN has a generic interface and supports nearly all tech-
1120 niques implemented in ray[tune]. The initial hyperparameters were found using grid
1121 search. The loss used for the hyperparameter selection is the classification perfor-
1122 mance of a Random Forest classifier trying to classify real vs. auto-encoded cells.
1123 Classification performance was measured using the area under the receiver operating
1124 characteristic curve (AUC / AUROC).

1125 4.7. Competing algorithms and methods

1126 We briefly discuss competing methods and have compared their performance to
1127 DISCERN in the results section. These algorithms can be grouped into two categories,
1128 i) imputation algorithms that were developed to estimate drop-out gene expression
1129 based on dataset inherent information (MAGIC, DCA, scImpute) and ii) algorithms
1130 designed for batch correction that we have modified or extended to reconstruct gene
1131 expression, although this is not their intended use (Seurat, scGen). Given the latter,
1132 it is clear that DISCERN could be used purely for batch correction in latent space, a
1133 subject beyond the scope of this manuscript.

1134 *MAGIC.* [13] Markov affinity-based graph imputation of cells (MAGIC) denoises and
1135 imputes the single-cell count matrix using data diffusion-based information sharing.
1136 The construction of a good similarity metric is challenging for finding biologically
1137 similar cells due to high sparsity. MAGIC finds a good similarity metric using a so-
1138 phisticated graph-based approach that builds less-noisy cell-cell affinities and informa-
1139 tion sharing across cells. A particular focus of MAGIC was to understand gene-gene
1140 relationships and to characterize other dynamics in biological systems. MAGIC is
1141 provided as a Python package.

1142 *DCA.* [11] is a deep learning-based method for denoising single-cell count matrices.
1143 DCA is implemented in Python and uses an autoencoder with a Zero-Inflated Negative
1144 Binomial (ZINB) loss function. For each gene, DCA computes gene-specific param-
1145 eters of ZINB distribution, namely dropout, dispersion and mean. By modeling gene
1146 distributions as a noise model and also computing dropout probabilities of each gene,
1147 DCA is able to denoise and impute the missing counts by identifying and correcting
1148 dropout events.

1149 *scImpute*. [12] Similarly to MAGIC, *scImpute* focuses on identifying cells that are
1150 similar, which is challenging due to the high sparsity of single-cell count matrices.
1151 *scImpute* is a statistical model using a three step process to impute scRNA-seq data. In
1152 the first step spectral clustering is applied on principal components to find neighbors,
1153 which later can be used to detect and impute dropout values. In the second step
1154 *scImpute* fits a mixture model of a Gamma distribution and a Normal distribution
1155 to distinguish technical and biological dropouts. In the last step, the model uses a
1156 regression model for each cell to impute the expression of genes with high probability
1157 of dropout. With this approach, *scImpute* avoids hallucinations and keeps the gene
1158 expression distribution. *scImpute* is provided as an R package.

1159 *Seurat*. [26] is an open-source toolkit for the analysis of single cell RNA-sequencing
1160 data. In addition to general analysis functions, *Seurat* offers batch-correction function-
1161 ality. *Seurat* uses canonical correlation analysis to construct this lower dimensional
1162 representation and tries to find neighbors between batches in this shared space. These
1163 anchors are filtered considering the local neighborhood of the cell pairs and remain-
1164 ing anchors are finally used to construct correction vectors for all cells in this low
1165 dimensional representation. While *Seurats* is intended to work in a lower dimensional
1166 representation, it can also be used to reconstruct the expression information from this
1167 lower dimensional representation. *Seurat* is provided as an R package.

1168 *scGen*. [23] is a variational autoencoder based deep learning method with a focus on
1169 learning features that help distinguish responding and non-responding genes and cells.
1170 *scGen* constructs a latent space in which it estimates perturbation vectors associated
1171 with a change between different conditions. Since *scGen* models the perturbation and
1172 infection responses in single cells, it is focused on in-silico screening with the use of
1173 cells coming from healthy samples. It can also be used for batch correction. For batch
1174 correction, and unlike DISCERN or *Seurat*, *scGen* uses both batch and cell type labels.
1175 *scGen* is built using the *scvi-tools* toolbox and implemented in python and pytorch.

1176 *Multigrade*. [19] *multigrade* is an autoencoder based deep learning method developed
1177 for the integration of different modalities to improve single cell RNA-seq downstream
1178 analysis, mainly clustering. The main focus is the integration of CITE-seq protein
1179 abundance since it is often available together with scRNA-seq. They use individual
1180 encoders for each modality and build a shared latent representation by partially sharing
1181 the decoder. *Multigrade* is built using the *scvi-tools* toolbox and implemented in
1182 python and pytorch.

1183 *scVI*. [36] is a variational autoencoder-based deep learning method developed for sev-
1184 eral single cell analysis approaches like batch correction, clustering, and differential
1185 expression analysis. It models expression data using a zero-inflated negative binomial
1186 loss during the training. For comparison of *scVI* to other models, only the batch
1187 correction functionality was used. For the differential expression analysis we used
1188 the same workflow as for the other methods to allow for a fair comparison. *scVI* is
1189 implemented in python and pytorch.

1190 *CarDEC*. [14] is an autoencoder-based learning method developed for batch effect
1191 correction, denoising of expression data and cell clustering. The *CarDEC* pipeline
1192 computes highly variable genes across all batches and pre-trains an autoencoder to
1193 reconstruct the expression of these genes. In a second step, the weights are transferred

1194 to a bigger network, which is trained jointly on the highly variable and lowly variable
1195 genes using two reconstruction losses. Additionally, they include a self-supervised
1196 clustering loss in the latent space to improve batch mixing. CarDEC is implemented
1197 in python and Tensorflow.

1198 *DeepImpute*. [15] is an ensemble method consisting of multiple autoencoder-like deep
1199 neural networks, where each network is trained to learn the relationship between a
1200 set of input genes and a set of target genes. Input and target gene sets are selected
1201 based on correlation of gene expression values. The estimated expression values from
1202 each of the networks is combined to yield the final imputed dataset. DeepImpute is
1203 implemented in python and Tensorflow.

1204 *trVAE*. [35] is a variational autoencoder based deep learning method developed for
1205 the generation of unseen samples or conditions of single cell RNA-seq data. It uses an
1206 encoder with additional inputs for encoding the condition and a decoder which gets,
1207 together with the latent code, the target condition as input. To achieve a condition
1208 independence the first layer is regularized using maximum mean discrepancy. trVAE
1209 is implemented in python and Tensorflow.

1210 4.8. Evaluation metrics

1211 *t-SNE & UMAP*. For visualization of the datasets and to qualitatively assess the in-
1212 tegration performance tSNE and UMAP were used. Both methods are based on PCA
1213 representation and use non-linear representations to create a 2D representation of the
1214 data. We used the scanpy [70] implementation. Default settings were used in nearly
1215 all cases except: In the combined COVID-19 dataset analogue to Kobak *et al.*[84] the
1216 dataset was subset to 25 000 cells and tSNE was computed using a perplexity of 250,
1217 and a learning rate of 25 000/12. These positions were taken and used as input to tSNE
1218 of all cells using a perplexity of 30 a learning rate of (number of observations)/12 and
1219 a late exaggeration of 4.0 using Fit-SNE [85]. Clustering was performed using PARC
1220 [86] with default parameters except `dist_std_local=1.5` and `small_pop=300`. Meth-
1221 ods were changed here due to computation time issues for 350 000 cells. covid-blood
1222 data was analyzed using a learning rate of (number of observations)/6 a perplexity
1223 of (number of observations)/120 and `early_exaggeration=4`. Clustering was performed
1224 using default parameters except `knn=100` and `small_pop=100` to reduce the number
1225 of clusters with limited cell number. Clustering of the T helper cells in healthy blood
1226 was performed using coarse clustering with 30 nearest neighbors and leiden cluster-
1227 ing (<https://github.com/vtraag/leidenalg>) with a resolution of 0.6. Afterwards a
1228 combined cluster of IFN-regulated and TREG was reclustered using a resolution of 0.4
1229 and effector T cells were reclustered using a resolution of 0.8. Resolution was chosen
1230 to dissect the raw gene expression changes of known cell types.

1231 *Mean gene expression*. Mean gene expression was calculated as average over log-
1232 normalized expression over all cells, usually stratified by celltype. This evaluation
1233 of expression data consists of many data points where several have values close to
1234 zero, but could have a high weight on rank-based correlation methods. Thus Pearson
1235 correlation was used to evaluate the performance.

1236 *Differential gene expression.* Differential gene expression was performed using the
1237 scanpy [70] rank_gene_groups function using the t-test method for calculating sta-
1238 tistical significance on log-normalized expression data. Differential gene expression
1239 analysis was always performed under consideration of the cell type information. For
1240 comparison of differential gene expression analysis between conditions, the Pearson
1241 correlation was used. It is calculated either on the log₂ fold-change or in most cases
1242 on the t-statistics, computed during significance estimation. The data was compared
1243 using the t-statistics, because it aggregates information on both the variance and the
1244 change in mean expression. Thus it allows, roughly speaking, for simultaneously eval-
1245 uating the significance and the log₂ fold change. Usually all available genes were used
1246 for correlation, except in the in-silico gene removal experiment, where only the re-
1247 moved genes were considered. We used spearman rank correlation when all genes were
1248 available and pearson correlation otherwise.

1249 *Pathway analysis.* Pathway analysis or gene set enrichment analysis was done using
1250 the prerank function from gseapy [87] on the t-statistics, computed as described in the
1251 ‘Differential gene expression’ section of the methods. To this end, the gene set library
1252 “KEGG_2019_Human” provided by enrichr [88] was used. Top pathways were selected
1253 using the normalized enrichment score as previously described [87].

1254 *Gene regulation.* [48] The python implementation of the SCENIC (pySENIC) was
1255 used to infer regulons specific for CD4⁺ T helper cells. SCENIC infers a gene regula-
1256 tory network using GRNBoost2 and creates co-expression modules. The co-expression
1257 modules get associated with transcription factors using the transcription factor motif
1258 discovery tool RcisTarget. A pair of transcription factor and associated gene set is
1259 called a regulon. For each cell, the regulons get scored using the AUCell algorithm
1260 to examine if a cell is affected by the regulon. We used default parameters of the
1261 pySENIC implementation.

1262 *Silhouette Score.* [89] - is a measure to evaluate clustering performance by comparing
1263 the mean intra-cluster distance to the mean nearest-cluster distance. The Silhouette
1264 score is computed for batch and cell type labels on the scaled and PCA-transformed
1265 data using a varying number of principal components (interval [10, 50]). The score is
1266 defined in the interval [-1, 1], where a positive value indicates separated clusters, a
1267 value of zero signifies cluster overlap, and a negative value when the closest cluster is
1268 not the wrong cluster. For accessing batch mixing a low, close to zero, value is best,
1269 while for cell type clusters a value close to 1 is best. The scikit-learn implementation
1270 was used.

1271 *Adjusted Rand Index.* [90] - The Rand index estimates the similarity between two
1272 clusterings by comparing all possible pairings of samples. The Adjusted Rand Index
1273 is adjusted for chance, such that a random labeling would result in a value close to 0,
1274 while a perfect clustering yields a score of 1. The Adjusted Rand Index is computed
1275 on the result of the leiden clustering algorithm using 20 different resolution parameters
1276 in the interval of [0.1, 30]. The best value (lowest for batch mixing, highest for cell
1277 type clustering) was used as the final score. The neighborhood graph for the leiden
1278 clustering algorithm is computed on scaled and PCA-transformed values, similar to
1279 the silhouette score, for a varying number of principal components (interval [10, 5]).
1280 The scikit-learn implementation was used.

1281 *Adjusted Mutual Information.* [91] - Mutual Information measures the similarity between two clusterings by computing the sizes of the intersection of all possible cluster
1282 label pairs. The Adjusted Mutual Information is adjusted for chance, such that a
1283 random labeling would result in a value close to 0, while a perfect clustering yields a
1284 score of 1. Additionally, this accounts for the fact that Mutual Information is generally
1285 higher for clusterings with larger numbers of clusters. The AMI was computed on
1286 clustering results as described for the Adjusted Rand Index. The scikit-learn implementation
1287 was used.
1288

1289 *COVID-19 classification*

1290 To evaluate the importance of the cell types found in the covid-blood-severity-hq
1291 dataset after reconstruction with DISCERN, the fraction for all T cell subtypes was
1292 used to predict the disease severity, as provided in [53]. The data was classified using a
1293 Gradient boosting classifier ([92], implemented in scikit-learn v1.0.2, default settings)
1294 using 25 rounds of leave-one-out cross-validation (LOOCV). Each round consists of n
1295 training-prediction iterations with $n - 1$ samples for training and 1 sample for testing,
1296 such that after one round prediction results for all n samples could be evaluated.
1297 We chose LOOCV over k -fold cross-validation and testing due to the limited size
1298 of the dataset, consisting of only 71 patients. We used pycm ([93], v3.3) for the
1299 performance evaluation. The final evaluation was done using the accuracy and F1 score
1300 as provided by pycm. The area under the receiver operating characteristic (AUROC)
1301 curve is computed with scikit-learn. Before training the classifiers a forward feature
1302 selection was performed using the SequentialFeatureSelector implemented in scikit-
1303 learn with default parameters. In total four experiments were performed. In the
1304 first experiment, classification with three disease categories (mild, moderate, severe)
1305 was used. Patients who died were excluded. For the other two experiments only
1306 patients with asymptomatic, mild, severe and critical symptoms were included. In all
1307 experiments the asymptomatic and mild category was merged to mild and severe and
1308 critical to severe.

1309 **References**

- 1310 [1] N. Editorial, Method of the year 2013, *Nat. Methods* 11 (1) (2014) 1.
- 1311 [2] Y. Zhao, U. Panzer, S. Bonn, C. F. Krebs, Single-cell biology to decode the
1312 immune cellular composition of kidney inflammation, *Cell and tissue research*
1313 385 (2) (2021) 435–443.
- 1314 [3] M. Stoeckius, C. Hafemeister, W. Stephenson, B. Houck-Loomis, P. K. Chat-
1315 topadhyay, H. Swerdlow, R. Satija, P. Smibert, Simultaneous epitope and tran-
1316 scriptome measurement in single cells, *Nature methods* 14 (9) (2017) 865–868.
- 1317 [4] A. A. Tu, T. M. Gierahn, B. Monian, D. M. Morgan, N. K. Mehta, B. Ruiter,
1318 W. G. Shreffler, A. K. Shalek, J. C. Love, Tcr sequencing paired with massively
1319 parallel 3' rna-seq reveals clonotypic t cell signatures, *Nature immunology* 20 (12)
1320 (2019) 1692–1699.
- 1321 [5] J. A. Pai, A. T. Satpathy, High-throughput and single-cell t cell receptor sequenc-
1322 ing technologies, *Nature Methods* 18 (8) (2021) 881–892.

- 1323 [6] S. Oller-Moreno, K. Kloiber, P. Machart, S. Bonn, Algorithmic advances in ma-
1324 chine learning for single-cell expression analysis, *Current Opinion in Systems*
1325 *Biology* 25 (2021) 27–33.
- 1326 [7] D. Lähnemann, J. Köster, E. Szczurek, D. J. McCarthy, S. C. Hicks, M. D.
1327 Robinson, C. A. Vallejos, K. R. Campbell, N. Beerenwinkel, A. Mahfouz, et al.,
1328 Eleven grand challenges in single-cell data science, *Genome biology* 21 (1) (2020)
1329 1–35.
- 1330 [8] X. Wang, Y. He, Q. Zhang, X. Ren, Z. Zhang, Direct comparative analyses of
1331 10x genomics chromium and smart-seq2, *Genomics, proteomics & bioinformatics*
1332 19 (2) (2021) 253–266.
- 1333 [9] A. K. Shalek, R. Satija, J. Shuga, J. J. Trombetta, D. Gennert, D. Lu, P. Chen,
1334 R. S. Gertner, J. T. Gaublotte, N. Yosef, et al., Single-cell rna-seq reveals
1335 dynamic paracrine control of cellular variation, *Nature* 510 (7505) (2014) 363–
1336 369.
- 1337 [10] W. Hou, Z. Ji, H. Ji, S. C. Hicks, A systematic evaluation of single-cell rna-
1338 sequencing imputation methods, *Genome biology* 21 (1) (2020) 1–30.
- 1339 [11] G. Eraslan, L. M. Simon, M. Mircea, N. S. Mueller, F. J. Theis, Single-cell rna-seq
1340 denoising using a deep count autoencoder, *Nature communications* 10 (1) (2019)
1341 1–14.
- 1342 [12] W. V. Li, J. J. Li, An accurate and robust imputation method scimpute for
1343 single-cell rna-seq data, *Nature communications* 9 (1) (2018) 1–9.
- 1344 [13] D. Van Dijk, R. Sharma, J. Nainys, K. Yim, P. Kathail, A. J. Carr, C. Burdziak,
1345 K. R. Moon, C. L. Chaffer, D. Pattabiraman, et al., Recovering gene interactions
1346 from single-cell data using data diffusion, *Cell* 174 (3) (2018) 716–729.
- 1347 [14] J. Lakkis, D. Wang, Y. Zhang, G. Hu, K. Wang, H. Pan, L. Ungar, M. P. Reilly,
1348 X. Li, M. Li, A joint deep learning model for simultaneous batch effect correction,
1349 denoising and clustering in single-cell transcriptomics, *bioRxiv* (2020).
- 1350 [15] C. Arisdakessian, O. Poirion, B. Yunits, X. Zhu, L. X. Garmire, Deepimpute:
1351 an accurate, fast, and scalable deep neural network method to impute single-cell
1352 rna-seq data, *Genome biology* 20 (1) (2019) 1–14.
- 1353 [16] Z.-H. Wen, J. L. Langsam, L. Zhang, W. Shen, X. Zhou, A bayesian factorization
1354 method to recover single-cell rna sequencing data, *Cell reports methods* 2 (1)
1355 (2022) 100133.
- 1356 [17] T. Peng, Q. Zhu, P. Yin, K. Tan, Scrabble: single-cell rna-seq imputation con-
1357 strained by bulk rna-seq data, *Genome biology* 20 (1) (2019) 1–12.
- 1358 [18] Z. Hu, S. Zu, J. S. Liu, Simples: a single-cell rna sequencing imputation strategy
1359 preserving gene modules and cell clusters variation, *NAR genomics and bioinfor-*
1360 *matics* 2 (4) (2020) lqaa077.
- 1361 [19] M. Lotfollahi, A. Litinetskaya, F. J. Theis, Multigrade: single-cell multi-omic data
1362 integration, *bioRxiv* (2022).

- 1363 [20] T. H. Kim, X. Zhou, M. Chen, Demystifying “drop-outs” in single-cell umi data,
1364 *Genome biology* 21 (1) (2020) 1–19.
- 1365 [21] R. Jiang, T. Sun, D. Song, J. J. Li, Statistics or biology: the zero-inflation con-
1366 troversy about scrna-seq data, *Genome biology* 23 (1) (2022) 1–24.
- 1367 [22] M. Marouf, P. Machart, V. Bansal, C. Kilian, D. S. Magruder, C. F. Krebs,
1368 S. Bonn, Realistic in silico generation and augmentation of single-cell rna-seq
1369 data using generative adversarial networks, *Nature communications* 11 (1) (2020)
1370 1–12.
- 1371 [23] M. Lotfollahi, F. A. Wolf, F. J. Theis, scgen predicts single-cell perturbation
1372 responses, *Nature methods* 16 (8) (2019) 715–721.
- 1373 [24] Y. Zhao, C. Kilian, J.-E. Turner, L. Bosurgi, K. Roedl, P. Bartsch, A.-C. Gnirck,
1374 F. Cortesi, C. Schultheiß, M. Hellmig, et al., Clonal expansion and activation of
1375 tissue-resident memory-like th17 cells expressing gm-csf in the lungs of patients
1376 with severe covid-19, *Science Immunology* 6 (56) (2021) eabf6692.
- 1377 [25] I. Tolstikhin, O. Bousquet, S. Gelly, B. Schoelkopf, Wasserstein auto-encoders,
1378 arXiv preprint arXiv:1711.01558 (2017).
- 1379 [26] T. Stuart, A. Butler, P. Hoffman, C. Hafemeister, E. Papalexi, W. M. Mauck III,
1380 Y. Hao, M. Stoeckius, P. Smibert, R. Satija, Comprehensive integration of single-
1381 cell data, *Cell* 177 (7) (2019) 1888–1902.
- 1382 [27] V. Dumoulin, J. Shlens, M. Kudlur, A learned representation for artistic style,
1383 arXiv preprint arXiv:1610.07629 (2016).
- 1384 [28] J. L. Ba, J. R. Kiros, G. E. Hinton, Layer normalization, arXiv preprint
1385 arXiv:1607.06450 (2016).
- 1386 [29] Å. Segerstolpe, A. Palasantza, P. Eliasson, E.-M. Andersson, A.-C. Andréasson,
1387 X. Sun, S. Picelli, A. Sabirsh, M. Clausen, M. K. Bjursell, et al., Single-cell
1388 transcriptome profiling of human pancreatic islets in health and type 2 diabetes,
1389 *Cell metabolism* 24 (4) (2016) 593–607.
- 1390 [30] D. Grün, M. J. Muraro, J.-C. Boisset, K. Wiebrands, A. Lyubimova, G. Dharm-
1391 madhikari, M. van den Born, J. Van Es, E. Jansen, H. Clevers, et al., De novo
1392 prediction of stem cell identity using single-cell transcriptome data, *Cell stem cell*
1393 19 (2) (2016) 266–277.
- 1394 [31] N. Lawlor, J. George, M. Bolisetty, R. Kursawe, L. Sun, V. Sivakamasundari,
1395 I. Kycia, P. Robson, M. L. Stitzel, Single-cell transcriptomes identify human islet
1396 cell signatures and reveal cell-type-specific expression changes in type 2 diabetes,
1397 *Genome research* 27 (2) (2017) 208–222.
- 1398 [32] M. J. Muraro, G. Dharmadhikari, D. Grün, N. Groen, T. Dielen, E. Jansen,
1399 L. Van Gurp, M. A. Engelse, F. Carlotti, E. J. De Koning, et al., A single-cell
1400 transcriptome atlas of the human pancreas, *Cell systems* 3 (4) (2016) 385–394.
- 1401 [33] M. Baron, A. Veres, S. L. Wolock, A. L. Faust, R. Gaujoux, A. Vetere, J. H. Ryu,
1402 B. K. Wagner, S. S. Shen-Orr, A. M. Klein, et al., A single-cell transcriptomic
1403 map of the human and mouse pancreas reveals inter-and intra-cell population
1404 structure, *Cell systems* 3 (4) (2016) 346–360.

- 1405 [34] Z. Fu, E. R Gilbert, D. Liu, Regulation of insulin synthesis and secretion and
1406 pancreatic beta-cell dysfunction in diabetes, *Current diabetes reviews* 9 (1) (2013)
1407 25–53.
- 1408 [35] M. Lotfollahi, M. Naghipourfar, F. J. Theis, F. A. Wolf, Conditional out-
1409 of-distribution generation for unpaired data using transfer vae, *Bioinformatics*
1410 36 (Supplement_2) (2020) i610–i617.
- 1411 [36] A. Gayoso, R. Lopez, G. Xing, P. Boyeau, V. Valiollah Pour Amiri, J. Hong,
1412 K. Wu, M. Jayasuriya, E. Mehlman, M. Langevin, et al., A python library for
1413 probabilistic analysis of single-cell omics data, *Nature Biotechnology* 40 (2) (2022)
1414 163–166.
- 1415 [37] G. A. Bouland, A. Mahfouz, M. J. Reinders, Differential dropout analysis captures
1416 biological variation in single-cell rna sequencing data, *Biorxiv* (2021).
- 1417 [38] M. Slyper, C. Porter, O. Ashenberg, J. Waldman, E. Drokhyansky, I. Wakiro,
1418 C. Smillie, G. Smith-Rosario, J. Wu, D. Dionne, et al., A single-cell and single-
1419 nucleus rna-seq toolbox for fresh and frozen human tumors, *Nature medicine*
1420 26 (5) (2020) 792–802.
- 1421 [39] G. C. Linderman, J. Zhao, M. Roulis, P. Bielecki, R. A. Flavell, B. Nadler,
1422 Y. Kluger, Zero-preserving imputation of single-cell rna-seq data, *Nature Com-
1423 munications* 13 (1) (2022) 1–11.
- 1424 [40] E. Bakos, C. A. Thaiss, M. P. Kramer, S. Cohen, L. Radomir, I. Orr, N. Kaushan-
1425 sky, A. Ben-Nun, S. Becker-Herman, I. Shachar, *Ccr2* regulates the immune re-
1426 sponse by modulating the interconversion and function of effector and regulatory
1427 t cells, *The Journal of Immunology* 198 (12) (2017) 4659–4671.
- 1428 [41] G. Monaco, B. Lee, W. Xu, S. Mustafah, Y. Y. Hwang, C. Carré, N. Burdin,
1429 L. Visan, M. Ceccarelli, M. Poidinger, et al., Rna-seq signatures normalized by
1430 mrna abundance allow absolute deconvolution of human immune cell types, *Cell
1431 reports* 26 (6) (2019) 1627–1640.
- 1432 [42] V. A. Traag, L. Waltman, N. J. Van Eck, From louvain to leiden: guaranteeing
1433 well-connected communities, *Scientific reports* 9 (1) (2019) 1–12.
- 1434 [43] M. Croft, Control of immunity by the tnfr-related molecule ox40 (cd134), *Annual
1435 review of immunology* 28 (2009) 57–78.
- 1436 [44] T. Riaz, L. M. Sollid, I. Olsen, G. A. de Souza, Quantitative proteomics of gut-
1437 derived th1 and th1/th17 clones reveal the presence of cd28+ nkg2d-th1 cytotoxic
1438 cd4+ t cells, *Molecular & Cellular Proteomics* 15 (3) (2016) 1007–1016.
- 1439 [45] L. Peng, Y. Chen, Q. Ou, X. Wang, N. Tang, Lncrna miat correlates with im-
1440 mune infiltrates and drug reactions in hepatocellular carcinoma, *International
1441 immunopharmacology* 89 (2020) 107071.
- 1442 [46] D. P. Saraiva, A. Jacinto, P. Borralho, S. Braga, M. G. Cabral, Hla-dr in cy-
1443 totoxic t lymphocytes predicts breast cancer patients’ response to neoadjuvant
1444 chemotherapy, *Frontiers in immunology* (2018) 2605.

- 1445 [47] M. S. Lee, K. Hanspers, C. S. Barker, A. P. Korn, J. M. McCune, Gene expression
1446 profiles during human cd4+ t cell differentiation, *International immunology* 16 (8)
1447 (2004) 1109–1124.
- 1448 [48] S. Aibar, C. B. González-Blas, T. Moerman, V. A. Huynh-Thu, H. Imrichova,
1449 G. Hulselmans, F. Rambow, J.-C. Marine, P. Geurts, J. Aerts, et al., Scenic:
1450 single-cell regulatory network inference and clustering, *Nature methods* 14 (11)
1451 (2017) 1083–1086.
- 1452 [49] A. M. Thornton, J. Lu, P. E. Korty, Y. C. Kim, C. Martens, P. D. Sun, E. M. She-
1453 vach, Helios+ and helios- treg subpopulations are phenotypically and functionally
1454 distinct and express dissimilar tcr repertoires, *European journal of immunology*
1455 49 (3) (2019) 398–412.
- 1456 [50] C. Imbratta, H. Hussein, F. Andris, G. Verdeil, c-maf, a swiss army knife for
1457 tolerance in lymphocytes, *Frontiers in immunology* 11 (2020) 206.
- 1458 [51] X. O. Yang, B. P. Pappu, R. Nurieva, A. Akimzhanov, H. S. Kang, Y. Chung,
1459 L. Ma, B. Shah, A. D. Panopoulos, K. S. Schluns, et al., T helper 17 lineage dif-
1460 ferentiation is programmed by orphan nuclear receptors ror α and ror γ , *Immunity*
1461 28 (1) (2008) 29–39.
- 1462 [52] K. Street, D. Risso, R. B. Fletcher, D. Das, J. Ngai, N. Yosef, E. Purdom, S. Du-
1463 doit, Slingshot: cell lineage and pseudotime inference for single-cell transcrip-
1464 tomics, *BMC genomics* 19 (1) (2018) 1–16.
- 1465 [53] E. Stephenson, G. Reynolds, R. A. Botting, F. J. Calero-Nieto, M. D. Morgan,
1466 Z. K. Tuong, K. Bach, W. Sunagak, K. B. Worlock, M. Yoshida, et al., Single-cell
1467 multi-omics analysis of the immune response in covid-19, *Nature medicine* 27 (5)
1468 (2021) 904–916.
- 1469 [54] M. Ota, Y. Nagafuchi, H. Hatano, K. Ishigaki, C. Terao, Y. Takeshima,
1470 H. Yanaoka, S. Kobayashi, M. Okubo, H. Shirai, et al., Dynamic landscape of
1471 immune cell-specific gene regulation in immune-mediated diseases, *Cell* 184 (11)
1472 (2021) 3006–3021.
- 1473 [55] G. Meyer Zu Horste, C. Wu, C. Wang, L. Cong, M. Pawlak, Y. Lee, W. Elyaman,
1474 S. Xiao, A. Regev, V. Kuchroo, Rbpj controls development of pathogenic th17
1475 cells by regulating il-23 receptor expression. *cell rep* 16 (2): 392–404 (2016).
- 1476 [56] S. De Biasi, M. Meschiari, L. Gibellini, C. Bellinazzi, R. Borella, L. Fidanza,
1477 L. Gozzi, A. Iannone, D. Lo Tartaro, M. Mattioli, et al., Marked t cell activa-
1478 tion, senescence, exhaustion and skewing towards th17 in patients with covid-19
1479 pneumonia, *Nature communications* 11 (1) (2020) 1–17.
- 1480 [57] B. J. Meckiff, C. Ramírez-Suástegui, V. Fajardo, S. J. Chee, A. Kusnadi, H. Si-
1481 mon, S. Eschweiler, A. Grifoni, E. Pelosi, D. Weiskopf, et al., Imbalance of reg-
1482 ulatory and cytotoxic sars-cov-2-reactive cd4+ t cells in covid-19, *Cell* 183 (5)
1483 (2020) 1340–1353.
- 1484 [58] L. Loyal, S. Warth, K. Jürchott, F. Mölder, C. Nikolaou, N. Babel, M. Nienen,
1485 S. Durlanik, R. Stark, B. Kruse, et al., Slamf7 and il-6r define distinct cytotoxic
1486 versus helper memory cd8+ t cells, *Nature communications* 11 (1) (2020) 1–12.

- 1487 [59] J. Yang, M. Zhong, E. Zhang, K. Hong, Q. Yang, D. Zhou, J. Xia, Y.-Q. Chen,
1488 M. Sun, B. Zhao, et al., Broad phenotypic alterations and potential dysfunc-
1489 tion of lymphocytes in individuals clinically recovered from covid-19, *Journal of*
1490 *Molecular Cell Biology* 13 (3) (2021) 197–209.
- 1491 [60] T. S. Andrews, M. Hemberg, False signals induced by single-cell imputation,
1492 *F1000Research* 7 (2018).
- 1493 [61] M. Lotfollahi, M. Naghipourfar, M. D. Luecken, M. Khajavi, M. Büttner, M. Wa-
1494 genstetter, Ž. Avsec, A. Gayoso, N. Yosef, M. Interlandi, et al., Mapping single-
1495 cell data to reference atlases by transfer learning, *Nature Biotechnology* 40 (1)
1496 (2022) 121–130.
- 1497 [62] M. D. Luecken, M. Büttner, K. Chaichoompu, A. Danese, M. Interlandi, M. F.
1498 Müller, D. C. Strobl, L. Zappia, M. Dugas, M. Colomé-Tatché, et al., Benchmark-
1499 ing atlas-level data integration in single-cell genomics, *Nature methods* 19 (1)
1500 (2022) 41–50.
- 1501 [63] C. Wagner, M. Griesel, A. Mikolajewska, A. Mueller, M. Nothacker, K. Kley,
1502 M.-I. Metzendorf, A.-L. Fischer, M. Kopp, M. Stegemann, et al., Systemic cor-
1503 ticosteroids for the treatment of covid-19, *Cochrane Database of Systematic Re-*
1504 *views* (8) (2021).
- 1505 [64] W. Chen, J. Luo, Y. Ye, R. Hoyle, W. Liu, R. Borst, S. Kazani, E. A. Shikatani,
1506 V. J. Erpenbeck, I. D. Pavord, et al., The roles of type 2 cytotoxic t cells in
1507 inflammation, tissue remodeling, and prostaglandin (pg) d2 production are at-
1508 tenuated by pgd2 receptor 2 antagonism, *The Journal of Immunology* 206 (11)
1509 (2021) 2714–2724.
- 1510 [65] S. Lab, *panc8.SeuratData: Eight Pancreas Datasets Across Five Technologies*, r
1511 package version 3.0.2 (2019).
- 1512 [66] J. Ding, X. Adiconis, S. K. Simmons, M. S. Kowalczyk, C. C. Hession, N. D.
1513 Marjanovic, T. K. Hughes, M. H. Wadsworth, T. Burks, L. T. Nguyen, et al.,
1514 Systematic comparison of single-cell and single-nucleus rna-sequencing methods,
1515 *Nature biotechnology* 38 (6) (2020) 737–746.
- 1516 [67] A. Gayoso, R. Lopez, G. Xing, P. Boyeau, K. Wu, M. Jayasuriya, E. Melhman,
1517 M. Langevin, Y. Liu, J. Samaran, et al., *Scvi-tools: A library for deep probabilis-*
1518 *tic analysis of single-cell omics data*, *bioRxiv* (2021).
- 1519 [68] R. Menon, A. S. Bomback, B. B. Lake, C. Stutzke, S. M. Grewenow, S. Menez,
1520 V. D. D’Agati, S. Jain, R. Knight, S. H. Lecker, et al., Integrated single-cell se-
1521 quencing and histopathological analyses reveal diverse injury and repair responses
1522 in a participant with acute kidney injury: a clinical-molecular-pathologic corre-
1523 lation, *Kidney International* 101 (6) (2022) 1116–1125.
- 1524 [69] G. X. Zheng, J. M. Terry, P. Belgrader, P. Ryvkin, Z. W. Bent, R. Wilson, S. B.
1525 Ziraldo, T. D. Wheeler, G. P. McDermott, J. Zhu, et al., Massively parallel digital
1526 transcriptional profiling of single cells, *Nature communications* 8 (1) (2017) 1–12.
- 1527 [70] F. A. Wolf, P. Angerer, F. J. Theis, *Scanpy: large-scale single-cell gene expression*
1528 *data analysis*, *Genome biology* 19 (1) (2018) 1–5.

- 1529 [71] Y. Le Cun, F. Fogelman-Soulié, Modèles connexionnistes de l'apprentissage, *Intellectica* 2 (1) (1987) 114–143.
1530
- 1531 [72] G. E. Hinton, R. Zemel, Autoencoders, minimum description length and
1532 helmholtz free energy, *Advances in neural information processing systems* 6
1533 (1993).
- 1534 [73] D. P. Kingma, M. Welling, Auto-encoding variational bayes, *arXiv preprint*
1535 *arXiv:1312.6114* (2013).
- 1536 [74] C. Villani, *Optimal transport: old and new*, Vol. 338, Springer, 2009.
- 1537 [75] M. Arjovsky, S. Chintala, L. Bottou, Wasserstein generative adversarial networks,
1538 in: *International conference on machine learning*, PMLR, 2017, pp. 214–223.
- 1539 [76] P. K. Rubenstein, B. Schoelkopf, I. Tolstikhin, On the latent space of wasserstein
1540 auto-encoders, *arXiv preprint arXiv:1802.03761* (2018).
- 1541 [77] D. Risso, F. Perraudeau, S. Gribkova, S. Dudoit, J.-P. Vert, A general and flexible
1542 method for signal extraction from single-cell rna-seq data, *Nature communications*
1543 9 (1) (2018) 1–17.
- 1544 [78] P. J. Huber, Robust estimation of a location parameter, *Annals Mathematics*
1545 *Statistics* (1964).
- 1546 [79] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, A. Smola, A kernel
1547 two-sample test, *Journal of Machine Learning Research* 13 (25) (2012) 723–773.
1548 URL <http://jmlr.org/papers/v13/gretton12a.html>
- 1549 [80] D. Misra, Mish: A self regularized non-monotonic activation function, *arXiv*
1550 *preprint arXiv:1908.08681* (2019).
- 1551 [81] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, J. Han, On the variance of the
1552 adaptive learning rate and beyond, *arXiv preprint arXiv:1908.03265* (2019).
- 1553 [82] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, *arXiv preprint*
1554 *arXiv:1412.6980* (2014).
- 1555 [83] R. Liaw, E. Liang, R. Nishihara, P. Moritz, J. E. Gonzalez, I. Stoica, Tune: A
1556 research platform for distributed model selection and training, *arXiv preprint*
1557 *arXiv:1807.05118* (2018).
- 1558 [84] D. Kobak, P. Berens, The art of using t-sne for single-cell transcriptomics, *Nature*
1559 *communications* 10 (1) (2019) 1–14.
- 1560 [85] G. C. Linderman, M. Rachh, J. G. Hoskins, S. Steinerberger, Y. Kluger, Fast
1561 interpolation-based t-sne for improved visualization of single-cell rna-seq data,
1562 *Nature methods* 16 (3) (2019) 243–245.
- 1563 [86] S. V. Stassen, D. M. Siu, K. C. Lee, J. W. Ho, H. K. So, K. K. Tsia, Parc:
1564 ultrafast and accurate clustering of phenotypic data of millions of single cells,
1565 *Bioinformatics* 36 (9) (2020) 2778–2786.

- 1566 [87] A. Subramanian, P. Tamayo, V. K. Mootha, S. Mukherjee, B. L. Ebert, M. A.
1567 Gillette, A. Paulovich, S. L. Pomeroy, T. R. Golub, E. S. Lander, et al., Gene set
1568 enrichment analysis: a knowledge-based approach for interpreting genome-wide
1569 expression profiles, *Proceedings of the National Academy of Sciences* 102 (43)
1570 (2005) 15545–15550.
- 1571 [88] E. Y. Chen, C. M. Tan, Y. Kou, Q. Duan, Z. Wang, G. V. Meirelles, N. R. Clark,
1572 A. Ma’ayan, *Enrichr*: interactive and collaborative html5 gene list enrichment
1573 analysis tool, *BMC bioinformatics* 14 (1) (2013) 1–14.
- 1574 [89] P. J. Rousseeuw, *Silhouettes*: a graphical aid to the interpretation and validation
1575 of cluster analysis, *Journal of computational and applied mathematics* 20 (1987)
1576 53–65.
- 1577 [90] L. Hubert, P. Arabie, *Comparing partitions*, *Journal of classification* 2 (1) (1985)
1578 193–218.
- 1579 [91] N. X. Vinh, J. Epps, Bailey, j2738784: Information theoretic measures for clus-
1580 terings comparison: variants, properties, normalization and correction for chance.
1581 vol. 11, *J Mach Learn Res* (2010) 2837–2854.
- 1582 [92] J. H. Friedman, *Greedy function approximation: a gradient boosting machine*,
1583 *Annals of statistics* (2001) 1189–1232.
- 1584 [93] S. Haghghi, M. Jasemi, S. Hessabi, A. Zolanvari, *Pycm*: Multiclass confusion
1585 matrix library in python, *Journal of Open Source Software* 3 (25) (2018) 729.