

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,f}, 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*

^f*Immunology and Allergy Unit, Department of Medicine Solna, Karolinska Institute, 17176 Stockholm, Sweden*

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 novel 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,

*Corresponding authors

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

¹Authors contributed equally

cell type identification, expression reconstruction, Deep Learning, Machine 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, and scImpute [11, 12, 13]. DCA is an autoencoder-based method for denoising and imputation of scRNAseq 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. All of these algorithms use information from similar cells with measured expression of the same dataset for imputation.

While current imputation methods provide improved gene expression information, they still rely on the comparison of similar cells with largely absent gene

39 expression information. Genes that are not expressed in neighboring cells can-
40 not be imputed, limiting the value of classical imputation. In an ideal case, it
41 would be possible to obtain information of the expected true gene expression per
42 cell, or at least expression information with less technical noise, to reconstruct
43 the true expression at single cell level.

44 Recent work has shown the effectiveness of deep generative models (e.g. Au-
45 toencoders and Generative Adversarial Networks) to infer realistic scRNA-seq
46 data and augment scarce cell populations using Generative Adversarial Net-
47 works [14] or the prediction of perturbation response using Autoencoders [15].
48 We hypothesized that a deep generative model could allow for the reconstruc-
49 tion of missing single cell gene expression information (low quality - lq) by using
50 related data with more genes expressed (high-quality - hq) as a reference, a com-
51 pletely novel approach to gene expression inference (Figure 1A). In other words,
52 lq data with many missing gene expression values and bad clustering could be
53 transformed into data with few missing genes and improved clustering if the
54 “style” of a related hq dataset could be transferred to it. In the best case,
55 it would be possible to infer gene expression information for single cell data
56 (lq) by using purified bulk RNA-seq data (hq), obtaining over ten thousand
57 genes expressed per cell. We envision that this novel approach, when properly
58 calibrated, is transformative for the analysis of single cell data, gaining deep
59 mechanistic insights into data beyond what is currently measurable. It is im-
60 portant to note that the concept of using hq data to reconstruct gene expression
61 in lq data is fundamentally different from classical imputation algorithms that
62 infer gene expression based on nearby cells from the same dataset, as outlined
63 above.

64 Based on the above considerations, we developed DISCERN, a novel deep
65 generative neural network for directed single cell expression reconstruction. DIS-
66 CERN allows for the realistic reconstruction of gene expression information by
67 transferring the style of hq data onto lq data, in latent and gene space. Our ex-
68 periments on real and simulated data show that DISCERN outperforms several
69 existing algorithms in gene expression inference across a wide array of single
70 cell datasets and technologies, improving cell clustering, cell type and activity
71 detection, and pathway and gene regulation identification. To obtain deep in-
72 sights into the cellular changes underlying COVID-19, we reconstructed single
73 cell expression data of patient blood and lung immune data. While in our ini-
74 tial analysis [16] of blood data we detected few immune cell types, expression
75 reconstruction with DISCERN resulted in the detection of 28 cell types and
76 states in blood, including two novel disease-associated T cell types, cytotoxic
77 CD4⁺ and CD8⁺ Tc2 T helper cells. Reconstructing a second COVID-19 blood
78 dataset with disease severity information, we were able to classify mild and se-
79 vere COVID-19 with an AUROC of 81 %, obtaining a potential biomarker of
80 disease stage. DISCERN can be easily integrated into existing workflows, as an
81 additional step after count mapping. Given that DISCERN is not limited by
82 a predefined distribution of data, we believe that it can be readily adapted to
83 enhance various other biomedical data types, especially other omics data such
84 as proteomics and spatial transcriptomics.

85 2. Results

86 2.1. The DISCERN algorithm for directed expression reconstruction

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

97 With these prerequisites in mind, we designed a deep neural network for
98 directed single cell expression reconstruction (DISCERN) (Figure S1B) that is
99 based on a modified Wasserstein Autoencoder [17]. A unique feature of DIS-
100 CERN is that it transfers the "style" of hq onto lq data to reconstruct missing
101 gene expression, which sets it apart from other batch correction methods such
102 as [18], which operate in a lower dimensional representation of the data (e.g.
103 PCA, CCA). To allow DISCERN to accurately reconstruct single cell RNA-
104 seq expression based on reference data, the structure of the network had to be
105 adapted in several ways. First, we implemented Conditional Layer Normaliza-
106 tion (CLN) [19, 20, 14] to allow for directed expression reconstruction of lq data
107 based on reference hq data (Figure S1B & S2). Second, we added two decoder
108 heads to the network to enable it to model dataset-specific dropout rates and
109 gene expression separately. Lastly, we extended DISCERN's loss function with
110 a binary cross-entropy term for learning the probability of dropouts to increase
111 general inference fidelity. Further algorithmic details of DISCERN can be found
112 in the methods and figs. S1 and S2.

113 We first demonstrate DISCERN's capabilities to faithfully reconstruct gene
114 expression using five pancreas single cell expression datasets of varying quality
115 (Tables S1 and S2). The pancreas data is widely used for benchmarking and it
116 is ideal to evaluate expression reconstruction for many cell types and sequencing
117 technologies. We consider a dataset as hq when the average number of genes
118 detected per cell (GDC) (e.g. smartseq2, GDC 6214) is much higher than in a
119 comparable lq dataset (Table S2). Conversely, a dataset is lq when the average
120 cell has lower counts and fewer genes expressed than a comparable hq dataset
121 (e.g. indrop, GDC 1887). Throughout this text, we will name sequencing tech-
122 nologies with capital (e.g. Smart-Seq2, InDrop) and datasets with lower case
123 first letters (smartseq2, indrop). We trained DISCERN on these five pancre-
124 atic single cell datasets and assessed the integration of data in gene space and
125 the average expression reconstruction per cell type. While uncorrected data
126 cluster by batch and not by cell type, DISCERN-integrated data show good
127 batch mixing and clustering of cells by cell type across all five datasets (fig. 1B
128 & fig. S2). To get a clearer picture of DISCERN's expression reconstruction
129 capabilities we next calculated correlation coefficients of measured expression

130 between the lowest quality inDrop and highest quality Smart-Seq2 data, before
131 and after expression reconstruction using DISCERN. The mean expression re-
132 construction of indrop-lq to smartseq2-hq and smartseq2-hq to indrop-lq data
133 is very accurate, showing a Pearson correlation of $r = 0.95$ ($p < 0.001$), while
134 mean expression correlation between uncorrected indrop-lq and smartseq2-hq
135 data is only $r = 0.77$ due to strong batch effects (Figure 1C & D, Figures S3
136 and S4). The improved quality of indrop-lq data reconstructed to smartseq2-hq
137 level is validated by the strong increase of genes expressed per cell, ranging from
138 ≈ 2000 genes per cell in the uncorrected indrop-lq data to ≈ 6000 genes in the
139 indrop-lq data after reconstruction (Figure S5).

140 We next investigated the effect of reconstruction of three cell type-specific
141 genes, before and after correction across the five pancreas datasets (Figure S6).
142 Insulin expression in the pancreas should be largely restricted to beta cells [21],
143 which can be observed in the uncorrected smartseq2-hq and celseq2 datasets,
144 while the indrop-lq batch shows a diffuse pattern of insulin expression across
145 cell types (Figure S6A left panel). This diffuse insulin expression is corrected
146 by reconstructing the smartseq2-hq expression pattern from the indrop-lq data
147 (Figure S6A middle panel). In general, the expected specificity of insulin ex-
148 pression in beta cells can be recovered for all datasets when using DISCERN's
149 reconstruction using the smartseq2-hq reference. Conversely, the reconstruction
150 from hq to the indrop-lq reference results in diffuse insulin expression across all
151 reconstructed datasets (Figure S6A right panel). We obtained similar results for
152 the pancreatic acinar cell-specific gene REG1A and the delta cell-specific gene
153 SST, both of which show diffuse expression across cell types in the uncorrected
154 inDrop data and cell-specific expression after reconstruction using smartseq2-hq
155 reference (Figure S6B & C). Interestingly, DISCERN can not only recover bio-
156 logical expression information, but it is also able to apply sequencing method-
157 specific effects after reconstruction. The smartseq2-hq dataset, for instance,
158 displays nearly no ribosomal protein coding gene expression after sequencing as
159 previously reported by [8], while data sequenced using InDrop, Cel-Seq, or Cel-
160 Seq shows prominent ribosomal protein coding gene expression (Figure S6D, left
161 panel). When reconstructing smartseq2-hq data to indrop-lq data, ribosomal
162 protein coding gene expression is re-instantiated (Figure S6D, right panel).

163 We further corroborated DISCERN's capability to integrate and reconstruct
164 gene expression in the more complex difftec dataset (Tables S1 and S2), consist-
165 ing of 14 single cell peripheral blood mononuclear cell (PBMC) datasets across a
166 wide range of technologies. Similar to pancreas, the difftec dataset is widely used
167 for benchmarking and it is ideal to evaluate expression reconstruction for even
168 more cell types and sequencing technologies. The different single cell technolo-
169 gies show large variation in quality, with an GDC ranging from 422 in Seq-Well
170 to 2795 in Smart-seq2. We trained DISCERN on these 14 PBMC single cell
171 datasets and observed very good integration in gene space (Figure S7). We
172 then reconstructed chromium-v2-lq (GDC 795) using a chromium-v3-hq refer-
173 ence (GDC 1514) and observed high mean gene expression correlation between
174 the reconstructed and reference datasets (Figures S8 and S9). These results
175 across 19 single cell datasets provide first evidence for the high-quality data in-

176 tegration and expression reconstruction that can be obtained with DISCERN.

177 *2.2. Specific and robust gene expression inference*

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

181 Since expression reconstruction can be seen as a generalization of expression
182 imputation, we compared DISCERN to DCA, MAGIC, and scImpute, three
183 state-of-the-art imputation algorithms [11, 12, 13]. Expression reconstruction
184 can also be viewed as a batch correction task in gene space, which is why we ad-
185 ditionally compared DISCERN to scGEN and Seurat [15, 18]. It is important to
186 note, however, that neither Seurat nor scGEN were designed for the expression
187 reconstruction task. Seurat and scGEN use a lower dimensional representation
188 in which a linear transformation aligns different batches. Seurat uses canonical
189 correlation analysis and scGEN uses the bottleneck layer representation of an
190 autoencoder to calculate and apply linear transformations.

191 To investigate the precision of gene expression reconstruction, we created an
192 artificial dataset by dividing the smartseq2-hq pancreas data into two batches,
193 smartseq-lq and smartseq2-hq. In the smartseq-lq batch, the top one KEGG
194 pathways per cell type were removed by setting the expression of genes con-
195 tained in these pathways to zero, while the smartseq2-hq remained unaltered.
196 Therefore, a reconstruction of smartseq-lq data using smartseq2-hq reference
197 (reconstructed-hq) should ideally recover the smartseq-lq expression to its orig-
198 inal state, prior to the removal of the genes. DISCERN is able to reconstruct
199 the mean expression for all cell types, achieving a correlation $r = 0.99$ (Fig-
200 ure 2A). DCA ($r = 0.66$), MAGIC ($r = 0.34$), scImpute ($r = 0.80$), and Seurat
201 ($r = 0.76$) have significantly lower correlation between the smartseq2-hq and
202 reconstructed-hq gene expression (Figure 2A). scGen shows only slightly reduced
203 performance ($r = 0.98$) compared to DISCERN, especially in the reconstruction
204 of highly expressed genes (Figure 2A) and low abundant cell types (Figure S10,
205 Megakaryocytes). We obtained similar results on the difftec dataset, with DIS-
206 CERN ($r = 0.98$) outperforming DCA ($r = 0.47$), Magic ($r = 0.21$), scImpute
207 ($r = 0.04$), Seurat ($r = 0.92$), and scGEN ($r = 0.94$) (Figure S10). To further
208 investigate gene expression reconstruction specificity, we compared the correla-
209 tion of reconstructed-hq to smartseq2-hq data after performing differential gene
210 expression (DEG) for each cell type against all other cell types (Figure 2B, up-
211 per panel). DISCERN is able to recover the correct DEG t-statistics with a
212 median correlation of 0.92, improving over state-of-the-art tools by more than
213 15 percentage points. In the corresponding experiment using the difftec dataset,
214 DISCERN achieves a median correlation of 0.85, which is a 25 percentage point
215 improvement over competing methods (Figure S11).

216 Since the genes were initially selected using KEGG gene set enrichment
217 analysis, the reconstruction of the corresponding pathways was investigated by
218 performing KEGG gene set enrichment analysis on the DEG results. DISCERN
219 is able to recover the pathway expression enrichment scores with a median cor-
220 relation of 0.93, exceeding the performance of Seurat and scGEN by more than

221 11 percentage points on median (Figure 2B, lower panel). In the corresponding
222 experiment using the difftec dataset, DISCERN achieves a median correlation
223 of 0.77, outperforming Seurat and scGen by more than 16 percentage points
224 (Figure S12).

225 While DISCERN outperforms competing algorithms in expression and path-
226 way reconstruction correlation, it achieves the second-best correlation for the
227 DEG fold-change (FC) of reconstructed-hq to smartseq2-hq data for the pan-
228 creas (Figure S13) and reconstructed-hq to chromium-v3-hq difftec datasets
229 (Figure S14). In both cases Seurat achieves slightly better correlation, which is
230 due to the fact that DISCERN slightly underestimates FC in favor of superior
231 DEG variance estimation.

232 Next, we show DISCERN's expression reconstruction robustness with re-
233 spect to varying sizes of lq to hq data. It is conceivable to assume that a large
234 amount of hq data would benefit the expression reconstruction of the lq data,
235 which makes it important to understand at what ratio good results can be ex-
236 pected. Interestingly, DISCERN seems to be very robust across a wide range
237 of smartseq2-lq to smartseq2-hq ratios, with correlations of 0.98 (ratio of lq/hq
238 0.14) to 0.93 (ratio of lq/hq 18.4), while the second-best performing algorithm
239 scGen showed a 11 percentage point decrease in performance (0.82 for ratio of
240 lq/hq 18.4) (Figure 2C, Figure S15). We observed similar results for the correla-
241 tion of t-statistics, showing a slight dependence of DISCERN's performance on
242 the lq/hq ratio (Figure S16). In general, all methods show better performance
243 with a small ratio of lq/hq data, while DISCERN shows least dependence and
244 outperforms other algorithms in the correlation of expression and t-statistics,
245 especially in the case of high lq/hq ratio.

246 Another aspect of expression reconstruction robustness is the dependence of
247 the algorithm on the cell type or cell state similarity of the lq and hq datasets.
248 In the optimal case, DISCERN would not require that the lq and hq datasets
249 have overlapping cell types to perform an accurate expression reconstruction,
250 which is theoretically possible if the network learns the general gene-regulatory
251 expression logic of the hq data (see discussion). To understand the dependence
252 on dataset similarity, we removed a complete cell type, pancreas alpha cells, from
253 the smartseq2-hq data and left the alpha cells in the smartseq2-lq data. We then
254 additionally varied the number of common cells in the lq and hq data, starting
255 with no overlapping cells (only alpha cells in the lq and all cells except alpha in
256 the hq data) and ending with almost complete overlap (all cells overlap between
257 the smartseq2-hq and -lq data, except for the alpha cells only present in lq data)
258 (Figure 2D). When evaluating DEG correlation, DISCERN was the only method
259 consistently achieving better performance than uncorrected data, outperforming
260 Seurat and scGen by more than 15 percentage points (Figure 2D). Similarly,
261 DISCERN was the only method consistently achieving better performance than
262 uncorrected data in the FC correlation task (Figure S17).

263 We next took a closer look at the integration and expression reconstruction
264 performance when no cell types overlap between the lq (alpha cells only) and hq
265 (all other cells) data. Notably, Seurat seems to over-integrate cell types, mix-
266 ing smartseq2-hq beta and gamma cells with reconstructed-hq alpha cells from

267 other batches (Figure S18), while scGEN and DISCERN keep the smartseq2-hq
268 and reconstructed-hq exclusive cell types separate (Figure 2E & Figure S18).
269 This over-integration seems to be causal for Seurat’s poor DEG correlation per-
270 formance ($r = 0.28$), while DISCERN ($r = 0.55$) is the only method achieving
271 better performance than uncorrected cells ($r = 0.47$) (Figure 2F). Thus, DIS-
272 CERN is able to keep existing expression correlations and improves the detec-
273 tion of cell type specific genes by reconstruction using an hq batch as reference.
274 In conclusion, DISCERN is both a precise and robust method for expression
275 reconstruction that outperforms existing methods by a significant margin.

276 *2.3. Improving cell cluster, type, and trajectory identification*

277 The comparison to competing methods provided evidence for DISCERN’s
278 superior expression reconstruction. Now, we will delineate how DISCERN’s
279 expression reconstruction improves downstream cell clustering, cell type and
280 activity state identification, marker gene determination, and gene regulatory
281 network and cell trajectory analysis.

282 To understand if cell-determining gene expression and pathways could be
283 recovered with expression reconstruction, we used a single nuclear sequencing
284 (sn-lq) and scRNA-seq (sc-hq) data pair that was prepared from the same liver
285 metastasis biopsy [22]. We reconstructed sn-lq data using the sc-hq reference,
286 obtaining reconstructed-hq data. While single nuclear sequencing provides re-
287 duced expression information in the average counts per cell as compared to
288 scRNA-seq (Table S2) [22], it is still the method of choice to obtain cell-specific
289 expression information when intact single cells cannot be recovered from a tissue
290 (e.g. after tissue fixation or freezing). It is important to note that nuclear tran-
291 scripts reflect current gene activity, which in part might not correlate with tran-
292 scripts that have lifetimes of up to days. Before integration, the sn-lq and sc-hq
293 datasets cluster by batch and not by cell type, while after expression reconstruc-
294 tion with DISCERN cells cluster by type and not by batch (Figure S19). This
295 is reflected in an expression correlation of 0.49 (sc-hq vs. sn-lq) before and 0.93
296 after reconstruction (sc-hq vs. reconstructed-hq) (Figure S20). Reconstruction
297 resulted in the expression of T cell receptor signaling genes in reconstructed T
298 cells (Figure S21) and antigen presentation genes in macrophages (Figure S22),
299 providing evidence that DISCERN faithfully recreates cell-determining genes
300 and pathways based on the hq data.

301 It is intriguing to observe that many genes of the antigen presentation path-
302 way in macrophages are not expressed in the sc-hq reference, most probably due
303 to dropout (Figure S22). We rationalized that bulk RNA sequencing (RNA-seq)
304 data of purified cell types (e.g. FACS sorted immune cells) is a suitable hq proxy
305 for the expected gene expression per cell. RNA-seq data of purified cells is read-
306 ily available from public repositories, making it possible to obtain thousands of
307 purified immune cell RNA-seq samples (see methods). We therefore set out to
308 increase cluster, cell type, gene regulatory network, and trajectory identification
309 of scRNA-seq data by reconstructing gene expression using a related RNA-seq
310 reference (Figure S23). For the scRNA-seq data we chose a cord blood mononu-
311 clear citeseq dataset (cite-lq) that was labeled with 15 antibodies (Table S3) to

312 allow for surface protein-based cell type discovery [23]. The CITE-seq informa-
313 tion allowed us to confirm expression reconstruction by DISCERN in cases where
314 gene expression is absent but protein expression and cell identity are validated
315 via antibody labeling. For the RNA-seq data, we selected 9.852 purified immune
316 samples (bulk-hq) and proceeded to reconstruct cite-lq (GDC 798) using a bulk-
317 hq (GDC 13.104) reference to obtain reconstructed-hq data with DISCERN. We
318 first investigated the correspondence of gene expression prior (cite-lq) and post
319 reconstruction (bulk-hq) with antibody-based surface protein labeling of *CD3D*,
320 *CD4*, *CD8A*, *CD2*, *B3GAT1*, *FCGR3A*, *CD14*, *ITGAX* and *CD19* (Figure 3A,
321 Figure S24). For several proteins (*CD8A*, *B3GAT1*, *CD4*), the corresponding
322 cite-lq gene expression was absent and cell type-specifically re-instantiated in
323 the reconstructed-hq expression data with DISCERN (Figure 3A, Figure S24).
324 In cases where cell type-specific gene and protein expression matched cite-lq
325 data (*CD3D*, *CD14*) the expression in reconstructed-hq data was left unaltered
326 (Figure S24). In some instances, we observed low cell type-specific expression
327 in the cite-lq data (*CD8A*, *CD2*, *FCGR3A*, *CD19*) that matched protein ex-
328 pression (Figure S24). In these cases, gene expression was increased in the cor-
329 rect cell types in the reconstructed-hq data. In general, we observed increased
330 agreement between cell type-specific surface protein and gene expression af-
331 ter reconstruction, showing that DISCERN doesn't invent or 'hallucinate' cell
332 types but reconstructs the expected expression specific for each cell type. We
333 further corroborated these results by selecting eight known cell type-specific
334 cytosolic proteins and investigated their expression before and after expression
335 reconstruction. *MS4A1* (B cells), *IL7R* ($CD4^+$ T cells), *MS4A7* (Monocytes),
336 *GNLY* and *NKG7* (NK cells) showed consistent expression before and after
337 reconstruction (Figure S25). The chemokine receptors *CCR2* (Monocytes, ac-
338 tivated T cells), *CXCR1* (NK cells), and *CXCR6* ($CD8^+$ T cells) showed the
339 correct cell type-specific expression only after expression reconstruction (Fig-
340 ure S25) [24]. It is notoriously hard to obtain cell subtype-specific information
341 from blood mononuclear scRNA-seq data, especially for $CD4^+$ T helper cells due
342 to their limited activation status in healthy individuals. This doesn't mean that
343 polarized $CD4^+$ T helper cells do not exist in healthy blood, as they are com-
344 monly detected after stimulation using FACS (Table S3) [25]. This lack of reso-
345 lution in scRNA-seq impedes clustering, marker gene, and trajectory analyses, a
346 drawback that could be overcome using DISCERN's expression reconstruction.
347 We therefore compared $CD4^+$ T cell (gene expression of *CD4* > 1 and *CD3E*
348 > 2.5) clustering and subtype identification using cite-lq and reconstructed-
349 hq data. While clustering with the leiden algorithm [26] using highly variable
350 genes of cite-lq data resulted in an unstructured distribution of $CD4^+$ T cell
351 subtypes (Figure 3B), clustering of reconstructed-hq data yields detailed in-
352 sights into T helper cell subtypes of blood mononuclear data (Figure 3C). Af-
353 ter reconstruction, we were able to characterize TH17, TH2, TH1, HLA-DR
354 expressing TREG (Active_TREG), naive $CD4^+$ T cells (*CD4_naive*), effector-
355 memory $CD4^+$ T cells (*CD4_EM*), central-memory $CD4^+$ T cells (*CD4_CM*),
356 and effector cells expressing IFN-regulated genes (*IFN_regulated*) (Figure 3C).
357 We selected published cell-determining marker genes and observed that many of

358 them were dropped out in the uncorrected data but present after reconstruction
359 (Figure S26). The absence of marker genes in uncorrected data results in poor
360 clustering and cell type identification, while single positive cells are detectable
361 in the respective neighborhood identified by reconstructed counts (Figure S26).
362 Importantly, we observed that in all cases the DISCERN-estimated proportions
363 of T helper subsets fall within the range of expected proportions as assessed by
364 previous FACS studies (Table S3, Figure S27). These findings are important,
365 as they prove once more that DISCERN discovers the correct cell subtypes and
366 cell proportions, in this case substantially outperforming the available CITE-seq
367 information in cell subtype resolution.

368 To further verify the cell type annotations, we extracted the top cluster-
369 determining genes from the reconstructed-hq data. Members of the TNF-
370 receptor superfamily are known to be expressed in T helper cell subtypes [27],
371 which can be observed after reconstruction in TH17 cells and partially in TH1,
372 TH2, Active_TREG and IFN_regulated cells (Figure S28). Similarly, recon-
373 structed TH1 cells show the expected high expression of granzymes *GZMK* and
374 *GZMA* [28], while *MIAT* and *HLA* expression are found in activated TREG
375 cells after reconstruction (Active_TREG cluster, Figure S28) [29, 30]. *NOG* ex-
376 pression is detected in reconstructed CD4_naive cells, as previously described
377 [31]. In addition, reconstructed CD4_naive, CD4_EM and CD4_CM show low
378 expression of the genes important for the T helper subtypes TH1, TH2, TH17,
379 Active_TREG and IFN_regulated. We further corroborated our cell type anno-
380 tation of reconstructed-hq data by observing the expected expression of several
381 established T cell subtype markers (Figure S29).

382 Similar to improved clustering and cell subtype detection, DISCERN reconstructed-
383 hq data resulted in improved gene regulatory network inference with SCENIC
384 [32]. SCENIC infers transcription factor-regulated gene expression modules
385 of single cell data. While cite-lq data resulted in a scattered distribution of
386 transcription factor networks across several T helper cell subtypes, SCENIC
387 with reconstructed-hq data showed transcription factor regulation in the cor-
388 rect subtypes (Figure 3D). After expression reconstruction the IKZF2 regulon
389 is detected in activated TREG cells [33] and the MAF regulon is found in dif-
390 ferentiated CD4⁺ T cells but not in naive CD4⁺ T cells [34]. A weak signal
391 of the MAF regulon is already detectable in the cite-lq data, yet strongly in-
392 creased in reconstructed-lq, while maintaining differentiated T helper cell speci-
393 ficity (Figure 3D). Furthermore, after reconstruction with DISCERN we could
394 identify the TH17 associated master transcriptional regulators RORC(+) and
395 RORA(+) [35], which were scattered over all TH17 cells before reconstruction
396 (Figure S30).

397 Finally, we wanted to investigate if DISCERN could also enhance cell trajec-
398 tory analyses with Slingshot of the citeseq data [36]. We focused on the differen-
399 tiation of effector and other T helper cell subtypes and found five lineages that
400 either pass through or terminate in the effector cell cluster in reconstructed-hq
401 data (Figure 3C). Two trajectories were of special interest to us: Lineage1 from
402 CD4_naive to TH1 cells (Figure S31) and Lineage2 from CD4_naive to TH17
403 cells (Figure S32). While the expression change along the trajectory in uncor-

404 rected data (Figure S31A, Figure S32A) is hardly visible, cell type-specific clusters
405 can be easily observed after DISCERN reconstruction (for lineage details
406 see Figure S31B, Figure S32B). The detailed insights into cell differentiation
407 that we obtained with reconstructed data are in stark contrast to the Slingshot
408 results obtained with cite-lq data. While terminal effector molecules can be detected
409 with cite-lq data, intermediate stages remain hidden, which prohibits the
410 detection of trajectories and results in a shuffling of marker gene expression (Figure
411 S31). Taken together these results highlight how expression reconstruction
412 using DISCERN improves downstream analyses and yields deeper biological insights
413 into cell type and state identification, gene regulation, and developmental
414 trajectories of cells.

415 *2.4. Discovering COVID-19 disease-relevant cells in lung and blood*

416 The previous sections have demonstrated DISCERN's utility to reconstruct
417 single cell expression data based on an hq reference, vastly improving the detection
418 of cell (sub-) types and their signaling. Given these advantages, we wondered if
419 DISCERN's expression reconstruction could deepen our understanding of cell type-
420 composition and signaling changes of immune cells in COVID-19 disease (Figure S33),
421 using two published datasets [37, 16]. To obtain best reconstruction results,
422 we again resorted to using bulk-hq immune reference data (Table S1) [38], as
423 outlined in the previous section.

424 First, we used a COVID-19 blood dataset (covid-blood-lq) with limited cell
425 type resolution, which was originally analyzed by our group using Seurat (Table
426 S1) [16]. While CD4⁺, CD8⁺, and NK cells formed separate clusters we were
427 unable to visibly distinguish subpopulations of these cells in covid-blood-lq
428 data [16]. Reconstruction of gene expression using bulk-hq data led to the
429 identification of 24 subtypes of CD4⁺ and CD8⁺ T cells in covid-blood-hq data
430 (Figure S34). Several cell clusters identified in covid-blood-hq data showed the
431 correct cell type-specific marker gene expression in covid-blood-lq data, albeit
432 in fewer cells, reduced in magnitude, and in some cases less specific (Figures S35
433 and S36). Reconstruction also led to the identification of CD4⁺ TH17 helper
434 cells that express *RORC* (Figure 4A & B, Figure S37). Based on the molecular
435 footprint of these TH17 cells they were further subdivided into TH17_cluster1
436 that exhibits a memory T cell phenotype with elevated *IL7R* expression and
437 TH17_cluster2 that exhibits an activated T cell phenotype with elevated *MHC-II*,
438 *CCR4* and *RBPJ* expression (Figure 4B, Figure S37). The expression of *RBPJ*
439 is of particular interest, as it is linked to TH17 cell pathogenicity, suggesting
440 a role of pathogenic TH17 cells in COVID-19 [39]. It is common practice
441 to stimulate memory T cells in vitro to trigger IL-17A production and a shift
442 towards a TH17 phenotype was previously described in COVID-19 [40]. With
443 DISCERN we are able to distinguish these cells in COVID-19 patient blood
444 without stimulation, identifying cytokine producing memory cells with a TH17-
445 like phenotype (Figure S37).

446 To further validate the existence of activated TH17 cells in COVID-19 patient
447 blood, we next analyzed the corresponding lung data (covid-lung) of the

448 patients for shared T cell receptor clones (Figure S38). The underlying assump-
449 tion is that cells with the same T cell receptor in lung and blood originate
450 from the same progenitor and therefore have a high probability of belonging
451 to the same cell type. For this comparison we used the cell type annotation
452 and representation of our original analysis of the covid-lung data, in which
453 memory T and TH17 cells were readily observed without reconstruction [16].
454 TH17_cluster1 cells showed strong clonal overlap with covid-lung CD4⁺ memory
455 T cells (Figure S38) and expressed comparable levels of *RORC* to covid-lung
456 effector memory TH17 cells (Figure S39), indicating that these CD4⁺ central
457 memory T cells could be TH17 (-like) cells. TH17_cluster2 in blood exhibited
458 strong clonal overlap with effector memory and resident memory TH17 cells
459 in covid-lung data (Figure S38) that express *RORC* and *IL-17A* (Figure S39).
460 Using the clonotype information of resident memory cells producing *IL-17A* in
461 inflamed lung (TRM17), we further corroborated the existence of the newly
462 identified population of IL-17A-producing TH17 cells in reconstructed COVID-
463 19 blood data (Figure S38). In general, the T cell receptor clonal information in
464 blood and lung therefore corroborated our cell type annotation in covid-blood-
465 hq data.

466 To understand the role of T cell subtypes in COVID-19 disease progression
467 we analyzed a second blood single cell dataset (covid-blood-severity-lq) contain-
468 ing disease-severity information for 130 COVID-19 patients [37]. To obtain opti-
469 mal cell type resolution, we combined the covid-blood-severity-lq T cell data [37]
470 with CD3⁺ covid-blood-lq cells [16] and reconstructed gene expression for the
471 combined dataset using bulk T cell sequencing reference data [38], resulting in
472 covid-blood-severity-hq data. Many of the 15 CD4⁺ T cell clusters identified in
473 covid-blood-severity-hq data (Figure S40) were also present in the covid-blood-
474 hq data, further validating the consistency of our cell type identification. This is
475 also corroborated by the available surface protein data for covid-blood-severity
476 data, substantiating that naive cells are CD45RA, memory cells are CD45RO,
477 and effector cell types are CD45RO positive (further details in Figure S41). We
478 compared the clusters that we identified in the covid-blood-hq with clusters iden-
479 tified in the covid-blood-severity-hq data and found confined and overlapping
480 regions of TFH, TH17_cluster1, and TH17_cluster2 cells (Figure S42). We also
481 compared the identified clusters to clusters defined in the original publication
482 (Figure S43). Cells identified as TFH in the original publication show signif-
483 icant overlap with naive CD4⁺ T cells (defined on transcriptome and protein
484 level) and CD4⁺ IL22⁺ cells (CD4.IL22) show marked overlap with TREG cells.
485 These results confirm once more the precise and robust cell type identification
486 that can be achieved with DISCERN.

487 Interestingly, we also identified two rather unexpected cell types after re-
488 construction. One cluster is positive for *CD4* and negative for *CD8A* while
489 otherwise expressing a signature of CD8⁺ effector memory cells with high ex-
490 pression of *GZMB*, *GZMH* and *PRF1* (Figure 4D & 4E). This signature points
491 to a CD4⁺ cytotoxic phenotype and indeed virus-reactive CD4⁺ cytotoxic cells
492 were described to be increased in blood during COVID-19 [41]. The other cell
493 type expresses *CD8*, *IL6R*, and *GATA3*, while being negative for *SLAMF7* (Fig-

494 ure 4D & 4E). These cells were described in the literature to be CD8⁺ T helper
495 cells [42], exert T helper function, and have been shown to lack cytotoxicity.
496 They lack expression of a significant number of cytokines and key transcription
497 factors pointing to a TH17 or TH22 phenotype. On a protein level these cells
498 express *CCR4*, while being negative for *CCR6*, making them cytolytic CD8⁺ T
499 helper type 2 cells (Tc2) cells. Part of this cluster overlaps with CD4 single-
500 positive cells and might explain why T helper type 2 cells are missing in the
501 CD4 cell clustering.

502 Overall, the highly specific and sensitive cell type identification in covid-
503 blood-severity-hq data enabled us to correlate the five COVID-19 disease sever-
504 ity categories to shifts in cell type and activity information. We first validated
505 the decrease in TFH cells with increasing disease severity, as described in the
506 original work (Figure S44) [37]. TH17 cells have been extensively studied using
507 flow cytometry and in accordance with our results MHC-II positive as well as
508 *CCR4* positive cells were described in COVID-19 patients (Figure 4B) [40]. We
509 observed a strong decrease in naive T helper cells in severe disease, most pro-
510 nounced for naive TREGs, while the fraction of TH17 cells showed little correla-
511 tion with disease severity (Figure S44). Of the two mixed cell types we detected
512 in COVID-19 data, cytotoxic CD4⁺ cells were increased in moderate and severe
513 disease (Figure S45). A similar increase is visible in patients with severe respi-
514 ratory disease without COVID-19 (Figure S46) and these cells might therefore
515 be a general marker of severe respiratory illness. Cytolytic CD8⁺ Tc2 cells are
516 increased in patients with severe symptoms and in those who died from COVID-
517 19 (Figure S45) and are described to be reduced after recovery from COVID-19
518 [43]. This positive correlation and the known role of Tc2 cells in fibroblast
519 proliferation induction and tissue remodeling could pinpoint a mechanistic role
520 of these cells in lung fibrosis as witnessed in severe COVID-19 patients. The
521 possibility to observe these cells in reconstructed single cell data may pave the
522 way to study the functional role of these cells in adverse COVID-19 outcome.

523 The relatively strong correlation of some cell types with COVID-19 out-
524 come suggests that blood cell fraction information might be used for patient
525 severity prediction. We trained a Gradient Boosting Machine (GBM) using
526 leave-one-out-cross-validation (LOOCV) on the fractions of all T cell types and
527 performed a forward feature elimination, to obtain a sparse, optimal model for
528 patient blood-based severity prediction. We first classified patients into three
529 groups, mild (union of asymptomatic and mild, $n = 26$), moderate ($n = 26$),
530 and severe (union of severe and critical, $n = 19$), reaching an AUROC of 0.63
531 (Table S4). We noticed that the mild and moderate groups were indistinguish-
532 able for the classifier (Figure S47). Training a GBM classifier on mild and severe
533 cases substantially increased classification performance, reaching an AUROC of
534 0.81 and accuracy, and F1 score of 0.82 (Table S4, Figure 4F & G). Compared
535 to the original T cell types and fractions reported (accuracy 0.61) [37], DIS-
536 CERN reconstructed T cell fractions are 33 % more accurate in the prediction
537 of COVID-19 disease severity (Figure 4G, Table S4). This classification improve-
538 ment is remarkable, given that DISCERN has no notion of disease severity when
539 it reconstructs gene expression. These results further demonstrate DISCERN's

540 precise and robust expression reconstruction that enabled the discovery of a
541 potential new blood-based biomarker for COVID-19 severity prediction.

542 3. Discussion

543 The sparsity of gene expression information and high technical noise in sin-
544 gle cell sequencing technologies limits the resolution of cell clustering, cell type
545 identification, and many other analyses. Several algorithms such as scImpute,
546 MAGIC, and DCA have addressed this problem by imputing missing gene ex-
547 pression in single cell data by borrowing expression information from similar
548 cells within the same dataset. While gene imputation clearly improves gene
549 expression by inferring values for dropped out genes, this imputation relies on
550 the comparison of similar cells with largely absent gene expression information
551 in the same dataset. With DISCERN we take a completely novel approach
552 to gene expression inference of single cell data, by realistic reconstruction of
553 missing gene expression in scRNA-seq data using a related dataset with more
554 complete gene expression information. We thus propose to call this procedure
555 ‘expression reconstruction’ to highlight the fundamental difference to classical
556 imputation and refer to the dataset with missing gene expression information
557 as low quality (lq) and the reference dataset as high-quality (hq).

558 We provide compelling evidence that our reference-based reconstruction out-
559 performs classical expression imputation algorithms as well as batch correction
560 algorithms such as Seurat and scGen, when they are repurposed for expression
561 reconstruction. To obtain an objective and thorough performance evaluation
562 for expression inference, we used seven performance metrics on 19 datasets,
563 including 12 single cell sequencing technologies. We focused our performance
564 evaluation on three scenarios with available ground-truth information, i) the
565 in silico creation of defined gene and pathway drop out events in scRNA-seq
566 data, ii) published hq and lq data pairs from the same tissue (pancreas, difftec,
567 sn/scRNA-seq datasets), and iii) CITE-seq protein expression as ground-truth
568 for cell types (citeseq dataset). In total, DISCERN achieved best performance
569 in 13 out of 15 experiments and obtained second rank in the remaining 2 com-
570 parisons. While DISCERN yields first place to Seurat in two FC expression
571 correlation comparisons, it always obtains best results across all datasets in
572 gene expression, gene regulatory network analysis, pathway reconstruction, and
573 cell type and activity identification and is the most stable algorithm for different
574 lq to hq size ratios and cell type overlaps.

575 It is important to note that DISCERN is a **precise** network that models
576 gene expression values realistically while retaining prior and vital biological in-
577 formation of the lq dataset after reconstruction. The network is also **robust**
578 to the presence of different cell types in hq and lq data, or an imbalance in
579 their relative ratios, and is robust to ‘hallucinating’ hq-specific cells into the lq
580 data. Several algorithmic choices are the foundation of DISCERN’s precision
581 and robustness. The network was designed to model the sequencing-technology-
582 specific and the underlying biological signals in separate components of its ar-
583 chitecture. Disentanglement of those two components is necessary to accurately

584 reconstruct expression information in the case where lq and hq datasets have
585 different content, i.e. cell type compositions. If the component designed to
586 model the effect of sequencing technology also captures the difference in the
587 biological signal, the reconstruction will lead to a lack of integration across the
588 two datasets where some cell types are still clustered by dataset (similar to
589 scGen in fig. S19). On the contrary, if the component modeling the biological
590 signal captures sequencing-technology-specific features, the reconstruction
591 will lead to an over-integration of the datasets where cells of different types are
592 mixed together (similar to Seurat in fig. S19). The demonstrated ability of DIS-
593 CERN to avoid those shortcomings, even in scenarios where there is very little
594 to no overlap between cell types across datasets, lies in the carefully crafted
595 balance between the expressivity of its components. The representational capa-
596 bilities of DISCERN, achieved via batch normalization, five loss terms, and a
597 dual head decoder, would reduce DISCERN's usability, if they would require fre-
598 quent dataset-specific tuning. The stability and usability was therefore a central
599 concern in the design and evaluation phase of DISCERN, which resulted in an
600 algorithm that gave very good results with a single set of default (hyper-) param-
601 eters. All comparisons to other algorithms, for instance, were performed with
602 default settings. Only the expression reconstruction of the exceptionally large
603 COVID-19 datasets required the fine-tuning of the learning rate, cross entropy
604 term, sigma, and the MMD penalty term. Another important technical feature
605 of DISCERN is that it can easily be integrated into existing workflows. It takes
606 a normalized count matrix, as created by nearly all existing single cell analysis
607 workflows, as input and produces a reconstructed expression matrix. This can
608 be used for most downstream applications (i.e. cell clustering, cell type identifi-
609 cation, cell trajectory analysis, and differential gene expression). DISCERN can
610 be trained on standard processors (CPU) for small and medium-sized datasets
611 and requires graphical processing units (GPU) for the expression reconstruction
612 of large datasets. Altogether, the usability and robustness of DISCERN should
613 enable even non-expert users to perform gene expression reconstruction.

614 A unique feature of DISCERN is the use of an hq reference to infer bio-
615 logically meaningful gene expression. While we consider this a main strength
616 of DISCERN, the dependence on a suitable reference dataset might also limit
617 its application. We took great care in this manuscript to mitigate this con-
618 cern by showing how DISCERN is able to reconstruct gene expression for many
619 different types of lq and hq pairs, ranging from indrop - smartseq2 to single
620 nucleus - single cell data pairs. Remarkable in this context is DISCERN's ro-
621 bustness to differences between the cell type compositions of lq and hq data
622 pairs, with DISCERN being the only algorithm obtaining robust expression re-
623 construction when few cell types overlap. We have also shown that purified
624 bulk RNA-seq samples can be used as hq reference, as successfully applied to
625 PBMC and COVID-19 datasets in this study. We used 9852 FACS purified
626 immune cell bulk sequencing samples [38], comprising 27 cell types, to success-
627 fully reconstruct single cell expression data. This implies that most single cell
628 studies involving immune cells (with or without other cell types present) can be
629 reconstructed with DISCERN using a single published bulk RNA-seq dataset.

630 Furthermore, public RNA-seq repositories such as NCBI GEO contain tens of
631 thousands of samples of immune and non-immune cells that could serve as refer-
632 ence for most expression reconstruction experiments. In conclusion, we provide
633 strong evidence that DISCERN is widely and easily applicable to many single
634 cell experiments.

635 While DISCERN gave good reconstruction results using default parameters
636 for most datasets we analyzed, we would like to highlight that the immense
637 representational power of generative neural networks can remove or hallucinate
638 biological information if not properly handled [6]. This is true for data inte-
639 gration [44] as well as for expression reconstruction algorithms and we would
640 highlight two guiding principles for optimal results. For non-expert users, we
641 would recommend the use of default settings and a careful selection of a re-
642 lated hq dataset. When datasets are large and complex, with many cell types
643 in the lq and several non-overlapping cell types in the hq data, one should al-
644 ways ensure that training does not merge or mix non-overlapping cell types with
645 other cells, by investigating that these cells keep their cell type-specific marker
646 gene expression. Keeping these ‘checks and balances’ will usually result in good
647 reconstruction results even for complex datasets such as covid-blood-severity.

648 To obtain novel insights into COVID-19 disease mechanisms and a new blood-
649 based biomarker for disease severity we reconstructed two published datasets
650 with DISCERN, Hamburg COVID-19 patients (covid-lung, -blood) and the
651 COVID-19 cell atlas (covid-blood-severity). The application of DISCERN to
652 the covid-blood dataset (COVID-19 patient blood) enabled us to detect 24 dif-
653 ferent immune cell types and activity states, which is quite remarkable given
654 that we find these cells in blood. Two TH17 subtypes caught our attention, as
655 they share the TCR clonality with the lung data from the same patients (covid-
656 lung), suggesting bloodstream re-entry of lung TH17 cells. We linked these two
657 subclusters to their functional role by separating them into a memory-like and
658 activated-like phenotype. The clonal overlap of activated TH17 cells in blood
659 with previously discovered lung-resident cells suggests that activated TH17 cells
660 in blood are resident T cells from the lung reentering circulation. These cells
661 might in part explain the multi-organ pathology observed in COVID-19, as
662 activated T cells might travel via the blood to secondary organs and cause in-
663 flammation and tissue damage. Future work might demonstrate the effect of
664 these activated T cells on tissue inflammation.

665 Given the detailed cell type and activity information we reached with gene
666 expression reconstruction, we wondered if changes in blood immune cell popu-
667 lations might be useful as a biomarker for disease severity prediction. We used
668 DISCERN to reconstruct the covid-blood and the covid-blood-severity datasets
669 and again identified a plethora of different T cell subtypes in the blood of pa-
670 tients with COVID-19. Using these cell proportions, we were able to classify
671 mild and severe disease using a GBM machine learning algorithm with 82 %
672 accuracy, outperforming classification with the originally published T cell types
673 by 21 percent points. This improvement is absolutely striking, as DISCERN
674 has no notion of the classification groups. It simply reconstructs gene expres-
675 sion and thereby improves cell type detection. These results are a convincing

676 implicit proof not only of the usefulness of DISCERN but more importantly of
677 its precision and robustness. While the use of this scRNA-seq-based biomarker
678 would be too expensive and time-consuming for clinical care, it strongly suggests
679 that FACS-based T cell fraction or count information from blood could be used
680 to trace and predict the severity state and potentially the disease trajectory of
681 COVID-19 patients.

682 Interestingly, we also discovered two atypical T cell types in reconstructed
683 COVID-19 patient blood single cell data. While cytotoxic CD4⁺ T cells have
684 been observed in COVID-19, we can show that this increase is not COVID-19
685 specific and is also observed in other types of pneumonia. Interestingly, we also
686 detected cytolytic CD8⁺ Tc2 cells that express *CD8A*, *GATA3*, *IL6R* and are
687 negative for *SLAMF6*. This cell type is linked to tissue fibrosis and steroid
688 refractory disease in asthma [45]. The increase in CD8⁺ Tc2 cells that we ob-
689 serve specifically in COVID-related death could be associated with COVID-19
690 patients that do not respond to steroids. Demonstration of increase of this cell
691 type in patients dying of COVID-19 points to a potential therapeutic inter-
692 vention with the drug Fevipiprant, which blocks CD8⁺ Tc2 cell activation and
693 its pro-fibrotic effects by inhibiting prostaglandin D2 signaling [46]. Functional
694 analysis of these cells has to demonstrate whether these cells are an early marker
695 of later death or whether it is a marker of already escalated treatment.

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

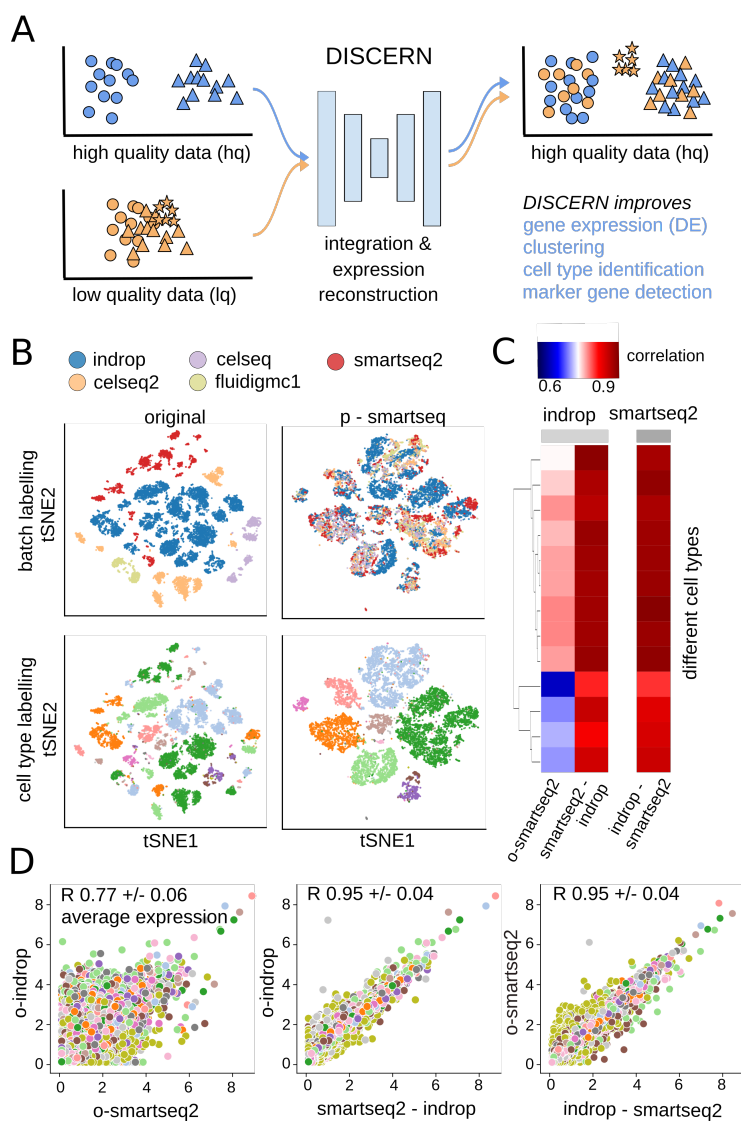


Figure 1: *Integration and expression reconstruction of single cell sequencing data.* **A:** DISCERN transfers the style of a high-quality (hq) dataset to a related low quality (lq) dataset, enabling gene expression reconstruction that results in improved clustering, cell type identification, marker 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 reconstruction 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 reconstruction (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, 'o' refers to original data. The mean Pearson correlation with one standard deviation over all cell types is shown in the figure title.

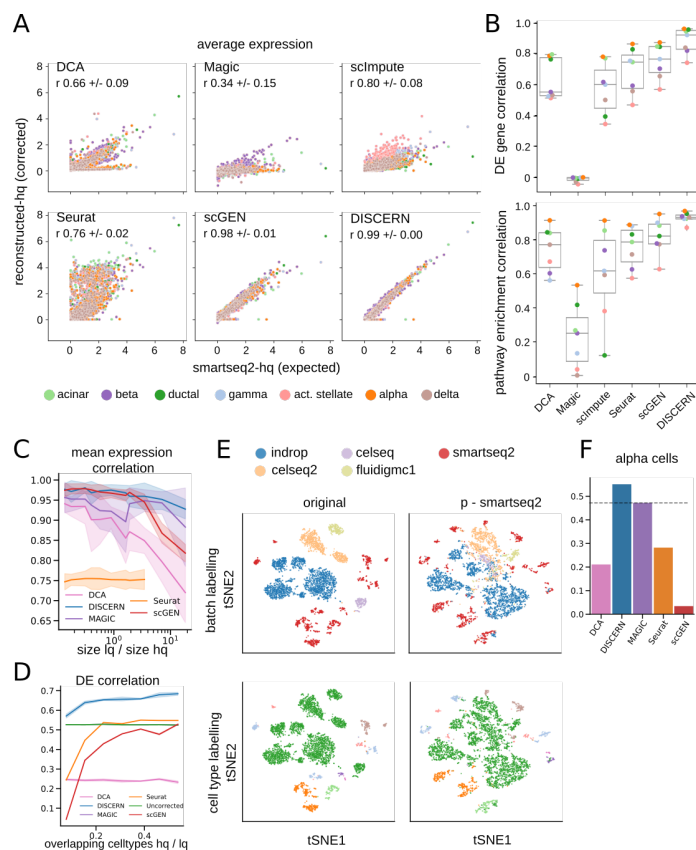


Figure 2: *Expression reconstruction benchmark of DISCERN and five state-of-the-art batch correction and imputation algorithms.* **A:** Comparison of the expression reconstruction performance of Seurat, scGEN, Magic, scImpute, DCA, and DISCERN using smartseq2 data. The smartseq2 data was split into a 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 smartseq-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 smartseq-hq batch and left in the smartseq-lq batch. The number of other overlapping cell types between the hq and lq data was then altered by removing cell types from the lq data before expression reconstruction (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. **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:** Pearson 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.

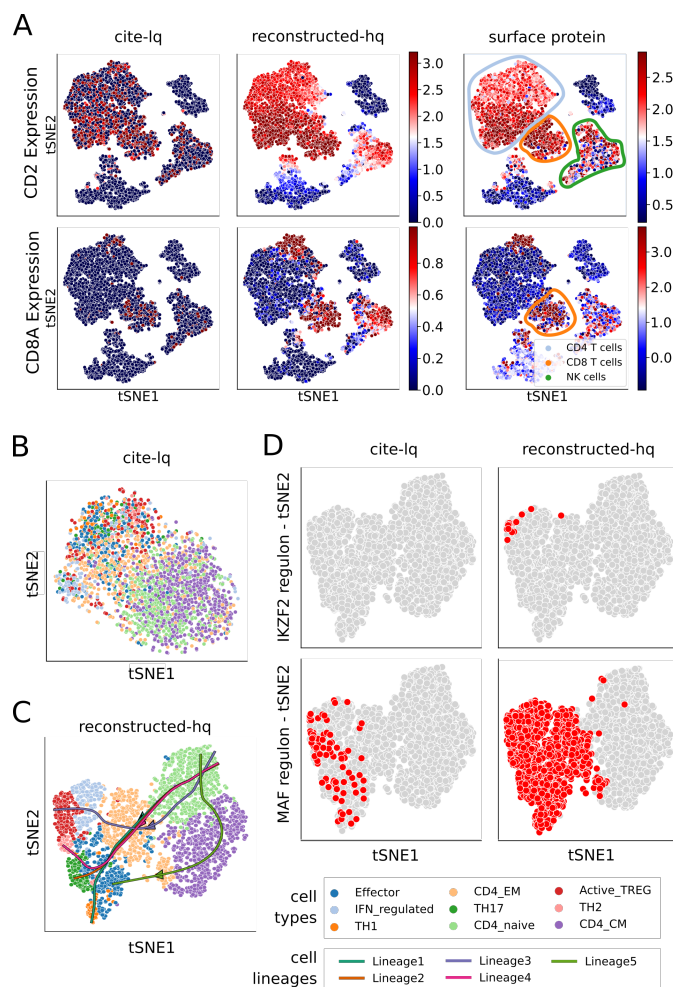


Figure 3: *Expression reconstruction improves downstream analyses including cell identification, gene regulation, and trajectory inference.* The cite-lq dataset was reconstructed using bulk-hq data and compared to ground truth CITE-seq (surface protein) information. The CITE-seq information was not used 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.

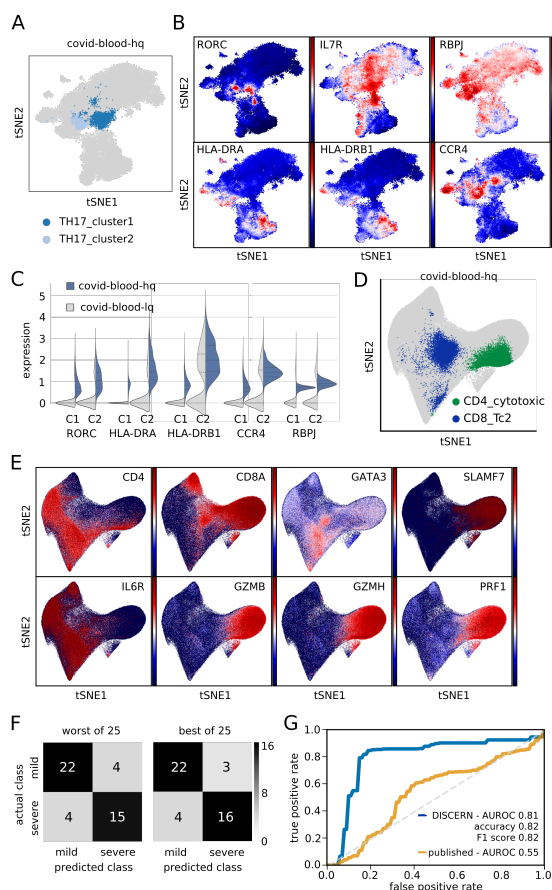


Figure 4: *Expression reconstruction improves COVID-19 cell type identification and allows for efficient disease severity prediction.* Two COVID-19 blood datasets were reconstructed and analyzed. Hamburg covid-blood-lq and covid-lung-lq data was reconstructed using bulk-hq data, resulting in the respective -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.

701
702

703 4. Methods

704 4.1. Data availability

705 In this manuscript many different scRNA-seq and RNA-seq datasets were
706 used. A comprehensive overview of dataset, method, cell type, origin, size, and
707 naming convention can be found in tables S1 to S3. All datasets are publicly
708 available as listed in table S1.

709 4.2. Dataset description

710 *Pancreas.* - The pancreas dataset is a collection of different scRNA-seq datasets,
711 profiling pancreas cells in the context of diabetes [47]. The pancreas data set is
712 a widely used data set for batch correction benchmark experiments and due to
713 its high number of cell types and sequencing technologies it allows to evaluate
714 differences between cells and sequencing technologies at the same time. The ex-
715 pression table, including the annotation, is available from SeuratData ([https://](https://github.com/satijalab/seurat-data)
716 github.com/satijalab/seurat-data) as `panc8.SeuratData` (v3.0.2) [47]. The
717 dataset was sequenced using five sequencing technologies (Smart-Seq2, Flu-
718 idigm C1, CelSeq, CEL-Seq2, inDrop) and consists of 13 cell types (alpha, beta
719 ,ductal, acinar, delta, gamma, activated_stellate, endothelial, quiescent_stellate,
720 macrophage, mast, epsilon, schwann). In total, before preprocessing, the data
721 set contains 14 890 cells.

722 *difftec.* - The difftec dataset was created for a systematic comparative anal-
723 ysis of scRNA-seq methods [48]. Similar to pancreas, the difftec dataset is
724 ideal for the evaluation of expression reconstruction across many cell types and
725 sequencing technologies. Seven sequencing technologies (10x Chromium v2,
726 10x Chromium v3, Smart-Seq2, Seq-Well, inDrop, Drop-seq, CEL-Seq2) were
727 used with at least two replicates each. In this dataset 10 different cell types
728 (Cytotoxic T cell, CD4⁺ T cell, CD14⁺ monocyte, B cell, Natural killer cell,
729 Megakaryocyte, CD16⁺ monocyte, Dendritic cell, Plasmacytoid dendritic cell,
730 Unassigned) were annotated, and make up for 31 021 cells in total before filter-
731 ing. The expression table including the annotation is available from SeuratData
732 as `pbmcsca.SeuratData` (v3.0.0).

733 *snRNA & scRNA.* - The dataset was created for the validation of a single
734 cell and single nuclei analysis toolbox [22]. Since snRNA-seq and scRNA-seq
735 data varies in the amount of counts per cell and the genes detected, we tested
736 if DISCERN could reconstruct snRNA-seq expression so that it would closely
737 resemble scRNA-seq expression, providing a biological ground-truth. While we
738 label snRNA-seq data as `lq` and scRNA-seq as `hq`, this distinction is incorrect
739 from a biological perspective, as gene expression should be in part different
740 between the nucleus and the cytosol. The dataset consists of a liver biopsy

741 sample (HTAPP-963) of metastatic breast cancer with single cell sequencing
742 and single nuclei sequencing. Eight cell types (Epithelial cells, Macrophages,
743 Hepatocytes, T cells, Endothelial cells, Fibroblasts, B cells, NK cells) were found
744 in the original publication in a total of 12 423 cells. The data was sequenced
745 using the Chromium V3 technology on a Illumina HiSeq X sequencer.

746 *covid-lung & covid-blood.* - The COVID-19 data set we have previously pub-
747 lished consists of blood and bronchoalveolar lavage (BAL) samples from four pa-
748 tients with bacterial pneumonia and eight patients with SARS-CoV-2 infection[16].
749 In total 155 706 cells were sequenced using TCR-seq technology, which allows
750 for the comparison of clonal expansion in both tissues. While we investigated
751 the lung data in detail in the original publication, the analysis of the blood was
752 largely limited to cell type identification. Using DISCERN, we use the blood
753 data to find previously unobserved cell types, link them to cell clones found in
754 the lung, and derive a biomarker based on cell fractions (see also covid-blood-
755 severity data). Cell type annotations for the BAL samples were used as in the
756 original publication.

757 *citeseq.* - This dataset contains CITE-seq information of healthy human PBMCs
758 for 6 cell types (B cells, CD4 T cells, NK cells, CD14⁺ Monocytes, FCGR3A⁺
759 Monocytes, CD8 T cells) [23]. In our analyses we used the cell type information
760 provided in the original publication [49]. The CITE-seq data is ideal to bench-
761 mark DISCERN, as the information of 13 surface proteins offers ground-truth
762 information on the cell types and a good proxy for the expression of the 13
763 corresponding genes.

764 *bulk.* - We used this large dataset of 28 FACS sorted and bulk sequenced immune
765 cell types as ‘ultimate’ hq reference data for lq immune single cell sequencing
766 data. Each of the 9852 samples provides an average expression information for
767 13 104 genes for a specific immune cell type, providing a hq reference for e.g.
768 lq single cell PBMC CiteSeq data with only 798 expressed genes per cell. We
769 further assume that this dataset is large enough to provide enough per cell type
770 variability for our deep neural network to faithfully learn and represent its gene
771 expression. In more detail, the dataset consists of 28 sorted immune cell types
772 (Naive CD4, Memory CD4, TH1, TH2, TH17, Tfh, Fr. I nTreg, Fr. II eTreg,
773 Fr. III T, Naive CD8, Memory CD8, CM CD8, EM CD8, TEMRA CD8, NK,
774 Naive B, USM B, SM B, Plasmablast, DN B, CL Monocytes, Int Monocytes,
775 NC Monocytes, mDC, pDC, Neutrophils, LDG) with \geq 99% purity [38]. Total
776 RNA was extracted using RNeasy Micro Kits (QIAGEN). Libraries for RNA-seq
777 were prepared using SMART-seq v4 Ultra Low Input RNA Kit (Takara Bio).
778 In total, the dataset contains 9852 samples collected in two phases from 416
779 donors, out of which 79 are healthy. For training DISCERN, bulk TPM counts
780 and all cell types were used if not stated otherwise.

781 *covid-blood-severity.* - This dataset is an aggregation of three COVID-19 se-
782 quencing studies using the 10X Genomics Chromium Single Cell 5’ v1.1 tech-

783 nology. It contains a large number of cell types with fine-grained cell type an-
784 notations that are complemented with information on COVID-19 disease sever-
785 ity for each patient sequenced. We used this dataset to obtain a blood-based
786 biomarker of COVID-19 disease severity, based on T cell fractions observed with
787 DISCERN. The data consists of PBMCs from 29 healthy, 89 COVID-19 and 12
788 LPS-treated patients. The authors detected 51 cell types in their original work
789 (see table S1) [37] and COVID-19 patients were classified by their disease sever-
790 ity (worst clinical outcome) into ‘asymptomatic’, ‘mild’, ‘moderate’, ‘severe’,
791 ‘critical’, and ‘death’. Count data together with CITE-seq information was
792 used as provided in the original publication ([https://covid19.cog.sanger.
793 ac.uk/submissions/release1/haniffa21.processed.h5ad](https://covid19.cog.sanger.ac.uk/submissions/release1/haniffa21.processed.h5ad)).

794 *4.3. Code availability*

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

799 *4.4. Preprocessing*

800 Raw expression data (Counts) preprocessing was performed as previously
801 described [50] using the scanpy (v1.6.1, [51]) implementation. In particular,
802 the intersection of genes between batches was used. The cells were filtered
803 to a minimum of 10 genes per cell and a minimum of 3 cells per gene. Li-
804 brary size normalization was performed to a value of 20 000 with subsequent
805 log-transformation. As model input for DISCERN the genes were scaled to
806 zero mean and unit variance. However, for all further evaluation the genes
807 were scaled to their uncorrected mean and variance not considering the batch
808 information.

809 *4.5. Description of DISCERN*

810 DISCERN is based on a Wasserstein Autoencoder with several added and
811 modified features. We will describe the details of DISCERN’s architecture in
812 the next paragraphs and a compact representation can be found in fig. S1B.

813 *Wasserstein Autoencoder.* - While neural network-based autoencoders have been
814 widely used for decades for dimensionality reduction [52, 53], recent advances
815 have also allowed their use to build a generative model of the distribution of
816 the data at hand[54]. More recently, leveraging results from optimal transport
817 [55], Wasserstein Generative Adversarial Networks (WGAN) [56] and Wasser-
818 stein Autoencoders (WAE) [17] have been designed to explicitly minimize the
819 (Wasserstein, or earth-mover) distance between the distribution of the input
820 data and their reconstruction. WGANs only implicitly encode their input into
821 a latent representation (called latent code), while WAE has the useful property
822 of using an explicit encoder, which makes it possible for the model to directly
823 manipulate the different representations of single-cell data. Finally, the WAE

824 framework, established in [17], allows the use of a wide range of architecture and
825 losses, which we are going to detail now. First of all, in order to effectively use a
826 number of latent dimensions that adaptively matches the intrinsic dimension of
827 the scRNA-seq data at hand, DISCERN uses a random encoder as prescribed
828 in [57].

829 *Architecture.* - Autoencoders widely used for transcriptomics applications are
830 shown to perform well on several tasks, like drug perturbation prediction [15]
831 or dropout imputation [12]. Since the ordering of the genes in scRNA-seq count
832 matrices is mostly arbitrary, fully-connected layers are usually used in this task.
833 In our case, DISCERN consists of three fully connected layers in the encoder
834 and the decoder. The bottleneck of the autoencoder (or latent space) contains
835 48 neurons, which is sufficient to accurately model all the datasets we used in
836 our experiments. Additionally, we exploit a finding from [57] to let the net-
837 work learn the appropriate amount of latent dimensions. While the encoder
838 will be tasked to transform the distribution of the input data into a fixed,
839 low-dimensional prior distribution (i.e. a standard Gaussian), the decoder will
840 perform the opposite, i.e. transforming the fixed, low-dimensional prior distri-
841 bution into gene space. scRNA-seq data is known to display a high level of zero
842 measurements, called dropout, which is essential to accurately model the count
843 distribution. To describe scRNAseq data in a parametric way, it is common to
844 model the expression level of a gene with zero-inflated negative binomial distri-
845 bution [58]. Despite the several non-linearities in the decoder architecture,
846 it is, however, difficult to learn an encoding function that maps a simple prior
847 to the distribution leading to low quality modeling of low expressed genes. To
848 address this issue, we scale the gene expression and attach a second head to the
849 decoder (i.e. a second decoder sharing all weights with the first, except for the
850 last layer). The task of the second decoder head is to predict, for each gene
851 of a cell, the probability of its expression to be dropped out, giving rise to a
852 random decoder. Thus, this second decoder head predicts dropout probabili-
853 ties and models the dropout probabilities for different batches. This additional
854 head allows modeling the dropout and the expression independently, to capture
855 the specific distribution of single cell data without the need for further explicit
856 assumption about the distribution.

857 *Loss function.* - The loss optimized during the training of DISCERN is com-
858 posed of four terms: a data-fitting (or reconstruction) loss, a dropout fitting
859 (cross entropy) loss, a prior-fitting term (ensuring that DISCERN approxi-
860 mately minimizes the Wasserstein distance) and a variance penalty term (that
861 controls the randomness of the encoder). Thus, DISCERN can be considered as
862 a Wasserstein Autoencoder as introduced in [17]. For the reconstruction term,
863 the framework introduced in [17] allows the use of any positive cost function.
864 We elected to use the Huber loss [59] as it is well suited for modeling scaled
865 scRNA-seq expression data, because it allows to select a threshold value to give
866 lower weight to high differences in highly expressed genes and thus allows the
867 model to learn a more robust expression estimate without focusing too much

868 on outlier values. For the prior-fitting term, following [17], DISCERN uses the
869 Maximum Mean Discrepancy (MMD) [60] between the aggregate posterior (i.e.
870 the distribution of the input single-cells after encoding) and a standard Gaus-
871 sian. We use the sum over an inverse multiquadratic kernel with different sizes
872 for this task. Then, to prevent the random encoder (with diagonal covariance)
873 from collapsing to a deterministic one, a penalty term that enforces that some
874 components of the variance are close to 1. Intuitively, that means that the su-
875 perfluous latent dimensions will only contain random noise (see [57] for more
876 details). Another loss term, namely the binary cross-entropy loss, on the second
877 decoder head is used to enable the model to learn a dropout probability for each
878 gene and sample. The loss on the dropout layer enables the model to capture
879 the bimodal distribution of single cell data. Additionally, activity regularization
880 is applied on the Conditional Layer Normalization (CLN), such that the weights
881 of the conditional layers are only regularized in a batch-specific manner and the
882 regularization is not applied for batches, which are not present in the current
883 mini-batch. This has the advantage that the batch dependent weights are not
884 influenced too much by different batch sizes. The four loss terms are added (and
885 weighed) together to form the loss that DISCERN minimizes during training
886 (see also fig. S1 for loss terms).

887 *Conditional Layer Normalization.* - The weights of those fully-connected lay-
888 ers are shared for all the batches that DISCERN is trained on. However, to
889 model the batch-specific differences, we use a Conditional Layer Normalization
890 (CLN) that applies the idea proposed in [19] to Layer Normalization [20]. In
891 essence, for each batch, different sets of shifting factors are learned. Note that
892 in DISCERN, no scaling factors are used to limit the expressivity of the con-
893 ditioning and therefore reduce the chance of over integration. This allows not
894 only to accurately model the batch-specific differences between batches, but also
895 to transfer the batch effect from one dataset onto another, in the spirit of the
896 style-transfer approach developed in [19]. To make things clear, DISCERN does
897 not explicitly train to integrate datasets. Instead, it trains to accurately model
898 the input data, capturing the batch-specific differences with the weights of the
899 CLN layers (i.e. conditioning), and the biological signal (which is mostly shared
900 across the batches to integrate) with the weights of the fully-connected layers.
901 After training, we encode all the cells we want to reconstruct, conditioning the
902 process on their batch of origin. Then, we take the batch chosen by the user and
903 proceed to decode all the cells conditioning on that specific batch, effectively
904 transferring the batch effect of one specific batch onto all of the batches we want
905 to integrate and reconstruct.

906 *Activations & dropout.* - With the exception of the output layer, every other
907 fully-connected layer of the encoder and the decoder was followed by a CLN,
908 a Mish ([61] activation function, and dropout during model training to reduce
909 overfitting.

910 *Optimization.* - To optimize the weights of our model, DISCERN uses Recti-
911 fied Adam ([62], which addresses some of the shortcomings of the widely used

912 Adam [63] and generally yields more stable training. To prevent overfitting, the
913 optimization is stopped early. It is implemented as a modification of the Keras
914 EarlyStopping (with parameter minDelta set to 0.01 and the patience to 30)
915 where the callback is delayed by a fixed number of 5 epochs. The delay was
916 implemented to prevent too early stopping due to the optimization procedure.

917 *4.6. Hyperparameters*

918 As outlined in the architecture section of the methods and depicted in fig. S1,
919 DISCERN features several learnable hyperparameters. The complexity of the
920 hyperparameter search space is a potential downside of DISCERN, if these hy-
921 perparameters would be unstable across different datasets or in other words,
922 would require constant tuning. Fortunately, DISCERN's hyperparameters are
923 very stable across the multitude of datasets tested in this manuscript, which we
924 will outline in this paragraph. Naturally, there is no rule without an exception,
925 which in this manuscript are the COVID-19 datasets that required optimization
926 for several hyperparameters.

927 *Constant hyperparameters.* - DISCERN features a number of hyper-parameters
928 that can be tuned through hyperparameter optimization (see below for details).
929 Most of them have default values that yield reasonable performance across the
930 different datasets we used and are being kept constant across experiments, in-
931 cluding the COVID-19 dataset. Those constant hyperparameters are: the choice
932 of the reconstruction loss (Huber loss), activation functions (Mish), CLN for the
933 conditioning, number of fully-connected layers (3) and their size (1024, 512, 256
934 and 256, 512, 1024 neurons for the encoder and the decoder respectively), num-
935 ber of latent dimensions (48), learning rate (1×10^{-3}), decay rates β_1 and β_2 of
936 Rectified Adam (0.85 and 0.95 respectively), batch size (192), label smoothing
937 for our custom cross entropy loss (0.1), dropout rates (0.4 in the encoder and 0
938 in the decoder), delta parameter of the Huber loss (9.0), weight on the penalty
939 on the randomness of the encoder λ_{σ} (1×10^{-8}), weight on the cross entropy
940 loss term $\lambda_{dropout}$ (1×10^5), weight on the MMD penalty term λ_{prior} (1500).

941 *Dataset-specific hyperparameters.* - The optimal value of the L2 regularization
942 applied on the weights of our custom CLN highly depends on the dataset at hand
943 and thus requires dataset-specific tuning. For datasets with a very small vari-
944 ance in cell compositions the L2 CLN regularization can be turned off (weight
945 set to 0). When datasets have different compositions the L2 CLN regularization
946 requires higher values (typically between 1×10^{-3} and 0.2).

947 *COVID-19 hyperparameters.* - For the experiments with COVID-19 datasets
948 slightly adjusted hyperparameters were used: learning rate of 6e-3, label smooth-
949 ing for our custom crossentropy loss of 0.05, weight on the penalty on the ran-
950 domness of the encoder λ_{σ} (1e-4), weight on the cross entropy loss term
951 $\lambda_{dropout}$ (2e3), weight on the MMD penalty term λ_{prior} (2000).

952 *Hyperparameter optimization.* - DISCERN implements different techniques for
953 hyperparameter optimization by using the ray[tune] library [64]. For most use
954 cases the model does not require hyperparameter tuning and the default pa-
955 rameter should be sufficient. However, DISCERN has a generic interface and
956 supports nearly all techniques implemented in ray[tune]. The initial hyperpa-
957 rameters were found using grid search. The loss used for the hyperparameter
958 selection is the classification performance of a Random Forest classifier trying
959 to classify real vs. auto-encoded cells. Classification performance was mea-
960 sured using the area under the receiver operating characteristic curve (AUC /
961 AUROC).

962 4.7. Competing algorithms and methods

963 We briefly discuss competing methods and have compared their performance
964 to DISCERN in the results section. These algorithms can be grouped into two
965 categories, i) imputation algorithms that were developed to estimate drop-out
966 gene expression based on dataset inherent information (MAGIC, DCA, scIm-
967 pute) and ii) algorithms designed for batch correction that we have modified
968 or extended to reconstruct gene expression, although this is not their intended
969 use (Seurat, scGen). Given the latter, it is clear that DISCERN could be used
970 purely for batch correction in latent space, a subject beyond the scope of this
971 manuscript.

972 *MAGIC.* [13] - Markov affinity-based graph imputation of cells (MAGIC) de-
973 noises and imputes the single-cell count matrix using data diffusion-based in-
974 formation sharing. The construction of a good similarity metric is challenging
975 for finding biologically similar cells due to high sparsity. MAGIC finds a good
976 similarity metric using a sophisticated graph-based approach that builds less-
977 noisy cell-cell affinities and information sharing across cells. A particular focus
978 of MAGIC was to understand gene-gene relationships and to characterize other
979 dynamics in biological systems. MAGIC is provided as a Python package.

980 *DCA.* [11] - Deep count autoencoder (DCA) is a deep learning-based method for
981 denoising single-cell count matrices. DCA is implemented in Python and uses
982 an autoencoder with a Zero-Inflated Negative Binomial (ZINB) loss function.
983 For each gene, DCA computes gene-specific parameters of ZINB distribution,
984 namely dropout, dispersion and mean. By modeling gene distributions as a noise
985 model and also computing dropout probabilities of each gene, DCA is able to
986 denoise and impute the missing counts by identifying and correcting dropout
987 events.

988 *scImpute.* [12] - Similarly to MAGIC, scImpute focuses on identifying cells that
989 are similar, which is challenging due to the high sparsity of single-cell count
990 matrices. scImpute is a statistical model using a three step process to impute
991 scRNA-seq data. In the first step spectral clustering is applied on principal com-
992 ponents to find neighbors, which later can be used to detect and impute dropout

993 values. In the second step scImpute fits a mixture model of a Gamma distribu-
994 tion and a Normal distribution to distinguish technical and biological dropouts.
995 In the last step, the model uses a regression model for each cell to impute the
996 expression of genes with high probability of dropout. With this approach, scIm-
997 pute avoids hallucinations and keeps the gene expression distribution. scImpute
998 is provided as an R package.

999 *Seurat*. [18] - Seurat is an open-source toolkit for the analysis of single cell
1000 RNA-sequencing data. In addition to general analysis functions, Seurat of-
1001 fers batch-correction functionality. Seurat uses canonical correlation analysis
1002 to construct this lower dimensional representation and tries to find neighbors
1003 between batches in this shared space. These anchors are filtered considering
1004 the local neighborhood of the cell pairs and remaining anchors are finally used
1005 to construct correction vectors for all cells in this low dimensional representa-
1006 tion. While Seurats is intended to work in a lower dimensional representation,
1007 it can also be used to reconstruct the expression information from this lower
1008 dimensional representation. Seurat is provided as an R package.

1009 *scGen*. [15] - scGen is a variational autoencoder based deep learning method
1010 with a focus on learning features that help distinguish responding and non-
1011 responding genes and cells. scGen constructs a latent space in which it es-
1012 timates perturbation vectors associated with a change between different con-
1013 ditions. Since scGen models the perturbation and infection responses in single
1014 cells, it is focused on in-silico screening with the use of cells coming from healthy
1015 samples. It can also be used for batch correction. For batch correction, and unlike
1016 DISCERN or Seurat, scGen uses both batch and cell type labels.

1017 4.8. Evaluation metrics

1018 *t-SNE & UMAP*. - For visualization of the data sets and to qualitatively assess
1019 the integration performance tSNE and UMAP were used. Both methods are
1020 based on PCA representation and use non-linear representations to create a 2D
1021 representation of the data. We used the scanpy [51] implementation. Default
1022 settings were used in nearly all cases except: In the combined COVID-19 dataset
1023 analogue to Kobak *et al.*[65] the dataset was subset to 25 000 cells and tSNE
1024 was computed using a perplexity of 250, and a learning rate of 25 000/12. These
1025 positions were taken and used as input to tSNE of all cells using a perplexity
1026 of 30 a learning rate of (number of observations)/12 and a late exaggeration of
1027 4.0 using FIt-SNE [66]. Clustering was performed using PARC [67] with de-
1028 fault parameters except `dist_std_local=1.5` and `small_pop=300`. Methods were
1029 changed here due to computation time issues for 350 000 cells. covid-blood data
1030 was analyzed using a learning rate of (number of observations)/6 a perplexity
1031 of (number of observations)/120 and `early_exaggeration=4`. Clustering was per-
1032 formed using default parameters except `knn=100` and `small_pop=100` to reduce
1033 the number of clusters with limited cell number. Clustering of the T helper cells
1034 in healthy blood was performed using coarse clustering with 30 nearest neigh-
1035 bors and leiden clustering (<https://github.com/vtraag/leidenalg>) with a

1036 resolution of 0.6. Afterwards a combined cluster of IFN-regulated and TREG
1037 was reclustered using a resolution of 0.4 and effector T cells were reclustered us-
1038 ing a resolution of 0.8. Resolution was chosen to dissect the raw gene expression
1039 changes of known cell types.

1040 *Mean gene expression.* - Mean gene expression was calculated as average over
1041 log-normalized expression over all cells, usually stratified by celltype. This eval-
1042 uation of expression data consists of many data points where several have values
1043 close to zero, but could have a high weight on rank-based correlation methods.
1044 Thus Pearson correlation was used to evaluate the performance.

1045 *Differential gene expression.* - Differential gene expression was performed using
1046 the scanpy [51] rank_gene_groups function using the t-test method for calculat-
1047 ing statistical significance on log-normalized expression data. Differential gene
1048 expression analysis was always performed under consideration of the cell type
1049 information. For comparison of differential gene expression analysis between
1050 conditions, the Pearson correlation was used. It is calculated either on the log2
1051 fold-change or in most cases on the t-statistics, computed during significance
1052 estimation. The data was compared using the t-statistics, because it aggregates
1053 information on both the variance and the change in mean expression. Thus it
1054 allows, roughly speaking, for simultaneously evaluating the significance and the
1055 log2 fold change.

1056 *Pathway analysis.* - Pathway analysis or gene set enrichment analysis was done
1057 using the prerank function from gseapy [68] on the t-statistics, computed as
1058 described in the ‘Differential gene expression’ section of the methods. To this
1059 end, the gene set library “KEGG_2019_Human” provided by enrichr [69] was
1060 used. Top pathways were selected using the normalized enrichment score as
1061 previously described [68].

1062 *Gene regulation.* [32] - The python implementation of the SCENIC (pySENIC)
1063 was used to infer regulons specific for CD4⁺ T helper cells. SCENIC infers a
1064 gene regulatory network using GRNBoost2 and creates co-expression modules.
1065 The co-expression modules get associated with transcription factors using the
1066 transcription factor motif discovery tool RcisTarget. A pair of transcription
1067 factor and associated gene set is called a regulon. For each cell, the regulons
1068 get scored using the AUCell algorithm to examine if a cell is affected by the
1069 regulon. We used default parameters of the pySENIC implementation.

1070 4.9. COVID-19 classification

1071 To evaluate the importance of the cell types found in the covid-blood-
1072 severity-hq data set after reconstruction with DISCERN, the fraction for all
1073 T cell subtypes was used to predict the disease severity, as provided in [37].
1074 The data was classified using a Gradient boosting classifier ([70], implemented
1075 in scikit-learn v1.0.2, default settings) using 25 rounds of leave-one-out cross-
1076 validation (LOOCV). Each round consists of n training-prediction iterations

1077 with $n - 1$ samples for training and 1 sample for testing, such that after one
1078 round prediction results for all n samples could be evaluated. We chose LOOCV
1079 over k -fold cross-validation and testing due to the limited size of the dataset,
1080 consisting of only 71 patients. We used `pymc` ([71], v3.3) for the performance
1081 evaluation. The final evaluation was done using the accuracy and F1 score
1082 as provided by `pymc`. The area under the receiver operating characteristic
1083 (AUROC) curve is computed with `scikit-learn`. Before training the classifiers
1084 a forward feature selection was performed using the `SequentialFeatureSelector`
1085 implemented in `scikit-learn` with default parameters. In total four experiments
1086 were performed. In the first experiment, classification with three disease cat-
1087 egories (mild, moderate, severe) was used. Patients who died were excluded.
1088 For the other two experiments only patients with asymptomatic, mild, severe
1089 and critical symptoms were included. In all experiments the asymptomatic and
1090 mild category was merged to mild and severe and critical to severe.

1091 5. Acknowledgements

1092 We thank Immo Prinz, Manuel Friese, Johannes Soeding, Robert Zinzen, Yu
1093 Zhao and Stefan Kurtz for their helpful comments and suggestions. FH, RK,
1094 MM, PM, & SB were supported by the LFF-FV 78, EU ERare-3 Maxomod,
1095 SFB 1286 Z2, FOR 296, and FOR 5068 research grants. Further support was
1096 obtained from the UKE R3 reduction of animal testing grant. CE was sup-
1097 ported by DFG ER 981/1-1 and the clinician scientist programme of university
1098 Hamburg, NG was supported by ERC StG-715271, and SHu was supported by
1099 ERC CoG-865466 and has an endowed Heisenberg-Professorship awarded by
1100 the Deutsche Forschungsgemeinschaft. SHa was funded by the BMBF STOP-
1101 FSGS-01GM1518C and SFB 1192 B08 research grants.

1102 6. Competing interests

1103 The authors declare no competing interests.

1104 7. Author contributions

1105 SB initiated and SB, PM, FH, and CE conceptualized the study with help
1106 from MM. FH and CE implemented DISCERN, MM refactored the code, and
1107 PM reviewed the DISCERN implementation. FH, CE, and RK performed the
1108 analyses. SB, PM, NG, and SHu supervised the study. SB, FH, and CE wrote
1109 the manuscript. SHu, PM, NG, RK and SHa provided ideas, contributed to
1110 manuscript text and critically reviewed the manuscript. All authors read and
1111 approved the final manuscript.

1112 **References**

- 1113 [1] N. Editorial, Method of the year 2013, *Nat. Methods* 11 (1) (2014) 1.
- 1114 [2] Y. Zhao, U. Panzer, S. Bonn, C. F. Krebs, Single-cell biology to decode
1115 the immune cellular composition of kidney inflammation, *Cell and tissue*
1116 *research* 385 (2) (2021) 435–443.
- 1117 [3] M. Stoeckius, C. Hafemeister, W. Stephenson, B. Houck-Loomis, P. K.
1118 Chattopadhyay, H. Swerdlow, R. Satija, P. Smibert, Simultaneous epi-
1119 tope and transcriptome measurement in single cells, *Nature methods* 14 (9)
1120 (2017) 865–868.
- 1121 [4] A. A. Tu, T. M. Gierahn, B. Monian, D. M. Morgan, N. K. Mehta,
1122 B. Ruitter, W. G. Shreffler, A. K. Shalek, J. C. Love, Tcr sequencing paired
1123 with massively parallel 3' rna-seq reveals clonotypic t cell signatures, *Nature*
1124 *immunology* 20 (12) (2019) 1692–1699.
- 1125 [5] J. A. Pai, A. T. Satpathy, High-throughput and single-cell t cell receptor
1126 sequencing technologies, *Nature Methods* 18 (8) (2021) 881–892.
- 1127 [6] S. Oller-Moreno, K. Kloiber, P. Machart, S. Bonn, Algorithmic advances
1128 in machine learning for single-cell expression analysis, *Current Opinion in*
1129 *Systems Biology* 25 (2021) 27–33.
- 1130 [7] D. Lähnemann, J. Köster, E. Szczurek, D. J. McCarthy, S. C. Hicks, M. D.
1131 Robinson, C. A. Vallejos, K. R. Campbell, N. Beerenwinkel, A. Mahfouz,
1132 et al., Eleven grand challenges in single-cell data science, *Genome biology*
1133 21 (1) (2020) 1–35.
- 1134 [8] X. Wang, Y. He, Q. Zhang, X. Ren, Z. Zhang, Direct comparative anal-
1135 yses of 10x genomics chromium and smart-seq2, *Genomics, proteomics &*
1136 *bioinformatics* 19 (2) (2021) 253–266.
- 1137 [9] A. K. Shalek, R. Satija, J. Shuga, J. J. Trombetta, D. Gennert, D. Lu,
1138 P. Chen, R. S. Gertner, J. T. Gaublotte, N. Yosef, et al., Single-cell
1139 rna-seq reveals dynamic paracrine control of cellular variation, *Nature*
1140 510 (7505) (2014) 363–369.
- 1141 [10] W. Hou, Z. Ji, H. Ji, S. C. Hicks, A systematic evaluation of single-cell
1142 rna-sequencing imputation methods, *Genome biology* 21 (1) (2020) 1–30.
- 1143 [11] G. Eraslan, L. M. Simon, M. Mircea, N. S. Mueller, F. J. Theis, Single-cell
1144 rna-seq denoising using a deep count autoencoder, *Nature communications*
1145 10 (1) (2019) 1–14.
- 1146 [12] W. V. Li, J. J. Li, An accurate and robust imputation method scimpute
1147 for single-cell rna-seq data, *Nature communications* 9 (1) (2018) 1–9.

- 1148 [13] D. Van Dijk, R. Sharma, J. Nainys, K. Yim, P. Kathail, A. J. Carr, C. Bur-
1149 dziak, K. R. Moon, C. L. Chaffer, D. Pattabiraman, et al., Recovering gene
1150 interactions from single-cell data using data diffusion, *Cell* 174 (3) (2018)
1151 716–729.
- 1152 [14] M. Marouf, P. Machart, V. Bansal, C. Kilian, D. S. Magruder, C. F. Krebs,
1153 S. Bonn, Realistic in silico generation and augmentation of single-cell rna-
1154 seq data using generative adversarial networks, *Nature communications*
1155 11 (1) (2020) 1–12.
- 1156 [15] M. Lotfollahi, F. A. Wolf, F. J. Theis, scgen predicts single-cell perturbation
1157 responses, *Nature methods* 16 (8) (2019) 715–721.
- 1158 [16] Y. Zhao, C. Kilian, J.-E. Turner, L. Bosurgi, K. Roedl, P. Bartsch, A.-C.
1159 Gnirck, F. Cortesi, C. Schultheiß, M. Hellmig, et al., Clonal expansion and
1160 activation of tissue-resident memory-like th17 cells expressing gm-csf in the
1161 lungs of patients with severe covid-19, *Science Immunology* 6 (56) (2021)
1162 eabf6692.
- 1163 [17] I. Tolstikhin, O. Bousquet, S. Gelly, B. Schoelkopf, Wasserstein auto-
1164 encoders, arXiv preprint arXiv:1711.01558 (2017).
- 1165 [18] T. Stuart, A. Butler, P. Hoffman, C. Hafemeister, E. Papalexi, W. M.
1166 Mauck III, Y. Hao, M. Stoeckius, P. Smibert, R. Satija, Comprehensive
1167 integration of single-cell data, *Cell* 177 (7) (2019) 1888–1902.
- 1168 [19] V. Dumoulin, J. Shlens, M. Kudlur, A learned representation for artistic
1169 style, arXiv preprint arXiv:1610.07629 (2016).
- 1170 [20] J. L. Ba, J. R. Kiros, G. E. Hinton, Layer normalization, arXiv preprint
1171 arXiv:1607.06450 (2016).
- 1172 [21] Z. Fu, E. R. Gilbert, D. Liu, Regulation of insulin synthesis and secretion
1173 and pancreatic beta-cell dysfunction in diabetes, *Current diabetes reviews*
1174 9 (1) (2013) 25–53.
- 1175 [22] M. Slyper, C. Porter, O. Ashenberg, J. Waldman, E. Drokhlyansky,
1176 I. Wakiro, C. Smillie, G. Smith-Rosario, J. Wu, D. Dionne, et al., A single-
1177 cell and single-nucleus rna-seq toolbox for fresh and frozen human tumors,
1178 *Nature medicine* 26 (5) (2020) 792–802.
- 1179 [23] G. C. Linderman, J. Zhao, M. Roulis, P. Bielecki, R. A. Flavell, B. Nadler,
1180 Y. Kluger, Zero-preserving imputation of single-cell rna-seq data, *Nature*
1181 *Communications* 13 (1) (2022) 1–11.
- 1182 [24] E. Bakos, C. A. Thaiss, M. P. Kramer, S. Cohen, L. Radomir, I. Orr,
1183 N. Kaushansky, A. Ben-Nun, S. Becker-Herman, I. Shachar, Ccr2 regulates
1184 the immune response by modulating the interconversion and function of
1185 effector and regulatory t cells, *The Journal of Immunology* 198 (12) (2017)
1186 4659–4671.

- 1187 [25] G. Monaco, B. Lee, W. Xu, S. Mustafah, Y. Y. Hwang, C. Carré, N. Burdin,
1188 L. Visan, M. Ceccarelli, M. Poidinger, et al., Rna-seq signatures normalized
1189 by mrna abundance allow absolute deconvolution of human immune cell
1190 types, *Cell reports* 26 (6) (2019) 1627–1640.
- 1191 [26] V. A. Traag, L. Waltman, N. J. Van Eck, From louvain to leiden: guaran-
1192 teeing well-connected communities, *Scientific reports* 9 (1) (2019) 1–12.
- 1193 [27] M. Croft, Control of immunity by the tnfr-related molecule ox40 (cd134),
1194 *Annual review of immunology* 28 (2009) 57–78.
- 1195 [28] T. Riaz, L. M. Sollid, I. Olsen, G. A. de Souza, Quantitative proteomics of
1196 gut-derived th1 and th1/th17 clones reveal the presence of cd28+ nkg2d-
1197 th1 cytotoxic cd4+ t cells, *Molecular & Cellular Proteomics* 15 (3) (2016)
1198 1007–1016.
- 1199 [29] L. Peng, Y. Chen, Q. Ou, X. Wang, N. Tang, Lncrna miat correlates with
1200 immune infiltrates and drug reactions in hepatocellular carcinoma, *Inter-
1201 national immunopharmacology* 89 (2020) 107071.
- 1202 [30] D. P. Saraiva, A. Jacinto, P. Borralho, S. Braga, M. G. Cabral, Hla-dr in
1203 cytotoxic t lymphocytes predicts breast cancer patients’ response to neoad-
1204 jvant chemotherapy, *Frontiers in immunology* (2018) 2605.
- 1205 [31] M. S. Lee, K. Hanspers, C. S. Barker, A. P. Korn, J. M. McCune, Gene
1206 expression profiles during human cd4+ t cell differentiation, *International
1207 immunology* 16 (8) (2004) 1109–1124.
- 1208 [32] S. Aibar, C. B. González-Blas, T. Moerman, V. A. Huynh-Thu, H. Im-
1209 richova, G. Hulselmans, F. Rambow, J.-C. Marine, P. Geurts, J. Aerts,
1210 et al., Scenic: single-cell regulatory network inference and clustering, *Nat-
1211 ure methods* 14 (11) (2017) 1083–1086.
- 1212 [33] A. M. Thornton, J. Lu, P. E. Korty, Y. C. Kim, C. Martens, P. D. Sun,
1213 E. M. Shevach, Helios+ and helios- treg subpopulations are phenotypically
1214 and functionally distinct and express dissimilar tcr repertoires, *European
1215 journal of immunology* 49 (3) (2019) 398–412.
- 1216 [34] C. Imbratta, H. Hussein, F. Andris, G. Verdeil, c-maf, a swiss army knife
1217 for tolerance in lymphocytes, *Frontiers in immunology* 11 (2020) 206.
- 1218 [35] X. O. Yang, B. P. Pappu, R. Nurieva, A. Akimzhanov, H. S. Kang,
1219 Y. Chung, L. Ma, B. Shah, A. D. Panopoulos, K. S. Schluns, et al., T
1220 helper 17 lineage differentiation is programmed by orphan nuclear recep-
1221 tors ror α and ror γ , *Immunity* 28 (1) (2008) 29–39.
- 1222 [36] K. Street, D. Risso, R. B. Fletcher, D. Das, J. Ngai, N. Yosef, E. Purdom,
1223 S. Dudoit, Slingshot: cell lineage and pseudotime inference for single-cell
1224 transcriptomics, *BMC genomics* 19 (1) (2018) 1–16.

- 1225 [37] E. Stephenson, G. Reynolds, R. A. Botting, F. J. Calero-Nieto, M. D.
1226 Morgan, Z. K. Tuong, K. Bach, W. Sungnak, K. B. Worlock, M. Yoshida,
1227 et al., Single-cell multi-omics analysis of the immune response in covid-19,
1228 *Nature medicine* 27 (5) (2021) 904–916.
- 1229 [38] M. Ota, Y. Nagafuchi, H. Hatano, K. Ishigaki, C. Terao, Y. Takeshima,
1230 H. Yanaoka, S. Kobayashi, M. Okubo, H. Shirai, et al., Dynamic landscape
1231 of immune cell-specific gene regulation in immune-mediated diseases, *Cell*
1232 184 (11) (2021) 3006–3021.
- 1233 [39] G. Meyer Zu Horste, C. Wu, C. Wang, L. Cong, M. Pawlak, Y. Lee,
1234 W. Elyaman, S. Xiao, A. Regev, V. Kuchroo, Rbpj controls development
1235 of pathogenic th17 cells by regulating il-23 receptor expression. *cell rep* 16
1236 (2): 392–404 (2016).
- 1237 [40] S. De Biasi, M. Meschiari, L. Gibellini, C. Bellinazzi, R. Borella, L. Fianza,
1238 L. Gozzi, A. Iannone, D. Lo Tartaro, M. Mattioli, et al., Marked t cell
1239 activation, senescence, exhaustion and skewing towards th17 in patients
1240 with covid-19 pneumonia, *Nature communications* 11 (1) (2020) 1–17.
- 1241 [41] B. J. Meckiff, C. Ramírez-Suástegui, V. Fajardo, S. J. Chee, A. Kusanadi,
1242 H. Simon, S. Eschweiler, A. Grifoni, E. Pelosi, D. Weiskopf, et al., Imbal-
1243 ance of regulatory and cytotoxic sars-cov-2-reactive cd4+ t cells in covid-19,
1244 *Cell* 183 (5) (2020) 1340–1353.
- 1245 [42] L. Loyal, S. Warth, K. Jürchott, F. Mölder, C. Nikolaou, N. Babel,
1246 M. Nienen, S. Durlanik, R. Stark, B. Kruse, et al., Slamf7 and il-6r define
1247 distinct cytotoxic versus helper memory cd8+ t cells, *Nature communica-*
1248 *tions* 11 (1) (2020) 1–12.
- 1249 [43] J. Yang, M. Zhong, E. Zhang, K. Hong, Q. Yang, D. Zhou, J. Xia, Y.-Q.
1250 Chen, M. Sun, B. Zhao, et al., Broad phenotypic alterations and potential
1251 dysfunction of lymphocytes in individuals clinically recovered from covid-
1252 19, *Journal of Molecular Cell Biology* 13 (3) (2021) 197–209.
- 1253 [44] M. Lotfollahi, M. Naghipourfar, M. D. Luecken, M. Khajavi, M. Büttner,
1254 M. Wagenstetter, Ž. Avsec, A. Gayoso, N. Yosef, M. Interlandi, et al.,
1255 Mapping single-cell data to reference atlases by transfer learning, *Nature*
1256 *Biotechnology* 40 (1) (2022) 121–130.
- 1257 [45] C. Wagner, M. Griesel, A. Mikolajewska, A. Mueller, M. Nothacker,
1258 K. Kley, M.-I. Metzendorf, A.-L. Fischer, M. Kopp, M. Stegemann, et al.,
1259 Systemic corticosteroids for the treatment of covid-19, *Cochrane Database*
1260 *of Systematic Reviews* (8) (2021).
- 1261 [46] W. Chen, J. Luo, Y. Ye, R. Hoyle, W. Liu, R. Borst, S. Kazani, E. A.
1262 Shikatani, V. J. Erpenbeck, I. D. Pavord, et al., The roles of type 2 cyto-
1263 toxic t cells in inflammation, tissue remodeling, and prostaglandin (pg) d2
1264 production are attenuated by pgd2 receptor 2 antagonism, *The Journal of*
1265 *Immunology* 206 (11) (2021) 2714–2724.

- 1266 [47] S. Lab, `panc8.SeuratData`: Eight Pancreas Datasets Across Five Technolo-
1267 gies, `r` package version 3.0.2 (2019).
- 1268 [48] J. Ding, X. Adiconis, S. K. Simmons, M. S. Kowalczyk, C. C. Hession,
1269 N. D. Marjanovic, T. K. Hughes, M. H. Wadsworth, T. Burks, L. T.
1270 Nguyen, et al., Systematic comparison of single-cell and single-nucleus rna-
1271 sequencing methods, *Nature biotechnology* 38 (6) (2020) 737–746.
- 1272 [49] A. Gayoso, R. Lopez, G. Xing, P. Boyeau, K. Wu, M. Jayasuriya, E. Melh-
1273 man, M. Langevin, Y. Liu, J. Samarán, et al., `Scvi-tools`: A library for
1274 deep probabilistic analysis of single-cell omics data, *bioRxiv* (2021).
- 1275 [50] G. X. Zheng, J. M. Terry, P. Belgrader, P. Ryvkin, Z. W. Bent, R. Wilson,
1276 S. B. Ziraldo, T. D. Wheeler, G. P. McDermott, J. Zhu, et al., Massively
1277 parallel digital transcriptional profiling of single cells, *Nature communica-*
1278 *tions* 8 (1) (2017) 1–12.
- 1279 [51] F. A. Wolf, P. Angerer, F. J. Theis, `Scanpy`: large-scale single-cell gene
1280 expression data analysis, *Genome biology* 19 (1) (2018) 1–5.
- 1281 [52] Y. Le Cun, F. Fogelman-Soulié, Modèles connexionnistes de
1282 l’apprentissage, *Intellectica* 2 (1) (1987) 114–143.
- 1283 [53] G. E. Hinton, R. Zemel, Autoencoders, minimum description length and
1284 helmholtz free energy, *Advances in neural information processing systems*
1285 6 (1993).
- 1286 [54] D. P. Kingma, M. Welling, Auto-encoding variational bayes, *arXiv preprint*
1287 *arXiv:1312.6114* (2013).
- 1288 [55] C. Villani, *Optimal transport: old and new*, Vol. 338, Springer, 2009.
- 1289 [56] M. Arjovsky, S. Chintala, L. Bottou, Wasserstein generative adversarial
1290 networks, in: *International conference on machine learning*, PMLR, 2017,
1291 pp. 214–223.
- 1292 [57] P. K. Rubenstein, B. Schoelkopf, I. Tolstikhin, On the latent space of
1293 wasserstein auto-encoders, *arXiv preprint arXiv:1802.03761* (2018).
- 1294 [58] D. Risso, F. Perraudeau, S. Gribkova, S. Dudoit, J.-P. Vert, A general and
1295 flexible method for signal extraction from single-cell rna-seq data, *Nature*
1296 *communications* 9 (1) (2018) 1–17.
- 1297 [59] P. J. Huber, Robust estimation of a location parameter: *Annals mathe-*
1298 *matics statistics*, 35 (1964).
- 1299 [60] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, A. Smola, A ker-
1300 nel two-sample test, *Journal of Machine Learning Research* 13 (25) (2012)
1301 723–773.
1302 URL <http://jmlr.org/papers/v13/gretton12a.html>

- 1303 [61] D. Misra, Mish: A self regularized non-monotonic activation function,
1304 arXiv preprint arXiv:1908.08681 (2019).
- 1305 [62] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, J. Han, On the variance
1306 of the adaptive learning rate and beyond, arXiv preprint arXiv:1908.03265
1307 (2019).
- 1308 [63] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv
1309 preprint arXiv:1412.6980 (2014).
- 1310 [64] R. Liaw, E. Liang, R. Nishihara, P. Moritz, J. E. Gonzalez, I. Stoica, Tune:
1311 A research platform for distributed model selection and training, arXiv
1312 preprint arXiv:1807.05118 (2018).
- 1313 [65] D. Kobak, P. Berens, The art of using t-sne for single-cell transcriptomics,
1314 Nature communications 10 (1) (2019) 1–14.
- 1315 [66] G. C. Linderman, M. Rachh, J. G. Hoskins, S. Steinerberger, Y. Kluger,
1316 Fast interpolation-based t-sne for improved visualization of single-cell rna-
1317 seq data, Nature methods 16 (3) (2019) 243–245.
- 1318 [67] S. V. Stassen, D. M. Siu, K. C. Lee, J. W. Ho, H. K. So, K. K. Tsia, Parc:
1319 ultrafast and accurate clustering of phenotypic data of millions of single
1320 cells, Bioinformatics 36 (9) (2020) 2778–2786.
- 1321 [68] A. Subramanian, P. Tamayo, V. K. Mootha, S. Mukherjee, B. L. Ebert,
1322 M. A. Gillette, A. Paulovich, S. L. Pomeroy, T. R. Golub, E. S. Lander,
1323 et al., Gene set enrichment analysis: a knowledge-based approach for in-
1324 terpreting genome-wide expression profiles, Proceedings of the National
1325 Academy of Sciences 102 (43) (2005) 15545–15550.
- 1326 [69] E. Y. Chen, C. M. Tan, Y. Kou, Q. Duan, Z. Wang, G. V. Meirelles, N. R.
1327 Clark, A. Ma’ayan, Enrichr: interactive and collaborative html5 gene list
1328 enrichment analysis tool, BMC bioinformatics 14 (1) (2013) 1–14.
- 1329 [70] J. H. Friedman, Greedy function approximation: a gradient boosting ma-
1330 chine, Annals of statistics (2001) 1189–1232.
- 1331 [71] S. Haghghi, M. Jasemi, S. Hessabi, A. Zolanvari, Pycm: Multiclass con-
1332 fusion matrix library in python, Journal of Open Source Software 3 (25)
1333 (2018) 729.