

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,f}, Nicola Gagliani^{c,d,e}, Samuel Huber^{c,d}, Pierre Machart^{a,b,*}, Stefan Bonn^{a,b,*}

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

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

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

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

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

^f*III. Department of Medicine, University Medical Center Hamburg-Eppendorf, 20246 Hamburg, Germany*

Abstract

Single cell sequencing provides detailed insights into biological processes including cell differentiation and identity. While providing deep cell-specific information, the method suffers from technical constraints, most notably a limited number of expressed genes per cell, which leads to suboptimal clustering and cell type identification. Here we present DISCERN, a novel deep generative network that reconstructs missing single cell gene expression using a reference dataset. DISCERN outperforms competing algorithms in expression inference resulting in greatly improved cell clustering, cell type and activity detection, and insights into the cellular regulation of disease. We used DISCERN to detect two 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 scRNA-seq data using a zero-inflated negative binomial model of the gene expression. MAGIC uses a nearest neighbor diffusion graph to impute gene expression and scImpute estimates gene expression and drop-out probabilities using linear regression. All of these algorithms use information from similar cells with measured expression of the same dataset for imputation. Another class of imputation algorithms use bulk RNA-seq data to constrain scRNA-seq expression imputation. BfImpute [14] uses Bayesian factorization, SCRABBLE [15] matrix regularization, and

39 SIMPLEs [16] a prior distribution on the bulk data to impute scRNA-seq ex-
40 pression. Unfortunately, SCRABBLE and Bfimpute do not scale beyond small
41 single cell datasets and few genes (3000 cells and genes in our hands), and
42 SIMPLEs requires matching single cell and bulk RNA-seq samples, severely
43 constraining their usability.

44 Similarly, multigrade[17], BABEL [18], and Cross Modal Autoencoders[19]
45 use scRNA-seq in combination with complementary, matching data (e.g. CITE-
46 seq, ATAC-seq) to improve imputation. While complementary CITE-seq infor-
47 mation is available for many scRNA-seq datasets, other information such as
48 ATAC-seq data of the same sample is usually missing, which severely constrains
49 the usability of BABEL and Cross Modal Autoencoder.

50 While current imputation methods provide improved gene expression infor-
51 mation, they still rely on the comparison of similar cells with largely absent gene
52 expression information, for example by using clustering approaches. Genes that
53 are not expressed in neighboring cells cannot be imputed, limiting the value of
54 classical imputation. In an ideal case, it would be possible to obtain information
55 of the expected true gene expression per cell, or at least expression information
56 with less technical noise, to reconstruct the true expression at single cell level.

57 Recent work has shown the effectiveness of deep generative models (e.g. Au-
58 toencoders and Generative Adversarial Networks) to infer realistic scRNA-seq
59 data and augment scarce cell populations using Generative Adversarial Net-
60 works [20] or the prediction of perturbation response using Autoencoders [21].
61 We hypothesized that a deep generative model could allow for the reconstruc-
62 tion of missing single cell gene expression information (low quality - lq) by using
63 related data with more genes expressed (high-quality - hq) as a reference, a com-
64 pletely novel approach to gene expression inference (Figure 1A). In other words,
65 lq data with many missing gene expression values and bad clustering could be
66 transformed into data with few missing genes and improved clustering if the
67 “style” of a related hq dataset could be transferred to it. In the best case,
68 it would be possible to infer gene expression information for single cell data
69 (lq) by using purified bulk RNA-seq data (hq), obtaining over ten thousand
70 genes expressed per cell. We envision that this novel approach, when properly
71 calibrated, is transformative for the analysis of single cell data, gaining deep
72 mechanistic insights into data beyond what is currently measurable. It is im-
73 portant to note that the concept of using hq data to reconstruct gene expression
74 in lq data is fundamentally different from classical imputation algorithms that
75 infer gene expression based on nearby cells from the same dataset, as outlined
76 above.

77 Based on the above considerations, we developed DISCERN, a novel deep
78 generative neural network for directed single cell expression reconstruction. DIS-
79 CERN allows for the realistic reconstruction of gene expression information by
80 transferring the style of hq data onto lq data, in latent and gene space. Our ex-
81 periments on real and simulated data show that DISCERN outperforms several
82 existing algorithms in gene expression inference across a wide array of single
83 cell datasets and technologies, improving cell clustering, cell type and activity
84 detection, and pathway and gene regulation identification. To obtain deep in-

85 sights into the cellular changes underlying COVID-19, we reconstructed single
86 cell expression data of patient blood and lung immune data. While in our ini-
87 tial analysis [22] of blood data we detected few immune cell types, expression
88 reconstruction with DISCERN resulted in the detection of 28 cell types and
89 states in blood, including two novel disease-associated T cell types, cytotoxic
90 CD4⁺ and CD8⁺ Tc2 T helper cells. Reconstructing a second COVID-19 blood
91 dataset with disease severity information, we were able to classify mild and se-
92 vere COVID-19 with an AUROC of 81 %, obtaining a potential biomarker of
93 disease stage. DISCERN can be easily integrated into existing workflows, as an
94 additional step after count mapping. Given that DISCERN is not limited by
95 a predefined distribution of data, we believe that it can be readily adapted to
96 enhance various other biomedical data types, especially other omics data such
97 as proteomics and spatial transcriptomics.

98 2. Results

99 2.1. The DISCERN algorithm for directed expression reconstruction

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

110 With these prerequisites in mind, we designed a deep neural network for
111 directed single cell expression reconstruction (DISCERN) (Figure S1B) that is
112 based on a modified Wasserstein Autoencoder [23]. A unique feature of DIS-
113 CERN is that it transfers the "style" of hq onto lq data to reconstruct missing
114 gene expression, which sets it apart from other batch correction methods such
115 as [24], which operate in a lower dimensional representation of the data (e.g.
116 PCA, CCA). To allow DISCERN to accurately reconstruct single cell RNA-
117 seq expression based on reference data, the structure of the network had to be
118 adapted in several ways. First, we implemented Conditional Layer Normaliza-
119 tion (CLN) [25, 26, 20] to allow for directed expression reconstruction of lq data
120 based on reference hq data (Figure S1B & S2). Second, we added two decoder
121 heads to the network to enable it to model dataset-specific dropout rates and
122 gene expression separately. Lastly, we extended DISCERN's loss function with
123 a binary cross-entropy term for learning the probability of dropouts to increase
124 general inference fidelity. Further algorithmic details of DISCERN can be found
125 in the methods and Figures S1 and S2.

126 We first demonstrate DISCERN's capabilities to faithfully reconstruct gene
127 expression using five pancreas single cell expression datasets of varying quality

128 (Tables S1 and S2). The pancreas data is widely used for benchmarking and it
129 is ideal to evaluate expression reconstruction for many cell types and sequencing
130 technologies. We consider a dataset as hq when the average number of genes
131 detected per cell (GDC) (e.g. smartseq2, GDC 6214) is much higher than in a
132 comparable lq dataset (Table S2). Conversely, a dataset is lq when the average
133 cell has lower counts and fewer genes expressed than a comparable hq dataset
134 (e.g. indrop, GDC 1887). Throughout this text, we will name sequencing tech-
135 nologies with capital (e.g. Smart-Seq2, InDrop) and datasets with lower case
136 first letters (smartseq2, indrop). We trained DISCERN on these five pancreatic
137 single cell datasets and assessed the integration of data in gene space and the
138 average expression reconstruction per cell type. While uncorrected data cluster
139 by batch and not by cell type, DISCERN-integrated data show good batch
140 mixing and clustering of cells by cell type across all five datasets (Figure 1B &
141 Figure S2). To get a clearer picture of DISCERN's expression reconstruction
142 capabilities we next calculated correlation coefficients of measured expression
143 between the lowest quality inDrop and highest quality Smart-Seq2 data, before
144 and after expression reconstruction using DISCERN. The mean expression re-
145 construction of indrop-lq to smartseq2-hq and smartseq2-hq to indrop-lq data
146 is very accurate, showing a Pearson correlation of $r = 0.95$ ($p < 0.001$), while
147 mean expression correlation between uncorrected indrop-lq and smartseq2-hq
148 data is only $r = 0.77$ due to strong batch effects (Figure 1C & D, Figures S3
149 and S4). The improved quality of indrop-lq data reconstructed to smartseq2-hq
150 level is validated by the strong increase of genes expressed per cell, ranging from
151 ≈ 2000 genes per cell in the uncorrected indrop-lq data to ≈ 6000 genes in the
152 indrop-lq data after reconstruction (Figure S5).

153 We next investigated the effect of reconstruction of three cell type-specific
154 genes, before and after correction across the five pancreas datasets (Figure S6).
155 Insulin expression in the pancreas should be largely restricted to beta cells [27],
156 which can be observed in the uncorrected smartseq2-hq and celseq2 datasets,
157 while the indrop-lq batch shows a diffuse pattern of insulin expression across
158 cell types (Figure S6A left panel). This diffuse insulin expression is corrected
159 by reconstructing the smartseq2-hq expression pattern from the indrop-lq data
160 (Figure S6A middle panel). In general, the expected specificity of insulin ex-
161 pression in beta cells can be recovered for all datasets when using DISCERN's
162 reconstruction using the smartseq2-hq reference. Conversely, the reconstruction
163 from hq to the indrop-lq reference results in diffuse insulin expression across all
164 reconstructed datasets (Figure S6A right panel). We obtained similar results for
165 the pancreatic acinar cell-specific gene REG1A and the delta cell-specific gene
166 SST, both of which show diffuse expression across cell types in the uncorrected
167 inDrop data and cell-specific expression after reconstruction using smartseq2-hq
168 reference (Figure S6B & C). Interestingly, DISCERN can not only recover bio-
169 logical expression information, but it is also able to apply sequencing method-
170 specific effects after reconstruction. The smartseq2-hq dataset, for instance,
171 displays nearly no ribosomal protein coding gene expression after sequencing as
172 previously reported by [8], while data sequenced using InDrop, Cel-Seq, or Cel-
173 Seq shows prominent ribosomal protein coding gene expression (Figure S6D, left

174 panel). When reconstructing smartseq2-hq data to indrop-lq data, ribosomal
175 protein coding gene expression is re-instantiated (Figure S6D, right panel).

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

190 2.2. Specific and robust gene expression inference

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

194 Since expression reconstruction can be seen as a generalization of expression
195 imputation, we compared DISCERN to DCA, MAGIC, and scImpute, three
196 state-of-the-art imputation algorithms [11, 12, 13]. Expression reconstruction
197 can also be viewed as a batch correction task in gene space, which is why we ad-
198 ditionally compared DISCERN to scGEN and Seurat [21, 24]. It is important to
199 note, however, that neither Seurat nor scGEN were designed for the expression
200 reconstruction task. Seurat and scGEN use a lower dimensional representation
201 in which a linear transformation aligns different batches. Seurat uses canonical
202 correlation analysis and scGEN uses the bottleneck layer representation of an
203 autoencoder to calculate and apply linear transformations.

204 To investigate the precision of gene expression reconstruction, we created an
205 artificial dataset by dividing the smartseq2-hq pancreas data into two batches,
206 smartseq-lq and smartseq2-hq. In the smartseq-lq batch, the top one KEGG
207 pathways per cell type were removed by setting the expression of genes con-
208 tained in these pathways to zero, while the smartseq2-hq remained unaltered.
209 Therefore, a reconstruction of smartseq-lq data using smartseq2-hq reference
210 (reconstructed-hq) should ideally recover the smartseq-lq expression to its orig-
211 inal state, prior to the removal of the genes. DISCERN is able to reconstruct
212 the mean expression for all cell types, achieving a correlation $r = 0.99$ (Fig-
213 ure 2A). DCA ($r = 0.66$), MAGIC ($r = 0.34$), scImpute ($r = 0.80$), and Seurat
214 ($r = 0.76$) have significantly lower correlation between the smartseq2-hq and
215 reconstructed-hq gene expression (Figure 2A). scGen shows only slightly reduced
216 performance ($r = 0.98$) compared to DISCERN, especially in the reconstruction
217 of highly expressed genes (Figure 2A) and low abundant cell types (Figure S10,

218 Megakaryocytes). We obtained similar results on the difftec dataset, with DIS-
219 CERN ($r = 0.98$) outperforming DCA ($r = 0.47$), Magic ($r = 0.21$), scImpute
220 ($r = 0.04$), Seurat ($r = 0.92$), and scGEN ($r = 0.94$) (Figure S10). To further
221 investigate gene expression reconstruction specificity, we compared the correla-
222 tion of reconstructed-hq to smartseq2-hq data after performing differential gene
223 expression (DEG) for each cell type against all other cell types (Figure 2B, up-
224 per panel). DISCERN is able to recover the correct DEG t-statistics with a
225 median correlation of 0.92, improving over state-of-the-art tools by more than
226 15 percentage points. In the corresponding experiment using the difftec dataset,
227 DISCERN achieves a median correlation of 0.85, which is a 25 percentage point
228 improvement over competing methods (Figure S11).

229 Since the genes were initially selected using KEGG gene set enrichment
230 analysis, the reconstruction of the corresponding pathways was investigated by
231 performing KEGG gene set enrichment analysis on the DEG results. DISCERN
232 is able to recover the pathway expression enrichment scores with a median cor-
233 relation of 0.93, exceeding the performance of Seurat and scGEN by more than
234 11 percentage points on median (Figure 2B, lower panel). In the corresponding
235 experiment using the difftec dataset, DISCERN achieves a median correlation
236 of 0.77, outperforming Seurat and scGen by more than 16 percentage points
237 (Figure S12).

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

245 Next, we show DISCERN's expression reconstruction robustness with re-
246 spect to varying sizes of lq to hq data. It is conceivable to assume that a large
247 amount of hq data would benefit the expression reconstruction of the lq data,
248 which makes it important to understand at what ratio good results can be ex-
249 pected. Interestingly, DISCERN seems to be very robust across a wide range
250 of smartseq2-lq to smartseq2-hq ratios, with correlations of 0.98 (ratio of lq/hq
251 0.14) to 0.93 (ratio of lq/hq 18.4), while the second-best performing algorithm
252 scGen showed a 11 percentage point decrease in performance (0.82 for ratio of
253 lq/hq 18.4) (Figure 2C, Figure S15). We observed similar results for the correla-
254 tion of t-statistics, showing a slight dependence of DISCERN's performance on
255 the lq/hq ratio (Figure S16). In general, all methods show better performance
256 with a small ratio of lq/hq data, while DISCERN shows least dependence and
257 outperforms other algorithms in the correlation of expression and t-statistics,
258 especially in the case of high lq/hq ratio.

259 Another aspect of expression reconstruction robustness is the dependence of
260 the algorithm on the cell type or cell state similarity of the lq and hq datasets.
261 In the optimal case, DISCERN would not require that the lq and hq datasets
262 have overlapping cell types to perform an accurate expression reconstruction,
263 which is theoretically possible if the network learns the general gene-regulatory

264 expression logic of the hq data (see discussion). To understand the dependence
265 on dataset similarity, we removed a complete cell type, pancreas alpha cells, from
266 the smartseq2-hq data and left the alpha cells in the smartseq2-lq data. We then
267 additionally varied the number of common cells in the lq and hq data, starting
268 with no overlapping cells (only alpha cells in the lq and all cells except alpha in
269 the hq data) and ending with almost complete overlap (all cells overlap between
270 the smartseq2-hq and -lq data, except for the alpha cells only present in lq data)
271 (Figure 2D). When evaluating DEG correlation, DISCERN was the only method
272 consistently achieving better performance than uncorrected data, outperforming
273 Seurat and scGen by more than 15 percentage points (Figure 2D). Similarly,
274 DISCERN was the only method consistently achieving better performance than
275 uncorrected data in the FC correlation task (Figure S17).

276 We next took a closer look at the integration and expression reconstruction
277 performance when no cell types overlap between the lq (alpha cells only) and hq
278 (all other cells) data. Notably, Seurat seems to over-integrate cell types, mix-
279 ing smartseq2-hq beta and gamma cells with reconstructed-hq alpha cells from
280 other batches (Figure S18), while scGEN and DISCERN keep the smartseq2-hq
281 and reconstructed-hq exclusive cell types separate (Figure 2E & Figure S18).
282 This over-integration seems to be causal for Seurat's poor DEG correlation per-
283 formance ($r = 0.28$), while DISCERN ($r = 0.55$) is the only method achieving
284 better performance than uncorrected cells ($r = 0.47$) (Figure 2F). Thus, DIS-
285 CERN is able to keep existing expression correlations and improves the detec-
286 tion of cell type specific genes by reconstruction using an hq batch as reference.
287 In conclusion, DISCERN is both a precise and robust method for expression
288 reconstruction that outperforms existing methods by a significant margin.

289 *2.3. Improving cell cluster, type, and trajectory identification*

290 The comparison to competing methods provided evidence for DISCERN's
291 superior expression reconstruction. Now, we will delineate how DISCERN's
292 expression reconstruction improves downstream cell clustering, cell type and
293 activity state identification, marker gene determination, and gene regulatory
294 network and cell trajectory analysis.

295 To understand if cell-determining gene expression and pathways could be
296 recovered with expression reconstruction, we used a single nuclear sequencing
297 (sn-lq) and scRNA-seq (sc-hq) data pair that was prepared from the same liver
298 metastasis biopsy [28]. We reconstructed sn-lq data using the sc-hq reference,
299 obtaining reconstructed-hq data. While single nuclear sequencing provides re-
300 duced expression information in the average counts per cell as compared to
301 scRNA-seq (Table S2) [28], it is still the method of choice to obtain cell-specific
302 expression information when intact single cells cannot be recovered from a tis-
303 sue (e.g. after tissue fixation or freezing). It is important to note that nuclear
304 transcripts reflect current gene activity, which in part might not correlate with
305 transcripts that have lifetimes of up to days. Before integration, the sn-lq and
306 sc-hq datasets cluster by batch and not by cell type, while after expression re-
307 construction with DISCERN cells cluster by type and not by batch (Figure S19).
308 This is reflected in an expression correlation of 0.49 (sc-hq vs. sn-lq) before and

0.93 after reconstruction (sc-hq vs. reconstructed-hq) (Figure S20). Seurat reconstructed expression, on the other hand, is barely different from uncorrected sn-lq data. This is reflected in a similar UMAP representation (Figure S19) and an identical expression correlation of 0.49 with uncorrected sn-lq data (Figure S20). DISCERN reconstruction resulted in the expression of T cell receptor signaling genes in reconstructed T cells (Figure S21) and antigen presentation genes in macrophages (Figure S22), providing evidence that DISCERN faithfully recreates cell-determining genes and pathways based on the hq data. Seurat is not able to reconstruct the expression information and shows a similar expression pattern as the uncorrected sn-lq dataset. In both datasets (seurat-hq and sn-lq) the expression of important T cell marker genes such as *CD3E*, *CD3D* and *CD8A* is largely absent, while in sc-hq and reconstructed-hq the expression is easily detectable (Figure S21). To further corroborate the advantage of single nuclear expression reconstruction, we next aimed to increase the T cell subtype resolution of human single nucleus acute kidney injury data (kidney-lq) by using matching single cell data (kidney-hq). Only 1% of kidney-lq nuclei show *CD3D*, *CD3E* or *CD3G* expression, compared to 7% of the cells in the kidney-hq dataset. Seurat and DISCERN were able to detect T cells in the reconstructed kidney-lq (reconstructed-hq) and the kidney-hq data with notable *CD3D* expression in this cluster (Figure S23). The reconstructed-hq and the kidney-hq T cells were further classified into T cell subtypes and activation states (Figure S23C). While a large proportion of T cells detected in Seurat reconstructed data could not be annotated due to missing *CD3D*, *CD4*, and *CD8A* expression, DISCERN reconstructed data does not present these limitations.

It is intriguing to observe that many marker genes are still hard to detect in kidney single cell RNA-seq data but also in the antigen presentation pathway in macrophages (Figure S22). This is most probably due to dropout. Thus, we rationalized that bulk RNA sequencing (RNA-seq) data of purified cell types (e.g. FACS sorted immune cells) is a suitable hq proxy for the expected gene expression per cell. RNA-seq data of purified cells is readily available from public repositories, making it possible to obtain thousands of purified immune cell RNA-seq samples (see methods). We therefore set out to increase cluster, cell type, gene regulatory network, and trajectory identification of scRNA-seq data by reconstructing gene expression using a related RNA-seq reference (Figure S24). For the scRNA-seq data we chose a cord blood mononuclear citeseq dataset (cite-lq) that was labeled with 15 antibodies (Table S3) to allow for surface protein-based cell type discovery [29]. The CITE-seq information allowed us to confirm expression reconstruction by DISCERN in cases where gene expression is absent but protein expression and cell identity are validated via antibody labeling. For the RNA-seq data, we selected 9.852 purified immune samples (bulk-hq) and proceeded to reconstruct cite-lq (GDC 798) using a bulk-hq (GDC 13.104) reference to obtain reconstructed-hq data with DISCERN. We first investigated the correspondence of gene expression prior (cite-lq) and post reconstruction (bulk-hq) with antibody-based surface protein labeling of *CD3D*, *CD4*, *CD8A*, *CD2*, *B3GAT1*, *FCGR3A*, *CD14*, *ITGAX* and *CD19* (Figure 3A, Figure S25). For several proteins (CD8A, B3GAT1, CD4), the corresponding

355 cite-lq gene expression was absent and cell type-specifically re-instantiated in
356 the reconstructed-hq expression data with DISCERN (Figure 3A, Figure S25).
357 In cases where cell type-specific gene and protein expression matched cite-lq
358 data (*CD3D*, *CD14*) the expression in reconstructed-hq data was left unaltered
359 (Figure S25). In some instances, we observed low cell type-specific expression
360 in the cite-lq data (*CD8A*, *CD2*, *FCGR3A*, *CD19*) that matched protein ex-
361 pression (Figure S25). In these cases, gene expression was increased in the cor-
362 rect cell types in the reconstructed-hq data. In general, we observed increased
363 agreement between cell type-specific surface protein and gene expression af-
364 ter reconstruction, showing that DISCERN doesn't invent or 'hallucinate' cell
365 types but reconstructs the expected expression specific for each cell type. We
366 further corroborated these results by selecting eight known cell type-specific
367 cytosolic proteins and investigated their expression before and after expression
368 reconstruction. *MS4A1* (B cells), *IL7R* (CD4⁺ T cells), *MS4A7* (Monocytes),
369 *GNLY* and *NKG7* (NK cells) showed consistent expression before and after
370 reconstruction (Figure S26). The chemokine receptors *CCR2* (Monocytes, ac-
371 tivated T cells), *CXCR1* (NK cells), and *CXCR6* (CD8⁺ T cells) showed the
372 correct cell type-specific expression only after expression reconstruction (Fig-
373 ure S26) [30]. It is notoriously hard to obtain cell subtype-specific information
374 from blood mononuclear scRNA-seq data, especially for CD4⁺ T helper cells due
375 to their limited activation status in healthy individuals. This doesn't mean that
376 polarized CD4⁺ T helper cells do not exist in healthy blood, as they are com-
377 monly detected after stimulation using FACS (Table S3) [31]. This lack of reso-
378 lution in scRNA-seq impedes clustering, marker gene, and trajectory analyses, a
379 drawback that could be overcome using DISCERN's expression reconstruction.
380 We therefore compared CD4⁺ T cell (gene expression of *CD4* > 1 and *CD3E*
381 > 2.5) clustering and subtype identification using cite-lq and reconstructed-
382 hq data. While clustering with the leiden algorithm [32] using highly variable
383 genes of cite-lq data resulted in an unstructured distribution of CD4⁺ T cell
384 subtypes (Figure 3B), clustering of reconstructed-hq data yields detailed in-
385 sights into T helper cell subtypes of blood mononuclear data (Figure 3C). Af-
386 ter reconstruction, we were able to characterize TH17, TH2, TH1, HLA-DR
387 expressing TREG (Active_TREG), naive CD4⁺ T cells (CD4_naive), effector-
388 memory CD4⁺ T cells (CD4_EM), central-memory CD4⁺ T cells (CD4_CM),
389 and effector cells expressing IFN-regulated genes (IFN_regulated) (Figure 3C).
390 We selected published cell-determining marker genes and observed that many of
391 them were dropped out in the uncorrected data but present after reconstruction
392 (Figure S27). The absence of marker genes in uncorrected data results in poor
393 clustering and cell type identification, while single positive cells are detectable
394 in the respective neighborhood identified by reconstructed counts (Figure S27).
395 Importantly, we observed that in all cases the DISCERN-estimated proportions
396 of T helper subsets fall within the range of expected proportions as assessed by
397 previous FACS studies (Table S3, Figure S28). These findings are important,
398 as they prove once more that DISCERN discovers the correct cell subtypes and
399 cell proportions, in this case substantially outperforming the available CITE-seq
400 information in cell subtype resolution.

401 To further verify the cell type annotations, we extracted the top cluster-
402 determining genes from the reconstructed-hq data. Members of the TNF-
403 receptor superfamily are known to be expressed in T helper cell subtypes [33],
404 which can be observed after reconstruction in TH17 cells and partially in TH1,
405 TH2, Active_TREG and IFN_regulated cells (Figure S29). Similarly, recon-
406 structed TH1 cells show the expected high expression of granzymes *GZMK* and
407 *GZMA* [34], while *MIAT* and *HLA* expression are found in activated TREG
408 cells after reconstruction (Active_TREG cluster, Figure S29) [35, 36]. *NOG* ex-
409 pression is detected in reconstructed CD4_naive cells, as previously described
410 [37]. In addition, reconstructed CD4_naive, CD4_EM and CD4_CM show low
411 expression of the genes important for the T helper subtypes TH1, TH2, TH17,
412 Active_TREG and IFN_regulated. We further corroborated our cell type anno-
413 tation of reconstructed-hq data by observing the expected expression of several
414 established T cell subtype markers (Figure S30). We compared these newly
415 found clusters to representations found with Seurat, multigrade, and in uncor-
416 rected cite-lq data. The uncorrected cite-lq data manifests cluster separation
417 for some cell types, most notably IFN_regulated and Active_TREG cells (Fig-
418 ure S31A). Seurat reconstruction and multigrade imputation with CITE-seq
419 information results in the mixing of cell types and clusters (Figure S31B & C).
420 A further comparison to Bfimpute and SCRABBLE was impossible due to the
421 dataset size, as outlined in the introduction.

422 Similar to improved clustering and cell subtype detection, DISCERN reconstructed-
423 hq data resulted in improved gene regulatory network inference with SCENIC
424 [38]. SCENIC infers transcription factor-regulated gene expression modules
425 of single cell data. While cite-lq data resulted in a scattered distribution of
426 transcription factor networks across several T helper cell subtypes, SCENIC
427 with reconstructed-hq data showed transcription factor regulation in the cor-
428 rect subtypes (Figure 3D). After expression reconstruction the IKZF2 regulon
429 is detected in activated TREG cells [39] and the MAF regulon is found in differ-
430 entiated CD4⁺ T cells but not in naive CD4⁺ T cells [40]. A weak signal of the
431 MAF regulon is already detectable in the cite-lq data, yet strongly increased in
432 reconstructed-lq, while maintaining differentiated T helper cell specificity (Fig-
433 ure 3D). Furthermore, after reconstruction with DISCERN we could identify
434 the TH17 associated master transcriptional regulators RORC(+) and RORA(+)
435 [41], which were scattered over all TH17 cells before reconstruction (Figure S32).
436 Seurat is able to partially reconstruct the expression of the RORC(+) regulon
437 but fails to detect the more specific RORA(+) expression (Figure S32).

438 Finally, we wanted to investigate if DISCERN could also enhance cell trajec-
439 tory analyses with Slingshot of the citeseq data [42]. We focused on the differen-
440 tiation of effector and other T helper cell subtypes and found five lineages that
441 either pass through or terminate in the effector cell cluster in reconstructed-hq
442 data (Figure 3C). Two trajectories were of special interest to us: Lineage1 from
443 CD4_naive to TH1 cells (Figure S33) and Lineage2 from CD4_naive to TH17
444 cells (Figure S34). While the expression change along the trajectory in uncor-
445 rected data (Figure S33A, Figure S34A) is hardly visible, cell type-specific clus-
446 ters can be easily observed after DISCERN reconstruction (for lineage details

447 see Figure S33B, Figure S34B). The detailed insights into cell differentiation
448 that we obtained with reconstructed data are in stark contrast to the Slingshot
449 results obtained with cite-lq data. While terminal effector molecules can be de-
450 tected with cite-lq data and seurat-hq data, intermediate stages remain hidden,
451 which prohibits the detection of trajectories and results in a shuffling of marker
452 gene expression (Figures S33 and S34). Taken together these results highlight
453 how expression reconstruction using DISCERN improves downstream analyses
454 and yields deeper biological insights into cell type and state identification, gene
455 regulation, and developmental trajectories of cells.

456 *2.4. Discovering COVID-19 disease-relevant cells in lung and blood*

457 The previous sections have demonstrated DISCERN's utility to reconstruct
458 single cell expression data based on an hq reference, vastly improving the detec-
459 tion of cell (sub-) types and their signaling. Given these advantages, we won-
460 dered if DISCERN's expression reconstruction could deepen our understanding
461 of cell type-composition and signaling changes of immune cells in COVID-19
462 disease (Figure S35), using two published datasets [43, 22]. To obtain best re-
463 construction results, we again resorted to using bulk-hq immune reference data
464 (Table S1) [44], as outlined in the previous section.

465 First, we used a COVID-19 blood dataset (covid-blood-lq) with limited cell
466 type resolution, which was originally analyzed by our group using Seurat (Ta-
467 ble S1) [22]. While $CD4^+$, $CD8^+$, and NK cells formed separate clusters we
468 were unable to visibly distinguish subpopulations of these cells in covid-blood-
469 lq data [22]. Reconstruction of gene expression using bulk-hq data led to the
470 identification of 24 subtypes of $CD4^+$ and $CD8^+$ T cells in covid-blood-hq data
471 (Figure S36). Several cell clusters identified in covid-blood-hq data showed the
472 correct cell type-specific marker gene expression in covid-blood-lq data, albeit
473 in fewer cells, reduced in magnitude, and in some cases less specific (Figures S37
474 and S38). Reconstruction also led to the identification of $CD4^+$ TH17 helper
475 cells that express *RORC* (Figure 4A & B, Figure S39). Based on the molecular
476 footprint of these TH17 cells they were further subdivided into TH17_cluster1
477 that exhibits a memory T cell phenotype with elevated *IL7R* expression and
478 TH17_cluster2 that exhibits an activated T cell phenotype with elevated *MHC-*
479 *II*, *CCR4* and *RBPJ* expression (Figure 4B, Figure S39). The expression of
480 *RBPJ* is of particular interest, as it is linked to TH17 cell pathogenicity, sug-
481 gesting a role of pathogenic TH17 cells in COVID-19 [45]. It is common practice
482 to stimulate memory T cells in vitro to trigger IL-17A production and a shift
483 towards a TH17 phenotype was previously described in COVID-19 [46]. With
484 DISCERN we are able to distinguish these cells in COVID-19 patient blood
485 without stimulation, identifying cytokine producing memory cells with a TH17-
486 like phenotype (Figure S39).

487 To further validate the existence of activated TH17 cells in COVID-19 pa-
488 tient blood, we next analyzed the corresponding lung data (covid-lung) of the
489 patients for shared T cell receptor clones (Figure S40). The underlying assump-
490 tion is that cells with the same T cell receptor in lung and blood originate
491 from the same progenitor and therefore have a high probability of belonging

492 to the same cell type. For this comparison we used the cell type annotation
493 and representation of our original analysis of the covid-lung data, in which
494 memory T and TH17 cells were readily observed without reconstruction [22].
495 TH17_cluster1 cells showed strong clonal overlap with covid-lung CD4⁺ memory
496 T cells (Figure S40) and expressed comparable levels of *RORC* to covid-lung
497 effector memory TH17 cells (Figure S41), indicating that these CD4⁺ central
498 memory T cells could be TH17 (-like) cells. TH17_cluster2 in blood exhibited
499 strong clonal overlap with effector memory and resident memory TH17 cells
500 in covid-lung data (Figure S40) that express *RORC* and *IL-17A* (Figure S41).
501 Using the clonotype information of resident memory cells producing *IL-17A* in
502 inflamed lung (TRM17), we further corroborated the existence of the newly
503 identified population of IL-17A-producing TH17 cells in reconstructed COVID-
504 19 blood data (Figure S40). In general, the T cell receptor clonal information in
505 blood and lung therefore corroborated our cell type annotation in covid-blood-
506 hq data.

507 To understand the role of T cell subtypes in COVID-19 disease progression
508 we analyzed a second blood single cell dataset (covid-blood-severity-lq) contain-
509 ing disease-severity information for 130 COVID-19 patients [43]. To obtain opti-
510 mal cell type resolution, we combined the covid-blood-severity-lq T cell data[43]
511 with CD3⁺ covid-blood-lq cells [22] and reconstructed gene expression for the
512 combined dataset using bulk T cell sequencing reference data[44], resulting in
513 covid-blood-severity-hq data. Many of the 15 CD4⁺ T cell clusters identified in
514 covid-blood-severity-hq data (Figure S42) were also present in the covid-blood-
515 hq data, further validating the consistency of our cell type identification. This is
516 also corroborated by the available surface protein data for covid-blood-severity
517 data, substantiating that naive cells are CD45RA, memory cells are CD45RO,
518 and effector cell types are CD45RO positive (further details in Figure S43). We
519 compared the clusters that we identified in the covid-blood-hq with clusters iden-
520 tified in the covid-blood-severity-hq data and found confined and overlapping
521 regions of TFH, TH17_cluster1, and TH17_cluster2 cells (Figure S44). We also
522 compared the identified clusters to clusters defined in the original publication
523 (Figure S45). Cells identified as TFH in the original publication show signif-
524 icant overlap with naive CD4⁺ T cells (defined on transcriptome and protein
525 level) and CD4⁺ IL22⁺ cells (CD4.IL22) show marked overlap with TREG cells.
526 These results confirm once more the precise and robust cell type identification
527 that can be achieved with DISCERN.

528 Interestingly, we also identified two rather unexpected cell types after re-
529 construction. One cluster is positive for *CD4* and negative for *CD8A* while
530 otherwise expressing a signature of CD8⁺ effector memory cells with high ex-
531 pression of *GZMB*, *GZMH* and *PRF1* (Figure 4D & 4E). This signature points
532 to a CD4⁺ cytotoxic phenotype and indeed virus-reactive CD4⁺ cytotoxic cells
533 were described to be increased in blood during COVID-19 [47]. The other cell
534 type expresses *CD8*, *IL6R*, and *GATA3*, while being negative for *SLAMF7* (Fig-
535 ure 4D & 4E). These cells were described in the literature to be CD8⁺ T helper
536 cells [48], exert T helper function, and have been shown to lack cytotoxicity.
537 They lack expression of a significant number of cytokines and key transcription

538 factors pointing to a TH17 or TH22 phenotype. On a protein level these cells
539 express *CCR4*, while being negative for CCR6, making them cytolytic CD8⁺ T
540 helper type 2 cells (Tc2) cells. Part of this cluster overlaps with CD4 single-
541 positive cells and might explain why T helper type 2 cells are missing in the
542 CD4 cell clustering.

543 Overall, the highly specific and sensitive cell type identification in covid-
544 blood-severity-hq data enabled us to correlate the five COVID-19 disease sever-
545 ity categories to shifts in cell type and activity information. We first validated
546 the decrease in TFH cells with increasing disease severity, as described in the
547 original work (Figure S46) [43]. TH17 cells have been extensively studied using
548 flow cytometry and in accordance with our results MHC-II positive as well as
549 *CCR4* positive cells were described in COVID-19 patients (Figure 4B) [46]. We
550 observed a strong decrease in naive T helper cells in severe disease, most pro-
551 nounced for naive TREGs, while the fraction of TH17 cells showed little correla-
552 tion with disease severity (Figure S46). Of the two mixed cell types we detected
553 in COVID-19 data, cytotoxic CD4⁺ cells were increased in moderate and severe
554 disease (Figure S47). A similar increase is visible in patients with severe respi-
555 ratory disease without COVID-19 (Figure S48) and these cells might therefore
556 be a general marker of severe respiratory illness. Cytolytic CD8⁺ Tc2 cells are
557 increased in patients with severe symptoms and in those who died from COVID-
558 19 (Figure S47) and are described to be reduced after recovery from COVID-19
559 [49]. This positive correlation and the known role of Tc2 cells in fibroblast
560 proliferation induction and tissue remodeling could pinpoint a mechanistic role
561 of these cells in lung fibrosis as witnessed in severe COVID-19 patients. The
562 possibility to observe these cells in reconstructed single cell data may pave the
563 way to study the functional role of these cells in adverse COVID-19 outcome.

564 The relatively strong correlation of some cell types with COVID-19 out-
565 come suggests that blood cell fraction information might be used for patient
566 severity prediction. We trained a Gradient Boosting Machine (GBM) using
567 leave-one-out-cross-validation (LOOCV) on the fractions of all T cell types and
568 performed a forward feature elimination, to obtain a sparse, optimal model for
569 patient blood-based severity prediction. We first classified patients into three
570 groups, mild (union of asymptomatic and mild, $n = 26$), moderate ($n = 26$),
571 and severe (union of severe and critical, $n = 19$), reaching an AUROC of 0.63
572 (Table S4). We noticed that the mild and moderate groups were indistinguish-
573 able for the classifier (Figure S49). Training a GBM classifier on mild and severe
574 cases substantially increased classification performance, reaching an AUROC of
575 0.81 and accuracy, and F1 score of 0.82 (Table S4, Figure 4F & G). Compared
576 to the original T cell types and fractions reported (accuracy 0.61) [43], DIS-
577 CERN reconstructed T cell fractions are 33 % more accurate in the prediction
578 of COVID-19 disease severity (Figure 4G, Table S4). This classification improve-
579 ment is remarkable, given that DISCERN has no notion of disease severity when
580 it reconstructs gene expression. These results further demonstrate DISCERN's
581 precise and robust expression reconstruction that enabled the discovery of a
582 potential new blood-based biomarker for COVID-19 severity prediction.

583 3. Discussion

584 The sparsity of gene expression information and high technical noise in sin-
585 gle cell sequencing technologies limits the resolution of cell clustering, cell type
586 identification, and many other analyses. Several algorithms such as scImpute,
587 MAGIC, and DCA have addressed this problem by imputing missing gene ex-
588 pression in single cell data by borrowing expression information from similar
589 cells within the same dataset. While gene imputation clearly improves gene
590 expression by inferring values for dropped out genes, this imputation relies on
591 the comparison of similar cells with largely absent gene expression information
592 in the same dataset. With DISCERN we take a completely novel approach
593 to gene expression inference of single cell data, by realistic reconstruction of
594 missing gene expression in scRNA-seq data using a related dataset with more
595 complete gene expression information. We thus propose to call this procedure
596 ‘expression reconstruction’ to highlight the fundamental difference to classical
597 imputation and refer to the dataset with missing gene expression information
598 as low quality (lq) and the reference dataset as high-quality (hq).

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

616 It is important to note that DISCERN is a **precise** network that models
617 gene expression values realistically while retaining prior and vital biological in-
618 formation of the lq dataset after reconstruction. The network is also **robust**
619 to the presence of different cell types in hq and lq data, or an imbalance in
620 their relative ratios, and is robust to ‘hallucinating’ hq-specific cells into the lq
621 data. Several algorithmic choices are the foundation of DISCERN’s precision
622 and robustness. The network was designed to model the sequencing-technology-
623 specific and the underlying biological signals in separate components of its ar-
624 chitecture. Disentanglement of those two components is necessary to accurately
625 reconstruct expression information in the case where lq and hq datasets have
626 different content, i.e. cell type compositions. If the component designed to
627 model the effect of sequencing technology also captures the difference in the

628 biological signal, the reconstruction will lead to a lack of integration across the
629 two datasets where some cell types are still clustered by dataset (similar to
630 scGen in Figure S18). On the contrary, if the component modeling the biolog-
631 ical signal captures sequencing-technology-specific features, the reconstruction
632 will lead to an over-integration of the datasets where cells of different types are
633 mixed together (similar to Seurat in Figure S18). The demonstrated ability of
634 DISCERN to avoid those shortcomings, even in scenarios where there is very
635 little to no overlap between cell types across datasets, lies in the carefully crafted
636 balance between the expressivity of its components. The representational capa-
637 bilities of DISCERN, achieved via batch normalization, five loss terms, and a
638 dual head decoder, would reduce DISCERN's usability, if they would require fre-
639 quent dataset-specific tuning. The stability and usability was therefore a central
640 concern in the design and evaluation phase of DISCERN, which resulted in an
641 algorithm that gave very good results with a single set of default (hyper-) param-
642 eters. All comparisons to other algorithms, for instance, were performed with
643 default settings. Only the expression reconstruction of the exceptionally large
644 COVID-19 datasets required the fine-tuning of the learning rate, cross entropy
645 term, sigma, and the MMD penalty term. Another important technical feature
646 of DISCERN is that it can easily be integrated into existing workflows. It takes
647 a normalized count matrix, as created by nearly all existing single cell analysis
648 workflows, as input and produces a reconstructed expression matrix. This can
649 be used for most downstream applications (i.e. cell clustering, cell type identifi-
650 cation, cell trajectory analysis, and differential gene expression). DISCERN can
651 be trained on standard processors (CPU) for small and medium-sized datasets
652 and requires graphical processing units (GPU) for the expression reconstruction
653 of large datasets. Altogether, the usability and robustness of DISCERN should
654 enable even non-expert users to perform gene expression reconstruction.

655 A unique feature of DISCERN is the use of an hq reference to infer biolog-
656 ically meaningful gene expression. While we consider this a main strength
657 of DISCERN, the dependence on a suitable reference dataset might also limit
658 its application. We took great care in this manuscript to mitigate this con-
659 cern by showing how DISCERN is able to reconstruct gene expression for many
660 different types of lq and hq pairs, ranging from indrop - smartseq2 to single
661 nucleus - single cell data pairs. Remarkable in this context is DISCERN's ro-
662 bustness to differences between the cell type compositions of lq and hq data
663 pairs, with DISCERN being the only algorithm obtaining robust expression re-
664 construction when few cell types overlap. We have also shown that purified
665 bulk RNA-seq samples can be used as hq reference, as successfully applied to
666 PBMC and COVID-19 datasets in this study. We used 9852 FACS purified
667 immune cell bulk sequencing samples [44], comprising 27 cell types, to success-
668 fully reconstruct single cell expression data. This implies that most single cell
669 studies involving immune cells (with or without other cell types present) can be
670 reconstructed with DISCERN using a single published bulk RNA-seq dataset.
671 Furthermore, public RNA-seq repositories such as NCBI GEO contain tens of
672 thousands of samples of immune and non-immune cells that could serve as refer-
673 ence for most expression reconstruction experiments. Conversely, pure cell type

674 or subtype bulk RNA-seq data could be hard to obtain as the sorting of cells
675 might have limited resolution or might be partially impure. In consequence,
676 the usage of bulk RNA-seq data as reference for expression reconstruction could
677 lead to a grouping or averaging of cell subtypes. While these potential caveats
678 might adversely affect expression reconstruction, we have not observed merging
679 or averaging effects of single cell subtypes when corresponding bulk RNA-seq
680 cell type information was not present or present at different proportions (Fig-
681 ure 3B & 3C, Figure S28). Importantly, cells do not necessarily cluster into
682 distinct classes but can build cell continua, as shown in the trajectory analy-
683 sis in Figure 3B & 3C, where T cells seem to differentiate into each other and
684 do not form clearly separable clusters. In general, handling continua of cell
685 types is challenging for imputation and batch correction algorithms, as many of
686 them, including for instance scGEN, Bfimpute, SIMPLEs, and cscGAN, require
687 or recommend cluster or cell type annotation. This might lead to under- or
688 over-integration of cell continua. DISCERN does not rely on cluster (or cell
689 type) information and seamlessly integrates and reconstructs cell clusters and
690 continua (Figure 3C, Figure S36). In conclusion, we provide strong evidence that
691 DISCERN is widely and easily applicable to many single cell experiments.

692 While DISCERN gave good reconstruction results using default parameters
693 for most datasets we analyzed, we would like to highlight that the immense
694 representational power of generative neural networks can remove or hallucinate
695 biological information if not properly handled [6]. This is true for data inte-
696 gration [50] as well as for expression reconstruction algorithms and we would
697 highlight two guiding principles for optimal results. For non-expert users, we
698 would recommend the use of default settings and a careful selection of a re-
699 lated hq dataset. When datasets are large and complex, with many cell types
700 in the lq and several non-overlapping cell types in the hq data, one should al-
701 ways ensure that training does not merge or mix non-overlapping cell types with
702 other cells, by investigating that these cells keep their cell type-specific marker
703 gene expression. Keeping these ‘checks and balances’ will usually result in good
704 reconstruction results even for complex datasets such as covid-blood-severity.

705 To obtain novel insights into COVID-19 disease mechanisms and a new blood-
706 based biomarker for disease severity we reconstructed two published datasets
707 with DISCERN, Hamburg COVID-19 patients (covid-lung, -blood) and the
708 COVID-19 cell atlas (covid-blood-severity). The application of DISCERN to
709 the covid-blood dataset (COVID-19 patient blood) enabled us to detect 24 dif-
710 ferent immune cell types and activity states, which is quite remarkable given
711 that we find these cells in blood. Two TH17 subtypes caught our attention, as
712 they share the TCR clonality with the lung data from the same patients (covid-
713 lung), suggesting bloodstream re-entry of lung TH17 cells. We linked these two
714 subclusters to their functional role by separating them into a memory-like and
715 activated-like phenotype. The clonal overlap of activated TH17 cells in blood
716 with previously discovered lung-resident cells suggests that activated TH17 cells
717 in blood are resident T cells from the lung reentering circulation. These cells
718 might in part explain the multi-organ pathology observed in COVID-19, as
719 activated T cells might travel via the blood to secondary organs and cause in-

720 inflammation and tissue damage. Future work might demonstrate the effect of
721 these activated T cells on tissue inflammation.

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

739 Interestingly, we also discovered two atypical T cell types in reconstructed
740 COVID-19 patient blood single cell data. While cytotoxic CD4⁺ T cells have
741 been observed in COVID-19, we can show that this increase is not COVID-19
742 specific and is also observed in other types of pneumonia. Interestingly, we also
743 detected cytolytic CD8⁺ Tc2 cells that express *CD8A*, *GATA3*, *IL6R* and are
744 negative for *SLAMF6*. This cell type is linked to tissue fibrosis and steroid
745 refractory disease in asthma [51]. The increase in CD8⁺ Tc2 cells that we ob-
746 serve specifically in COVID-related death could be associated with COVID-19
747 patients that do not respond to steroids. Demonstration of increase of this cell
748 type in patients dying of COVID-19 points to a potential therapeutic inter-
749 vention with the drug Fevipiprant, which blocks CD8⁺ Tc2 cell activation and
750 its pro-fibrotic effects by inhibiting prostaglandin D2 signaling [52]. Functional
751 analysis of these cells has to demonstrate whether these cells are an early marker
752 of later death or whether it is a marker of already escalated treatment.

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

758 Acknowledgements

759 We thank Immo Prinz, Manuel Friese, Johannes Soeding, Robert Zinzen, Yu
760 Zhao and Stefan Kurtz for their helpful comments and suggestions. We thank
761 Rajasree Menon and Matthias Kretzler for providing kidney single cell and sin-
762 gle nuclear RNA-seq data and their support for corresponding analysis. FH,

763 RK, MM, PM, & SB were supported by the LFF-FV 78, EU ERare-3 Maxo-
764 mod, SFB 1286 Z2, FOR 296, and FOR 5068 research grants. Further support
765 was obtained from the UKE R3 reduction of animal testing grant. CE was sup-
766 ported by DFG ER 981/1-1 and the clinician scientist programme of university
767 Hamburg, NG was supported by ERC StG-715271, and SHu was supported by
768 ERC CoG-865466 and has an endowed Heisenberg-Professorship awarded by
769 the Deutsche Forschungsgemeinschaft. SHa was funded by the BMBF STOP-
770 FSGS-01GM1518C and SFB 1192 B08 research grants.

771 **Competing interests**

772 The authors declare no competing interests.

773 **Author contributions**

774 SB initiated and SB, PM, FH, and CE conceptualized the study with help
775 from MM. FH and CE implemented DISCERN, MM refactored the code, and
776 PM reviewed the DISCERN implementation. FH, CE, and RK performed the
777 analyses. SB, PM, NG, and SHu supervised the study. SB, FH, and CE wrote
778 the manuscript. SHu, PM, NG, RK and SHa provided ideas, contributed to the
779 manuscript text and critically reviewed the manuscript. All authors read and
780 approved the final manuscript.

781 Main figures

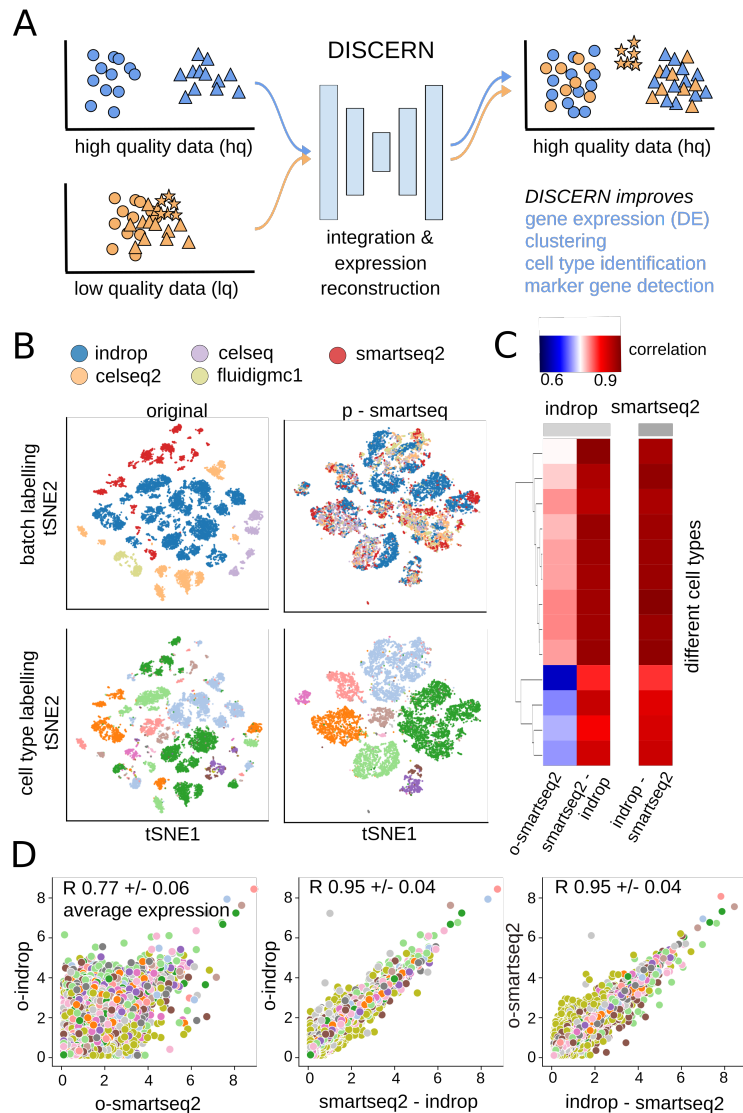


Figure 1: *Integration and expression reconstruction of single cell sequencing data.* **A:** DISCERN transfers the style of a high-quality (hq) dataset to a related low quality (lq) dataset, enabling gene expression reconstruction that results in improved clustering, cell type identification, marker gene detection, and mechanistic insights into cell function. The hq and lq datasets have to be related but not identical, containing for example several overlapping cell types but also exclusive cell types of cell activity states for one or the other dataset. **B:** t-SNE visualization of the pancreas dataset before reconstruction (original) and after transferring the style of the smartseq2 dataset using DISCERN (p-smartseq2). The upper row shows the dataset of origin before and after projection colored by batch and the lower row colored by cell type annotation (details of 13 cell types in supplements). **C** and **D:** Average gene expression (over all the cells of a given type) of the pancreas indrop and smartseq2 datasets before (first column and panel) and after smartseq2 to indrop (second column and panel), and after indrop to smartseq2 projection (third column and panel). **C:** Gene correlation by cell type shown in colored heatmap. **D:** Each colored point represents a single gene colored by the cell type, 'o' refers to original data. The mean Pearson correlation with one standard deviation over all cell types is shown in the figure title.

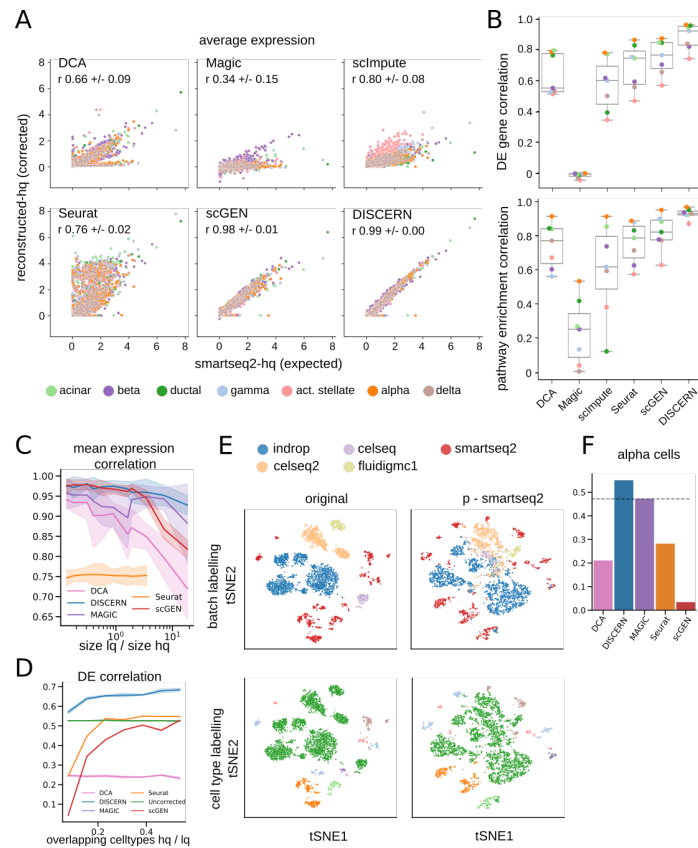


Figure 2: *Expression reconstruction benchmark of DISCERN and five state-of-the-art batch correction and imputation algorithms.* **A:** Comparison of the expression reconstruction performance of Seurat, scGEN, Magic, scImpute, DCA, and DISCERN using smartseq2 data. The smartseq2 data was split into a 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 smartseq2-hq data for different ratios of lq to hq data. The standard deviation indicates the deviation in correlation of the cell types. The datasets were created as described in **A**. **D:** Alpha cells were removed from the smartseq2-hq batch and left in the smartseq2-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.

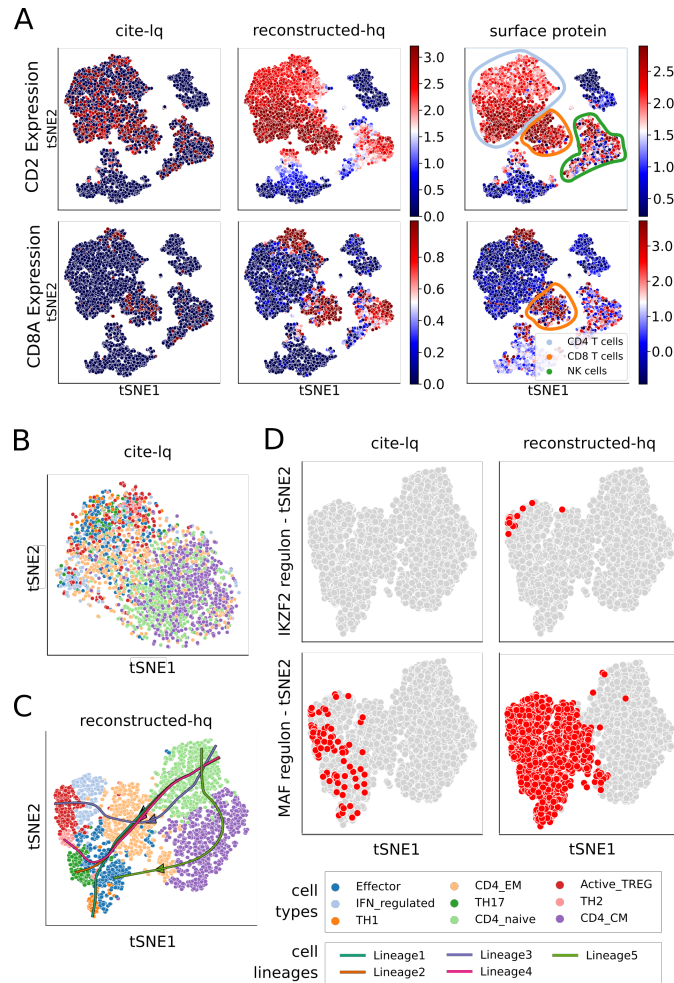


Figure 3: *Expression reconstruction improves downstream analyses including cell identification, gene regulation, and trajectory inference.* The cite-lq dataset was reconstructed using bulk-hq data and compared to ground truth CITE-seq (surface protein) information. The CITE-seq information was not used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

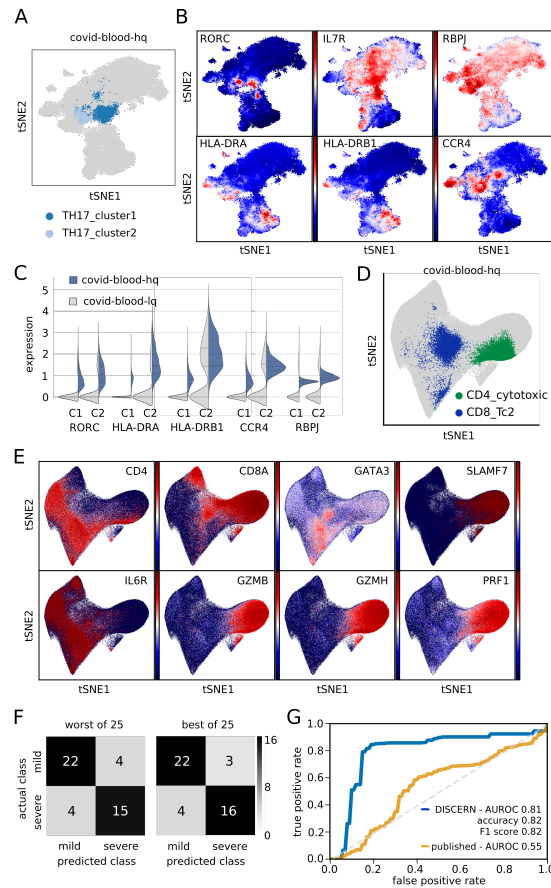


Figure 4: *Expression reconstruction improves COVID-19 cell type identification and allows for efficient disease severity prediction.* Two COVID-19 blood datasets were reconstructed and analyzed. Hamburg covid-blood-lq and covid-lung-lq data was reconstructed using bulk-hq data, resulting in the respective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

782

783

784 4. Methods

785 4.1. Data availability

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

790 4.2. Dataset description

791 *Pancreas.* The pancreas dataset is a collection of different scRNA-seq datasets,
792 profiling pancreas cells in the context of diabetes [53]. The pancreas dataset is
793 a widely used dataset for batch correction benchmark experiments and due to
794 its high number of cell types and sequencing technologies it allows to evaluate
795 differences between cells and sequencing technologies at the same time. The ex-
796 pression table, including the annotation, is available from SeuratData ([https://](https://github.com/satijalab/seurat-data)
797 github.com/satijalab/seurat-data) as `panc8.SeuratData` (v3.0.2) [53]. The
798 dataset was sequenced using five sequencing technologies (Smart-Seq2, Flu-
799 idigm C1, CelSeq, CEL-Seq2, inDrop) and consists of 13 cell types (alpha, beta
800 ,ductal, acinar, delta, gamma, activated_stellate, endothelial, quiescent_stellate,
801 macrophage, mast, epsilon, schwann). In total, before preprocessing, the dataset
802 contains 14 890 cells.

803 *difftec.* The difftec dataset was created for a systematic comparative analysis
804 of scRNA-seq methods [54]. Similar to pancreas, the difftec dataset is ideal
805 for the evaluation of expression reconstruction across many cell types and se-
806 quencing technologies. Seven sequencing technologies (10x Chromium v2, 10x
807 Chromium v3, Smart-Seq2, Seq-Well, inDrop, Drop-seq, CEL-Seq2) were used
808 with at least two replicates each. In this dataset 10 different cell types (Cy-
809 tototoxic T cell, CD4⁺ T cell, CD14⁺ monocyte, B cell, Natural killer cell,
810 Megakaryocyte, CD16⁺ monocyte, Dendritic cell, Plasmacytoid dendritic cell,
811 Unassigned) were annotated, and make up for 31 021 cells in total before filter-
812 ing. The expression table including the annotation is available from SeuratData
813 as `pbmcsca.SeuratData` (v3.0.0).

814 *snRNA & scRNA.* The dataset was created for the validation of a single cell
815 and single nuclei analysis toolbox [28]. Since snRNA-seq and scRNA-seq data
816 varies in the amount of counts per cell and the genes detected, we tested if
817 DISCERN could reconstruct snRNA-seq expression so that it would closely
818 resemble scRNA-seq expression, providing a biological ground-truth. While we
819 label snRNA-seq data as `lq` and scRNA-seq as `hq`, this distinction is incorrect
820 from a biological perspective, as gene expression should be in part different
821 between the nucleus and the cytosol. The dataset consists of a liver biopsy

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

827 *covid-lung & covid-blood.* The COVID-19 dataset we have previously published
828 consists of blood and bronchoalveolar lavage (BAL) samples from four patients
829 with bacterial pneumonia and eight patients with SARS-CoV-2 infection[22].
830 In total 155 706 cells were sequenced using TCR-seq technology, which allows
831 for the comparison of clonal expansion in both tissues. While we investigated
832 the lung data in detail in the original publication, the analysis of the blood was
833 largely limited to cell type identification. Using DISCERN, we use the blood
834 data to find previously unobserved cell types, link them to cell clones found in
835 the lung, and derive a biomarker based on cell fractions (see also covid-blood-
836 severity data). Cell type annotations for the BAL samples were used as in the
837 original publication.

838 *citeseq.* This dataset contains CITE-seq information of healthy human PBMCs
839 for 6 cell types (B cells, CD4 T cells, NK cells, CD14⁺ Monocytes, FCGR3A⁺
840 Monocytes, CD8 T cells) [29]. In our analyses we used the cell type information
841 provided in the original publication [55]. The CITE-seq data is ideal to bench-
842 mark DISCERN, as the information of 13 surface proteins offers ground-truth
843 information on the cell types and a good proxy for the expression of the 13
844 corresponding genes.

845 *bulk.* We used this large dataset of 28 FACS sorted and bulk sequenced immune
846 cell types as ‘ultimate’ hq reference data for lq immune single cell sequencing
847 data. Each of the 9852 samples provides an average expression information for
848 13 104 genes for a specific immune cell type, providing a hq reference for e.g. lq
849 single cell PBMC CITE-seq data with only 798 expressed genes per cell. We
850 further assume that this dataset is large enough to provide enough per cell type
851 variability for our deep neural network to faithfully learn and represent its gene
852 expression. In more detail, the dataset consists of 28 sorted immune cell types
853 (Naive CD4, Memory CD4, TH1, TH2, TH17, Tfh, Fr. I nTreg, Fr. II eTreg,
854 Fr. III T, Naive CD8, Memory CD8, CM CD8, EM CD8, TEMRA CD8, NK,
855 Naive B, USM B, SM B, Plasmablast, DN B, CL Monocytes, Int Monocytes,
856 NC Monocytes, mDC, pDC, Neutrophils, LDG) with \geq 99% purity [44]. Total
857 RNA was extracted using RNeasy Micro Kits (QIAGEN). Libraries for RNA-seq
858 were prepared using SMART-seq v4 Ultra Low Input RNA Kit (Takara Bio).
859 In total, the dataset contains 9852 samples collected in two phases from 416
860 donors, out of which 79 are healthy. For training DISCERN, bulk TPM counts
861 and all cell types were used if not stated otherwise.

862 *covid-blood-severity.* This dataset is an aggregation of three COVID-19 sequenc-
863 ing studies using the 10X Genomics Chromium Single Cell 5’ v1.1 technology.

864 It contains a large number of cell types with fine-grained cell type annotations
865 that are complemented with information on COVID-19 disease severity for each
866 patient sequenced. We used this dataset to obtain a blood-based biomarker of
867 COVID-19 disease severity, based on T cell fractions observed with DISCERN.
868 The data consists of PBMCs from 29 healthy, 89 COVID-19 and 12 LPS-treated
869 patients. The authors detected 51 cell types in their original work (see Ta-
870 ble S1) [43] and COVID-19 patients were classified by their disease severity
871 (worst clinical outcome) into ‘asymptomatic’, ‘mild’, ‘moderate’, ‘severe’, ‘crit-
872 ical’, and ‘death’. Count data together with CITE-seq information was used
873 as provided in the original publication ([https://covid19.cog.sanger.ac.uk/
874 submissions/release1/haniffa21.processed.h5ad](https://covid19.cog.sanger.ac.uk/submissions/release1/haniffa21.processed.h5ad)).

875 *kidney-lq (snRNA-seq) & kidney-hq (scRNA-seq)*. The kidney dataset consists
876 of single cell RNA-seq and single nuclei RNA-seq data of 9 patients with acute
877 kidney injury sequenced using 10X Genomics Chromium technology. It contains
878 in total 82 701 cells with 52 934 cells sequenced using snRNA-seq and 29 767 cells
879 sequenced using scRNA-seq. The dataset does not contain cell type annotation,
880 but in initial analysis using a different subset [56] suggested that identification
881 of T cells in the snRNA-seq data is challenging. For this reason, the analysis
882 was focused on the detection of T cells and their subtypes.

883 4.3. Code availability

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

888 4.4. Preprocessing

889 Raw expression data (Counts) preprocessing was performed as previously
890 described [57] using the scanpy (v1.6.1, [58]) implementation. In particular,
891 the intersection of genes between batches was used. The cells were filtered
892 to a minimum of 10 genes per cell and a minimum of 3 cells per gene. Li-
893 brary size normalization was performed to a value of 20 000 with subsequent
894 log-transformation. As model input for DISCERN the genes were scaled to
895 zero mean and unit variance. However, for all further evaluation the genes
896 were scaled to their uncorrected mean and variance not considering the batch
897 information.

898 4.5. Description of DISCERN

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

902 *Wasserstein Autoencoder.* While neural network-based autoencoders have been
903 widely used for decades for dimensionality reduction [59, 60], recent advances
904 have also allowed their use to build a generative model of the distribution of
905 the data at hand[61]. More recently, leveraging results from optimal transport
906 [62], Wasserstein Generative Adversarial Networks (WGAN) [63] and Wasser-
907 stein Autoencoders (WAE) [23] have been designed to explicitly minimize the
908 (Wasserstein, or earth-mover) distance between the distribution of the input
909 data and their reconstruction. WGANs only implicitly encode their input into
910 a latent representation (called latent code), while WAE has the useful property
911 of using an explicit encoder, which makes it possible for the model to directly
912 manipulate the different representations of single-cell data. Finally, the WAE
913 framework, established in [23], allows the use of a wide range of architecture and
914 losses, which we are going to detail now. First of all, in order to effectively use a
915 number of latent dimensions that adaptively matches the intrinsic dimension of
916 the scRNA-seq data at hand, DISCERN uses a random encoder as prescribed
917 in [64].

918 *Architecture.* Autoencoders widely used for transcriptomics applications are
919 shown to perform well on several tasks, like drug perturbation prediction [21]
920 or dropout imputation [12]. Since the ordering of the genes in scRNA-seq
921 matrices is mostly arbitrary, fully-connected layers are usually used in this task.
922 In our case, DISCERN consists of three fully connected layers in the encoder
923 and the decoder. The bottleneck of the autoencoder (or latent space) contains
924 48 neurons, which is sufficient to accurately model all the datasets we used in
925 our experiments. Additionally, we exploit a finding from [64] to let the net-
926 work learn the appropriate amount of latent dimensions. While the encoder
927 will be tasked to transform the distribution of the input data into a fixed,
928 low-dimensional prior distribution (i.e. a standard Gaussian), the decoder will
929 perform the opposite, i.e. transforming the fixed, low-dimensional prior distri-
930 bution into gene space. scRNA-seq data is known to display a high level of zero
931 measurements, called dropout, which is essential to accurately model the count
932 distribution. To describe scRNA-seq data in a parametric way, it is common to
933 model the expression level of a gene with zero-inflated negative binomial distri-
934 bution [65]. Despite the several non-linearities in the decoder architecture,
935 it is, however, difficult to learn an encoding function that maps a simple prior
936 to the distribution leading to low quality modeling of low expressed genes. To
937 address this issue, we scale the gene expression and attach a second head to the
938 decoder (i.e. a second decoder sharing all weights with the first, except for the
939 last layer). The task of the second decoder head is to predict, for each gene
940 of a cell, the probability of its expression to be dropped out, giving rise to a
941 random decoder. Thus, this second decoder head predicts dropout probabili-
942 ties and models the dropout probabilities for different batches. This additional
943 head allows modeling the dropout and the expression independently, to capture
944 the specific distribution of single cell data without the need for further explicit
945 assumption about the distribution.

946 *Loss function.* The loss optimized during the training of DISCERN is composed of four terms: a data-fitting (or reconstruction) loss, a dropout fitting (cross entropy) loss, a prior-fitting term (ensuring that DISCERN approximately minimizes the Wasserstein distance) and a variance penalty term (that controls the randomness of the encoder). Thus, DISCERN can be considered as a Wasserstein Autoencoder as introduced in [23]. For the reconstruction term, the framework introduced in [23] allows the use of any positive cost function. We elected to use the Huber loss [66] as it is well suited for modeling scaled scRNA-seq expression data, because it allows to select a threshold value to give lower weight to high differences in highly expressed genes and thus allows the model to learn a more robust expression estimate without focusing too much on outlier values. For the prior-fitting term, following [23], DISCERN uses the Maximum Mean Discrepancy (MMD) [67] between the aggregate posterior (i.e. the distribution of the input single-cells after encoding) and a standard Gaussian. We use the sum over an inverse multiquadratic kernel with different sizes for this task. Then, to prevent the random encoder (with diagonal covariance) from collapsing to a deterministic one, a penalty term that enforces that some components of the variance are close to 1. Intuitively, that means that the superfluous latent dimensions will only contain random noise (see [64] for more details). Another loss term, namely the binary cross-entropy loss, on the second decoder head is used to enable the model to learn a dropout probability for each gene and sample. The loss on the dropout layer enables the model to capture the bimodal distribution of single cell data. Additionally, activity regularization is applied on the Conditional Layer Normalization (CLN), such that the weights of the conditional layers are only regularized in a batch-specific manner and the regularization is not applied for batches, which are not present in the current mini-batch. This has the advantage that the batch dependent weights are not influenced too much by different batch sizes. The four loss terms are added (and weighed) together to form the loss that DISCERN minimizes during training (see also Figure S1 for loss terms).

976 *Conditional Layer Normalization.* The weights of those fully-connected layers are shared for all the batches that DISCERN is trained on. However, to model the batch-specific differences, we use a Conditional Layer Normalization (CLN) that applies the idea proposed in [25] to Layer Normalization [26]. In essence, for each batch, different sets of shifting factors are learned. Note that in DISCERN, no scaling factors are used to limit the expressivity of the conditioning and therefore reduce the chance of over integration. This allows not only to accurately model the batch-specific differences between batches, but also to transfer the batch effect from one dataset onto another, in the spirit of the style-transfer approach developed in [25]. To make things clear, DISCERN does not explicitly train to integrate datasets. Instead, it trains to accurately model the input data, capturing the batch-specific differences with the weights of the CLN layers (i.e. conditioning), and the biological signal (which is mostly shared across the batches to integrate) with the weights of the fully-connected layers. After training, we encode all the cells we want to reconstruct, conditioning the process on

991 their batch of origin. Then, we take the batch chosen by the user and proceed to
992 decode all the cells conditioning on that specific batch, effectively transferring
993 the batch effect of one specific batch onto all of the batches we want to integrate
994 and reconstruct.

995 *Activations & dropout.* With the exception of the output layer, every other
996 fully-connected layer of the encoder and the decoder was followed by a CLN,
997 a Mish ([68]) activation function, and dropout during model training to reduce
998 overfitting.

999 *Optimization.* To optimize the weights of our model, DISCERN uses Rectified
1000 Adam ([69]), which addresses some of the shortcomings of the widely used Adam
1001 [70] and generally yields more stable training. To prevent overfitting, the op-
1002 timization is stopped early. It is implemented as a modification of the Keras
1003 EarlyStopping (with parameter minDelta set to 0.01 and the patience to 30)
1004 where the callback is delayed by a fixed number of 5 epochs. The delay was
1005 implemented to prevent too early stopping due to the optimization procedure.

1006 4.6. Hyperparameters

1007 As outlined in the architecture section of the methods and depicted in Fig-
1008 ure S1, DISCERN features several learnable hyperparameters. The complexity
1009 of the hyperparameter search space is a potential downside of DISCERN, if
1010 these hyperparameters would be unstable across different datasets or in other
1011 words, would require constant tuning. Fortunately, DISCERN’s hyperparame-
1012 ters are very stable across the multitude of datasets tested in this manuscript,
1013 which we will outline in this paragraph. Naturally, there is no rule without an
1014 exception, which in this manuscript are the COVID-19 datasets that required
1015 optimization for several hyperparameters.

1016 *Constant hyperparameters.* DISCERN features a number of hyper-parameters
1017 that can be tuned through hyperparameter optimization (see below for details).
1018 Most of them have default values that yield reasonable performance across the
1019 different datasets we used and are being kept constant across experiments, in-
1020 cluding the COVID-19 dataset. Those constant hyperparameters are: the choice
1021 of the reconstruction loss (Huber loss), activation functions (Mish), CLN for the
1022 conditioning, number of fully-connected layers (3) and their size (1024, 512, 256
1023 and 256, 512, 1024 neurons for the encoder and the decoder respectively), num-
1024 ber of latent dimensions (48), learning rate (1×10^{-3}), decay rates β_1 and β_2 of
1025 Rectified Adam (0.85 and 0.95 respectively), batch size (192), label smoothing
1026 for our custom cross entropy loss (0.1), dropout rates (0.4 in the encoder and 0
1027 in the decoder), delta parameter of the Huber loss (9.0), weight on the penalty
1028 on the randomness of the encoder λ_{sigma} (1×10^{-8}), weight on the cross entropy
1029 loss term $\lambda_{dropout}$ (1×10^5), weight on the MMD penalty term λ_{prior} (1500).

1030 *Dataset-specific hyperparameters.* The optimal value of the L2 regularization
1031 applied on the weights of our custom CLN highly depends on the dataset at hand
1032 and thus requires dataset-specific tuning. For datasets with a very small vari-
1033 ance in cell compositions the L2 CLN regularization can be turned off (weight
1034 set to 0). When datasets have different compositions the L2 CLN regularization
1035 requires higher values (typically between 1×10^{-3} and 0.2).

1036 *COVID-19 hyperparameters.* For the experiments with COVID-19 datasets slightly
1037 adjusted hyperparameters were used: learning rate of 6e-3, label smoothing for
1038 our custom crossentropy loss of 0.05, weight on the penalty on the randomness
1039 of the encoder λ_{sigma} (1e-4), weight on the cross entropy loss term $\lambda_{dropout}$
1040 (2e3), weight on the MMD penalty term λ_{prior} (2000).

1041 *Hyperparameter optimization.* DISCERN implements different techniques for
1042 hyperparameter optimization by using the ray[tune] library [71]. For most use
1043 cases the model does not require hyperparameter tuning and the default pa-
1044 rameter should be sufficient. However, DISCERN has a generic interface and
1045 supports nearly all techniques implemented in ray[tune]. The initial hyperpa-
1046 rameters were found using grid search. The loss used for the hyperparameter
1047 selection is the classification performance of a Random Forest classifier trying
1048 to classify real vs. auto-encoded cells. Classification performance was mea-
1049 sured using the area under the receiver operating characteristic curve (AUC /
1050 AUROC).

1051 4.7. Competing algorithms and methods

1052 We briefly discuss competing methods and have compared their performance
1053 to DISCERN in the results section. These algorithms can be grouped into two
1054 categories, i) imputation algorithms that were developed to estimate drop-out
1055 gene expression based on dataset inherent information (MAGIC, DCA, scIm-
1056 pute) and ii) algorithms designed for batch correction that we have modified
1057 or extended to reconstruct gene expression, although this is not their intended
1058 use (Seurat, scGen). Given the latter, it is clear that DISCERN could be used
1059 purely for batch correction in latent space, a subject beyond the scope of this
1060 manuscript.

1061 *MAGIC.* [13] Markov affinity-based graph imputation of cells (MAGIC) de-
1062 noises and imputes the single-cell count matrix using data diffusion-based in-
1063 formation sharing. The construction of a good similarity metric is challenging
1064 for finding biologically similar cells due to high sparsity. MAGIC finds a good
1065 similarity metric using a sophisticated graph-based approach that builds less-
1066 noisy cell-cell affinities and information sharing across cells. A particular focus
1067 of MAGIC was to understand gene-gene relationships and to characterize other
1068 dynamics in biological systems. MAGIC is provided as a Python package.

1069 *DCA*. [11] Deep count autoencoder (DCA) is a deep learning-based method for
1070 denoising single-cell count matrices. DCA is implemented in Python and uses
1071 an autoencoder with a Zero-Inflated Negative Binomial (ZINB) loss function.
1072 For each gene, DCA computes gene-specific parameters of ZINB distribution,
1073 namely dropout, dispersion and mean. By modeling gene distributions as a noise
1074 model and also computing dropout probabilities of each gene, DCA is able to
1075 denoise and impute the missing counts by identifying and correcting dropout
1076 events.

1077 *scImpute*. [12] Similarly to MAGIC, scImpute focuses on identifying cells that
1078 are similar, which is challenging due to the high sparsity of single-cell count
1079 matrices. scImpute is a statistical model using a three step process to impute
1080 scRNA-seq data. In the first step spectral clustering is applied on principal com-
1081 ponents to find neighbors, which later can be used to detect and impute dropout
1082 values. In the second step scImpute fits a mixture model of a Gamma distribu-
1083 tion and a Normal distribution to distinguish technical and biological dropouts.
1084 In the last step, the model uses a regression model for each cell to impute the
1085 expression of genes with high probability of dropout. With this approach, scIm-
1086 pute avoids hallucinations and keeps the gene expression distribution. scImpute
1087 is provided as an R package.

1088 *Seurat*. [24] Seurat is an open-source toolkit for the analysis of single cell
1089 RNA-sequencing data. In addition to general analysis functions, Seurat of-
1090 fers batch-correction functionality. Seurat uses canonical correlation analysis
1091 to construct this lower dimensional representation and tries to find neighbors
1092 between batches in this shared space. These anchors are filtered considering
1093 the local neighborhood of the cell pairs and remaining anchors are finally used
1094 to construct correction vectors for all cells in this low dimensional representa-
1095 tion. While Seurats is intended to work in a lower dimensional representation,
1096 it can also be used to reconstruct the expression information from this lower
1097 dimensional representation. Seurat is provided as an R package.

1098 *scGen*. [21] scGen is a variational autoencoder based deep learning method with
1099 a focus on learning features that help distinguish responding and non-responding
1100 genes and cells. scGen constructs a latent space in which it estimates perturba-
1101 tion vectors associated with a change between different conditions. Since scGen
1102 models the perturbation and infection responses in single cells, it is focused on
1103 in-silico screening with the use of cells coming from healthy samples. It can also
1104 be used for batch correction. For batch correction, and unlike DISCERN or
1105 Seurat, scGen uses both batch and cell type labels.

1106 *Multigrade*. [17] multigrade is an autoencoder based deep learning method de-
1107 veloped for the integration of different modalities to improve single cell RNA-seq
1108 downstream analysis, mainly clustering. The main focus is the integration of
1109 CITE-seq protein abundance since it is often available together with scRNA-
1110 seq. They use individual encoders for each modality and build a shared latent

1111 representation by partially sharing the decoder. Multigrade is built using the
1112 scvi-tools toolbox and implemented in python and pytorch.

1113 *4.8. Evaluation metrics*

1114 *t-SNE & UMAP.* For visualization of the datasets and to qualitatively assess
1115 the integration performance tSNE and UMAP were used. Both methods are
1116 based on PCA representation and use non-linear representations to create a 2D
1117 representation of the data. We used the scanpy [58] implementation. Default
1118 settings were used in nearly all cases except: In the combined COVID-19 dataset
1119 analogue to Kobak *et al.*[72] the dataset was subset to 25 000 cells and tSNE
1120 was computed using a perplexity of 250, and a learning rate of 25 000/12. These
1121 positions were taken and used as input to tSNE of all cells using a perplexity
1122 of 30 a learning rate of (number of observations)/12 and a late exaggeration of
1123 4.0 using FIt-SNE [73]. Clustering was performed using PARC [74] with de-
1124 fault parameters except `dist_std_local=1.5` and `small_pop=300`. Methods were
1125 changed here due to computation time issues for 350 000 cells. covid-blood data
1126 was analyzed using a learning rate of (number of observations)/6 a perplexity
1127 of (number of observations)/120 and `early_exaggeration=4`. Clustering was per-
1128 formed using default parameters except `knn=100` and `small_pop=100` to reduce
1129 the number of clusters with limited cell number. Clustering of the T helper cells
1130 in healthy blood was performed using coarse clustering with 30 nearest neigh-
1131 bors and leiden clustering (<https://github.com/vtraag/leidenalg>) with a
1132 resolution of 0.6. Afterwards a combined cluster of IFN-regulated and TREG
1133 was reclustered using a resolution of 0.4 and effector T cells were reclustered us-
1134 ing a resolution of 0.8. Resolution was chosen to dissect the raw gene expression
1135 changes of known cell types.

1136 *Mean gene expression.* Mean gene expression was calculated as average over
1137 log-normalized expression over all cells, usually stratified by celltype. This eval-
1138 uation of expression data consists of many data points where several have values
1139 close to zero, but could have a high weight on rank-based correlation methods.
1140 Thus Pearson correlation was used to evaluate the performance.

1141 *Differential gene expression.* Differential gene expression was performed using
1142 the scanpy [58] `rank_gene_groups` function using the t-test method for calculat-
1143 ing statistical significance on log-normalized expression data. Differential gene
1144 expression analysis was always performed under consideration of the cell type
1145 information. For comparison of differential gene expression analysis between
1146 conditions, the Pearson correlation was used. It is calculated either on the log2
1147 fold-change or in most cases on the t-statistics, computed during significance
1148 estimation. The data was compared using the t-statistics, because it aggregates
1149 information on both the variance and the change in mean expression. Thus it
1150 allows, roughly speaking, for simultaneously evaluating the significance and the
1151 log2 fold change.

1152 *Pathway analysis.* Pathway analysis or gene set enrichment analysis was done
1153 using the prerank function from gseapy [75] on the t-statistics, computed as
1154 described in the ‘Differential gene expression’ section of the methods. To this
1155 end, the gene set library “KEGG_2019_Human” provided by enrichr [76] was
1156 used. Top pathways were selected using the normalized enrichment score as
1157 previously described [75].

1158 *Gene regulation.* [38] The python implementation of the SCENIC (pySENIC)
1159 was used to infer regulons specific for CD4⁺ T helper cells. SCENIC infers a
1160 gene regulatory network using GRNBoost2 and creates co-expression modules.
1161 The co-expression modules get associated with transcription factors using the
1162 transcription factor motif discovery tool RcisTarget. A pair of transcription
1163 factor and associated gene set is called a regulon. For each cell, the regulons
1164 get scored using the AUCell algorithm to examine if a cell is affected by the
1165 regulon. We used default parameters of the pySENIC implementation.

1166 *COVID-19 classification*

1167 To evaluate the importance of the cell types found in the covid-blood-
1168 severity-hq dataset after reconstruction with DISCERN, the fraction for all
1169 T cell subtypes was used to predict the disease severity, as provided in [43].
1170 The data was classified using a Gradient boosting classifier ([77], implemented
1171 in scikit-learn v1.0.2, default settings) using 25 rounds of leave-one-out cross-
1172 validation (LOOCV). Each round consists of n training-prediction iterations
1173 with $n - 1$ samples for training and 1 sample for testing, such that after one
1174 round prediction results for all n samples could be evaluated. We chose LOOCV
1175 over k-fold cross-validation and testing due to the limited size of the dataset,
1176 consisting of only 71 patients. We used pycm ([78], v3.3) for the performance
1177 evaluation. The final evaluation was done using the accuracy and F1 score
1178 as provided by pycm. The area under the receiver operating characteristic
1179 (AUROC) curve is computed with scikit-learn. Before training the classifiers
1180 a forward feature selection was performed using the SequentialFeatureSelector
1181 implemented in scikit-learn with default parameters. In total four experiments
1182 were performed. In the first experiment, classification with three disease cat-
1183 egories (mild, moderate, severe) was used. Patients who died were excluded.
1184 For the other two experiments only patients with asymptomatic, mild, severe
1185 and critical symptoms were included. In all experiments the asymptomatic and
1186 mild category was merged to mild and severe and critical to severe.

1187 **References**

- 1188 [1] N. Editorial, Method of the year 2013, Nat. Methods 11 (1) (2014) 1.
- 1189 [2] Y. Zhao, U. Panzer, S. Bonn, C. F. Krebs, Single-cell biology to decode
1190 the immune cellular composition of kidney inflammation, Cell and tissue
1191 research 385 (2) (2021) 435–443.

- 1192 [3] M. Stoeckius, C. Hafemeister, W. Stephenson, B. Houck-Loomis, P. K.
1193 Chattopadhyay, H. Swerdlow, R. Satija, P. Smibert, Simultaneous epi-
1194 tope and transcriptome measurement in single cells, *Nature methods* 14 (9)
1195 (2017) 865–868.
- 1196 [4] A. A. Tu, T. M. Gierahn, B. Monian, D. M. Morgan, N. K. Mehta,
1197 B. Ruiters, W. G. Shreffler, A. K. Shalek, J. C. Love, Tcr sequencing paired
1198 with massively parallel 3' rna-seq reveals clonotypic t cell signatures, *Nature*
1199 *immunology* 20 (12) (2019) 1692–1699.
- 1200 [5] J. A. Pai, A. T. Satpathy, High-throughput and single-cell t cell receptor
1201 sequencing technologies, *Nature Methods* 18 (8) (2021) 881–892.
- 1202 [6] S. Oller-Moreno, K. Kloiber, P. Machart, S. Bonn, Algorithmic advances
1203 in machine learning for single-cell expression analysis, *Current Opinion in*
1204 *Systems Biology* 25 (2021) 27–33.
- 1205 [7] D. Lähnemann, J. Köster, E. Szczurek, D. J. McCarthy, S. C. Hicks, M. D.
1206 Robinson, C. A. Vallejos, K. R. Campbell, N. Beerenwinkel, A. Mahfouz,
1207 et al., Eleven grand challenges in single-cell data science, *Genome biology*
1208 21 (1) (2020) 1–35.
- 1209 [8] X. Wang, Y. He, Q. Zhang, X. Ren, Z. Zhang, Direct comparative anal-
1210 yses of 10x genomics chromium and smart-seq2, *Genomics, proteomics &*
1211 *bioinformatics* 19 (2) (2021) 253–266.
- 1212 [9] A. K. Shalek, R. Satija, J. Shuga, J. J. Trombetta, D. Gennert, D. Lu,
1213 P. Chen, R. S. Gertner, J. T. Gaublomme, N. Yosef, et al., Single-cell
1214 rna-seq reveals dynamic paracrine control of cellular variation, *Nature*
1215 510 (7505) (2014) 363–369.
- 1216 [10] W. Hou, Z. Ji, H. Ji, S. C. Hicks, A systematic evaluation of single-cell
1217 rna-sequencing imputation methods, *Genome biology* 21 (1) (2020) 1–30.
- 1218 [11] G. Eraslan, L. M. Simon, M. Mircea, N. S. Mueller, F. J. Theis, Single-cell
1219 rna-seq denoising using a deep count autoencoder, *Nature communications*
1220 10 (1) (2019) 1–14.
- 1221 [12] W. V. Li, J. J. Li, An accurate and robust imputation method scimpute
1222 for single-cell rna-seq data, *Nature communications* 9 (1) (2018) 1–9.
- 1223 [13] D. Van Dijk, R. Sharma, J. Nainys, K. Yim, P. Kathail, A. J. Carr, C. Bur-
1224 dziak, K. R. Moon, C. L. Chaffer, D. Pattabiraman, et al., Recovering gene
1225 interactions from single-cell data using data diffusion, *Cell* 174 (3) (2018)
1226 716–729.
- 1227 [14] Z.-H. Wen, J. L. Langsam, L. Zhang, W. Shen, X. Zhou, A bayesian fac-
1228 torization method to recover single-cell rna sequencing data, *Cell reports*
1229 *methods* 2 (1) (2022) 100133.

- 1230 [15] T. Peng, Q. Zhu, P. Yin, K. Tan, Scrabble: single-cell rna-seq imputation
1231 constrained by bulk rna-seq data, *Genome biology* 20 (1) (2019) 1–12.
- 1232 [16] Z. Hu, S. Zu, J. S. Liu, Simples: a single-cell rna sequencing imputation
1233 strategy preserving gene modules and cell clusters variation, *NAR genomics
1234 and bioinformatics* 2 (4) (2020) lqaa077.
- 1235 [17] M. Lotfollahi, A. Litinetskaya, F. J. Theis, Multigrate: single-cell multi-
1236 omic data integration, *bioRxiv* (2022).
- 1237 [18] K. E. Wu, K. E. Yost, H. Y. Chang, J. Zou, Babel enables cross-modality
1238 translation between multiomic profiles at single-cell resolution, *Proceedings
1239 of the National Academy of Sciences* 118 (15) (2021).
- 1240 [19] K. D. Yang, A. Belyaeva, S. Venkatachalapathy, K. Damodaran, A. Kat-
1241 coff, A. Radhakrishnan, G. Shivashankar, C. Uhler, Multi-domain transla-
1242 tion between single-cell imaging and sequencing data using autoencoders,
1243 *Nature Communications* 12 (1) (2021) 1–10.
- 1244 [20] M. Marouf, P. Machart, V. Bansal, C. Kilian, D. S. Magruder, C. F. Krebs,
1245 S. Bonn, Realistic in silico generation and augmentation of single-cell rna-
1246 seq data using generative adversarial networks, *Nature communications*
1247 11 (1) (2020) 1–12.
- 1248 [21] M. Lotfollahi, F. A. Wolf, F. J. Theis, scgen predicts single-cell perturbation
1249 responses, *Nature methods* 16 (8) (2019) 715–721.
- 1250 [22] Y. Zhao, C. Kilian, J.-E. Turner, L. Bosurgi, K. Roedl, P. Bartsch, A.-C.
1251 Gnirck, F. Cortesi, C. Schultheiß, M. Hellmig, et al., Clonal expansion and
1252 activation of tissue-resident memory-like th17 cells expressing gm-csf in the
1253 lungs of patients with severe covid-19, *Science Immunology* 6 (56) (2021)
1254 eabf6692.
- 1255 [23] I. Tolstikhin, O. Bousquet, S. Gelly, B. Schoelkopf, Wasserstein auto-
1256 encoders, *arXiv preprint arXiv:1711.01558* (2017).
- 1257 [24] T. Stuart, A. Butler, P. Hoffman, C. Hafemeister, E. Papalexi, W. M.
1258 Mauck III, Y. Hao, M. Stoeckius, P. Smibert, R. Satija, Comprehensive
1259 integration of single-cell data, *Cell* 177 (7) (2019) 1888–1902.
- 1260 [25] V. Dumoulin, J. Shlens, M. Kudlur, A learned representation for artistic
1261 style, *arXiv preprint arXiv:1610.07629* (2016).
- 1262 [26] J. L. Ba, J. R. Kiros, G. E. Hinton, Layer normalization, *arXiv preprint
1263 arXiv:1607.06450* (2016).
- 1264 [27] Z. Fu, E. R. Gilbert, D. Liu, Regulation of insulin synthesis and secretion
1265 and pancreatic beta-cell dysfunction in diabetes, *Current diabetes reviews*
1266 9 (1) (2013) 25–53.

- 1267 [28] M. Slyper, C. Porter, O. Ashenberg, J. Waldman, E. Drokhyansky,
1268 I. Wakiro, C. Smillie, G. Smith-Rosario, J. Wu, D. Dionne, et al., A single-
1269 cell and single-nucleus rna-seq toolbox for fresh and frozen human tumors,
1270 *Nature medicine* 26 (5) (2020) 792–802.
- 1271 [29] G. C. Linderman, J. Zhao, M. Roulis, P. Bielecki, R. A. Flavell, B. Nadler,
1272 Y. Kluger, Zero-preserving imputation of single-cell rna-seq data, *Nature*
1273 *Communications* 13 (1) (2022) 1–11.
- 1274 [30] E. Bakos, C. A. Thaiss, M. P. Kramer, S. Cohen, L. Radomir, I. Orr,
1275 N. Kaushansky, A. Ben-Nun, S. Becker-Herman, I. Shachar, Ccr2 regulates
1276 the immune response by modulating the interconversion and function of
1277 effector and regulatory t cells, *The Journal of Immunology* 198 (12) (2017)
1278 4659–4671.
- 1279 [31] G. Monaco, B. Lee, W. Xu, S. Mustafah, Y. Y. Hwang, C. Carré, N. Burdin,
1280 L. Visan, M. Ceccarelli, M. Poidinger, et al., Rna-seq signatures normalized
1281 by mrna abundance allow absolute deconvolution of human immune cell
1282 types, *Cell reports* 26 (6) (2019) 1627–1640.
- 1283 [32] V. A. Traag, L. Waltman, N. J. Van Eck, From louvain to leiden: guaran-
1284 teeing well-connected communities, *Scientific reports* 9 (1) (2019) 1–12.
- 1285 [33] M. Croft, Control of immunity by the tnfr-related molecule ox40 (cd134),
1286 *Annual review of immunology* 28 (2009) 57–78.
- 1287 [34] T. Riaz, L. M. Sollid, I. Olsen, G. A. de Souza, Quantitative proteomics of
1288 gut-derived th1 and th1/th17 clones reveal the presence of cd28+ nkg2d-
1289 th1 cytotoxic cd4+ t cells, *Molecular & Cellular Proteomics* 15 (3) (2016)
1290 1007–1016.
- 1291 [35] L. Peng, Y. Chen, Q. Ou, X. Wang, N. Tang, Lncrna miat correlates with
1292 immune infiltrates and drug reactions in hepatocellular carcinoma, *Inter-
1293 national immunopharmacology* 89 (2020) 107071.
- 1294 [36] D. P. Saraiva, A. Jacinto, P. Borralho, S. Braga, M. G. Cabral, Hla-dr in
1295 cytotoxic t lymphocytes predicts breast cancer patients’ response to neoad-
1296 jvant chemotherapy, *Frontiers in immunology* (2018) 2605.
- 1297 [37] M. S. Lee, K. Hanspers, C. S. Barker, A. P. Korn, J. M. McCune, Gene
1298 expression profiles during human cd4+ t cell differentiation, *International
1299 immunology* 16 (8) (2004) 1109–1124.
- 1300 [38] S. Aibar, C. B. González-Blas, T. Moerman, V. A. Huynh-Thu, H. Im-
1301 richova, G. Hulselmans, F. Rambow, J.-C. Marine, P. Geurts, J. Aerts,
1302 et al., Scenic: single-cell regulatory network inference and clustering, *Nat-
1303 ure methods* 14 (11) (2017) 1083–1086.

- 1304 [39] A. M. Thornton, J. Lu, P. E. Korty, Y. C. Kim, C. Martens, P. D. Sun,
1305 E. M. Shevach, Helios+ and helios- treg subpopulations are phenotypically
1306 and functionally distinct and express dissimilar tcr repertoires, *European*
1307 *journal of immunology* 49 (3) (2019) 398–412.
- 1308 [40] C. Imbratta, H. Hussein, F. Andris, G. Verdeil, c-maf, a swiss army knife
1309 for tolerance in lymphocytes, *Frontiers in immunology* 11 (2020) 206.
- 1310 [41] X. O. Yang, B. P. Pappu, R. Nurieva, A. Akimzhanov, H. S. Kang,
1311 Y. Chung, L. Ma, B. Shah, A. D. Panopoulos, K. S. Schluns, et al., T
1312 helper 17 lineage differentiation is programmed by orphan nuclear recep-
1313 tors ror α and ror γ , *Immunity* 28 (1) (2008) 29–39.
- 1314 [42] K. Street, D. Risso, R. B. Fletcher, D. Das, J. Ngai, N. Yosef, E. Purdom,
1315 S. Dudoit, Slingshot: cell lineage and pseudotime inference for single-cell
1316 transcriptomics, *BMC genomics* 19 (1) (2018) 1–16.
- 1317 [43] E. Stephenson, G. Reynolds, R. A. Botting, F. J. Calero-Nieto, M. D.
1318 Morgan, Z. K. Tuong, K. Bach, W. Sungnak, K. B. Worlock, M. Yoshida,
1319 et al., Single-cell multi-omics analysis of the immune response in covid-19,
1320 *Nature medicine* 27 (5) (2021) 904–916.
- 1321 [44] M. Ota, Y. Nagafuchi, H. Hatano, K. Ishigaki, C. Terao, Y. Takeshima,
1322 H. Yanaoka, S. Kobayashi, M. Okubo, H. Shirai, et al., Dynamic landscape
1323 of immune cell-specific gene regulation in immune-mediated diseases, *Cell*
1324 184 (11) (2021) 3006–3021.
- 1325 [45] G. Meyer Zu Horste, C. Wu, C. Wang, L. Cong, M. Pawlak, Y. Lee,
1326 W. Elyaman, S. Xiao, A. Regev, V. Kuchroo, Rbpj controls development
1327 of pathogenic th17 cells by regulating il-23 receptor expression. *cell rep* 16
1328 (2): 392–404 (2016).
- 1329 [46] S. De Biasi, M. Meschiari, L. Gibellini, C. Bellinazzi, R. Borella, L. Fidanza,
1330 L. Gozzi, A. Iannone, D. Lo Tartaro, M. Mattioli, et al., Marked t cell
1331 activation, senescence, exhaustion and skewing towards th17 in patients
1332 with covid-19 pneumonia, *Nature communications* 11 (1) (2020) 1–17.
- 1333 [47] B. J. Meckiff, C. Ramírez-Suástegui, V. Fajardo, S. J. Chee, A. Kurnadi,
1334 H. Simon, S. Eschweiler, A. Grifoni, E. Pelosi, D. Weiskopf, et al., Imbal-
1335 ance of regulatory and cytotoxic sars-cov-2-reactive cd4+ t cells in covid-19,
1336 *Cell* 183 (5) (2020) 1340–1353.
- 1337 [48] L. Loyal, S. Warth, K. Jürchott, F. Mölder, C. Nikolaou, N. Babel,
1338 M. Nienen, S. Durlanik, R. Stark, B. Kruse, et al., Slamf7 and il-6r define
1339 distinct cytotoxic versus helper memory cd8+ t cells, *Nature communica-*
1340 *tions* 11 (1) (2020) 1–12.
- 1341 [49] J. Yang, M. Zhong, E. Zhang, K. Hong, Q. Yang, D. Zhou, J. Xia, Y.-Q.
1342 Chen, M. Sun, B. Zhao, et al., Broad phenotypic alterations and potential

- 1343 dysfunction of lymphocytes in individuals clinically recovered from covid-
1344 19, *Journal of Molecular Cell Biology* 13 (3) (2021) 197–209.
- 1345 [50] M. Lotfollahi, M. Naghipourfar, M. D. Luecken, M. Khajavi, M. Büttner,
1346 M. Wagenstetter, Ž. Avsec, A. Gayoso, N. Yosef, M. Interlandi, et al.,
1347 Mapping single-cell data to reference atlases by transfer learning, *Nature*
1348 *Biotechnology* 40 (1) (2022) 121–130.
- 1349 [51] C. Wagner, M. Griesel, A. Mikolajewska, A. Mueller, M. Nothacker,
1350 K. Kley, M.-I. Metzendorf, A.-L. Fischer, M. Kopp, M. Stegemann, et al.,
1351 Systemic corticosteroids for the treatment of covid-19, *Cochrane Database*
1352 *of Systematic Reviews* (8) (2021).
- 1353 [52] W. Chen, J. Luo, Y. Ye, R. Hoyle, W. Liu, R. Borst, S. Kazani, E. A.
1354 Shikatani, V. J. Erpenbeck, I. D. Pavord, et al., The roles of type 2 cyto-
1355 toxic t cells in inflammation, tissue remodeling, and prostaglandin (pg) d2
1356 production are attenuated by pgd2 receptor 2 antagonism, *The Journal of*
1357 *Immunology* 206 (11) (2021) 2714–2724.
- 1358 [53] S. Lab, *panc8.SeuratData: Eight Pancreas Datasets Across Five Technolo-*
1359 *gies, r package version 3.0.2* (2019).
- 1360 [54] J. Ding, X. Adiconis, S. K. Simmons, M. S. Kowalczyk, C. C. Hession,
1361 N. D. Marjanovic, T. K. Hughes, M. H. Wadsworth, T. Burks, L. T.
1362 Nguyen, et al., Systematic comparison of single-cell and single-nucleus rna-
1363 sequencing methods, *Nature biotechnology* 38 (6) (2020) 737–746.
- 1364 [55] A. Gayoso, R. Lopez, G. Xing, P. Boyeau, K. Wu, M. Jayasuriya, E. Melh-
1365 man, M. Langevin, Y. Liu, J. Samarán, et al., *Scvi-tools: A library for*
1366 *deep probabilistic analysis of single-cell omics data*, *bioRxiv* (2021).
- 1367 [56] R. Menon, A. S. Bomback, B. B. Lake, C. Stutzke, S. M. Grewenow,
1368 S. Menez, V. D. D’Agati, S. Jain, R. Knight, S. H. Lecker, et al., Inte-
1369 grated single-cell sequencing and histopathological analyses reveal diverse
1370 injury and repair responses in a participant with acute kidney injury:
1371 a clinical-molecular-pathologic correlation, *Kidney International* 101 (6)
1372 (2022) 1116–1125.
- 1373 [57] G. X. Zheng, J. M. Terry, P. Belgrader, P. Ryvkin, Z. W. Bent, R. Wilson,
1374 S. B. Ziraldo, T. D. Wheeler, G. P. McDermott, J. Zhu, et al., Massively
1375 parallel digital transcriptional profiling of single cells, *Nature communica-*
1376 *tions* 8 (1) (2017) 1–12.
- 1377 [58] F. A. Wolf, P. Angerer, F. J. Theis, *Scanpy: large-scale single-cell gene*
1378 *expression data analysis*, *Genome biology* 19 (1) (2018) 1–5.
- 1379 [59] Y. Le Cun, F. Fogelman-Soulié, *Modèles connexionnistes de*
1380 *l’apprentissage*, *Intellectica* 2 (1) (1987) 114–143.

- 1381 [60] G. E. Hinton, R. Zemel, Autoencoders, minimum description length and
1382 helmholtz free energy, *Advances in neural information processing systems*
1383 6 (1993).
- 1384 [61] D. P. Kingma, M. Welling, Auto-encoding variational bayes, arXiv preprint
1385 arXiv:1312.6114 (2013).
- 1386 [62] C. Villani, *Optimal transport: old and new*, Vol. 338, Springer, 2009.
- 1387 [63] M. Arjovsky, S. Chintala, L. Bottou, Wasserstein generative adversarial
1388 networks, in: *International conference on machine learning*, PMLR, 2017,
1389 pp. 214–223.
- 1390 [64] P. K. Rubenstein, B. Schoelkopf, I. Tolstikhin, On the latent space of
1391 wasserstein auto-encoders, arXiv preprint arXiv:1802.03761 (2018).
- 1392 [65] D. Risso, F. Perraudeau, S. Gribkova, S. Dudoit, J.-P. Vert, A general and
1393 flexible method for signal extraction from single-cell rna-seq data, *Nature*
1394 *communications* 9 (1) (2018) 1–17.
- 1395 [66] P. J. Huber, Robust estimation of a location parameter, *Annals Mathemat-*
1396 *ics Statistics* (1964).
- 1397 [67] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, A. Smola, A ker-
1398 nel two-sample test, *Journal of Machine Learning Research* 13 (25) (2012)
1399 723–773.
1400 URL <http://jmlr.org/papers/v13/gretton12a.html>
- 1401 [68] D. Misra, Mish: A self regularized non-monotonic activation function,
1402 arXiv preprint arXiv:1908.08681 (2019).
- 1403 [69] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, J. Han, On the variance
1404 of the adaptive learning rate and beyond, arXiv preprint arXiv:1908.03265
1405 (2019).
- 1406 [70] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv
1407 preprint arXiv:1412.6980 (2014).
- 1408 [71] R. Liaw, E. Liang, R. Nishihara, P. Moritz, J. E. Gonzalez, I. Stoica, Tune:
1409 A research platform for distributed model selection and training, arXiv
1410 preprint arXiv:1807.05118 (2018).
- 1411 [72] D. Kobak, P. Berens, The art of using t-sne for single-cell transcriptomics,
1412 *Nature communications* 10 (1) (2019) 1–14.
- 1413 [73] G. C. Linderman, M. Rachh, J. G. Hoskins, S. Steinerberger, Y. Kluger,
1414 Fast interpolation-based t-sne for improved visualization of single-cell rna-
1415 seq data, *Nature methods* 16 (3) (2019) 243–245.

- 1416 [74] S. V. Stassen, D. M. Siu, K. C. Lee, J. W. Ho, H. K. So, K. K. Tsia, Parc:
1417 ultrafast and accurate clustering of phenotypic data of millions of single
1418 cells, *Bioinformatics* 36 (9) (2020) 2778–2786.
- 1419 [75] A. Subramanian, P. Tamayo, V. K. Mootha, S. Mukherjee, B. L. Ebert,
1420 M. A. Gillette, A. Paulovich, S. L. Pomeroy, T. R. Golub, E. S. Lander,
1421 et al., Gene set enrichment analysis: a knowledge-based approach for in-
1422 terpreting genome-wide expression profiles, *Proceedings of the National*
1423 *Academy of Sciences* 102 (43) (2005) 15545–15550.
- 1424 [76] E. Y. Chen, C. M. Tan, Y. Kou, Q. Duan, Z. Wang, G. V. Meirelles, N. R.
1425 Clark, A. Ma’ayan, Enrichr: interactive and collaborative html5 gene list
1426 enrichment analysis tool, *BMC bioinformatics* 14 (1) (2013) 1–14.
- 1427 [77] J. H. Friedman, Greedy function approximation: a gradient boosting ma-
1428 chine, *Annals of statistics* (2001) 1189–1232.
- 1429 [78] S. Haghghi, M. Jasemi, S. Hessabi, A. Zolanvari, Pycm: Multiclass con-
1430 fusion matrix library in python, *Journal of Open Source Software* 3 (25)
1431 (2018) 729.