

# 1 Title

## 2 **Deep-learning-based cell composition analysis from tissue expression profiles.**

3 Kevin Menden<sup>§</sup>, Mohamed Marouf, Anupriya Dalmia, Peter Heutink, Stefan Bonn<sup>§</sup>

4

5 <sup>§</sup> Correspondence to sbonn@uke.de & Kevin.Menden@dzne.de

## 6 Abstract

7 We present Scaden, a deep neural network for cell deconvolution that uses gene  
8 expression information to infer the cellular composition of tissues. Scaden is trained  
9 on single cell RNA-seq data to engineer discriminative features that confer robustness  
10 to bias and noise, making complex data preprocessing and feature selection  
11 unnecessary. We demonstrate that Scaden outperforms existing deconvolution  
12 algorithms in both precision and robustness, across tissues and species. A single  
13 trained network reliably deconvolves bulk RNA-seq and microarray, human and  
14 mouse tissue expression data. Due to this stability and flexibility, we surmise that deep  
15 learning-based cell deconvolution will become a mainstay across data types and  
16 algorithmic approaches. Scaden's comprehensive software package is easy to use on  
17 novel as well as diverse existing expression datasets available in public resources,  
18 deepening the molecular and cellular understanding of developmental and disease  
19 processes.

## 20 Keywords

21 Cell Deconvolution, Deep Learning, Machine Learning, single cell RNA sequencing,  
22 RNA sequencing, Deep Sequencing, Source Separation.

## 23 Introduction

24 The analysis of tissue-specific gene expression using Next Generation Sequencing  
25 (RNA-seq) is a centerpiece of the molecular characterization of biological and medical  
26 processes<sup>1</sup>. A well-known limitation of tissue-based RNA-seq is that it typically  
27 measures average gene expression across many molecularly diverse cell types that  
28 can have distinct cellular states<sup>2</sup>. A change in gene expression between two conditions  
29 can therefore be attributed to a change in the cellular composition of the tissue or a  
30 change in gene expression in a specific cell population, or a mixture of the two. To  
31 deconvolve systematic differences in cell type composition is especially important in  
32 systems with cellular proliferation (e.g. cancer) or cellular death (e.g. neuronal loss in  
33 Neurodegenerative Diseases)<sup>3</sup>.

34 To account for this problem, several computational cell deconvolution methods have  
35 been proposed during the last years<sup>4,5</sup>. These algorithms attempt to calculate an  
36 approximation of the cell type composition of a given gene expression sample, such  
37 that systematic differences in cellular abundance between samples can be detected,  
38 interpreted, and possibly corrected for. Current algorithms utilize gene expression  
39 profiles (GEPs) of cell type-specifically expressed genes to estimate cellular fractions  
40 using linear regression<sup>4</sup>. While the best performing linear regression algorithms for  
41 deconvolution seem to be variations of Support Vector Regression (SVR)<sup>6-10</sup>, the  
42 selection of an optimal GEP is a field of active research<sup>10,11</sup>. Indeed, it has been  
43 recently shown that the design of the GEP is the most important factor in most  
44 deconvolution methods, as results from different algorithms strongly correlate given  
45 the same GEP<sup>11</sup>.

46 In theory, an optimal GEP should contain a set of genes that are predominantly  
47 expressed within each cell population of a complex sample<sup>12</sup>. They should be stably

48 expressed across experimental conditions, for example across health and disease,  
49 and resilient to experimental noise and bias. The negative impact of bias on  
50 deconvolution performance can be partly improved by using large, heterogeneous  
51 GEP matrices<sup>11</sup>. It is therefore not surprising that recent advancement in cell  
52 deconvolution relied almost exclusively on sophisticated algorithms to normalize the  
53 data and engineer optimal GEPs<sup>10</sup>.

54 While GEP-based approaches lay the foundational basis of modern cell deconvolution  
55 algorithms, we hypothesize that Deep Neural Networks (DNNs) could create optimal  
56 features for cell deconvolution, without relying on the complex generation of GEPs.  
57 DNNs such as multilayer perceptrons are universal function approximators that  
58 achieve state-of-the-art performance on classification and regression tasks. We  
59 theorize that by using gene expression information as network input, hidden layer  
60 nodes of the DNN would represent higher-order latent representations of cell types  
61 that are robust to input noise and technical bias.

62 An obvious limitation of DNNs is the requirement for large training data to avoid  
63 overfitting of the machine learning model. While ground truth information on tissue  
64 RNA-seq cell composition is scarce, one can use single cell RNA-seq (scRNA-seq)  
65 data to obtain virtually unlimited *in silico* tissue datasets of predefined cell  
66 composition<sup>7-9,13-15</sup>. This is achieved by sub-sampling and subsequently merging cells  
67 from scRNA-seq datasets and is limited only by the availability of tissue-specific  
68 scRNA-seq data. It is to be noted that scRNA-seq data suffers from known biases,  
69 such as drop-out, that RNA-seq data is not subject to<sup>16</sup>. While this complicates the use  
70 of scRNA-seq data for GEP design<sup>8</sup>, we surmise that latent network nodes could  
71 represent features that are robust to such biases.

72 Based on these assumptions we developed a single-cell-assisted deconvolutional  
73 DNN (Scaden) that uses simulated bulk RNA-seq samples for training and predicts  
74 cell type proportions for input expression samples of cell mixtures. Scaden is trained  
75 on publicly available scRNA- and RNA-seq data, does not rely on specific GEP  
76 matrices, and automatically infers informative features. Finally, we show that Scaden  
77 deconvolves expression data into cell types with higher precision and robustness than  
78 existing methods that rely on GEP matrices, across tissues, species, and data types.

## 79 Results

### 80 Scaden Overview, Model Selection, and Training

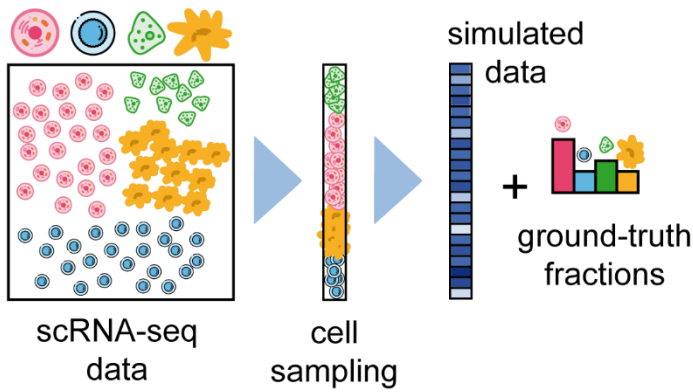
81 The basic architecture of Scaden is a DNN that takes gene counts of RNA-seq data  
82 as input and outputs predicted cell fractions (Fig. 1). To optimize the performance of  
83 the DNN, it is trained on data that contains both the gene expression and the real cell  
84 fraction information (Fig. 1A). The network then adjusts its weights to minimize the  
85 error between the predicted cell fractions and the real cell fractions (Fig. 1B).

86 For the model selection and training we made use of the virtually unlimited amount of  
87 artificial bulk RNA-seq datasets with defined composition that can be generated *in*  
88 *silico* from published scRNA-seq and RNA-seq datasets (simulated tissues) (Fig. 1,  
89 Tables S1 & S2). The only constraint being that the scRNA-seq and RNA-seq data  
90 must come from the same tissue as the bulk data subject to deconvolution.

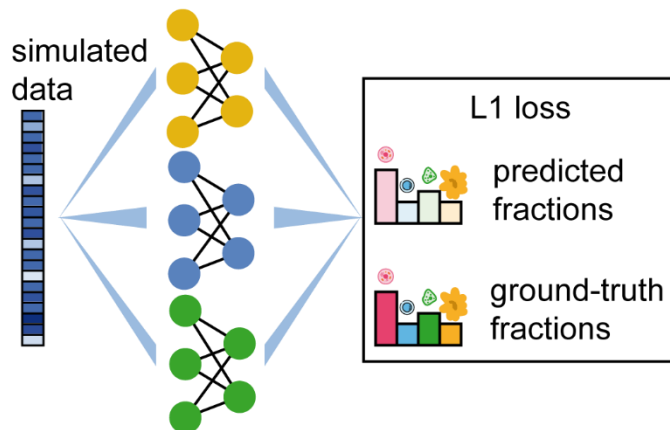
91 To find the optimal DNN architecture for cell deconvolution, we performed leave-one-  
92 dataset-out cross validation on simulated peripheral blood mononuclear cell (PBMC)  
93 tissue, training on mixtures of three scRNA-seq datasets and evaluating the  
94 performance on simulated tissue from a fourth scRNA-seq dataset (Table S1 & S3).

95 The final Scaden model is an ensemble of the three best performing models and the  
96 final cell type composition estimates are the averaged predictions of all three  
97 ensemble models (Fig. S1, Table S4).  
98

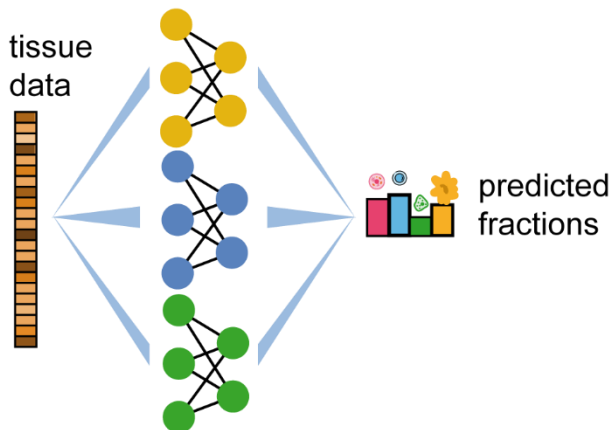
### A Simulated training data



### B Scaden training



### C Scaden predictions



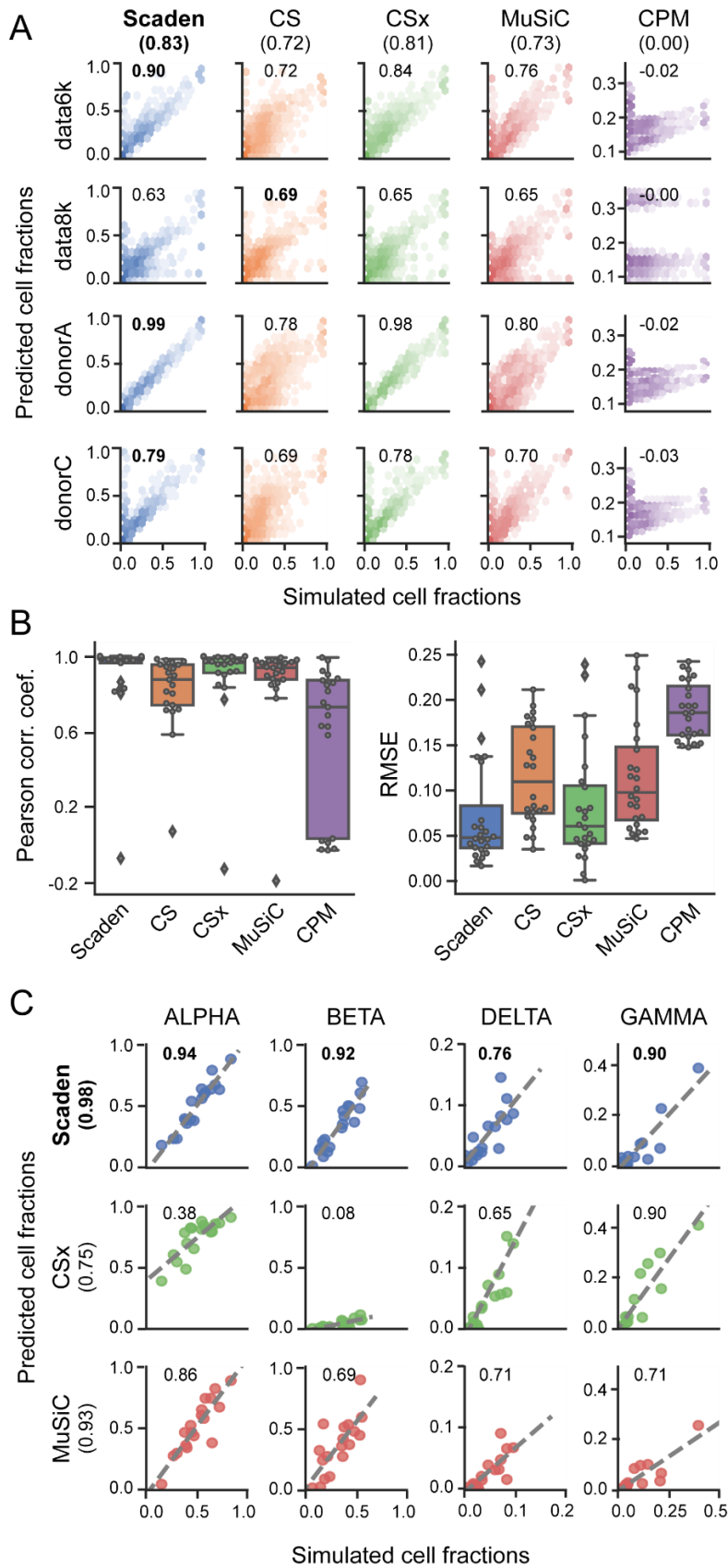
99

100 **Figure 1** Overview of training data generation and cell type deconvolution with Scaden. A:  
101 Artificial bulk samples are generated by subsampling random cells from a scRNA-seq datasets  
102 and merging their expression profiles. B: Model training and parameter optimization on  
103 simulated tissue RNA-seq data by comparing cell fraction predictions to ground-truth cell  
104 composition. C: Cell deconvolution of real tissue RNA-seq data using Scaden.

105

106 To get an initial estimate of Scaden's deconvolution fidelity we measured the root  
107 mean square error (RMSE), Lin's concordance correlation coefficient (CCC)<sup>17</sup>,  
108 Pearson's correlation coefficient (r), and the slope and intercept of the regression fitted  
109 for actual and predicted cell fractions. To this end, 32,000 human PBMC, 14,000  
110 human pancreas, 6,000 human ascites, and 30,000 mouse brain simulated tissue  
111 datasets were generated for network training and evaluation (Table S2). We then  
112 compared Scaden to four state-of-the-art GEP-based cell deconvolution algorithms,  
113 CIBERSORT (CS)<sup>6</sup>, CIBERSORTx (CSx)<sup>7</sup>, MuSiC<sup>8</sup>, and Cell Population Mapping  
114 (CPM)<sup>9</sup>. While CS relies on hand-curated GEP matrices, CSx, MuSiC, and CPM can  
115 generate GEPs using scRNA-seq data as input.

116 We first evaluated the deconvolution performance on simulated PBMC data, since  
117 curated GEP matrices and RNA-seq datasets with associated ground truth cell type  
118 compositions are available for human PBMCs, making this tissue uniquely suited  
119 toward deconvolution performance evaluation. Scaden was trained on simulated data  
120 from all datasets but a held-out dataset while CSx, MuSiC and CPM used a GEP  
121 generated from a scRNA-seq dataset excluding a held-out dataset (e.g. data6k,  
122 data8k, donorA). Subsequently the algorithms were tested on 500 simulated PBMC  
123 samples from a held-out scRNA-seq dataset (e.g. donorC) (Fig. 2A & B, Table S5).  
124 For CS we used the PBMC-optimized LM22 GEP matrix<sup>6</sup> and tested performance on  
125 the 500 simulated PBMC samples from a held-out scRNA-seq dataset (e.g. donorC).



126

127 **Figure 2** Deconvolution performance on simulated tissue data A: Ground truth values (x-axis)

128 plotted against cell type fraction estimates (y-axis) for predictions made on simulated data

129 from four PBMC scRNA-seq datasets. Darker color in a hexbin corresponds to more data  
130 points falling into this bin. Numbers inside the plotting area signify CCC values, the overall  
131 CCC is shown in parenthesis below the algorithm name. B: Boxplots of  $r$  and RMSE values  
132 for simulated PBMC data. C: Per-cell-type scatterplots of ground truth (x-axis) and predicted  
133 values (y-axis) for Scaden, CSx, and MuSiC on artificial pancreas data<sup>18</sup>. Numbers inside the  
134 plotting area signify CCC values.

135  
136 For three of four test datasets (data6k, donorA, donorC), Scaden obtained the highest  
137 CCC and lowest RMSE, followed by CSx, MuSiC, CS, and CPM (Fig. 2A, Table S5).  
138 For one test dataset, CS obtained the highest CCC and lowest RMSE, followed by  
139 CSx, MuSiC, Scaden and CPM. Overall, Scaden obtains the highest CCC and lowest  
140 RMSE (0.83, 0.09, respectively), followed by CSx (0.81, 0.10), MuSiC (0.73, 0.13), CS  
141 (0.72, 0.13), and CPM (0, 0.19) (Fig. 2A). As expected, all algorithms that use scRNA-  
142 seq data as reference perform good in this scenario with the notable exception of CPM,  
143 for which we could not generate reasonable predictions. We believe that CPM's  
144 modest performance might be due to the cell state mapping required by CPM. We  
145 created this mapping using UMAP, which might not be optimal for CPM. On average,  
146 Scaden also obtained the highest correlation and the best intercept and slope values  
147 on simulated PBMC data (Table S5).

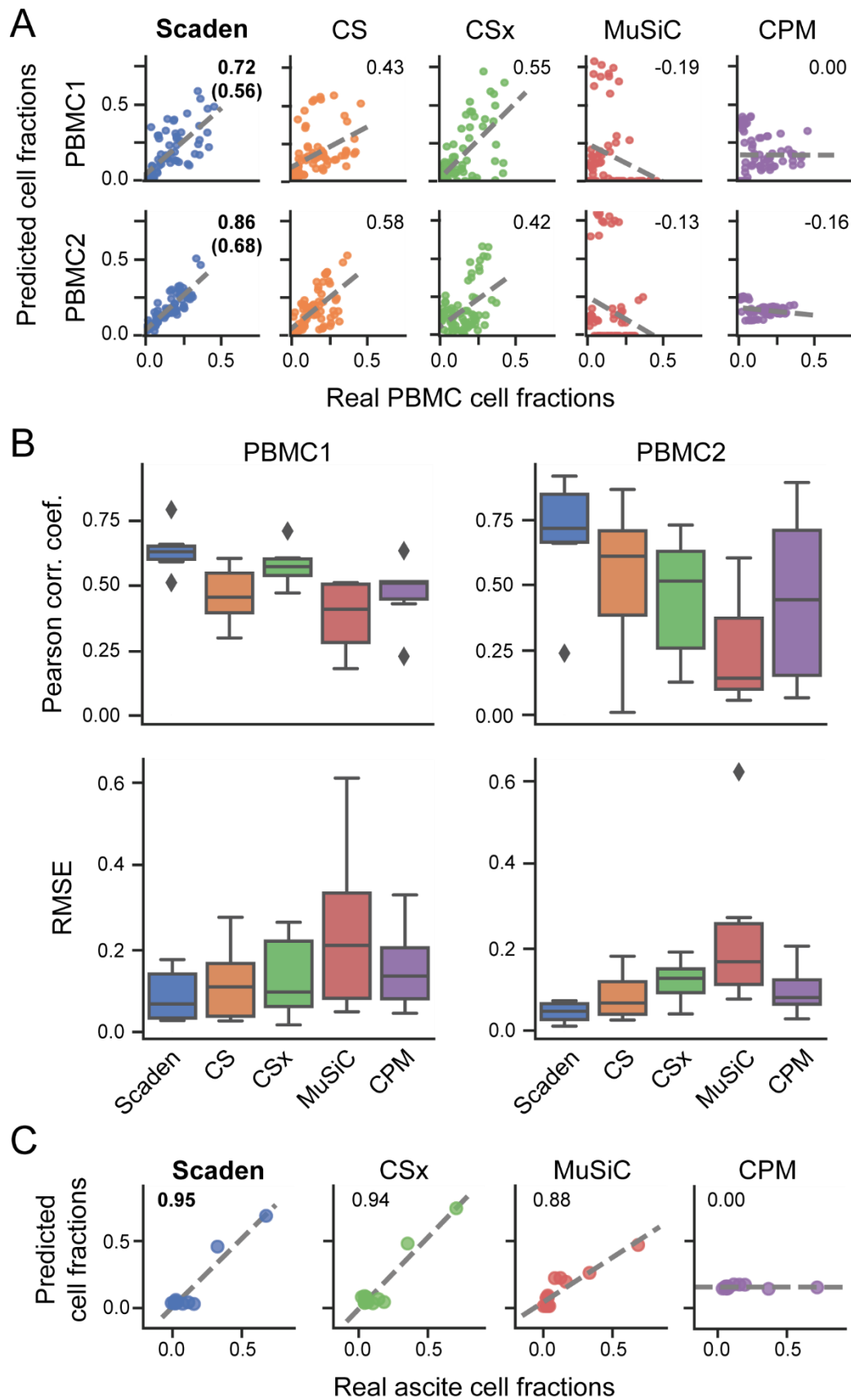
148 A specific feature of the MuSiC algorithm is that it preferentially selects genes with low  
149 inter-subject and intra-cell cluster variability for its GEP, which increases  
150 deconvolution robustness when high expression heterogeneity is observed between  
151 human subjects, for example<sup>8</sup>. To understand if Scaden can utilize multi-subject  
152 information to increase its deconvolution performance, we trained Scaden, CSx, and  
153 MuSiC on scRNA-seq pancreas data from several subjects<sup>19</sup> and assessed the  
154 performance on a separate simulated pancreas RNA-seq dataset<sup>18</sup> (Fig. 2C, Table



155 S6). To allow for direct comparison, we chose the same pancreas training and test  
156 datasets that were used in the original MuSiC publication (Table S1). To enable  
157 Scaden to leverage the heterogeneity of multi-subject data, training data was  
158 generated separately for every subject in the dataset (see Methods). CSx cannot profit  
159 from multi-subject data, but performed well on the artificial PBMC datasets and was  
160 therefore included in the comparison. The best performance is achieved by Scaden  
161 (CCC = 0.98), closely followed by MuSiC (CCC = 0.93), while CSx does not perform  
162 as well (CCC = 0.75) (Fig. 2C, Table S6). This provides strong evidence that Scaden,  
163 by separating training data generation for each subject, can learn inter-subject  
164 heterogeneity and outperform specialized multi-subject algorithms such as MuSiC on  
165 the cell-type deconvolution task.

## 166 Robust deconvolution of bulk expression data

167 The true use case of cell deconvolution algorithms is the cell fraction estimation of  
168 tissue RNA-seq data. We therefore assessed the performance of Scaden, CS, CSx,  
169 MuSiC, and CPM to deconvolve two publicly available human PBMC bulk RNA-seq  
170 datasets, for which ground-truth cell composition information was measured using flow  
171 cytometry (Fig. 3A, Tables S7 & S8). We will refer to these datasets that consists of  
172 12 samples each as PBMC1<sup>20</sup> and PBMC2<sup>10</sup>. Deconvolution for all methods was  
173 performed as described in the previous section, with the difference that data from all  
174 four PBMC scRNA-seq datasets was now deployed for Scaden training.



175

176 **Figure 3** Deconvolution of real tissue RNA-seq data A: Per-cell-type scatterplots of ground

177 truth (x-axis) and predicted values (y-axis) for Scaden, CS, CSx, MuSiC, and CPM on real

178 PBMC1 and PBMC2 cell fractions. Numbers inside the plotting area signify CCC values. For  
179 Scaden, the CCC using only scRNA-seq training data (in parenthesis) and the CCC using  
180 mixed scRNA-seq and RNA-seq training data is shown. B: Boxplots of  $r$  (first row) and RMSE  
181 (second row) values for real PBMC1 (first column) and PBMC2 (second column) data. C: Per-  
182 cell-type scatterplots of ground truth (x-axis) and predicted values (y-axis) for Scaden, CSx,  
183 MuSiC, and CPM on real ascite cell fractions. Numbers inside the plotting area signify CCC  
184 values.

185  
186 On the PBMC1 dataset, Scaden obtained the highest CCC and lowest RMSE (0.56,  
187 0.13), while CSx (0.55, 0.16) and CS (0.43, 0.15) performed well yet significantly worse  
188 than Scaden (Fig. 3A, Tables S8 & S9). CPM (0, 0.18) and MuSiC (-0.19, 0.32) both  
189 failed to deconvolve the cell fractions of the PBMC1 data. Scaden also obtained the  
190 best CCC and RMSE (0.68, 0.08) on the PBMC2 dataset, while CS (0.58, 0.10) and  
191 CSx (0.42, 0.13) obtained good deconvolution results. Similar to the PBMC1 data  
192 deconvolution results, CPM (-0.16, 0.11) as well as MuSiC (-0.13, 0.30) did not  
193 perform well on the PBMC2 deconvolution task. In addition to CCC and RMSE metrics,  
194 Scaden achieves the best correlation, intercept and slope on both PBMC datasets  
195 (Tables S8 & S9).

196 An additional algorithmic feature of Scaden is that it seamlessly integrates increasing  
197 amounts of training data, which can be of different types, such as a combination of  
198 simulated tissue and real tissue data with cell fraction information. In theory, even  
199 limited real tissue training data could make Scaden robust to data type bias and  
200 consequently improve Scaden's deconvolution performance on real tissue data. We  
201 therefore trained Scaden on a mix of simulated PBMC (500 samples) and real PBMC2  
202 (12 samples) data and evaluated its performance on real PBMC1 data (Fig. 3A, S2,  
203 Table S9). While the training contained only ~2% real data, Scaden's CCC increased

204 from 0.56 to 0.72 and the RMSE decreased from 0.13 to 0.10. We observed similar  
205 performance increases when Scaden was trained on simulated PBMC and real  
206 PBMC1 data and evaluated on real PBMC2 data (Fig. 3A, S2, Table S9).

207 We next evaluated Scaden's performance real ascites RNA-seq data, for which  
208 scRNA-seq and FACS cell proportion data is available<sup>21</sup> (Table S7). It is noteworthy  
209 that RNA-seq, scRNA-seq, and FACS data was generated for the same samples,  
210 which potentially entails reduced experimental and technical bias and consequently  
211 higher deconvolution fidelity for the ascites data as compared to the PBMC data. We  
212 did not evaluate CS's performance on the ascites data as there was no optimized  
213 ascites GEP available. For Scaden, CSx, and MuSiC we used scRNA-seq data to  
214 generate simulated tissue ascites data for training. Scaden, CSx, CPM, and MuSiC all  
215 accurately predict the cell type compositions for the three real ascites samples, while  
216 CPM does not perform well (Fig. 3C, Table S10). The highest CCC and lowest RMSE  
217 were achieved by Scaden (0.95, 0.06), followed by CSx (0.94, 0.07), MuSiC (0.88,  
218 0.08), and CPM (0, 0.18). This further validates that Scaden reliably deconvolves  
219 tissue RNA-seq data into the constituent cell fractions and that very accurate  
220 deconvolution results can be obtained if reference and target datasets are from the  
221 same experiment. Again, we surmise that CPM's drop in tissue RNA-seq  
222 deconvolution performance might be due to the required cell state space embedding,  
223 which makes deconvolution results not only depend on a good GEP matrix, but also  
224 on a good embedding. It might be that CPM's deconvolution suffers from the UMAP  
225 embedding used throughout this manuscript.

226 We next wanted to assess if Scaden's deconvolution performance is robust across  
227 species. We therefore tested whether a Scaden model trained on mouse brain scRNA-  
228 seq data could generate reasonable cell composition estimations for real human brain

229 RNA-seq data (Table S7). To this end, Scaden was trained on artificial data generated  
230 from five mouse brain scRNA-seq datasets and predicted the cell fractions on human  
231 post-mortem RNA-seq brain samples (390 prefrontal cortex samples) from the  
232 ROSMAP study<sup>22</sup>. Ground-truth cell fractions were not available for this data, which is  
233 why we used Braak stages<sup>23</sup> that correspond to Alzheimer's disease severity and  
234 correlate with the degree of neuronal loss. Overall, Scaden's cell fraction predictions  
235 capture the increased neuronal loss with increasing Braak stage (Fig. S3).  
236 Interestingly, the largest drop in neural percentage is observed at stage 5, when the  
237 neurodegeneration typically reaches the prefrontal cortex of the brain. By learning  
238 robust features, Scaden reliably deconvolves RNA-seq data in a cross-species  
239 comparison.

240 Given the robustness with which Scaden predicts tissue RNA-seq cell fractions using  
241 scRNA-seq training data, even across species, we next wanted to investigate if an  
242 scRNA-seq-trained Scaden model can also deconvolve other data types. To this end,  
243 we measured the deconvolution performance on a bulk PBMC microarray dataset (20  
244 samples)<sup>6</sup> of a Scaden model trained on scRNA-seq and RNA-seq PBMC data (see  
245 above). We compared Scaden to CS using the microarray-derived LM22 matrix. CS  
246 achieved a slightly higher CCC and slightly lower total RMSE (0.72, 0.11) than Scaden  
247 (0.71, 0.13), while Scaden obtained the highest average CCC (0.50) compared to CS  
248 (0.39) (Fig. S4, Table S11). Notably in this scenario, Scaden was trained entirely on  
249 simulated data and RNA-seq data, while CS's LM22 GEP was optimized on PBMC  
250 microarray data.

251 Overall, we provide strong evidence that Scaden robustly deconvolves tissue data  
252 across tissues, species, and even data types.

## 253 Discussion

254 Scaden is the first deep learning-based cell deconvolution algorithm. In many  
255 instances, it compares favorably in both prediction robustness and accuracy to existing  
256 deconvolution algorithms that rely on GEP design and linear regression. We believe  
257 that Scaden's performance relies to a large degree on the inherent feature engineering  
258 of the DNN. The network does not only select features (genes) for regression, it also  
259 creates novel features that are optimal for the regression task in the nodes of the  
260 hidden layers. These hidden features are non-linear combinations of the input features  
261 (gene expression), which makes it notoriously difficult to explain how a DNN works<sup>24</sup>.  
262 It is important to highlight that this feature creation is fundamentally different from all  
263 other existing cell deconvolution algorithms, which rely on heuristics that select a  
264 defined subset of genes as features for linear regression.

265 Another advantage of this inherent feature engineering is that Scaden can be trained  
266 to be robust to input noise and bias (e.g. batch effects). Noise and bias is all prevalent  
267 in experimental data, due to different sample quality, sample processing,  
268 experimenters, and instrumentation, for example. If the network is trained on different  
269 datasets of the same tissue, however, it learns to create hidden features that are  
270 robust to noise and bias, such as batch effects. This robustness is pivotal in real world  
271 cell deconvolution use cases, where the bulk RNA data for deconvolution and the  
272 training data (and therefore the network and GEP) contain different noise and biases.  
273 While especially recent cell deconvolution algorithms include batch correction  
274 heuristics prior to GEP construction, Scaden optimizes its hidden features  
275 automatically when trained on data from various batches.

276 The robustness to noise and bias, which might be due to hidden feature generation, is  
277 especially evident in Scaden's ability to deconvolve across data types. A network

278 trained on *in silico* bulk RNA-seq data can seamlessly deconvolve microarray data of  
279 the same tissue. This is quite noteworthy, as microarray data is known to have a  
280 reduced dynamic range and several hybridization-based biases compared to RNA-  
281 seq data. In other words, Scaden can deconvolve bulk data of types it has never been  
282 trained on, even in the face of strong data type bias. This raises the possibility that  
283 Scaden trained on scRNA-seq data might reliably deconvolve other bulk omics data  
284 as well, such as proteomic and metabolomic data. This assumption is strengthened  
285 by the fact that Scaden, trained on scRNA-seq data, attains state-of-the-art  
286 performance on the deconvolution of bulk RNA-seq data, two data types with very  
287 distinct biases<sup>16</sup>.

288 As highlighted in the introduction, a drawback for many DNNs is the large amount of  
289 training data required to obtain robust performance. Here, we used scRNA-seq data  
290 to create virtually unlimited amounts of *in silico* bulk RNA-seq data of predefined type  
291 (target tissue) with known composition, across datasets. This immediately highlights  
292 Scaden's biggest limitation, the dependency on scRNA-seq data of the target tissue.  
293 In this study we have shown that Scaden, trained solely on simulated data from  
294 scRNA-seq datasets, can outperform GEP-based deconvolution algorithms. We did  
295 observe, however, that the addition of labeled RNA-seq samples to the training data  
296 did significantly improve deconvolution performance in the case of PBMC data. We  
297 therefore believe that efforts to increase the similarity between simulated training data  
298 and the target bulk RNA-seq data could increase Scaden's performance further.  
299 Mixtures of *in silico* bulk RNA-seq data and publically available RNA-seq data, of  
300 purified cell types for example, could further increase the deconvolution performance  
301 of Scaden. Furthermore, domain adaptation methods can be used to improve  
302 performance of models that are trained on data (here, scRNA-seq data) that is similar

303 to the target data (here, RNA-seq data)<sup>25</sup>. In future versions, Scaden's simple  
304 multilayer perceptron architecture could leverage domain adaptation to further  
305 stabilize and improve its cell deconvolution performance.

306 Recent cell deconvolution algorithms have used cell fraction estimates to infer cell  
307 type-specific gene expression from bulk RNA-seq data. It is straightforward to use  
308 Scaden's cell fraction estimates to infer per group<sup>3</sup> and per sample<sup>7</sup> cell type-specific  
309 gene expression using simple regression or non-negative matrix factorization,  
310 respectively. We would like to add a note of caution, however, as the error of cell  
311 fraction estimates, which can be quite significant, is propagated into the gene  
312 expression calculations and will affect any downstream statistical analysis.

313 In summary, the deconvolution performance, robustness to noise and bias, the  
314 flexibility to learn from large numbers of *in silico* datasets, across data types (scRNA-  
315 seq and RNA-seq mixtures), and potentially even tissues makes us believe that DNN-  
316 based architectures will become an algorithmic mainstay of cell type deconvolution.

317



## 318 Methods

### 319 Datasets and pre-processing

#### 320 scRNA-seq datasets

321 The following human PBMC scRNA-seq datasets were downloaded from the 10X  
322 Genomics data download page: 6k PBMCs from a Healthy Donor, 8k PBMCs from a  
323 Healthy Donor, Frozen PBMCs (Donor A), Frozen PBMCs (Donor C){Zheng et al,  
324 2017}. Throughout this paper, these datasets are referred to with the handles data6k,  
325 data8k, donorA and donorC, respectively. These four datasets were chosen because  
326 of clearly identifiable cell types for the majority of cells. The Ascites scRNA-seq dataset  
327 was downloaded from <https://figshare.com> as provided by Schelker<sup>21</sup>. Pancreas and  
328 mouse brain datasets were downloaded from the scRNA-seq dataset collection of the  
329 Hemberg lab (<https://hemberg-lab.github.io/scRNA.seq.datasets/>). A table listing all  
330 datasets including references to the original publications can be found in Table S1.

#### 331 scRNA-seq preprocessing and analysis

332 All datasets were processed using the Python package Scanpy (v. 1.2.2)<sup>26</sup> following  
333 the Scanpy's reimplementation of the popular Seurat's clustering workflow. First, the  
334 corresponding cell-gene matrices were filtered for cells with less than 500 detected  
335 genes, and genes expressed in less than 5 cells. The resulting count matrix for each  
336 dataset was filtered for outliers with high or low numbers of counts. Gene expression  
337 was normalized to library size using the Scanpy function 'normalize\_per\_cell'. The  
338 normalized matrix of all filtered cells and genes was saved for the subsequent data  
339 generation step.

340 The following processing and analysis steps had the sole purpose of assigning cell  
341 type labels to every cell. All cells were clustered using the louvain clustering  
342 implementation of the Scanpy package. The louvain clustering resolution was chosen  
343 for each dataset, using the lowest possible resolution value (low resolution values lead  
344 to less clusters) for which the calculated clusters separated the cell types  
345 appropriately. The top 1000 highly variable genes were used for clustering, which were  
346 calculated using Scanpy's 'filter\_genes\_dispersion' function with parameters  
347 min\_mean=0.0125, max\_mean=3 and min\_disp=0.5. Principal Component Analysis  
348 (PCA) was used for dimensionality reduction.

349 To identify cell types, marker genes were investigated for all cell types in question. For  
350 PBMC datasets, useful marker genes were adopted from public resources such as the  
351 Seurat tutorial for 2700 PBMCs<sup>27</sup>. Briefly, IL7R was taken as marker for CD4 T-cells,  
352 LYZ for Monocytes, MS4A1 for B-cells, GNLY for Natural Killer cells, FCER1A for  
353 Dendritic cells and CD8A and CCL5 as markers for CD8 T-cells. For all other scRNA-  
354 seq datasets, marker genes and expected cell types were inferred from the original  
355 publication of the dataset. For instance, to annotate cell types of the mouse brain  
356 dataset from Zeisel et al.<sup>28</sup>, we used the same marker genes as Zeisel and colleagues.  
357 We did not use the same cell type labels from the original publications because a main  
358 objective was to assure that cell type labeling is consistent between all datasets of a  
359 certain tissue.

360 Cell type annotation was performed manually across all the clusters for each dataset,  
361 such that all cells belonging to the same cluster were labeled with the same cell type.  
362 The cell type identity of each cluster was chosen by crossing the cluster's highly  
363 differentially expressed genes with the curated cell type's marker genes. Clusters that

364 could not be clearly identified with a cell type were grouped into the ‘Unknown’  
365 category.

### 366 Tissue Datasets for Benchmarking

367 To assess the deconvolution performance on real tissue expression data, we used  
368 datasets for which the corresponding cell fractions were measured and published. The  
369 first dataset is the **PBMC1** dataset which was obtained from Zimmermann *et al.*<sup>20</sup>. The  
370 second dataset, **PBMC2**, was downloaded from GEO with accession code  
371 GSE107011<sup>10</sup>. This dataset contains both RNA-seq profiles of immune cells (S4  
372 cohort) and from bulk individuals (S13 cohort). As we were interested in the bulk  
373 profiles, we only used 12 samples from the S13 cohort from this data. Flow cytometry  
374 fractions were collected from the Monaco *et al.* publication<sup>10</sup>.

375 In addition to the above mentioned two PBMC datasets, we used Ascites RNA-seq  
376 data. This dataset was kindly provided by the authors and cell type fractions for this  
377 dataset were taken from the supplementary materials of the publication<sup>21</sup>.

378 For the evaluation on pancreas data, artificial bulk RNA-seq samples created from the  
379 scRNA-seq dataset of Xin *et al.*<sup>18</sup> were used. This dataset was downloaded from the  
380 resources of the MuSiC publication<sup>8</sup>. The artificial bulk RNA-seq samples used for  
381 evaluation were then created using the ‘bulk\_construct’ function of the MuSiC tool.

382 To assess how Scaden deals with unknown cell types in a bulk mixture, we used the  
383 whole blood dataset from Newman *et al.*<sup>7</sup>, which consists of 12 samples (GSE127813).

384 Cell type fractions were downloaded from the CSx website  
385 (<https://cibersortx.stanford.edu/download.php>).

386 The microarray dataset GSE65133 was downloaded from GEO, and cell type fractions  
387 taken from the original CS publication<sup>6</sup>.

388 Finally, we wanted to get insights into neurodegenerative cell fraction changes in the  
389 brain. While it is known that neurodegenerative diseases like Alzheimer's Disease are  
390 accompanied by a gradual loss of brain neurons, stage-specific cell type shifts are still  
391 hard to come by. Here we use the ROSMAP (Religious Orders Study and Memory and  
392 Aging Project Study) cortical RNA-seq dataset along with the corresponding clinical  
393 metadata, to infer cell type composition over six clinically relevant stages of  
394 neurodegeneration<sup>22</sup>.

#### 395 RNA-seq preprocessing and analysis

396 For the RNA-seq datasets analyzed in this study, we did not apply any additional  
397 processing steps, but used the obtained count or expression tables directly as  
398 downloaded for all dataset except the ROSMAP dataset. For the latter, we generated  
399 count tables from raw FastQ-files using Salmon<sup>29</sup> and the GRCh38 reference genome.  
400 FastQ-files from the ROSMAP study were downloaded from Synapse  
401 ([www.synapse.org](http://www.synapse.org)).

#### 402 Simulation of bulk RNA-seq samples from scRNA-seq data

403 Scadan's deep neural network requires large amounts of training RNA-seq samples  
404 with known cell fractions. This explains why the generation of artificial bulk RNA-seq  
405 data is one of the key elements of the Scaden workflow.

406 In order to generate the training data, preprocessed scRNA-seq datasets were used  
407 (see section 'Data Collection and Processing'), comprising the gene expression matrix  
408 and the cell type labels. Artificial RNA-seq samples were simulated by sub-sampling  
409 cells from individual scRNA-seq datasets - cells from different datasets were not  
410 merged into samples to preserve within-subject relationships. Datasets generated  
411 from multiple subjects were split according to subject and each sub-sampling was

412 constrained to cells from one subject in order to capture the cross-subject  
413 heterogeneity and keep subject-specific gene dependencies.

414 The exact sub-sampling procedure is described in the following. First, for every  
415 simulated sample, random fractions were created for all different cell types within each  
416 scRNA-seq dataset using the random module of the Python package NumPy. Briefly,  
417 a random number was chosen from a uniform distribution between 0 and 1 using the  
418 NumPy function 'random.rand()' for each cell type, and then this number was divided  
419 by the sum of all random numbers created to ensure the constraint of all fractions  
420 adding up to 1:

$$421 \quad f_c = \frac{r_c}{\sum_{C_{all}} r_c}$$

422 where  $r_c$  is the random number created for cell type  $c$ , and  $C_{all}$  is the set of all cell  
423 types. Here,  $f_c$  is the calculated random fraction for cell type  $c$ . Then, each fraction  
424 was multiplied with the total number of cells selected for each sample, yielding the  
425 number of cells to choose for a specific cell type:

$$426 \quad N_c = f_c * N_{total}$$

427  
428 where  $N_c$  is the number of cells to select for the cell type  $c$ , and  $N_{total}$  is the total  
429 number of cells contributing to one simulated RNA-seq sample (400, in this study).  
430 Next,  $N_c$  cells were randomly sampled from the scRNA-seq gene expression matrix  
431 for each cell type  $c$ . Afterwards, the randomly selected single-cell expression profiles  
432 for every cell type are then aggregated by summing their expression values, to yield  
433 the artificial bulk expression profile for this sample.  
434

435 Using the above described approach, cell compositions that are strongly biased  
436 toward a certain cell type or are missing specific cell types are rare among the  
437 generated training samples. To account for this and to simulate cell compositions with

438 a heavy bias to and the absence of certain cell types, a variation of the sub-sampling  
439 procedure was used to generate samples with sparse compositions, which we refer to  
440 as sparse samples. Before generating the random fractions for all cell types, a random  
441 number of cell types was selected to be absent from the sample, with the requirement  
442 of at least one cell type constituting the sample. After these leave-out cell types were  
443 chosen, random fractions were created and samples generated as described above.  
444 Using this procedure, we generated 32,000 samples for the human PBMC training  
445 dataset, 14,000 samples for the human pancreas training dataset and 30,000 samples  
446 for the mouse brain training dataset (Table S2).  
447 Artificial bulk RNA-seq datasets were stored in 'h5ad' format using the Anndata  
448 package<sup>26</sup>, which allows to store the samples together with their corresponding cell  
449 type ratios, while also keeping information about the scRNA-seq dataset of origin for  
450 each sample. This allowed to access samples from specific datasets, which is useful  
451 for cross validation.

## 452 Scaden Overview

453 The following section contains an overview of the input data preprocessing, the  
454 Scaden model, model selection, and how Scaden predictions are generated.

### 455 Input Data Preprocessing

456 The data preprocessing step is aimed to make the input data more suitable for  
457 machine learning algorithms. To achieve this, an optimal preprocessing procedure  
458 should transform any input data from the simulated samples or from the bulk RNA-seq  
459 to the same feature scale. Before any scaling procedure can be applied, it must be  
460 ensured that both the training data and the bulk RNA-seq data subject to prediction  
461 share the same features. Therefore, before scaling, both datasets are limited to

462 contain features (genes) that are available in both datasets.. The two-step processing  
463 procedure used for Scaden is described in the following:

464 First, to account for heteroscedasticity, a feature inherent to RNA-seq data, the data  
465 was transformed into logarithmic space by adding a pseudocount of 1 and then taking  
466 the Logarithm (base 2). Additional to stabilizing the variance, this transformation yields  
467 data that is approximately Gaussian.

468 Second, every sample was scaled to the range [0,1] using the MinMaxScaler() class  
469 from the Sklearn preprocessing module {ref}. Per sample scaling, unlike per feature  
470 scaling that is more common in machine learning, assures that inter-gene relative  
471 expression patterns in every sample are preserved. This is important, as our  
472 hypothesis was that a neural network could learn the deconvolution from these inter-  
473 gene expression patterns.

$$474 \quad x_{scaled,i} = (x_i - \min(X_i)) / (\max(X_i) - \min(X_i))$$

475 where  $x_{scaled,i}$  is the log2 expression value of gene x in sample i,  $X_i$  is the vector of  
476 log2 expression values for all genes of sample i,  $\min(X_i)$  is the minimum gene  
477 expression of vector  $X_i$ , and  $\max(X_i)$  the maximum gene expression of vector  $X_i$ .

478 Note that all training datasets are stored as expression values and are only processed  
479 as described above. In the deployment use-case the simulated training data should  
480 contain the same features as in the bulk RNA-seq sample that shall be deconvolved.

## 481 Model Selection

482 The goal of model selection was to find an architecture and hyperparameters that  
483 robustly deconvolve simulated tissue RNA-seq data and, more importantly, real bulk  
484 RNA-seq data. Due to the very limited availability of bulk RNA-seq datasets with known  
485 cell fractions, model selection was mainly optimized on the simulated PBMC datasets.

486 To capture inter-experimental variation, we used leave-one-dataset-out cross  
487 validation for model optimization: a model was trained on simulated data from all but  
488 one dataset, and performance was tested on simulated samples from the left-out  
489 dataset. This allows to simulate batch effects between datasets and helps to test the  
490 generalizability of the model. Model performance was evaluated based on pearson  
491 product moment correlation and absolute deviation between predicted and ground  
492 truth values. As averaging the predictions of models with different architectures  
493 increased performance, we decided to use an ensemble architecture for Scaden. For  
494 this ensemble, the three best performing architectures were chosen. Model training  
495 and prediction is done separately for each model, with the prediction averaging step  
496 combining all model predictions (Fig. S1). We provide a list of all tested parameters in  
497 the supplementary materials (Table S4).

#### 498 Final Scaden Model

499 The Scaden model learns cell type deconvolution through supervised training on  
500 datasets of simulated bulk RNA-seq samples simulated with scRNA-seq data. To  
501 account for model biases and to improve performance, Scaden consists of an  
502 ensemble of three deep neural networks with varying architectures and degrees of  
503 dropout regularization. All models of the ensemble use four layers of varying sizes  
504 between 32 and 1024 nodes, with dropout-regularization implemented in two of the  
505 three ensemble models. The exact layer sizes and dropout rates are listed in Table  
506 S3. The Rectified Linear Unit (ReLU) is used as activation function in every internal  
507 layer. We used a Softmax function to predict cell fractions, as we did not see any  
508 improvements in using a linear output function with consecutive non-negativity  
509 correction and sum-to-one scaling. Python (v. 3.6.6) and the TensorFlow library (v.



510 1.10.0) were used for implementation of Scaden. A complete list of all software used  
511 for the implementation of Scaden is provided in Table S12.

## 512 Training and Prediction

513 After the preprocessing of the data a Scaden ensemble can be trained on simulated  
514 tissue RNA-seq data or mixtures of simulated and real tissue RNA-seq data.  
515 Parameters are optimized using Adam with a learning rate of 0.0001 and a batch size  
516 of 128. We used an L1 loss as optimization objective:

$$517 \quad L1(y_i, \hat{y}_i) = |y_i - \hat{y}_i|$$

518 where  $y_i$  is the vector of ground truth fractions of sample  $i$  and  $\hat{y}_i$  is the vector of  
519 predicted fractions of sample  $i$ . Each of the three ensemble models is trained  
520 independently for 5,000 steps. This ‘early stopping’ serves to avoid domain overfitting  
521 on the simulated tissue data, which would decrease the model performance on the  
522 real tissue RNA-seq data. We observed that training for more steps lead to an average  
523 performance decrease on real tissue RNA-seq data. To perform deconvolution with  
524 Scaden, a bulk RNA-seq sample is fed into a trained Scaden ensemble and three  
525 independent predictions for the cell type fractions of this sample are generated by the  
526 trained deep neural networks. These three predictions are then averaged per cell type  
527 to yield the final cell type composition for the input bulk RNA-seq sample:

$$528 \quad \hat{y}_c = \frac{\hat{y}_c^1 + \hat{y}_c^2 + \hat{y}_c^3}{3}$$

529 where  $\hat{y}_c$  is the final predicted fraction for cell type  $c$  and  $\hat{y}_c^i$  is the predicted fraction for  
530 cell type  $c$  of model  $i$ .

## 531 Scaden Overview

532 We used several performance measures to compare Scaden to four existing cell  
533 deconvolution algorithms, CIBERSORT with LM22 GEP (CS), CIBERSORTx (CSx),  
534 MuSiC and CPM. To compare the performance of the five deconvolution algorithms  
535 we measured the root mean squared error (RMSE), Lin's concordance correlation  
536 coefficient  $CCC$ , Pearson product moment correlation coefficient  $r$ , and  $R^2$  values  
537 comparing real and predicted cell fractions estimates. Additionally, to identify  
538 systematic prediction errors and biases, slope and intercept for the regression lines  
539 were calculated. These metrics are defined as follows:

$$540 \quad RMSE(y, \hat{y}) = \sqrt{avg(y - \hat{y})^2}$$

$$541 \quad r(y, \hat{y}) = \frac{cov(y, \hat{y})}{\sigma_y \sigma_{\hat{y}}}$$

$$542 \quad R^2(y, \hat{y}) = r(y, \hat{y})^2$$

$$543 \quad slope(y, \hat{y}) = \frac{\Delta y}{\Delta \hat{y}}$$

$$544 \quad CCC(y, \hat{y}) = \frac{2r\sigma_y\sigma_{\hat{y}}}{\sigma_y^2 + \sigma_{\hat{y}}^2 + (\mu_x - \mu_{\hat{y}})}$$

545 where  $y$  are the ground truth fractions,  $\hat{y}$  are the prediction fractions,  $\sigma_x$  is the standard  
546 deviation of  $x$ ,  $cov(y, \hat{y})$  is the covariance of  $y$  and  $\hat{y}$ , and  $\mu_y, \mu_{\hat{y}}$  are the mean of the  
547 predicted and ground truth fractions, respectively.

548 All metrics were calculated for all data points of a dataset, and separately for all data  
549 points of a specific cell type. For the latter approach, we then averaged the resulting  
550 values to recover single values. While in general the metrics calculated on all data  
551 points are sufficient, good performance on cell type-level is important if one is to  
552 compare fractions of a specific cell type between samples.

553 CIBERSORT (CS)

554 CS is a cell convolution algorithm based on specialized GEPs and support vector  
555 regression. Cell composition estimations were obtained using the CS web application  
556 (<https://cibersort.stanford.edu/>). For all deconvolutions with CS, we used the LM22  
557 GEP, which was generated by the CS authors from 22 leukocyte subsets profiled on  
558 the HGU133A microarray platform.

559 Because the LM22 GEP matrix contains cell types at a finer granularity than what was  
560 used for this study, predicted fractions of sub-cell types were added together. For cell  
561 grouping, we used the mapping of sub-cell types to broader types given by Figure 6  
562 from Monaco *et al.*<sup>10</sup>. We provide a table with the exact mappings used here in the  
563 supplementary material (Table S13). The deconvolution was performed using 500  
564 permutations with quantile normalization disabled for all datasets but GSE65133  
565 (Microarray), as is recommended for RNA-seq data. We used default settings for all  
566 other CS parameters.

567 CIBERSORTx (CSx)

568 CSx is a recent variant of CS that can generate GEP matrices from scRNA-seq data  
569 and use these for deconvolution. For additional deconvolution robustness, it applies  
570 batch normalization to the data. All signature matrices were created by uploading the  
571 labeled scRNA-seq expression matrices and using the default options. Quantile  
572 normalization was disabled. For deconvolution on simulated data, no batch  
573 normalization was used. For all bulk RNA-seq datasets, the S-Mode batch  
574 normalization was chosen. All PBMC datasets were deconvolved using a GEP matrix  
575 generated from the data6k dataset (for simulated samples from data6k, a donorA GEP  
576 matrix was chosen).

## 577 MuSiC

578 MuSiC is a deconvolution algorithm that uses multi-subject scRNA-seq datasets as  
579 GEP matrices in an attempt to include heterogeneity in the matrices to improve  
580 generalization. While MuSiC tries to address similar issues of previous deconvolution  
581 algorithms by using scRNA-seq data, the approach is very different. For  
582 deconvolution, MuSiC applies a sophisticated GEP-based deconvolution algorithm  
583 that uses weighted non-negative least squares regression with an iterative estimation  
584 procedure that imposes more weight on informative genes and less weight on non-  
585 informative genes.

586 The MuSiC R package contains functionality to generate the necessary GEP matrix  
587 given a scRNA-seq dataset and cell type labels. To generate MuSiC deconvolution  
588 predictions on PBMC datasets, we used the data8k scRNA-seq dataset as reference  
589 data for MuSiC and follow the tutorial provided by the authors to perform the  
590 deconvolution. For deconvolution of artificial samples generated from the data8k  
591 dataset, we provided MuSiC with the data6k dataset as reference instead.

592 MuSiC was developed with a focus on multi-subject scRNA-seq datasets, in which the  
593 algorithm tries to take advantage from the added heterogeneity that these datasets  
594 contain, by calculating a measure of cross-subject consistency for marker genes. To  
595 assess how Scaden performs on multi-subject datasets compared to MuSiC, we  
596 evaluated both methods on artificial bulk RNA-seq samples from human pancreas  
597 {Xin et al}. We used the 'bulk\_construct' function from MuSiC to combine the cells from  
598 all 18 subjects contained in the scRNA-seq dataset from Xin et al to generate artificial  
599 bulk samples for evaluation. Next, as a multi-subject reference dataset, we used the  
600 pancreas scRNA-seq dataset from Segerstolpe *et al.*<sup>19</sup>, which contains single-cell  
601 expression data from 10 different subjects, 4 of which with type-2 Diabetes. For

602 Scaden, the Segerstolpe scRNA-seq dataset was split by subjects, and training  
603 datasets were generated for each subject, yielding in total 10,000 samples. For  
604 MuSiC, a processed version of this dataset was downloaded from the resources  
605 provided by the MuSiC authors<sup>8</sup> and used as input reference dataset for the MuSiC  
606 deconvolution. Deconvolution was then performed according to the MuSiC tutorial,  
607 and performance compared according to the above-defined metrics.

#### 608 Cell Population Mapping (CPM)

609 CPM is a deconvolution algorithm that uses single-cell expression profiles to identify  
610 a so-called 'cell population map' from bulk RNA-seq data<sup>9</sup>. In CPM, the cell population  
611 map is defined as composition of cells over a cell-state space, where a cell-state is  
612 defined as a current phenotype of a single cell. Contrary to other deconvolution  
613 methods, CPM tries to estimate the abundance of all cell-states and types for a given  
614 bulk mixture, instead of only deconvolving the cell types. As input, CPM requires a  
615 scRNA-seq dataset and a low-dimensional embedding of all cells in this dataset, which  
616 represents the cell-state map. As CPM estimates abundances of both cell-states and  
617 types, it can be used for cell type deconvolution by summing up all estimated fractions  
618 for all cell-states of a given cell type - a method that is implemented in the scBio R  
619 package, which contains the CPM method. To perform deconvolution with CPM, we  
620 used the data6k PBMC scRNA-seq dataset as input reference for all PBMC samples.  
621 For samples simulated from the data6k dataset, we used the data8k dataset as  
622 reference. According to the CPM paper, a dimension reduction method can be used  
623 to obtain the cell-state space. We therefore used UMAP, a dimension reduction  
624 method widely used for scRNA-seq data, to generate the cell-state space mapping for  
625 the input scRNA-seq data. Deconvolution was then performed using the CPM function

626 of the scBio package with a scRNA-seq and accompanying UMAP embedding as

627 input.

628

## 629 List of abbreviations

630 RNA-seq : Next Generation RNA Sequencing

631 GEP : gene expression profile matrix

632 SVR : Support Vector Regression

633 DNN : Deep Neural Network

634 scRNA-seq : single cell RNA-seq

635 simulated tissue : training data generated by mixing proportions of scRNA-seq data

636 PBMC : peripheral blood mononuclear cells

637 CCC : concordance correlation coefficient

638  $r$  : Pearson's correlation coefficient

639 CS : CIBERSORT

640 CSx : CIBERSORTx

641 CPM : Cell Population Mapping

## 642 **Author information**

### 643 **Affiliations**

644 **German Center for Neurodegenerative Diseases Tuebingen, Germany**

645 Kevin Menden, Anupriya Dalmia, Peter Heutink, Stefan Bonn

646 **Institute of Medical Systems Biology, University Medical Center Hamburg-**

647 **Eppendorf, Germany**

648 Mohamed Marouf, Stefan Bonn

### 649 **Contributions**

650 KM and SB initiated the project. KM, PH, and SB designed the study, deep learning  
651 models, and analysis. KM and MM built the deep learning models. KM, MM, and AD  
652 analyzed the data. KM and SB wrote and MM, AD, and PH contributed to the  
653 manuscript writing.

654

655 **Competing interests**

656 The authors have no competing interests.

### 657 **Acknowledgements**

658 We would like to thank the people of the Genome Biology of Neurodegenerative  
659 Diseases group and Institute of Medical Systems Biology for helpful discussions and  
660 suggestions.

### 661 **Funding**



662 This study was supported in part by RiMod-FTD a EU Joint Programme -  
663 Neurodegenerative Disease Research (JPND) to PH, KM and SFB 1286/Z2 and  
664 BMBF Integrative Data Semantics for Neurodegenerative research (IDSN) to MM.

665 **Corresponding author**

666 Correspondence to Stefan Bonn (sbonn@uke.de) and Kevin Menden  
667 (kevin.menden@dzne.de).

668

## 669 References

- 670 1. Hrdlickova, R., Toloue, M. & Tian, B. RNA-Seq methods for transcriptome  
671 analysis. *Wiley Interdiscip. Rev. RNA* **8**, (2017).
- 672 2. Egeblad, M., Nakasone, E. S. & Werb, Z. Tumors as organs: Complex tissues  
673 that interface with the entire organism. *Dev. Cell* **18**, 884–901 (2010).
- 674 3. Kuhn, A., Thu, D., Waldvogel, H. J., Faull, R. L. M. & Luthi-Carter, R.  
675 Population-specific expression analysis (PSEA) reveals molecular changes in  
676 diseased brain. *Nat. Methods* **8**, 945–947 (2011).
- 677 4. Avila Cobos, F., Vandesompele, J., Mestdagh, P. & De Preter, K.  
678 Computational deconvolution of transcriptomics data from mixed cell  
679 populations. *Bioinformatics* **34**, 1969–1979 (2018).
- 680 5. Mohammadi, S., Zuckerman, N., Goldsmith, A. & Grama, A. A Critical Survey  
681 of Deconvolution Methods for Separating Cell Types in Complex Tissues.  
682 *Proc. IEEE* **105**, 340–366 (2017).
- 683 6. Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue  
684 expression profiles. *Nat. Methods* **12**, 453–457 (2015).
- 685 7. Newman, A. M. *et al.* Determining cell type abundance and expression from  
686 bulk tissues with digital cytometry. *Nat. Biotechnol.* (2019).
- 687 8. Wang, X., Park, J., Susztak, K., Zhang, N. R. & Li, M. Bulk tissue cell type  
688 deconvolution with multi-subject single-cell expression reference. *Nat.*  
689 *Commun.* **10**, 380 (2019).
- 690 9. Frishberg, A. *et al.* Cell composition analysis of bulk genomics using single-  
691 cell data. *Nat. Methods* **16**, (2019).
- 692 10. Monaco, G. *et al.* RNA-Seq Signatures Normalized by mRNA Abundance  
693 Allow Absolute Deconvolution of Human Immune Cell Types. *Cell Rep.* **26**,

- 694 1627–1640.e7 (2019).
- 695 11. Vallania, F. *et al.* Leveraging heterogeneity across multiple datasets increases  
696 cell-mixture deconvolution accuracy and reduces biological and technical  
697 biases. *Nat. Commun.* **9**, (2018).
- 698 12. Venet, D., Pecasse, F., Maenhaut, C. & Bersini, H. Separation of samples into  
699 their constituents using gene expression data. *Bioinformatics* **17**, 279–287  
700 (2001).
- 701 13. Shapiro, E., Biezuner, T. & Linnarsson, S. Single-cell sequencing-based  
702 technologies will revolutionize whole-organism science. *Nat. Rev. Genet.* **14**,  
703 618–630 (2013).
- 704 14. Schaum, N. *et al.* Single-cell transcriptomics of 20 mouse organs creates a  
705 Tabula Muris. *Nature* **562**, 367–372 (2018).
- 706 15. Kelley, K. W., Nakao-Inoue, H., Molofsky, A. V. & Oldham, M. C. Variation  
707 among intact tissue samples reveals the core transcriptional features of  
708 human CNS cell classes. *Nat. Neurosci.* **21**, 1171–1184 (2018).
- 709 16. Hicks, S. C., Townes, F. W., Teng, M. & Irizarry, R. A. Missing data and  
710 technical variability in single-cell RNA-sequencing experiments. *Biostatistics*  
711 **19**, 562–578 (2018).
- 712 17. Lin, L. I. A Concordance Correlation Coefficient to Evaluate Reproducibility  
713 *Biometrics* **45**, 255–268 (1989).
- 714 18. Xin, Y. *et al.* RNA Sequencing of Single Human Islet Cells Reveals Type 2  
715 Diabetes Genes. *Cell Metab.* **24**, 608–615 (2016).
- 716 19. Segerstolpe, Å. *et al.* Single-Cell Transcriptome Profiling of Human Pancreatic  
717 Islets in Health and Type 2 Diabetes. *Cell Metab.* **24**, 593–607 (2016).
- 718 20. Zimmermann, M. T. *et al.* System-wide associations between DNA-

- 719 methylation, gene expression, and humoral immune response to influenza  
720 vaccination. *PLoS One* **11**, 1–21 (2016).
- 721 21. Schelker, M. *et al.* Estimation of immune cell content in tumour tissue using  
722 single-cell RNA-seq data. *Nat. Commun.* **8**, 2032 (2017).
- 723 22. Bennett, D. A. *et al.* Religious Orders Study and Rush Memory and Aging  
724 Project. *J. Alzheimer's Dis.* **64**, S161–S189 (2018).
- 725 23. Braak, H. & Braak, E. Neuropathological staging of Alzheimer-related  
726 changes. *Acta Neuropathol.* **82**, 239–59 (1991).
- 727 24. Yosinski, J., Clune, J., Nguyen, A., Fuchs, T. & Lipson, H. Understanding  
728 Neural Networks Through Deep Visualization. (2015).
- 729 25. Athiwaratkun, B., Finzi, M., Izmailov, P. & Wilson, A. G. Improving  
730 Consistency-Based Semi-Supervised Learning with Weight Averaging. *Jmlr*  
731 **17**, 1–35 (2018).
- 732 26. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: Large-scale single-cell gene  
733 expression data analysis. *Genome Biol.* **19**, 1–5 (2018).
- 734 27. Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial  
735 reconstruction of single-cell gene expression data. *Nat. Biotechnol.* **33**, 495–  
736 502 (2015).
- 737 28. Zeisel, A. *et al.* Molecular Architecture of the Mouse Nervous System. *Cell*  
738 **174**, 999–1014.e22 (2018).
- 739 29. Love, M. I., Soneson, C., Patro, R., Vitting-seerup, K. & Oshlack, A.  
740 Swimming downstream : statistical analysis of differential transcript usage  
741 following Salmon quantification. *F1000* 1–50 (2019).
- 742
- 743

744 **Supplementary Figures & Tables**

Tissue	Name	# cells	# Subjects	Source
PBMC	data6k	5,419	1	10X Genomics
PBMC	data8k	8,381	1	10X Genomics
PBMC	donorA	2,900	1	10X Genomics
PBMC	donorC	9,519	1	10X Genomics
Mouse Brain	Tasic	1,679	1	Tasic et al., Nat. Neurosci., 2016
Mouse Brain	Zeisel	3,005	1	Zeisel et al., Science, 2015
Mouse Brain	Romanov	2,881	1	Romanov et al., Nat. Neurosci., 2018
Mouse Brain	Campbell	21,086	1	Campbell et al, Nat. Neurosci., 2017
Mouse Brain	Chen	14,437	1	Chen et al., Cell Rep., 2017
Pancreas	Segerstolpe	3,514	10	Segerstolpe et al., Cell Metab., 2016
Pancreas	Baron	8,569	4	Baron et al., Cell Syst., 2016
Ascites	Ascites	3,114	3	Schelker et al, Nat. Comm., 2018

745 **Table S1** *scRNA-seq datasets used for the generation of simulated tissues for Scaden*  
746 *training.*

747

748

Tissue	# Samples	# Datasets	Size
PBMC	32,000	4	1.2 GB
Pancreas	14,000	2	0.6 GB
Mouse Brain	30,000	5	1.5 GB
Ascites	6,000	1	0.38 GB

749 **Table S2** *Number of samples, datasets, and size of the simulated training data.*

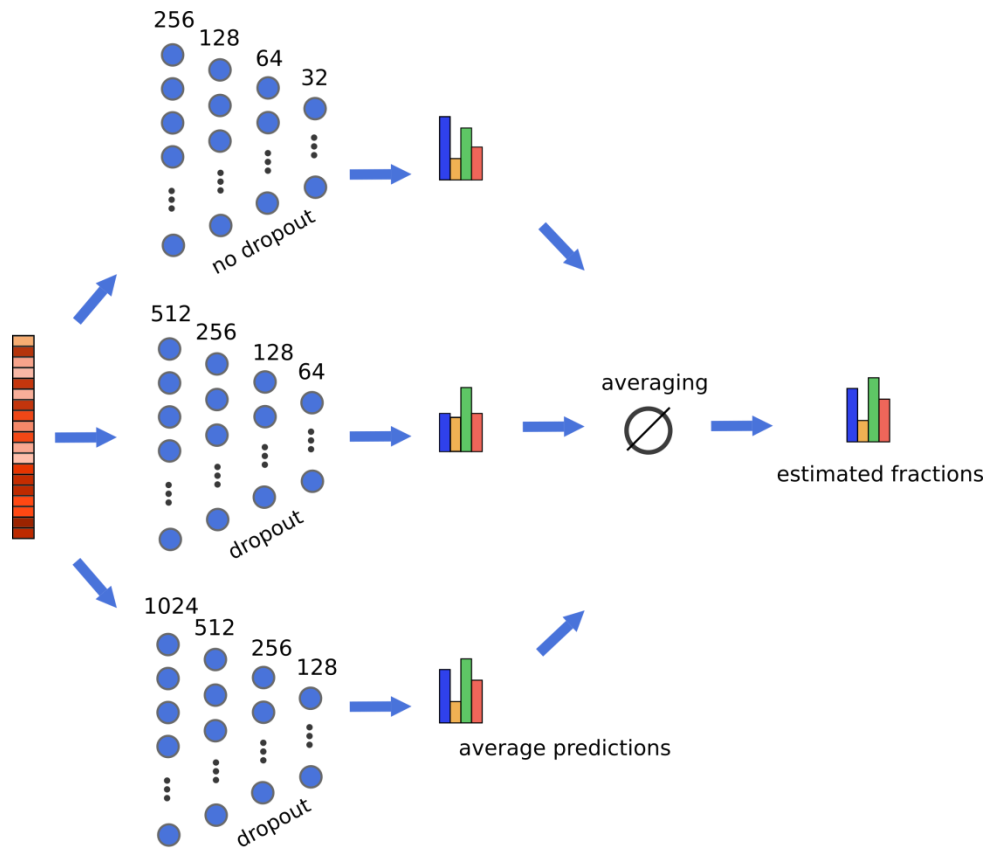
750

751

Parameter	Values tested
Batch size	32, 64, 128, 256, 512
# Layers	2, 3, 4
Layer sizes	2048, 1024, 512, 256, 128, 64, 32, 16
Dropout rate	[0, 0.8]
Loss function	L1, L2

752 **Table S3** *Hyperparameters used for model optimization.*

753



754

755 **Figure S1** Overview of the Scaden neural network ensemble model. A bulk RNA-seq  
756 sample is the input to three separate deep neural networks with varying layer sizes  
757 and dropout regularization. The predictions of all three models are subsequently  
758 averaged to obtain the final Scaden predictions. During training, predictions are not  
759 averaged and each model is trained separately.

760



Model	# Layers	Layer sizes	Dropout rates
M256	4	256, 128, 64, 32	0, 0, 0, 0
M512	4	512, 256, 128, 64	0, 0.3, 0.2., 0.1
M1024	4	1024, 512, 256, 128	0, 0.6, 0.3, 0.1

761 **Table S4** Architectures of deep neural network models used in Scaden ensemble.

762

763

Method	DS	RMSE	Slope	Correlation	Intercept	CCC
<b>CPM</b>	data6k	0.189	-0.011	-0.05	0.168	-0.021
<b>CPM</b>	data8k	0.183	-0.002	-0.005	0.167	-0.003
<b>CPM</b>	donorA	0.215	-0.008	-0.046	0.168	-0.016
<b>CPM</b>	donorC	0.174	-0.015	-0.062	0.168	-0.026
<b>CS</b>	data6k	0.131	0.664	0.728	0.074	0.717
<b>CS</b>	data8k	0.121	0.567	0.714	0.087	0.685
<b>CS</b>	donorA	0.131	0.708	0.788	0.049	0.784
<b>CS</b>	donorC	0.131	0.742	0.719	0.079	0.688
<b>CSx</b>	data6k	0.1	0.83	0.854	0.05	0.844
<b>CSx</b>	data8k	0.134	0.601	0.662	0.082	0.651
<b>CSx</b>	donorA	0.045	0.9	0.979	0.017	0.975
<b>CSx</b>	donorC	0.11	0.819	0.811	0.07	0.776
<b>MuSiC</b>	data6k	0.126	0.768	0.768	0.059	0.759
<b>MuSiC</b>	data8k	0.135	0.597	0.659	0.083	0.648
<b>MuSiC</b>	donorA	0.13	0.771	0.802	0.038	0.801
<b>MuSiC</b>	donorC	0.132	0.815	0.737	0.07	0.704
<b>Scaden</b>	data6k	0.077	0.855	0.917	0.047	0.904
<b>Scaden</b>	data8k	0.132	0.525	0.651	0.093	0.627
<b>Scaden</b>	donorA	0.034	0.924	0.989	0.013	0.986
<b>Scaden</b>	donorC	0.109	0.848	0.821	0.067	0.786

764 **Table S5** Deconvolution evaluation on simulated PBMC data.

765

766

Method	Celltype	RMSE	Correlation	Slope	Intercept	CCC
<b>CSx</b>	ALPHA	0.282	0.816	0.691	0.431	0.375
<b>CSx</b>	BETA	0.309	0.833	0.175	-0.017	0.078
<b>CSx</b>	DELTA	0.04	0.812	1.567	-0.013	0.647
<b>CSx</b>	GAMMA	0.052	0.921	1.131	0.0	0.897
<b>CSx</b>	Total	0.212	0.79	1.113	-0.028	0.746
<b>MuSiC</b>	ALPHA	0.11	0.887	1.108	-0.042	0.863
<b>MuSiC</b>	BETA	0.148	0.752	1.067	0.017	0.694
<b>MuSiC</b>	DELTA	0.023	0.817	0.716	-0.003	0.707
<b>MuSiC</b>	GAMMA	0.068	0.881	0.552	-0.003	0.711
<b>MuSiC</b>	Total	0.099	0.938	1.078	-0.019	0.929
<b>Scaden</b>	ALPHA	0.067	0.949	1.071	-0.034	0.942
<b>Scaden</b>	BETA	0.07	0.936	1.152	-0.045	0.916
<b>Scaden</b>	DELTA	0.024	0.807	1.012	0.008	0.764
<b>Scaden</b>	GAMMA	0.045	0.914	0.89	-0.008	0.901
<b>Scaden</b>	Total	0.055	0.978	1.033	-0.008	0.976

767 **Table S6** Deconvolution performance on simulated pancreas data from Xin et al..

768

769

770

Tissue	Name	# Samples	Reference
PBMC	PBMC1	12	Zimmermann et al., PLOS one, 2016
PBMC	PBMC2	12	Monaco et al., Cell Reports, 2019
Pancreas	Xin	18	Xin et al., Cell Metab., 2016
Human Brain	ROSMAP	390	Bennett et al., Curr Alzheimer Res., 2012
Ascites	Ascites	3	Schelker at al., Nat. Comm. 2018

771 **Table S7** *Tissue RNA-seq datasets used for performance evaluation.*

772

773

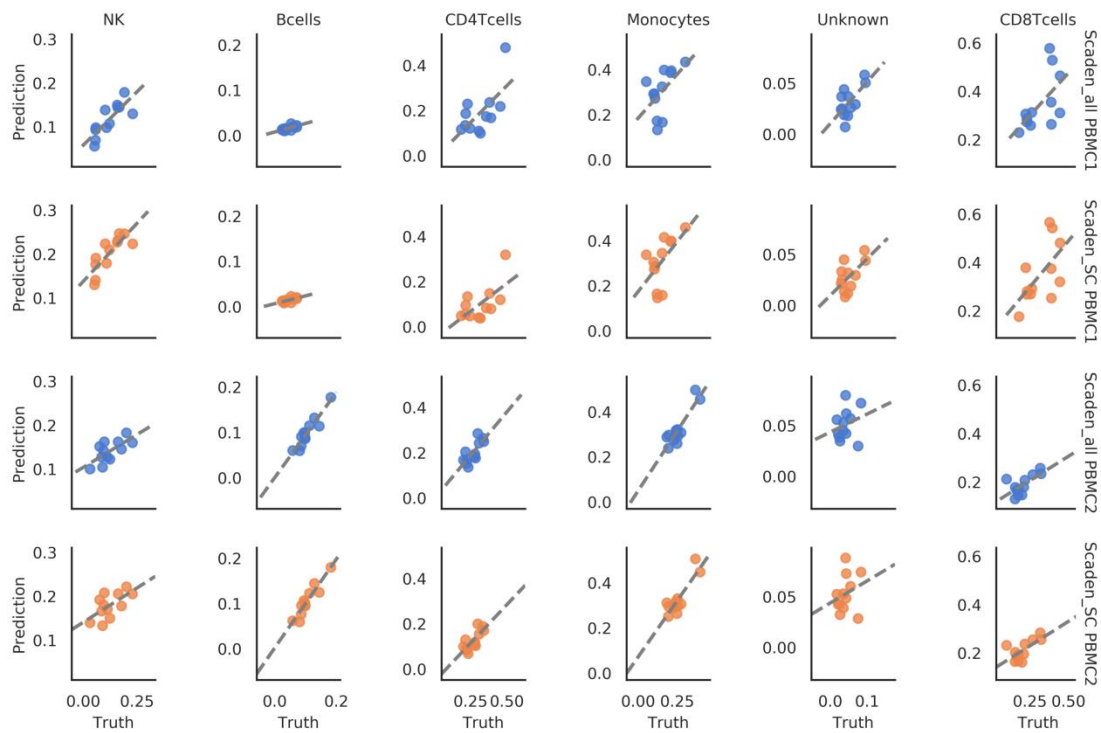
774

Method	Dataset	Celltype	RMSE	Correlation	Slope	Intercept	CCC
CPM	PBMC1	Total	0.18	-0.003	-0.003	0.167	-0.003
CPM	PBMC2	Total	0.114	-0.203	-0.094	0.182	-0.155
CS	PBMC1	Total	0.147	0.437	0.491	0.085	0.434
CS	PBMC2	Total	0.101	0.594	0.754	0.041	0.577
CSx	PBMC1	Total	0.16	0.603	0.925	0.012	0.552
CSx	PBMC2	Total	0.13	0.456	0.67	0.055	0.424
MuSiC	PBMC1	Total	0.316	-0.235	-0.468	0.245	-0.189
MuSiC	PBMC2	Total	0.299	-0.197	-0.542	0.257	-0.127
Scaden	PBMC1	Total	0.104	0.722	0.805	0.032	0.717
Scaden	PBMC2	Total	0.052	0.855	0.848	0.025	0.855

775 **Table S8** Deconvolution performance on real PBMC RNA-seq datasets PBMC1 and

776 PBMC2.

777



778

779 **Figure S2** Comparison of Scaden deconvolution results on PBMC1 and PBMC2  
780 datasets with and without (Scaden\_all, Scaden\_SC, respectively) bulk RNA-seq  
781 samples included in training data.

782

783

Method	Dataset	Celltype	RMSE	Correlation	Slope	Intercept	CCC
Scaden_SC	PBMC1	Total	0.131	0.564	0.644	0.059	0.559
Scaden_SC	PBMC2	Total	0.077	0.684	0.689	0.052	0.684
Scaden_all	PBMC1	Total	0.104	0.722	0.805	0.032	0.717
Scaden_all	PBMC2	Total	0.052	0.855	0.848	0.025	0.855
Scaden_SC	PBMC1	Bcells	0.033	0.648	0.172	0.006	0.083
Scaden_SC	PBMC1	CD4Tcells	0.228	0.633	0.492	-0.055	0.149
Scaden_SC	PBMC1	CD8Tcells	0.101	0.603	0.761	0.108	0.562
Scaden_SC	PBMC1	Monocytes	0.178	0.556	0.885	0.173	0.186
Scaden_SC	PBMC1	NK	0.087	0.81	0.531	0.137	0.312
Scaden_SC	PBMC1	Unknown	0.029	0.577	0.361	0.009	0.287
Scaden_SC	PBMC2	Bcells	0.012	0.936	0.977	0.002	0.935
Scaden_SC	PBMC2	CD4Tcells	0.145	0.767	0.682	-0.057	0.119
Scaden_SC	PBMC2	CD8Tcells	0.049	0.67	0.403	0.129	0.587
Scaden_SC	PBMC2	Monocytes	0.078	0.865	0.994	0.071	0.558
Scaden_SC	PBMC2	NK	0.071	0.629	0.314	0.14	0.276
Scaden_SC	PBMC2	Unknown	0.025	0.247	0.217	0.044	0.209
Scaden_all	PBMC1	Bcells	0.031	0.668	0.188	0.007	0.1
Scaden_all	PBMC1	CD4Tcells	0.151	0.638	0.652	-0.017	0.345
Scaden_all	PBMC1	CD8Tcells	0.096	0.6	0.704	0.123	0.569
Scaden_all	PBMC1	Monocytes	0.172	0.518	0.777	0.184	0.177
Scaden_all	PBMC1	NK	0.036	0.804	0.488	0.058	0.71
Scaden_all	PBMC1	Unknown	0.026	0.64	0.41	0.01	0.365
Scaden_all	PBMC2	Bcells	0.013	0.936	0.94	0.0	0.917
Scaden_all	PBMC2	CD4Tcells	0.074	0.772	0.769	-0.005	0.373
Scaden_all	PBMC2	CD8Tcells	0.051	0.672	0.398	0.106	0.562
Scaden_all	PBMC2	Monocytes	0.072	0.895	1.058	0.049	0.614
Scaden_all	PBMC2	NK	0.045	0.69	0.301	0.103	0.467
Scaden_all	PBMC2	Unknown	0.023	0.241	0.178	0.043	0.203

784 **Table S9** Deconvolution performance on real PBMC RNA-seq data for Scaden models trained  
785 only on scRNA-seq simulated tissues (Scaden\_SC) or on a mix of simulated and real tissue  
786 data (Scaden\_all).

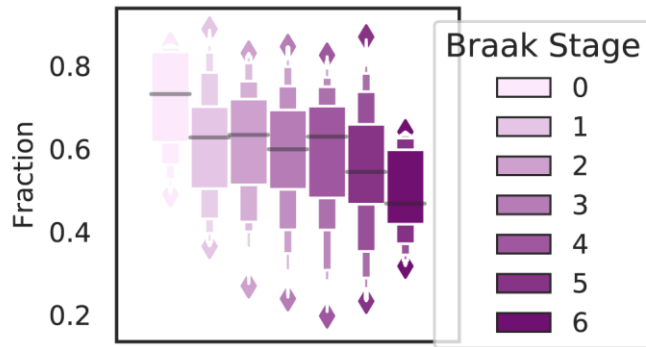
787

Method	Type	CCC	Correlation	Intercept	RMSE	Slope
<b>CPM</b>	Total	-0.0	0.004	0.153	0.183	-0.0
<b>CSx</b>	Total	0.938	0.952	0.002	0.069	1.115
<b>MuSiC</b>	Total	0.876	0.907	0.033	0.079	0.696
<b>Scaden</b>	Total	0.948	0.955	-0.030	0.061	1.066

788 **Table S10** *Deconvolution performance on real Ascites RNA-seq data.*

789

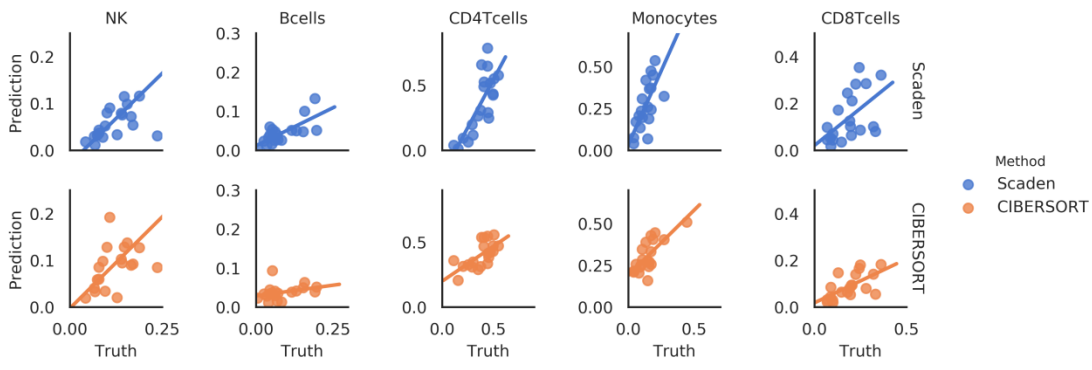




790

791 **Figure S3** *Deconvolution performance on real human brain RNA-seq data.* Scaden was  
792 trained on mouse scRNA-seq data and the trained model was used to deconvolve cell  
793 fractions of ROSMAP human brain RNA-seq data. This data does not contain cell fraction  
794 ground-truth information. Instead, the box plot shows the decrease of neuronal cell fractions  
795 with increasing Braak disease stage, a well-known phenomenon in AD.

796



797

798 **Figure S4** Deconvolution performance comparison of CS (LM22) and Scaden on the

799 GSE65133 PBMC microarray dataset.

800

801

Method	Celltype	CCC	Correlation	Intercept	RMSE	Slope
CS	Bcells	0.122	0.33	0.029	0.068	0.109
CS	CD4Tcells	0.629	0.658	0.199	0.095	0.537
CS	CD8Tcells	0.285	0.635	0.018	0.12	0.375
CS	Monocytes	0.295	0.741	0.19	0.17	0.779
CS	NK	0.623	0.698	-0.003	0.059	0.78
CS	Total	0.717	0.728	0.026	0.11	0.869
Scaden	Bcells	0.431	0.728	0.012	0.055	0.388
Scaden	CD4Tcells	0.64	0.778	-0.195	0.153	1.474
Scaden	CD8Tcells	0.474	0.543	0.02	0.104	0.635
Scaden	Monocytes	0.43	0.838	0.033	0.191	1.764
Scaden	NK	0.516	0.741	-0.029	0.074	0.77
Scaden	Total	0.705	0.749	-0.015	0.126	1.067

802 **Table S11** *Deconvolution performance on real PBMC microarray data.*

803

804

Software	Version
pandas	0.23.4
Python	3.6.8
Tensorflow	1.10.0
matplotlib	2.2.3
nb_conda	2.2.1
numpy	1.15.0
scipy	1.1.0
seaborn	0.9.0
anndata	0.6.9
scanpy	1.2.2
scikit-learn	0.20.0
ipython	6.5.0
python-igraph	0.7.1.post6
louvain	0.6.1
tqdm	4.7.2
igraph	0.7.1

805 **Table S12** *Software packages and versions used.*

806

807

Target Cell Type	LM22 Cell Types
B cells	B cells naive, B cells memory
CD8 T cells	T cells CD8, T cells follicular helper, T cells gamma delta
CD4 T cells	T cells CD4 naive, T cells regulatory (Tregs), T cells CD4 memory resting, T cells CD4 memory activated
NK	NK cells resting, NK cells activated
Dendritic	Dendritic cells resting, Dendritic cells activated
Monocytes	Monocytes, Macrophages M0, Macrophages M1, Macrophages M2
Unknown	Mast cells resting, Mast cells activated, Eosinophils, T cells follicular helper, T cells gamma delta, Plasma cells, Neutrophils, Dendritic

808 **Table S13** Mapping of the LM22 GEP to cell types.

809