# 1    Title

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

3    Kevin Menden[$], Mohamed Marouf, Anupriya Dalmia, Peter Heutink, Stefan Bonn[$]

4

5    [$] 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.

## Introduction

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

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

In theory, an optimal GEP should contain a set of genes that are predominantly expressed within each cell population of a complex sample[12]. They should be stably

2

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[11]. 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[10].

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[7–9,13–15]. 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[16]. While this complicates the use

70    of scRNA-seq data for GEP design[8], 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
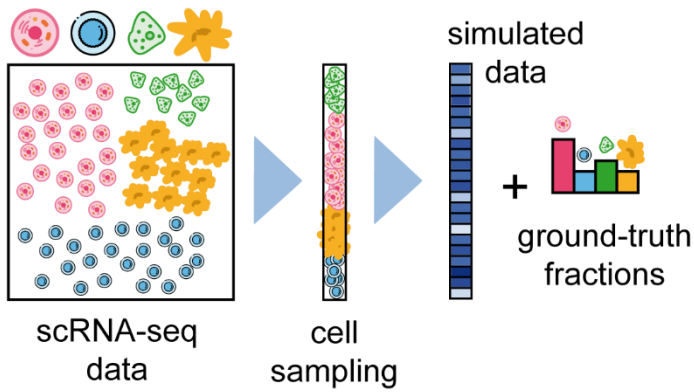

81    The basic architecture of Scaden is a DNN that takes gene counts of RNA-seq data

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.
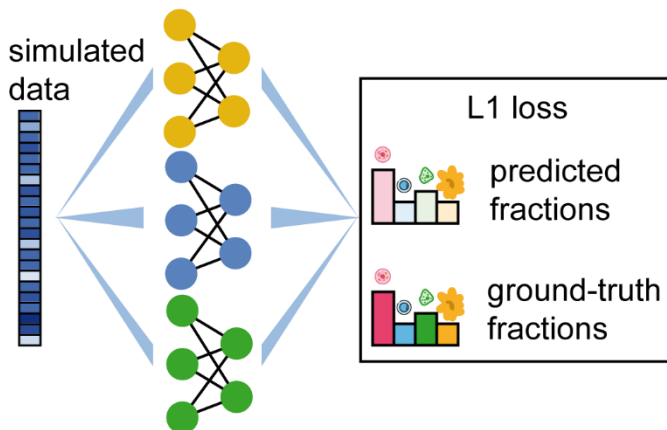
91    To find the optimal DNN architecture for cell deconvolution, we performed leave-one-

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

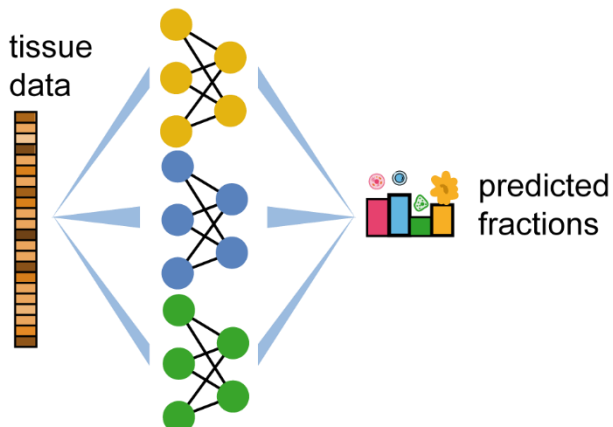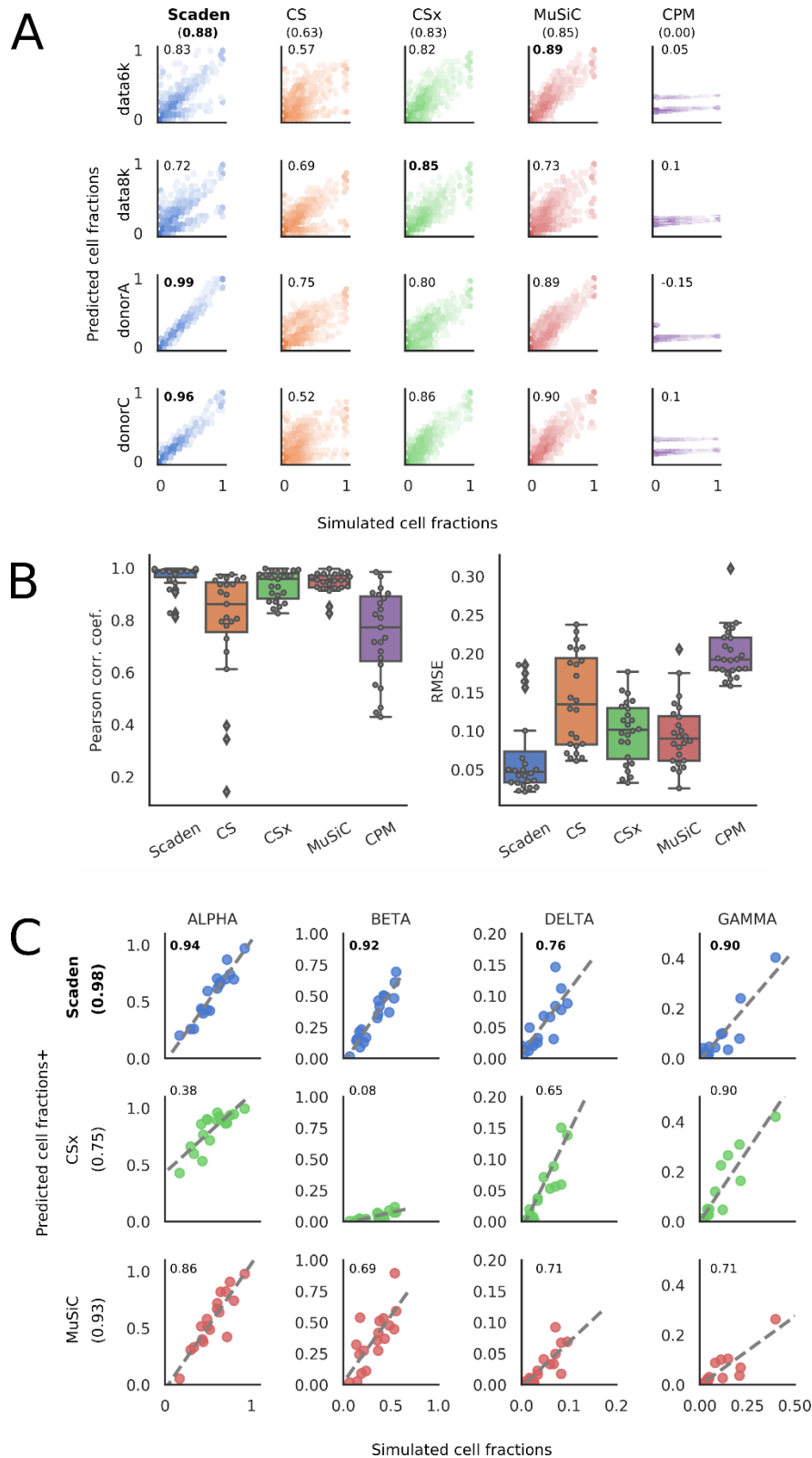97    ensemble models (Fig. S1, Table S4).

98



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)[17],

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  samples 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)[6], CIBERSORTx (CSx)[7], MuSiC[8], and Cell Population Mapping

114  (CPM)[9]. 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[6] and tested performance on

125  the 500 simulated PBMC samples from a held-out scRNA-seq dataset (e.g. donorC).

6

**Figure 2** *Deconvolution performance on simulated tissue data* A: Ground truth values (x-axis) plotted against cell type fraction estimates (y-axis) for predictions made on simulated data from four PBMC scRNA-seq datasets. Darker color in a hexbin corresponds to more data 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[18]. Numbers inside the

134    plotting area signify CCC values.

135

136    For two of four test datasets (donorA, donorC), Scaden obtained the highest CCC and

137    lowest RMSE, followed by CSx, MuSiC, CS, and CPM (Fig. 2A, Table S5). CSx and

138    MuSiC obtain the highest CCC values for the data8k and data6k datasets,

139    respectively. Overall, Scaden obtains the highest CCC and lowest RMSE (0.88, 0.08,

140    respectively), followed by MuSiC(0.85, 0.10), CSx(0.83, 0.11), CS (0.63 0.15), and

141    CPM (0, 0.20) (Fig. 2A). As expected, all algorithms that use scRNA-seq data as

142    reference perform good in this scenario with the notable exception of CPM. We want

143    to mention that CPM was not primarily developed for cell deconvolution, but merely

144    incorporates this as an additional feature. On average, Scaden also obtained the

145    highest correlation and the best intercept and slope values on simulated PBMC data

146    (Table S5).

147    A specific feature of the MuSiC algorithm is that it preferentially weighs genes

148    according to low inter-subject and intra-cell cluster variability for its GEP, which

149    increases deconvolution robustness when high expression heterogeneity is observed

150    between human subjects, for example[8]. To understand if Scaden can utilize multi-

151    subject information to increase its deconvolution performance, we trained Scaden,

152    CSx, and MuSiC on scRNA-seq pancreas data from several subjects[19] and assessed

153    the performance on a separate simulated pancreas RNA-seq dataset[18] (Fig. 2C, Table

154    S6). To allow for direct comparison, we chose the same pancreas training and test

155    datasets that were used in the original MuSiC publication (Table S1). To enable

156    Scaden to leverage the heterogeneity of multi-subject data, training data was
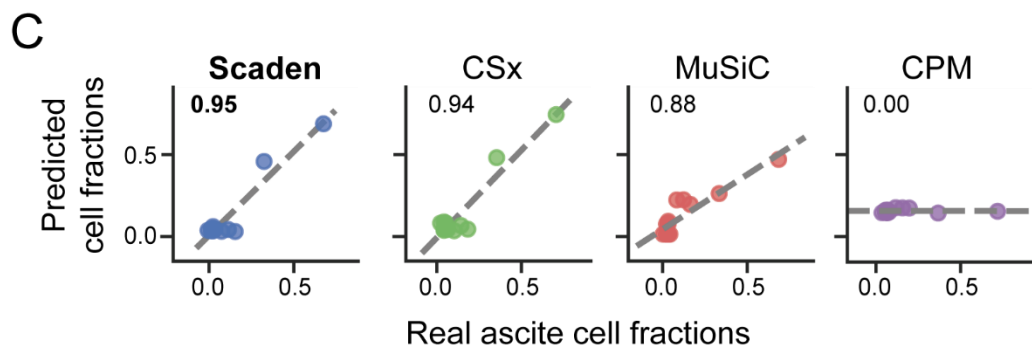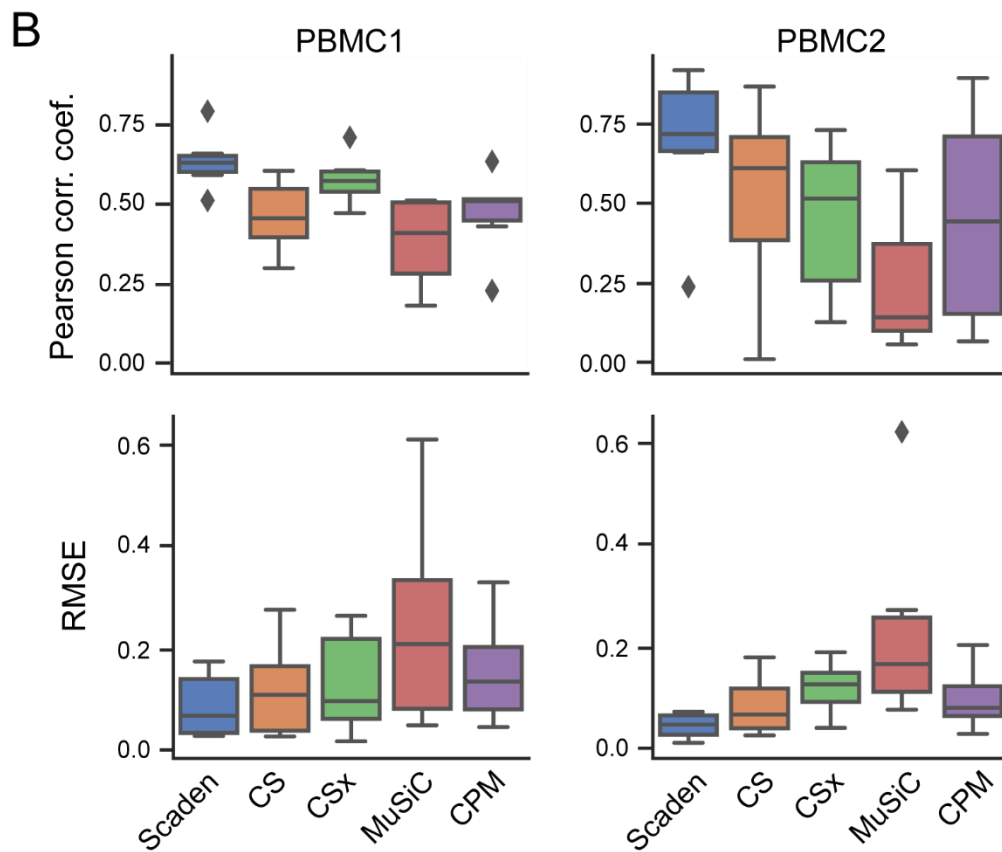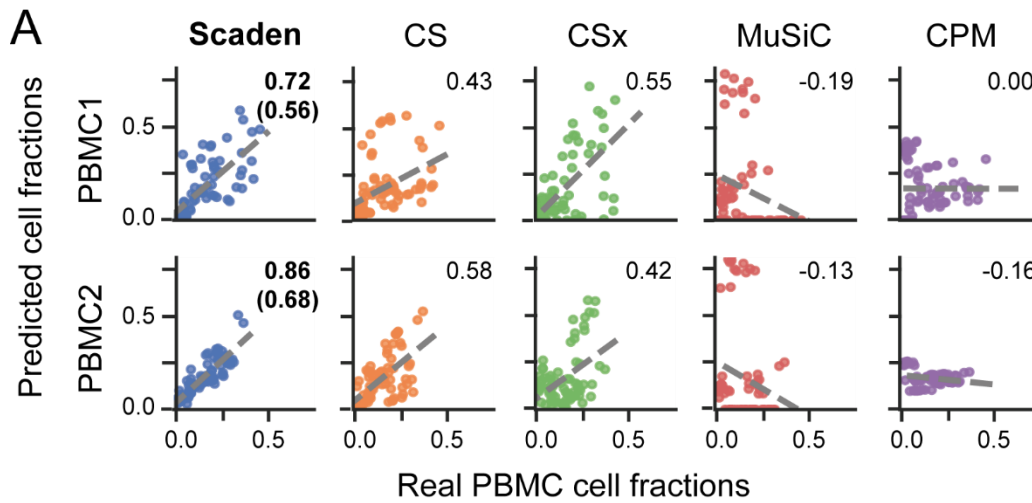
8

157  generated separately for every subject in the dataset (see Methods). CSx cannot profit

158  from multi-subject data but performed well on the artificial PBMC datasets and was

159  therefore included in the comparison. The best performance is achieved by Scaden

160  (CCC = 0.98), closely followed by MuSiC (CCC = 0.93), while CSx does not perform

161  as well (CCC = 0.75) (Fig. 2C, Table S6). This provides strong evidence that Scaden,

162  by separating training data generation for each subject, can learn inter-subject

163  heterogeneity and outperform specialized multi-subject algorithms such as MuSiC on

164  the cell-type deconvolution task.

165  Additionally, we wanted to test how the best performing deconvolution algorithms

166  Scaden, MuSiC, and CSx behave when unknown cell content is part of the mixture.

167  To test this, all cells falling into the 'Unknown' category were removed from the training

168  or reference datasets but added to the simulated mixture samples at fixed percentages

169  (5%, 10%, 20%, 30%) (see Methods). Scaden obtains the highest CCC for all tested

170  percentages of unknown cell content (Fig. S2, Table S8). The general deconvolution

171  performance declines linearly with increasing percentage of unknown content for all

172  tested algorithms (Fig. S2, Table S8), indicating that Scaden, MuSiC, and CSx have a

173  similar robustness against unknown mixture content.


174  Robust deconvolution of bulk expression data

175  The true use case of cell deconvolution algorithms is the cell fraction estimation of

176  tissue RNA-seq data. We therefore assessed the performance of Scaden, CS, CSx,

177  MuSiC, and CPM to deconvolve two publicly available human PBMC bulk RNA-seq

178  datasets, for which ground-truth cell composition information was measured using flow

179  cytometry (Fig. 3A, Tables S7 & S9). We will refer to these datasets that consists of

180  12 samples each as PBMC1 [20] and PBMC2 [10]. Deconvolution for all methods was

181    performed as described in the previous section, with the difference that data from all

182    four PBMC scRNA-seq datasets was now deployed for Scaden training.



183

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

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

186   PBMC1 and PBMC2 cell fractions. Numbers inside the plotting area signify CCC values. For

187   Scaden, the CCC using only scRNA-seq training data (in parenthesis) and the CCC using

188   mixed scRNA-seq and RNA-seq training data is shown. B: Boxplots of r (first row) and RMSE

189   (second row) values for real PBMC1 (first column) and PBMC2 (second column) data. C: Per-

190   cell-type scatterplots of ground truth (x-axis) and predicted values (y-axis) for Scaden, CSx,

191   MuSiC, and CPM on real ascite cell fractions. Numbers inside the plotting area signify CCC

192   values.

193

194   On the PBMC1 dataset, Scaden obtained the highest CCC and lowest RMSE (0.56,

195   0.13), while CSx (0.55, 0.16) and CS (0.43, 0.15) performed well yet significantly worse

196   than Scaden (Fig. 3A, Tables S8 & S9). CPM (0, 0.18) and MuSiC (-0.19, 0.32) both

197   failed to deconvolve the cell fractions of the PBMC1 data. Scaden also obtained the

198   best CCC and RMSE (0.68, 0.08) on the PBMC2 dataset, while CS (0.58, 0.10) and

199   CSx (0.42, 0.13) obtained good deconvolution results. Similar to the PBMC1 data

200   deconvolution results, CPM (-0.16, 0.11) as well as MuSiC (-0.13, 0.30) did not

201   perform well on the PBMC2 deconvolution task. In addition to CCC and RMSE metrics,

202   Scaden achieves the best correlation, intercept and slope on both PBMC datasets

203   (Tables S9 & S10).

204   An additional algorithmic feature of Scaden is that it seamlessly integrates increasing

205   amounts of training data, which can be of different types, such as a combination of

206   simulated tissue and real tissue data with cell fraction information. In theory, even

207   limited real tissue training data could make Scaden robust to data type bias and

208   consequently improve Scaden's deconvolution performance on real tissue data. We

209   therefore trained Scaden on a mix of simulated PBMC (500 samples) and real PBMC2

210   (12 samples) data and evaluated its performance on real PBMC1 data (Fig. 3A, S3,

211    Table S9). While the training contained only ~2% real data, Scaden's CCC increased

212    from 0.56 to 0.72 and the RMSE decreased from 0.13 to 0.10. We observed similar

213    performance increases when Scaden was trained on simulated PBMC and real

214    PBMC1 data and evaluated on real PBMC2 data (Fig. 3A, S3, Table S10).

215    We next evaluated Scaden's performance on real ascites RNA-seq data, for which

216    scRNA-seq and FACS cell proportion data is available[21] (Table S7). It is noteworthy

217    that RNA-seq, scRNA-seq, and FACS data was generated for the same samples,

218    which potentially entails reduced experimental and technical bias and consequently

219    higher deconvolution fidelity for the ascites data as compared to the PBMC data. We

220    did not evaluate CS's performance on the ascites data as there was no optimized

221    ascites GEP available. 'For Scaden, CSx, CPM and MuSiC we used scRNA-seq data

222    to generate either simulated tissue data for training (Scaden) or a reference GEP (CSx,

223    CPM, MuSiC). Scaden, CSx, CPM, and MuSiC all accurately predict the cell type

224    compositions for the three real ascites samples, while CPM does not perform well (Fig.

225    3C, Table S11). The highest CCC and lowest RMSE were achieved by Scaden (0.95,

226    0.06), followed by CSx (0.94, 0.07), MuSiC (0.88, 0.08), and CPM (0, 0.18). This

227    further validates that Scaden reliably deconvolves tissue RNA-seq data into the

228    constituent cell fractions and that very accurate deconvolution results can be obtained

229    if reference and target datasets are from the same experiment. Again, we stress that

230    CPM was not primarily developed for cell deconvolution, but mainly for a different

231    functionality.

232    We next wanted to assess if Scaden's deconvolution performance is robust across

233    species. We therefore tested whether a Scaden model trained on mouse brain scRNA-

234    seq data could generate reasonable cell composition estimations for real human brain

235    RNA-seq data (Table S7). To this end, Scaden was trained on artificial data generated

236   from five mouse brain scRNA-seq datasets and predicted the cell fractions on human

237   post-mortem RNA-seq brain samples (390 prefrontal cortex samples) from the

238   ROSMAP study[22]. Ground-truth cell fractions were not available for this data, which is

239   why we used Braak stages[23] that correspond to Alzheimer's disease severity and

240   correlate with the degree of neuronal loss. Overall, Scaden's cell fraction predictions

241   capture the increased neuronal loss with increasing Braak stage (Fig. S4).

242   Interestingly, the largest drop in neural percentage is observed at stage 5, when the

243   neurodegeneration typically reaches the prefrontal cortex of the brain. By learning

244   robust features, Scaden reliably deconvolves RNA-seq data in a cross-species

245   comparison.

246   Given the robustness with which Scaden predicts tissue RNA-seq cell fractions using

247   scRNA-seq training data, even across species, we next wanted to investigate if a

248   scRNA-seq-trained Scaden model can also deconvolve other data types. To this end,

249   we measured the deconvolution performance on a bulk PBMC microarray dataset (20

250   samples)[6] of a Scaden model trained on scRNA-seq and RNA-seq PBMC data (see

251   above). We compared Scaden to CS using the microarray-derived LM22 matrix. CS

252   achieved a slightly higher CCC and slightly lower total RMSE (0.72, 0.11) than Scaden

253   (0.71, 0.13), while Scaden obtained the highest average CCC (0.50) compared to CS

254   (0.39) (Fig. S5, Table S12). Notably in this scenario, Scaden was trained entirely on

255   simulated data and RNA-seq data, while CS's LM22 GEP was optimized on PBMC

256   microarray data.

257   Overall, we provide strong evidence that Scaden robustly deconvolves tissue data

258   across tissues, species, and even data types.

## Discussion

259

260    Scaden is the first deep learning-based cell deconvolution algorithm. In many

261    instances, it compares favorably in both prediction robustness and accuracy to existing

262    deconvolution algorithms that rely on GEP design and linear regression. We believe

263    that Scaden's performance relies to a large degree on the inherent feature engineering

264    of the DNN. The network does not only select features (genes) for regression, it also

265    creates novel features that are optimal for the regression task in the nodes of the

266    hidden layers. These hidden features are non-linear combinations of the input features

267    (gene expression), which makes it notoriously difficult to explain how a DNN works[24].

268    It is important to highlight that this feature creation is fundamentally different from all

269    other existing cell deconvolution algorithms, which rely on heuristics that select a

270    defined subset of genes as features for linear regression.

271    Another advantage of this inherent feature engineering is that Scaden can be trained

272    to be robust to input noise and bias (e.g. batch effects). Noise and bias are all prevalent

273    in experimental data, due to different sample quality, sample processing,

274    experimenters, and instrumentation, for example. If the network is trained on different

275    datasets of the same tissue, however, it learns to create hidden features that are

276    robust to noise and bias, such as batch effects. This robustness is pivotal in real world

277    cell deconvolution use cases, where the bulk RNA data for deconvolution and the

278    training data (and therefore the network and GEP) contain different noise and biases.

279    While especially recent cell deconvolution algorithms include batch correction

280    heuristics prior to GEP construction, Scaden optimizes its hidden features

281    automatically when trained on data from various batches.

282    The robustness to noise and bias, which might be due to hidden feature generation, is

283    especially evident in Scaden's ability to deconvolve across data types. A network

14

284    trained on *in silico* bulk RNA-seq data can seamlessly deconvolve microarray data of

285    the same tissue. This is quite noteworthy, as microarray data is known to have a

286    reduced dynamic range and several hybridization-based biases compared to RNA-

287    seq data. In other words, Scaden can deconvolve bulk data of types it has never been

288    trained on, even in the face of strong data type bias. This raises the possibility that

289    Scaden trained on scRNA-seq data might reliably deconvolve other bulk omics data

290    as well, such as proteomic and metabolomic data. This assumption is strengthened

291    by the fact that Scaden, trained on scRNA-seq data, attains state-of-the-art

292    performance on the deconvolution of bulk RNA-seq data, two data types with very

293    distinct biases[16].

294    As highlighted in the introduction, a drawback for many DNNs is the large amount of

295    training data required to obtain robust performance. Here, we used scRNA-seq data

296    to create virtually unlimited amounts of *in silico* bulk RNA-seq data of predefined type

297    (target tissue) with known composition, across datasets. This immediately highlights

298    Scaden's biggest limitation, the dependency on scRNA-seq data of the target tissue.

299    In this study we have shown that Scaden, trained solely on simulated data from

300    scRNA-seq datasets, can outperform GEP-based deconvolution algorithms. We did

301    observe, however, that the addition of labeled RNA-seq samples to the training data

302    did significantly improve deconvolution performance in the case of PBMC data. We

303    therefore believe that efforts to increase the similarity between simulated training data

304    and the target bulk RNA-seq data could increase Scaden's performance further.

305    Mixtures of *in silico* bulk RNA-seq data and publically available RNA-seq data, of

306    purified cell types for example, could further increase the deconvolution performance

307    of Scaden. Furthermore, domain adaptation methods can be used to improve

308    performance of models that are trained on data (here, scRNA-seq data) that is similar

309    to the target data (here, RNA-seq data)[25]. In future versions, Scaden's simple

310    multilayer perceptron architecture could leverage domain adaptation to further

311    stabilize and improve its cell deconvolution performance.

312    Recent cell deconvolution algorithms have used cell fraction estimates to infer cell

313    type-specific gene expression from bulk RNA-seq data. It is straightforward to use

314    Scaden's cell fraction estimates to infer per group[3] and per sample[7] cell type-specific

315    gene expression using simple regression or non-negative matrix factorization,

316    respectively. We would like to add a note of caution, however, as the error of cell

317    fraction estimates, which can be quite significant, is propagated into the gene

318    expression calculations and will affect any downstream statistical analysis.

319    In summary, the deconvolution performance, robustness to noise and bias, the

320    flexibility to learn from large numbers of *in silico* datasets, across data types (scRNA-

321    seq and RNA-seq mixtures), and potentially even tissues makes us believe that DNN-

322    based architectures will become an algorithmic mainstay of cell type deconvolution.

323

# 324 Methods

## 325 Datasets and pre-processing

### 326 scRNA-seq datasets

327 The following human PBMC scRNA-seq datasets were downloaded from the 10X

328 Genomics data download page: 6k PBMCs from a Healthy Donor, 8k PBMCs from a

329 Healthy Donor, Frozen PBMCs (Donor A), Frozen PBMCs (Donor C){Zheng et al,

330 2017}. Throughout this paper, these datasets are referred to with the handles data6k,

331 data8k, donorA and donorC, respectively. These four datasets were chosen because

332 of clearly identifiable cell types for the majority of cells. The Ascites scRNA-seq dataset

333 was downloaded from https://figshare.com as provided by Schelker[21]. Pancreas and

334 mouse brain datasets were downloaded from the scRNA-seq dataset collection of the

335 Hemberg lab (https://hemberg-lab.github.io/scRNA.seq.datasets/). A table listing all

336 datasets including references to the original publications can be found in Table S1.

### 337 scRNA-seq preprocessing and analysis

338 All datasets were processed using the Python package Scanpy (v. 1.2.2)[26] following

339 the Scanpy's reimplementation of the popular Seurat's clustering workflow. First, the

340 corresponding cell-gene matrices were filtered for cells with less than 500 detected

341 genes, and genes expressed in less than 5 cells. The resulting count matrix for each

342 dataset was filtered for outliers with high or low numbers of counts. Gene expression

343 was normalized to library size using the Scanpy function 'normalize_per_cell'. The

344 normalized matrix of all filtered cells and genes was saved for the subsequent data

345 generation step.

346      The following processing and analysis steps had the sole purpose of assigning cell

347      type labels to every cell. All cells were clustered using the louvain clustering

348      implementation of the Scanpy package. The louvain clustering resolution was chosen

349      for each dataset, using the lowest possible resolution value (low resolution values lead

350      to less clusters) for which the calculated clusters separated the cell types

351      appropriately. The top 1000 highly variable genes were used for clustering, which were

352      calculated using Scanpy's 'filter_genes_dispersion' function with parameters

353      min_mean=0.0125, max_mean=3 and min_disp=0.5. Principal Component Analysis

354      (PCA) was used for dimensionality reduction.

355      To identify cell types, marker genes were investigated for all cell types in question. For

356      PBMC datasets, useful marker genes were adopted from public resources such as the

357      Seurat tutorial for 2700 PBMCs[27]. Briefly, IL7R was taken as marker for CD4 T-cells,

358      LYZ for Monocytes, MS4A1 for B-cells, GNLY for Natural Killer cells, FCER1A for

359      Dendritic cells and CD8A and CCL5 as markers for CD8 T-cells. For all other scRNA-

360      seq datasets, marker genes and expected cell types were inferred from the original

361      publication of the dataset. For instance, to annotate cell types of the mouse brain

362      dataset from Zeisel et al.[28], we used the same marker genes as Zeisel and colleagues.

363      We did not use the same cell type labels from the original publications because a main

364      objective was to assure that cell type labeling is consistent between all datasets of a

365      certain tissue.

366      Cell type annotation was performed manually across all the clusters for each dataset,

367      such that all cells belonging to the same cluster were labeled with the same cell type.

368      The cell type identity of each cluster was chosen by crossing the cluster's highly

369      differentially expressed genes with the curated cell type's marker genes. Clusters that

370    could not be clearly identified with a cell type were grouped into the 'Unknown'

371    category.

372    Tissue Datasets for Benchmarking

373    To assess the deconvolution performance on real tissue expression data, we used

374    datasets for which the corresponding cell fractions were measured and published. The

375    first dataset is the **PBMC1** dataset which was obtained from Zimmermann *et al.*[20]. The

376    second dataset, **PBMC2**, was downloaded from GEO with accession code

377    GSE107011 [10]. This dataset contains both RNA-seq profiles of immune cells (S4

378    cohort) and from bulk individuals (S13 cohort). As we were interested in the bulk

379    profiles, we only used 12 samples from the S13 cohort from this data. Flow cytometry

380    fractions were collected from the Monaco *et al.* publication[10].

381    In addition to the above mentioned two PBMC datasets, we used Ascites RNA-seq

382    data. This dataset was kindly provided by the authors and cell type fractions for this

383    dataset were taken from the supplementary materials of the publication[21].

384    For the evaluation on pancreas data, artificial bulk RNA-seq samples created from the

385    scRNA-seq dataset of Xin *et al.*[18] were used. This dataset was downloaded from the

386    resources of the MuSiC publication[8]. The artificial bulk RNA-seq samples used for

387    evaluation were then created using the 'bulk_construct' function of the MuSiC tool.

388    To assess how Scaden deals with unknown cell types in a bulk mixture, we used the

389    whole blood dataset from Newman *et al.*[7], which consists of 12 samples (GSE127813).

390    Cell    type    fractions    were    downloaded    from    the    CSx    website

391    (https://cibersortx.stanford.edu/download.php).

392    To assess robustness against unknown mixture content, all cells classified as

393    'Unknown' were removed from the data6k, data8k, donorA, and donorC datasets to

394    generate training samples for Scaden and reference datasets for MuSiC and CSx.

19

395 Then, test datasets were generated with fixed content of 'Unknown' cells at 5%, 10%,

396 20% and 30%. Performance on these samples was then assessed to test robustness

397 against unseen cell types in the bulk mixture. Scaden was trained on samples from all

398 datasets but the test dataset, while CSx and MuSiC used data8k as a reference.

399 The microarray dataset GSE65133 was downloaded from GEO, and cell type fractions

400 taken from the original CS publication[6].

401 Finally, we wanted to get insights into neurodegenerative cell fraction changes in the

402 brain. While it is known that neurodegenerative diseases like Alzheimer's Disease are

403 accompanied by a gradual loss of brain neurons, stage-specific cell type shifts are still

404 hard to come by. Here we use the ROSMAP (Religious Orders Study and Memory and

405 Aging Project Study) cortical RNA-seq dataset along with the corresponding clinical

406 metadata, to infer cell type composition over six clinically relevant stages of

407 neurodegeneration[22].

408 RNA-seq preprocessing and analysis

409 For the RNA-seq datasets analyzed in this study, we did not apply any additional

410 processing steps, but used the obtained count or expression tables directly as

411 downloaded for all dataset except the ROSMAP dataset. For the latter, we generated

412 count tables from raw FastQ-files using Salmon[29] and the GRCh38 reference genome.

413 FastQ-files from the ROSMAP study were downloaded from Synapse

414 (www.synapse.org).

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

416 Scadan's deep neural network requires large amounts of training RNA-seq samples

417 with known cell fractions. This explains why the generation of artificial bulk RNA-seq

418 data is one of the key elements of the Scaden workflow.

419   In order to generate the training data, preprocessed scRNA-seq datasets were used

420   (see section 'Data Collection and Processing'), comprising the gene expression matrix

421   and the cell type labels. Artificial RNA-seq samples were simulated by sub-sampling

422   cells from individual scRNA-seq datasets - cells from different datasets were not

423   merged into samples to preserve within-subject relationships. Datasets generated

424   from multiple subjects were split according to subject and each sub-sampling was

425   constrained to cells from one subject in order to capture the cross-subject

426   heterogeneity and keep subject-specific gene dependencies.

427   The exact sub-sampling procedure is described in the following. First, for every

428   simulated sample, random fractions were created for all different cell types within each

429   scRNA-seq dataset using the random module of the Python package NumPy. Briefly,

430   a random number was chosen from a uniform distribution between 0 and 1 using the

431   NumPy function 'random.rand()' for each cell type, and then this number was divided

432   by the sum of all random numbers created to ensure the constraint of all fractions

433   adding up to 1:

434
$$f_c = \frac{r_c}{\sum_{C_{all}} r_c}$$

435   where $r_c$ is the random number created for cell type $c$, and $C_{all}$ is the set of all cell

436   types. Here, $f_c$ is the calculated random fraction for cell type $c$. Then, each fraction

437   was multiplied with the total number of cells selected for each sample, yielding the

438   number of cells to choose for a specific cell type:

439
440
$$N_c = f_c * N_{total}$$
441
442   where $N_c$ is the number of cells to select for the cell type $c$, and $N_{total}$ is the total

443   number of cells contributing to one simulated RNA-seq sample (400, in this study).

444   Next, $N_c$ cells were randomly sampled from the scRNA-seq gene expression matrix

21

445    for each cell type $c$. Afterwards, the randomly selected single-cell expression profiles

446    for every cell type are then aggregated by summing their expression values, to yield

447    the artificial bulk expression profile for this sample.

448    Using the above described approach, cell compositions that are strongly biased

449    toward a certain cell type or are missing specific cell types are rare among the

450    generated training samples. To account for this and to simulate cell compositions with

451    a heavy bias to and the absence of certain cell types, a variation of the sub-sampling

452    procedure was used to generate samples with sparse compositions, which we refer to

453    as sparse samples. Before generating the random fractions for all cell types, a random

454    number of cell types was selected to be absent from the sample, with the requirement

455    of at least one cell type constituting the sample. After these leave-out cell types were

456    chosen, random fractions were created and samples generated as described above.

457    Using this procedure, we generated 32,000 samples for the human PBMC training

458    dataset, 14,000 samples for the human pancreas training dataset and 30,000 samples

459    for the mouse brain training dataset (Table S2).

460    Artificial bulk RNA-seq datasets were stored in 'h5ad' format using the Anndata

461    package[26], which allows to store the samples together with their corresponding cell

462    type ratios, while also keeping information about the scRNA-seq dataset of origin for

463    each sample. This allowed to access samples from specific datasets, which is useful

464    for cross validation.


465    Scaden Overview

466    The following section contains an overview of the input data preprocessing, the

467    Scaden model, model selection, and how Scaden predictions are generated.

468    Input Data Preprocessing

469    The data preprocessing step is aimed to make the input data more suitable for

470    machine learning algorithms. To achieve this, an optimal preprocessing procedure

471    should transform any input data from the simulated samples or from the bulk RNA-seq

472    to the same feature scale. Before any scaling procedure can be applied, it must be

473    ensured that both the training data and the bulk RNA-seq data subject to prediction

474    share the same features. Therefore, before scaling, both datasets are limited to

475    contain features (genes) that are available in both datasets.. The two-step processing

476    procedure used for Scaden is described in the following:

477    First, to account for heteroscedasticity, a feature inherent to RNA-seq data, the data

478    was transformed into logarithmic space by adding a pseudocount of 1 and then taking

479    the Logarithm (base 2). Additional to stabilizing the variance, this transformation yields

480    data that is approximately Gaussian.

481    Second, every sample was scaled to the range [0,1] using the MinMaxScaler() class

482    from the Sklearn preprocessing module. Per sample scaling, unlike per feature scaling

483    that is more common in machine learning, assures that inter-gene relative expression

484    patterns in every sample are preserved. This is important, as our hypothesis was that

485    a neural network could learn the deconvolution from these inter-gene expression

486    patterns.

487    $$x_{scaled,i} = (x_i - min(\boldsymbol{X}_i)) / (max(X_i) - min(X_i))$$

488    where $x_{scaled,i}$ is the log2 expression value of gene x in sample i, $X_i$ is the vector of

489    log2 expression values for all genes of sample i, $min(\boldsymbol{X}_i)$ is the minimum gene

490    expression of vector $X_i$, and $max(X_i)$ the maximum gene expression of vector $X_i$.

23

491    Note that all training datasets are stored as expression values and are only processed

492    as described above. In the deployment use-case the simulated training data should

493    contain the same features as in the bulk RNA-seq sample that shall be deconvolved.

494    Model Selection

495    The goal of model selection was to find an architecture and hyperparameters that

496    robustly deconvolve simulated tissue RNA-seq data and, more importantly, real bulk

497    RNA-seq data. Due to the very limited availability of bulk RNA-seq datasets with known

498    cell fractions, model selection was mainly optimized on the simulated PBMC datasets.

499    To capture inter-experimental variation, we used leave-one-dataset-out cross

500    validation for model optimization: a model was trained on simulated data from all but

501    one dataset, and performance was tested on simulated samples from the left-out

502    dataset. This allows to simulate batch effects between datasets and helps to test the

503    generalizability of the model. Model performance was evaluated based on pearson

504    product moment correlation and absolute deviation between predicted and ground

505    truth values. As averaging the predictions of models with different architectures

506    increased performance, we decided to use an ensemble architecture for Scaden. For

507    this ensemble, the three best performing architectures were chosen. Model training

508    and prediction is done separately for each model, with the prediction averaging step

509    combining all model predictions (Fig. S1). We provide a list of all tested parameters in

510    the supplementary materials (Table S4).

511    Final Scaden Model

512    The Scaden model learns cell type deconvolution through supervised training on

513    datasets of simulated bulk RNA-seq samples simulated with scRNA-seq data. To

514    account for model biases and to improve performance, Scaden consists of an

515 ensemble of three deep neural networks with varying architectures and degrees of

516 dropout regularization. All models of the ensemble use four layers of varying sizes

517 between 32 and 1024 nodes, with dropout-regularization implemented in two of the

518 three ensemble models. The exact layer sizes and dropout rates are listed in Table

519 S3. The Rectified Linear Unit (ReLU) is used as activation function in every internal

520 layer. We used a Softmax function to predict cell fractions, as we did not see any

521 improvements in using a linear output function with consecutive non-negativity

522 correction and sum-to-one scaling. Python (v. 3.6.6) and the TensorFlow library (v.

523 1.10.0) were used for implementation of Scaden. A complete list of all software used

524 for the implementation of Scaden is provided in Table S12.

525 Training and Prediction

526 After the preprocessing of the data a Scaden ensemble can be trained on simulated

527 tissue RNA-seq data or mixtures of simulated and real tissue RNA-seq data.

528 Parameters are optimized using Adam with a learning rate of 0.0001 and a batch size

529 of 128. We used an L1 loss as optimization objective:

530
$$L1(y_i, \hat{y}_i) \ = \ |y_i - \hat{y}_i|$$

531 where $y_i$ is the vector of ground truth fractions of sample $i$ and $\hat{y}_i$ is the vector of

532 predicted fractions of sample $i$. Each of the three ensemble models is trained

533 independently for 5,000 steps. This 'early stopping' serves to avoid domain overfitting

534 on the simulated tissue data, which would decrease the model performance on the

535 real tissue RNA-seq data. We observed that training for more steps lead to an average

536 performance decrease on real tissue RNA-seq data. To perform deconvolution with

537 Scaden, a bulk RNA-seq sample is fed into a trained Scaden ensemble and three

538 independent predictions for the cell type fractions of this sample are generated by the

25

539    trained deep neural networks. These three predictions are then averaged per cell type

540    to yield the final cell type composition for the input bulk RNA-seq sample:

541
$$\hat{y}_c = \frac{\widehat{y_c^1} + \widehat{y_c^2} + \widehat{y_c^3}}{3}$$

542    where $\hat{y}_c$ is the final predicted fraction for cell type $c$ and $\widehat{y_c^i}$ is the predicted fraction for

543    cell type $c$ of model $i$.


## Algorithm Comparison

545    We used several performance measures to compare Scaden to four existing cell

546    deconvolution algorithms, CIBERSORT with LM22 GEP (CS), CIBERSORTx (CSx),

547    MuSiC and CPM. To compare the performance of the five deconvolution algorithms

548    we measured the root mean squared error (RMSE), Lin's concordance correlation

549    coefficient $CCC$, Pearson product moment correlation coefficient $r$, and $R^2$ values

550    comparing real and predicted cell fractions estimates. Additionally, to identify

551    systematic prediction errors and biases, slope and intercept for the regression lines

552    were calculated. These metrics are defined as follows:

553
$$RMSE(y, \hat{y}) = \sqrt{avg(y - \hat{y})^2}$$

554
$$r(y, \hat{y}) = \frac{cov(y, \hat{y})}{\sigma_y \sigma_{\hat{y}}}$$

555
$$R^2(y, \hat{y}) = r(y, \hat{y})^2$$

556
$$slope(y, \hat{y}) = \frac{\Delta y}{\Delta \hat{y}}$$

557
$$CCC(y, \hat{y}) = \frac{2r\sigma_y \sigma_{\hat{y}}}{\sigma_y{}^2 + \sigma_{\hat{y}}{}^2 + (\mu_x - \mu_{\hat{y}})}$$

558    where $y$ are the ground truth fractions, $\hat{y}$ are the prediction fractions, $\sigma_x$ is the standard

559    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

560    predicted and ground truth fractions, respectively.

26

561 All metrics were calculated for all data points of a dataset, and separately for all data

562 points of a specific cell type. For the latter approach, we then averaged the resulting

563 values to recover single values. While in general the metrics calculated on all data

564 points are sufficient, good performance on cell type-level is important if one is to

565 compare fractions of a specific cell type between samples.

## CIBERSORT (CS)

567 CS is a cell convolution algorithm based on specialized GEPs and support vector

568 regression. Cell composition estimations were obtained using the CS web application

569 (https://cibersort.stanford.edu/). For all deconvolutions with CS, we used the LM22

570 GEP, which was generated by the CS authors from 22 leukocyte subsets profiled on

571 the HGU133A microarray platform.

572 Because the LM22 GEP matrix contains cell types at a finer granularity than what was

573 used for this study, predicted fractions of sub-cell types were added together. For cell

574 grouping, we used the mapping of sub-cell types to broader types given by Figure 6

575 from Monaco *et al.*[10]. We provide a table with the exact mappings used here in the

576 supplementary material (Table S13). The deconvolution was performed using 500

577 permutations with quantile normalization disabled for all datasets but GSE65133

578 (Microarray), as is recommended for RNA-seq data. We used default settings for all

579 other CS parameters.

## CIBERSORTx (CSx)

581 CSx is a recent variant of CS that can generate GEP matrices from scRNA-seq data

582 and use these for deconvolution. For additional deconvolution robustness, it applies

583 batch normalization to the data. All signature matrices were created by uploading the

584 labeled scRNA-seq expression matrices and using the default options. Quantile

27

585 normalization was disabled. For deconvolution on simulated data, no batch

586 normalization was used. For all bulk RNA-seq datasets, the S-Mode batch

587 normalization was chosen. All PBMC datasets were deconvolved using a GEP matrix

588 generated from the data6k dataset (for simulated samples from data6k, a donorA GEP

589 matrix was chosen).

590 MuSiC

591 MuSiC is a deconvolution algorithm that uses multi-subject scRNA-seq datasets as

592 GEP matrices in an attempt to include heterogeneity in the matrices to improve

593 generalization. While MuSiC tries to address similar issues of previous deconvolution

594 algorithms by using scRNA-seq data, the approach is very different. For

595 deconvolution, MuSiC applies a sophisticated GEP-based deconvolution algorithm

596 that uses weighted non-negative least squares regression with an iterative estimation

597 procedure that imposes more weight on informative genes and less weight on non-

598 informative genes.

599 The MuSiC R package contains functionality to generate the necessary GEP matrix

600 given a scRNA-seq dataset and cell type labels. To generate MuSiC deconvolution

601 predictions on PBMC datasets, we used the data8k scRNA-seq dataset as reference

602 data for MuSiC and follow the tutorial provided by the authors to perform the

603 deconvolution. For deconvolution of artificial samples generated from the data8k

604 dataset, we provided MuSiC with the data6k dataset as reference instead.

605 MuSiC was developed with a focus on multi-subject scRNA-seq datasets, in which the

606 algorithm tries to take advantage from the added heterogeneity that these datasets

607 contain, by calculating a measure of cross-subject consistency for marker genes. To

608 assess how Scaden performs on multi-subject datasets compared to MuSiC, we

609 evaluated both methods on artificial bulk RNA-seq samples from human pancreas. We

28

610   used the 'bulk_construct' function from MuSiC to combine the cells from all 18 subjects

611   contained in the scRNA-seq dataset from Xin et al to generate artificial bulk samples

612   for evaluation. Next, as a multi-subject reference dataset, we used the pancreas

613   scRNA-seq dataset from Segerstolpe *et al*.[19], which contains single-cell expression

614   data from 10 different subjects, 4 of which with type-2 Diabetes. For Scaden, the

615   Segerstolpe scRNA-seq dataset was split by subjects, and training datasets were

616   generated for each subject, yielding in total 10,000 samples. For MuSiC, a processed

617   version of this dataset was downloaded from the resources provided by the MuSiC

618   authors[8] and used as input reference dataset for the MuSiC deconvolution.

619   Deconvolution was then performed according to the MuSiC tutorial, and performance

620   compared according to the above-defined metrics.

621   Cell Population Mapping (CPM)

622   CPM is a deconvolution algorithm that uses single-cell expression profiles to identify

623   a so-called 'cell population map' from bulk RNA-seq data[9]. In CPM, the cell population

624   map is defined as composition of cells over a cell-state space, where a cell-state is

625   defined as a current phenotype of a single cell. Contrary to other deconvolution

626   methods, CPM tries to estimate the abundance of all cell-states and types for a given

627   bulk mixture, instead of only deconvolving the cell types. As input, CPM requires a

628   scRNA-seq dataset and a low-dimensional embedding of all cells in this dataset, which

629   represents the cell-state map. As CPM estimates abundances of both cell-states and

630   types, it can be used for cell type deconvolution by summing up all estimated fractions

631   for all cell-states of a given cell type - a method that is implemented in the scBio R

632   package, which contains the CPM method. To perform deconvolution with CPM, we

633   used the data6k PBMC scRNA-seq dataset as input reference for all PBMC samples.

634   For samples simulated from the data6k dataset, we used the data8k dataset as

29

635    reference. According to the CPM paper, a dimension reduction method can be used

636    to obtain the cell-state space. We therefore used UMAP, a dimension reduction

637    method widely used for scRNA-seq data, to generate the cell-state space mapping for

638    the input scRNA-seq data. Deconvolution was then performed using the CPM function

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

640    input.

641

## Data Availability

Only publicly available datasets were used during this study. The scRNA-seq PBMC

datasetse donorA, donorC, data6k and data8k were all downloaded from 10X

Genomics (https://support.10xgenomics.com/single-cell-gene-expression/datasets),

were they are listed as 'Frozen PBMCs (Donor A)', 'Frozen PBMCs (Donor C)', '6k

PBMCs from a Healthy Donor' and '8k PBMCs from a Healthy Donor', respectively.

The Segerstolpe et al. scRNA-seq pancreas dataset was downloaded from

ArrayExpress with accession code E-MTAB-5061. The scRNA-seq datasets from

Baron et al. (pancreas), Tasic et al., Zeisel et al., Romanov et al., Campbell et al.

and Chen et al. (all mouse brain) were all downloaded from https://hemberg-

lab.github.io/scRNA.seq.datasets/. The ascites scRNA-seq dataset was downloaded

from https://figshare.com/s/711d3fb2bd3288c8483a. The bulk RNA-seq dataset

PBMC1 is accessible from ImmPort with accession code SDY67. The PBMC2

dataset was downloaded from GEO with accession code GSE107011. The

ROSMAP human brain RNA-seq dataset was downloaded from Synapse (ID:

syn3219045). The bulk RNA-seq data from ascites was kindly provided by Schelker

et al. The pancreas scRNA-seq dataset from Xin et al. was accessed from the

MuSiC tutorial site (https://xuranw.github.io/MuSiC/articles/pages/data.html).

660     ## Code Availability

661     The source code for Scaden is available at https://github.com/KevinMenden/scaden.

662     Documentation is published at https://scaden.readthedocs.io. Code to generate the

663     figures along with the training datasets used in this study is published at figshare:

664     https://figshare.com/projects/Scaden/62834.

665

## 666    List of abbreviations

667    RNA-seq : Next Generation RNA Sequencing

668    GEP : gene expression profile matrix

669    SVR : Support Vector Regression

670    DNN : Deep Neural Network

671    scRNA-seq : single cell RNA-seq

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

673    PBMC : peripheral blood mononuclear cells

674    CCC : concordance correlation coefficient

675    r : Pearson's correlation coefficient

676    CS : CIBERSORT

677    CSx : CIBERSORTx

678    CPM : Cell Population Mapping

679 # Author information

680 ## Affiliations

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

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

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

684 **Eppendorf, Germany**

685 Mohamed Marouf, Stefan Bonn

686 ## Contributions

687 KM and SB initiated the project. KM, PH, and SB designed the study, deep learning

688 models, and analysis. KM and MM built the deep learning models. KM, MM, and AD

689 analyzed the data. KM and SB wrote and MM, AD, and PH contributed to the

690 manuscript writing.

691

692 Competing interests

693 The authors have no competing interests.

702    **Corresponding author**

703    Correspondence    to    Stefan    Bonn    (sbonn@uke.de)    and    Kevin    Menden

704    (kevin.menden@dzne.de).

705

# References

707 1. Hrdlickova, R., Toloue, M. & Tian, B. RNA-Seq methods for transcriptome

708      analysis. *Wiley Interdiscip. Rev. RNA* **8,** (2017).

709 2. Egeblad, M., Nakasone, E. S. & Werb, Z. Tumors as organs: Complex tissues

710      that interface with the entire organism. *Dev. Cell* **18,** 884–901 (2010).

711 3. Kuhn, A., Thu, D., Waldvogel, H. J., Faull, R. L. M. & Luthi-Carter, R.

712      Population-specific expression analysis (PSEA) reveals molecular changes in

713      diseased brain. *Nat. Methods* **8,** 945–947 (2011).

714 4. Avila Cobos, F., Vandesompele, J., Mestdagh, P. & De Preter, K.

715      Computational deconvolution of transcriptomics data from mixed cell

716      populations. *Bioinformatics* **34,** 1969–1979 (2018).

717 5. Mohammadi, S., Zuckerman, N., Goldsmith, A. & Grama, A. A Critical Survey

718      of Deconvolution Methods for Separating Cell Types in Complex Tissues.

719      *Proc. IEEE* **105,** 340–366 (2017).

720 6. Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue

721      expression profiles. *Nat. Methods* **12,** 453–457 (2015).

722 7. Newman, A. M. *et al.* Determining cell type abundance and expression from

723      bulk tissues with digital cytometry. *Nat. Biotechnol.* (2019).

724      doi:10.1038/s41587-019-0114-2

725 8. Wang, X., Park, J., Susztak, K., Zhang, N. R. & Li, M. Bulk tissue cell type

726      deconvolution with multi-subject single-cell expression reference. *Nat.*

727      *Commun.* **10,** 380 (2019).

728 9. Frishberg, A. *et al.* Cell composition analysis of bulk genomics using single-cell

729      data. *Nat. Methods* **16,** (2019).

730 10. Monaco, G. *et al.* RNA-Seq Signatures Normalized by mRNA Abundance

731       Allow Absolute Deconvolution of Human Immune Cell Types. *Cell Rep.* **26,**

732       1627–1640.e7 (2019).

733   11.  Vallania, F. *et al.* Leveraging heterogeneity across multiple datasets increases

734       cell-mixture deconvolution accuracy and reduces biological and technical

735       biases. *Nat. Commun.* **9,** (2018).

736   12.  Venet, D., Pecasse, F., Maenhaut, C. & Bersini, H. Separation of samples into

737       their constituents using gene expression data. *Bioinformatics* **17,** 279–287

738       (2001).

739   13.  Shapiro, E., Biezuner, T. & Linnarsson, S. Single-cell sequencing-based

740       technologies will revolutionize whole-organism science. *Nat. Rev. Genet.* **14,**

741       618–630 (2013).

742   14.  Schaum, N. *et al.* Single-cell transcriptomics of 20 mouse organs creates a

743       Tabula Muris. *Nature* **562,** 367–372 (2018).

744   15.  Kelley, K. W., Nakao-Inoue, H., Molofsky, A. V. & Oldham, M. C. Variation

745       among intact tissue samples reveals the core transcriptional features of human

746       CNS cell classes. *Nat. Neurosci.* **21,** 1171–1184 (2018).

747   16.  Hicks, S. C., Townes, F. W., Teng, M. & Irizarry, R. A. Missing data and

748       technical variability in single-cell RNA-sequencing experiments. *Biostatistics*

749       **19,** 562–578 (2018).

750   17.  Lin, L. I. A Concordance Correlation Coefficient to Evaluate Reproducibility

751       Author ( s ): Lawrence I-Kuei Lin Published by : International Biometric Society

752       Stable URL : http://www.jstor.org/stable/2532051 REFERENCES Linked

753       references are available on JSTOR for thi. *Biomatrics* **45,** 255–268 (1989).

754   18.  Xin, Y. *et al.* RNA Sequencing of Single Human Islet Cells Reveals Type 2

755       Diabetes Genes. *Cell Metab.* **24,** 608–615 (2016).

756   19.   Segerstolpe, Å. *et al.* Single-Cell Transcriptome Profiling of Human Pancreatic

757         Islets in Health and Type 2 Diabetes. *Cell Metab.* **24,** 593–607 (2016).

758   20.   Zimmermann, M. T. *et al.* System-wide associations between DNA-

759         methylation, gene expression, and humoral immune response to influenza

760         vaccination. *PLoS One* **11,** 1–21 (2016).

761   21.   Schelker, M. *et al.* Estimation of immune cell content in tumour tissue using

762         single-cell RNA-seq data. *Nat. Commun.* **8,** 2032 (2017).

763   22.   Bennett, D. A. *et al.* Religious Orders Study and Rush Memory and Aging

764         Project. *J. Alzheimer's Dis.* **64,** S161–S189 (2018).

765   23.   Braak, H. & Braak, E. Neuropathological stageing of Alzheimer-related

766         changes. *Acta Neuropathol.* **82,** 239–59 (1991).

767   24.   Yosinski, J., Clune, J., Nguyen, A., Fuchs, T. & Lipson, H. Understanding

768         Neural Networks Through Deep Visualization. (2015).

769   25.   Athiwaratkun, B., Finzi, M., Izmailov, P. & Wilson, A. G. Improving

770         Consistency-Based Semi-Supervised Learning with Weight Averaging. *Jmlr*

771         **17,** 1–35 (2018).

772   26.   Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: Large-scale single-cell gene

773         expression data analysis. *Genome Biol.* **19,** 1–5 (2018).

774   27.   Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial

775         reconstruction of single-cell gene expression data. *Nat. Biotechnol.* **33,** 495–

776         502 (2015).

777   28.   Zeisel, A. *et al.* Molecular Architecture of the Mouse Nervous System. *Cell*

778         **174,** 999–1014.e22 (2018).

779   29.   Love, M. I., Soneson, C., Patro, R., Vitting-seerup, K. & Oshlack, A. Swimming

780         downstream : statistical analysis of differential transcript usage following

781        Salmon quantification. 1–50 (2019).

782

783

784  # 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 |

785  **Table S1** *scRNA-seq datasets used for the generation of simulated tissues for Scaden*

786  *training.*

787

788

| 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 |

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

790

791

| 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 |

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

793

794

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

795

796

797

798

799

800

43

| 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 |

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

802

| Model | # Layers | Layer sizes | Dropout rates |
|-------|----------|-------------|---------------|

| Method | DS | RMSE | Slope | Correlation | Intercept | CCC |
|---|---|---|---|---|---|---|
| **CPM** | data6k | 0.192 | 0.03 | 0.082 | 0.162 | 0.053 |
| **CPM** | data8k | 0.185 | 0.048 | 0.263 | 0.159 | 0.093 |
| **CPM** | donorA | 0.239 | -0.081 | -0.259 | 0.18 | -0.147 |
| **CPM** | donorC | 0.189 | 0.038 | 0.102 | 0.16 | 0.066 |
| **CS** | data6k | 0.163 | 0.508 | 0.57 | 0.082 | 0.566 |
| **CS** | data8k | 0.136 | 0.551 | 0.708 | 0.075 | 0.687 |
| **CS** | donorA | 0.137 | 0.605 | 0.767 | 0.066 | 0.746 |
| **CS** | donorC | 0.168 | 0.45 | 0.522 | 0.092 | 0.517 |
| **CSx** | data6k | 0.106 | 0.756 | 0.824 | 0.041 | 0.821 |
| **CSx** | data8k | 0.097 | 0.744 | 0.863 | 0.043 | 0.854 |
| **CSx** | donorA | 0.125 | 0.696 | 0.81 | 0.051 | 0.801 |
| **CSx** | donorC | 0.094 | 0.829 | 0.865 | 0.029 | 0.864 |
| **MuSiC** | data6k | 0.086 | 0.848 | 0.887 | 0.025 | 0.886 |
| **MuSiC** | data8k | 0.136 | 0.663 | 0.728 | 0.056 | 0.725 |
| **MuSiC** | donorA | 0.1 | 0.811 | 0.883 | 0.031 | 0.88 |
| **MuSiC** | donorC | 0.084 | 0.897 | 0.896 | 0.017 | 0.896 |
| **Scaden** | data6k | 0.104 | 0.747 | 0.83 | 0.042 | 0.825 |
| **Scaden** | data8k | 0.133 | 0.625 | 0.73 | 0.063 | 0.722 |
| **Scaden** | donorA | 0.035 | 0.92 | 0.988 | 0.013 | 0.985 |
| **Scaden** | donorC | 0.046 | 0.849 | 0.973 | 0.025 | 0.964 |

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

804

45

805

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

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

807

808

809

| 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 |

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

811

812

813

**Figure S2** *Deconvolution performance on datasets with added unknown mixture contents.*

815 Unknown cell type content was added to the simulated bulk mixture in fixed concentrations

816 (5%, 10%, 20%, 30%). The deconvolution performance was assessed on samples

817 generated from the data6k, donorA and donorC datasets.

818

| Method | Content | RMSE | CCC |
|---|---|---|---|
| CSx | 0.02 | 0.097 | 0.731 |
| CSx | 0.05 | 0.092 | 0.751 |
| CSx | 0.1 | 0.092 | 0.715 |
| CSx | 0.2 | 0.091 | 0.694 |
| CSx | 0.3 | 0.099 | 0.632 |
| MuSiC | 0.02 | 0.084 | 0.793 |
| MuSiC | 0.05 | 0.083 | 0.799 |
| MuSiC | 0.1 | 0.089 | 0.739 |
| MuSiC | 0.2 | 0.095 | 0.669 |
| MuSiC | 0.3 | 0.101 | 0.665 |
| Scaden | 0.02 | 0.041 | 0.934 |
| Scaden | 0.05 | 0.044 | 0.925 |
| Scaden | 0.1 | 0.046 | 0.902 |
| Scaden | 0.2 | 0.054 | 0.846 |
| Scaden | 0.3 | 0.063 | 0.798 |

819

820 **Table S8** *Deconvolution performance on datasets with added unknown mixture*
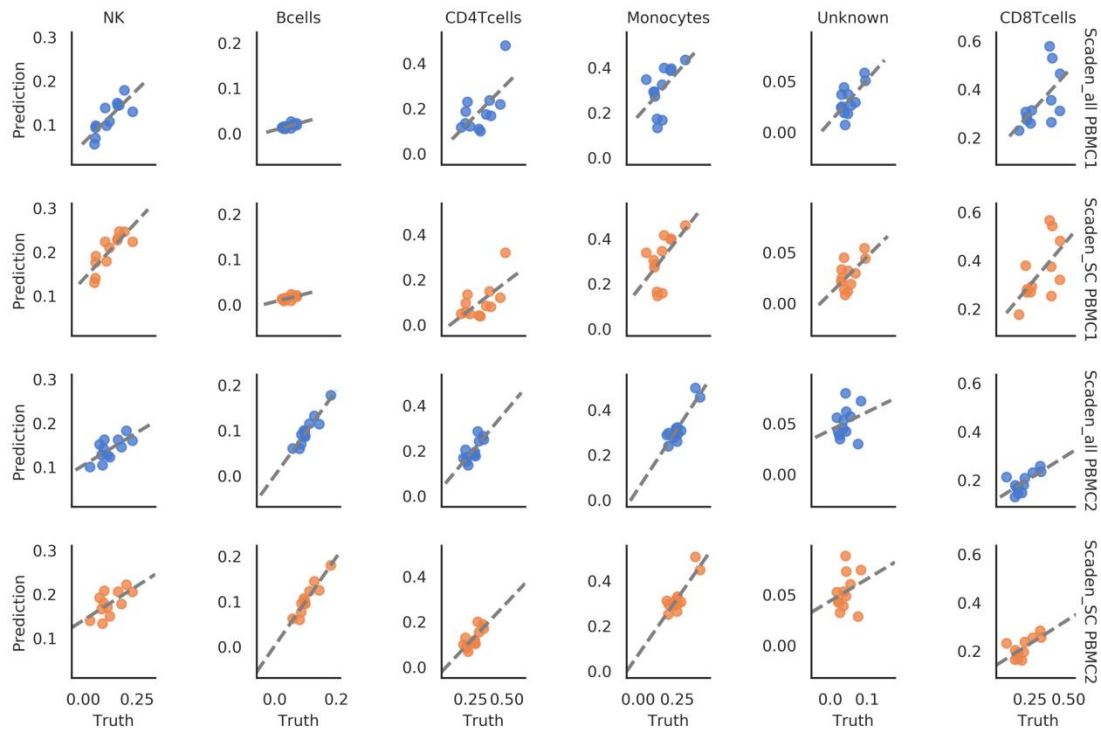
821 *contents.*

822

823

| 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 |

824     **Table S9** *Deconvolution performance on real PBMC RNA-seq datasets PBMC1 and*

825     *PBMC2.*

826

827

**Figure S3** *Comparison of Scaden deconvolution results on PBMC1 and PBMC2 datasets*

829 *with and withouth (Scaden_all, Scaden_SC, respectively) bulk RNA-seq samples included in*

830 *training data.*

831

832

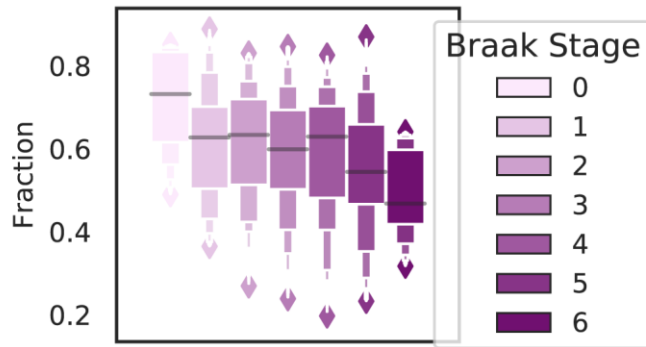| Method | Dataset | Celltype | RMSE | Correlation | Slope | Intercept | CCC |
|---|---|---|---|---|---|---|---|
| **Scaden_SC** | PBMC1 | Total | 0.131 | 0.564 | 0.644 | 0.059 | 0.559 |
| **Scaden_SC** | PBMC2 | Total | 0.077 | 0.684 | 0.689 | 0.052 | 0.684 |
| **Scaden_all** | PBMC1 | Total | 0.104 | 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 |

833  **Table S10** *Deconvolution performance on real PBMC RNA-seq data for Scaden models*

834  *trained only on scRNA-seq simulated tissues (Scaden_SC) or on a mix of simulated and real*

835  *tissue data (Scaden_all).*

836

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

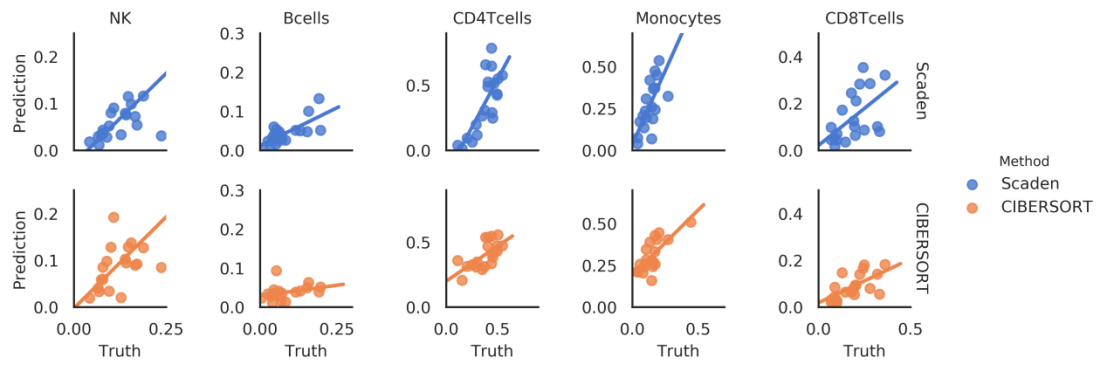837 **Table S11** *Deconvolution performance on real Ascites RNA-seq data.*

838

839

**Figure S4** *Deconvolution performance on real human brain RNA-seq data.* Scaden was trained on mouse scRNA-seq data and the trained model was used to deconvolve cell fractions of ROSMAP human brain RNA-seq data. This data does not contain cell fraction ground-truth information. Instead, the box plot shows the decrease of neuronal cell fractions with increasing Braak disease stage, a well-known phenomenon in AD.

845

**Figure S5** *Deconvolution performance comparison of CS (LM22) and Scaden on the*

*GSE65133 PBMC microarray dataset.*

850

| 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 |

851    **Table S12** *Deconvolution performance on real PBMC microarray data.*

852

853

| 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 |

854     **Table S13** *Software packages and versions used.*

855

856

| 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 folicular helper, T cells gamma delta, Plasma cells, Neutrophils, Dendritic |

857 **Table S14** *Mapping of the LM22 GEP to cell types.*

858