1  **Automated identification of Cell Types in Single Cell RNA Sequencing**

2  Feiyang Ma[1], Matteo Pellegrini[1,2,3*]

3

4  1. Molecular Biology Institute, University of California, Los Angeles, Los Angeles, California,

5  United States of America

6  2. Institute of Genomics and Proteomics, University of California, Los Angeles, Los Angeles,

7  California, United States of America

8  3. Department of Molecular, Cell, and Developmental Biology, University of California, Los

9  Angeles, Los Angeles, California, United States of America

10

11  * Corresponding author

12  E-mail: matteop@mcdb.ucla.edu

13

14

15

16

17

18

19

20

21

22

23

24 **Abstract**

25 Cell type identification is one of the major goals in single cell RNA sequencing (scRNA-seq).

26 Current methods for assigning cell types typically involve the use of unsupervised clustering, the

27 identification of signature genes in each cluster, followed by a manual lookup of these genes in

28 the literature and databases to assign cell types. However, there are several limitations associated

29 with these approaches, such as unwanted sources of variation that influence clustering and a lack

30 of canonical markers for certain cell types. Here, we present ACTINN (Automated Cell Type

31 Identification using Neural Networks), which employs a neural network with 3 hidden layers,

32 trains on datasets with predefined cell types, and predicts cell types for other datasets based on

33 the trained parameters. We trained the neural network on a mouse cell type atlas (Tabula Muris

34 Atlas) and a human immune cell dataset, and used it to predict cell types for mouse leukocytes,

35 human PBMCs and human T cell sub types. The results showed that our neural network is fast

36 and accurate, and should therefore be a useful tool to complement existing scRNA-seq pipelines.

37

38 **Author Summary**

39 Single cell RNA sequencing (scRNA-seq) provides high resolution profiling of the

40 transcriptomes of individual cells, which inevitably results in high volumes of data that require

41 complex data processing pipelines. Usually, one of the first steps in the analysis of scRNA-seq is

42 to assign individual cells to known cell types. To accomplish this, traditional methods first group

43 the cells into different clusters, then find marker genes, and finally use these to manually assign

44 cell types for each cluster. Thus these methods require prior knowledge of cell type canonical

45 markers, and some level of subjectivity to make the cell type assignments. As a result, the

46 process is often laborious and requires domain specific expertise, which is a barrier for

47    inexperienced users. By contrast, our neural network ACTINN automatically learns the features

48    for each predefined cell type and uses these features to predict cell types for individual cells.

49    This approach is computationally efficient and requires no domain expertise of the tissues being

50    studied. We believe ACTINN allows users to rapidly identify cell types in their datasets, thus

51    rendering the analysis of their scRNA-seq datasets more efficient.

52

53

54

## Introduction

56    Single cell RNA sequencing (scRNA-seq) enables the profiling of the transcriptomes of

57    individual cells, thus characterizing the heterogeneity of samples in manner that was not possible

58    using traditional bulk RNA-Seq[1]. However, scRNA-seq experiments typically yield high

59    volumes of data, especially when the number of cells is large (often many thousands). Thus, fast

60    and efficient computational methods are essential for scRNA-seq analyses.

61

62    One common goal of scRNA-seq analyses is to identify the cell type of each individual cell that

63    has been profiled. To accomplish this, typically cells are first grouped into different clusters in an

64    unsupervised way, and the number of clusters allows us to approximately determine how many

65    distinct cell types are present in the sample. Each cluster should contain cells with similar

66    expression profiles, and so the aggregated profile of a cluster increases the signal to noise of the

67    expression estimates. To attempt to interpret the identity of each cluster, marker genes are found

68    as those that are uniquely highly expressed in a cluster, compared to all the other clusters. These

69    canonical markers are then used to assign the cell types for the clusters, by cross referencing the

70   markers with lists of previously characterized cell type specific markers. While this process is

71   able to identify cell types, there are some limitations: 1. Since the clustering method is

72   unsupervised, all sources of variation influence the formation clusters, including effects that are

73   not directly related to cell types such as differential expression induced by cell cycles. 2. It is

74   often difficult to find an optimal match between the marker genes associated with each cluster

75   and the canonical markers for specific cell types. Moreover, depending on the clustering

76   parameters used,  one cluster might contain multiple cell types, or one cell type could be split

77   into multiple clusters. 3. Using canonical markers to assign cell types requires background

78   knowledge of cell type specific markers, and sometimes these are not well characterized or

79   difficult to find in the literature. Moreover, some canonical markers may be expressed by more

80   than one cell type, and some cell types may have no known markers. 4. The same types of cells

81   processed by two distinct scRNA-seq techniques tend to cluster separately due to technical batch

82   effects, which complicates cell type identification in composite datasets. 5. Cell subtypes are

83   often very similar to each other, which limits efforts to separate them accurately into different

84   clusters. To overcome many of the limitations of existing approaches, new methods need to be

85   developed.

86

87   Neural networks provide a popular framework for machine learning algorithms which can be

88   used to interpret complex datasets.  As a result, neural networks have been widely used in many

89   fields, including for the analysis of scRNA-seq data[2-5]. Since the output data from scRNA-seq is

90   feature-enriched and well-structured, it is well suited as an input for neural networks. Here, we

91   present  ACTINN (Automated Cell Type Identification using Neural Networks) for scRNA-seq

92   cell type identification. To overcome may of the limitations of traditional cell type identification

93 approaches described above, we used a neural network with 3 hidden layers, trained it on

94 scRNA-seq datasets with predefined cell types, and predicted cell types in other datasets based

95 on the trained parameters. We tested our neural network with several published datasets and

96 show that it is fast, efficient and accurate.

97

## 98 Results

### 99 Overview of the neural network

100 We used a neural network with 3 hidden layers, each containing 100, 50 and 25 nodes,

101 respectively (Fig 1). For the activation functions, we used the softmax function for the ouput

102 layer and the rectified linear unit (ReLU) function for the other layers. We used the cross-entropy

103 function as the loss function. The neural network model was implemented using TensorFlow

104 (https://www.tensorflow.org/), and the code was written in python. We trained the neural

105 network on 6 Intel(R) Xeon(R) CPU E5-2660 v3, and the training process took 0.5 minute to

106 complete with 1000 cells, 11 minutes with 32,000 cells and 21 minutes with 56,000 cells. The

107 maximum memory used in training with 56,000 cells was 18 GB. The code and datasets used in

108 this study are available at https://github.com/mafeiyang/ACTINN.

109

### 110 ACTINN model for murine cell types

111 We used 2 datasets from the Tabula Muris Consortium (The Tabula Muris Consortium. 2018) to

112 train and test our neural network. The datasets contain 100,605 cells from 20 mouse organs, and

113 were sequenced by two distinct techniques, 10X Genomics (10X) and Smart-seq2 (SS2). To

114 ensure we are using cells with high quality, we filtered out cells with less than 300 detected

115 genes, clustered the cells, and identified marker genes for each cluster using Seurat[6]. The details

116    of the Seurat analysis can be found in the methods section. We manually assigned cell types for

117    each cluster based on canonical markers (Fig 2A). We focused on 12 cell types and selected cells

118    that have the same labels between our analyses and the Tabula Muris Consortium's. This process

119    resulted in 56,112 cells (Fig 2B). Cells processed by 10X have a median of 4,787 unique

120    molecular identifiers (UMIs) and 1,558 genes detected, and cells processed by SS2 have a

121    median of 623,799 UMIs and 2,448 genes detected.

122

123    To test the robustness of our neural network's performance, we first trained and tested it on cells

124    processed by each scRNA-seq platform separately. To this end, we randomly sampled 3000 cells

125    for testing, and used the remainder of cells for training. We repeated this process 10 times, and

126    the average training accuracies for the 10X dataset and the SS2 dataset were 99.997% and

127    99.963%, respectively, and the average testing accuracies were 99.883% and 99.660%,

128    respectively (Fig 2D). These results show that our neural network can achieve very high

129    accuracy when training and testing on datasets generated by the same technique.

130

131    **ACTINN overcomes batch effects introduced by different techniques**

132    Different scRNA-seq techniques can introduce significant batch effects[7] with the same cell

133    types clustering separately due to technical artifacts (Fig 2C). To test our neural network's

134    performance accounting for the batch effects introduced by different techniques, we trained it on

135    cells processed by one platform and tested it on cells processed by the other. We first trained the

136    neural network on all the 10X cells and tested in on all the SS2 cells. The training accuracy was

137    99.997% and the testing accuracy was 98.625%. Among the 288 incorrectly predicted cells, 118

138    monocytes were predicted as B cells, 64 monocytes were predicted as epithelial cells, 47 NK

139      cells were predicted T cells (Supplementary File 1). We then trained the neural network on the

140      SS2 dataset and tested it on the 10X dataset. The training accuracy was 100% and the testing

141      accuracy was 99.195%. Among the 283 incorrectly predicted cells, 150 endothelial cells were

142      predicted as epidermis, 46 T cells were predicted as NK cells, and there were several other

143      mispredictions (Supplementary File 1).

144

145      **Early stopping prevents overfitting of the training set**

146      To prevent overfitting the parameters on the training set, we randomly sampled 5,000 cells from

147      the 10X dataset and 5,000 cells from the SS2 dataset. We trained the neural network on the 10X

148      cells and tested it on the SS2 cells. During the training process, we recorded the accuracy and the

149      cost after each epoch. The accuracy was defined as the percentage of cells whose cell type was

150      correctly predicted, and the cost was the output of the cost function after each epoch. We found

151      that the training accuracy saturated early (5 epochs), and the testing accuracy saturated at around

152      50 epochs (Fig 2E), and the cost decreased very slowly after 50 epochs (Fig 2F). These results

153      indicate that early stopping can be used to reduce training time and prevent overfitting.

154

155      **Cell type prediction using the mouse cell atlas**

156      Since the cell types from the two mouse cell atlas datasets can be accurately predicted, we

157      combined the two datasets and used the combined dataset as the reference to predict cell types

158      for other datasets. We first tried to predict cell types for a dataset that contains flow cytometry

159      sorted leukocytes from mouse aorta[8]. All cells were predicted as leukocytes except for 1

160      erythrocyte, which we think is a doublet of an erythrocyte and B cell as high expression of

161      hemoglobin genes was detected (Fig 3A). We also carried out unsupervised analysis on the

162    dataset and clustered the cells using Seurat. Then we used the canonical markers to assign the

163    cell types for each cluster (Fig 3B). Most cells had the same cell type assignment by the two

164    methods. However, our neural network detected some natural killer (NK) cells, which were in

165    the same cluster with the T cells, and were assigned as T cells in the unsupervised clustering. We

166    checked the expression of CD3D, CD8A and GZMA (Fig 3C), and found no expression of

167    CD3D and CD8A, but high expression of GZMA in the NK cells, which suggests that these are

168    likely NK cells.

169

170    It is generally thought that human and mouse share similar cell types, and the same cell type

171    from human and mouse share similar expression profiles. To test this, we trained our neural

172    network on the mouse cell atlas datasets and used the parameters to predict the cell types for a

173    human peripheral blood mononuclear cell (PBMC) dataset. We found 4 main populations in the

174    PBMC dataset, namely, B cells, monocytes, NK cells and T cells (Fig 3D). We plotted the

175    canonical markers for these 4 populations (Fig 3E) and found that the predicted cell types

176    matched the expected marker expression. These results suggest that the mouse cell atlas datasets

177    can be used as a reference to identify cell types for both human and mouse cells.

178

179    **ACTINN accurately predicts cell subtypes**

180    Although it is relatively easy to distinguish different cell types in scRNA-seq using the

181    unsupervised clustering methods, it is more difficult to further divide one cell type into cell

182    subtypes. Here, we collected 5 publicly available datasets[9], each containing one flow cytometry

183    sorted T cell subtype. We merged these datasets and selected the cells that have the same labels

184    between our analyses and the flow cytometry sorting, and then used these cells as a reference for

185    the neural network. We then clustered the selected cells and identified markers (Fig 4A and 4B)

186    for each sub cell type using Seurat. For the test set, we used the T cells from the human PBMC

187    datasets mentioned above.

188

189    To test our neural network's ability to predict cell subtypes, we trained it on the T cell subtype

190    reference, and predicted the subtypes for the T cells from the PBMC dataset (Fig 4D). We then

191    identified marker genes for each predicted subtype. As expected, the marker genes matched the

192    ones from the reference (Fig 4E). These results show that our neural network can be used to

193    accurately identify cell subtypes. We found that the subtypes predicted by the neural network did

194    not perfectly match the cell types associated with the Seurat clusters (Fig 4C). Some clusters

195    contained different subtypes and some subtypes were composed of several clusters. We think the

196    difference was influenced by two factors: 1. Unsupervised clustering considers all variance in the

197    data, while the neural network is trained to find the difference between the subtypes; 2. It is

198    difficult to   set the parameters optimally for the unsupervised analysis, which can result in

199    multiple cell types in one cluster or multiple clusters for one cell type.

200

## Discussion

202    scRNA-seq provides high resolution profiling of the transcriptomes of single cells. Typically, the

203    first step in scRNA-seq analysis is to assign each cell a cell type based on our prior knowledge of

204    marker genes. Current methods for cell type assignment first cluster the cells in an unsupervised

205    manner and rely on the canonical markers to identify the cell types for each cluster. However,

206    this approach has several limitations, including the fact that the clusters may not optimally

207    segregate single cell types, and certain cell types may not have previously characterized markers.

208 Moreover, these methods are computationally intensive, especially when the number of cells

209 becomes large. To render cell type identification in scRNA-seq more efficient, we employed a

210 neural network, trained it on cells with predefined cell types, and used it to predict cell types for

211 new datasets.

212

213 We first obtained and cleaned two datasets from the Tabula Muris Consortium, then trained and

214 tested our neural network on these datasets with or without batch effect introduced by different

215 scRNA-seq platforms. The training accuracy always approached 100%, and the testing accuracy

216 was around 99.8% within a platform and 99.0% when testing and training are performed across

217 different platforms. As the cell types in the two Tabula muris atlas datasets can be mutually

218 predicted using our neural network, we merged them and used the combined datasets as the

219 reference to predict cell types for other datasets. The predicted cell types were well matched with

220 the cell types assigned using the canonical markers for both the mouse and human datasets. We

221 also trained and tested the neural network on 5 T cell subtypes and found that the predicted

222 subtypes showed the same markers as the reference subtypes, which suggests that our neural

223 network can be used to predict sub cell types as well.

224

225 Compared to the traditional unsupervised methods used for cell type identification, our neural

226 network has the following advantages: 1. It uses all the genes to capture the features for each cell

227 type instead of relying on a limited number of canonical markers. 2. It focuses the analysis on the

228 signal associated with the variance between cell types, while unsupervised clustering tends to be

229 affected by other sources of cell type independent variation (i.e. platform or cell cycle). 3. It

230 requires no background knowledge of cell type markers, while the unsupervised method requires

231    users to have prior knowledge of canonical markers for each cell type in their data. 4. It is much

232    more computationally efficient than the traditional approach. Moreover, users can subsample the

233    reference cells to make the computation of the neural network less compute intensive and more

234    memory efficient.

235

236    There are some aspects of our approach that could be improved in the future. As the neural

237    network is supervised, the quantity and quality of the reference data are critical. We anticipate

238    that with time more cell types from larger atlases should be used to train a more comprehensive

239    neural network. Also, better pairing of reference and test sets will undoubtedly improve

240    performance. For example, the soon to be developed human cell atlas should be used to predict

241    human cell types instead of the mouse cell atlas. Nonetheless, we showed that even with the

242    current reference data our neural network is computationally efficient and accurate, and should

243    improve cell type identification pipelines.

244

245    **Materials and Methods**

246    **Data normalization**

247    We used several publicly available datasets in our analyses. The mouse cell atlas datasets were

248    collected from https://tabula-muris.ds.czbiohub.org/. The CD45 sorted leukocyte datasets were

249    published in Winkels et al[8] .The T cell subtypes and PBMC datasets were collected from

250    https://support.10xgenomics.com/single-cell-gene-expression/datasets.

251    To filter and normalize the data, we first identified genes that were detected in both training set

252    and test set.  The training set and the test set were then merged into one matrix based on the

253    common genes. Next, each cell's expression value was normalized to its total expression value

254 and multiplied by a scale factor of 10,000. The counts were increased by 1, and the log2 value

255 was calculated. To filter out outlier genes, the genes with the highest 1% and lowest 1%

256 expression were removed. The gene with the highest 1% and the lowest 1% standard deviation

257 were also removed. Finally, the matrix was split into the training set and the test set.

258

259 **Neural network configuration**

260 We used a neural network that contains an input layer, 3 hidden layers, and an output layer. The

261 input layer had a number of nodes equal to the number of genes in the training set. The 3 hidden

262 layers had 100, 50 and 25 nodes, respectively. The output layer had a number of nodes equal to

263 the number of cell types in the training set. Forward propagation was implemented as:

$$x^{[i]} = g\left(W^{[i]}x^{[i-1]} + b^{[i-1]}\right)$$

264 Where $x^{[i]}$ represents the output of the $i$th layer ($x^{[0]}$ represents the input layer), $b^{[i]}$ represents

265 the intercept of the $i$th layer, $W^{[i]}$ represents the weight matrix of the $i$th layer, and $g$ represents

266 the activation function used in the neural network. Specifically, for the activation function, a

267 rectified linear unit (ReLU) function was used for the input and hidden layers, which is defined

268 as:

$$ReLU(x) = max(0, x)$$

269 For the output layer, the softmax function was used, which is defined as:

$$softmax(x_{[j]}) = \frac{\exp{(x_{[j]})}}{\sum_{j=1}^{k} \exp{\left(x_{[j]}\right)}}$$

270 Where $x_{[j]}$ represents the $j$th element of the input vector for the output layer, which has $k$

271 elements, representing a total of $k$ cell types in the training set.

272 For the loss function, we used the cross-entropy function, which is defined as:

$$H(y', y) = \sum_{j=1}^{k} \left( y_{[j]} log(y'_{[j]}) + (1 - y_{[j]}) log(1 - y'_{[j]}) \right)$$

273    Where vector $y$ represents the true label for the cell, $y_{[j]}$ is defined to be 1 if the cell is the $j$th

274    cell type, and the other elements in $y$ are defined to be 0. $y'$ represents the output of the output

275    layer, and $y'_{[j]}$ represents the posterior probability that the cell is the $j$th cell type. L2

276    regularization was added to the loss function.

277

278    **Parameters used in the neural network**

279    The neural network model was implemented using TensorFlow (https://www.tensorflow.org/),

280    and the code was written in python. The parameters were initialized with Xavier initializer[10].

281    The starting learning rate was set to 0.0001 with staircase exponential decay, the decay rate was

282    set to 0.95, and the decay step was set to 1000. This means that after every 1000 global steps, the

283    learning rate would be the original learning rate multiplied by 0.95. 50 epochs were used to train

284    the neural network with a mini batch size of 128, which is the number of samples used in

285    training at every global step. The L2 regularization rate was set to 0.005.

286

287    **Unsupervised single cell analysis**

288    To identify different cell types and find signature genes for each cell type, Seurat[6] was used to

289    analyze the digital expression matrix generated by scRNA-seq. Specifically, in Seurat, cells with

290    less than 1000 unique molecular identifiers (UMIs) and genes detected in less than 10 cells were

291    first filtered out. Second, highly variable genes were detected and used for further analysis. Third,

292    the data was scaled for sequencing depth of each cell. Fourth, principle component analysis

293    (PCA) and t-distributed stochastic neighbor embedding (tSNE) were used to reduce the

294    dimension and plot the data on a two-dimensional graph. Lastly, a graph-based clustering

295    approach was used to cluster the cells, then signature genes were found and used to define cell

296    type for each cluster.

297

## Acknowledgments

301

## References

303    1.   Hwang B, Lee JH, Bang D. Single-cell RNA sequencing technologies and bioinformatics

304        pipelines. Exp Mol Med. 2018;50(8):96. Published 2018 Aug 7. doi:10.1038/s12276-018-

305        0071-8

306

307    2.   Lin C, Jain S, Kim H, Bar-Joseph Z. Using neural networks for reducing the dimensions

308        of single-cell RNA-Seq data. Nucleic Acids Res. 2017;45(17):e156.

309

310    3.   Shaham U, Stanton KP, Zhao J, et al. Removal of batch effects using distribution-

311        matching residual networks. Bioinformatics. 2017;33(16):2539-2546.

312

313    4.   Lopez, R., et al., Deep generative modeling for single-cell transcriptomics. Nature

314        Methods, 2018. 15(12): p. 1053-1058.

315

316     5.  Cho, H., Berger, B. & Peng, J. Generalizable and Scalable Visualization of SingleCell

317         Data Using Neural Networks. Cell Syst 7, 185-191 e184, doi:10.1016/j.cels.2018.05.017

318         (2018).

319

320     6.  Butler, A., Hoffman, P., Smibert, P., Papalexi, E.  Satija, R. Integrating single-cell

321         transcriptomic data across different conditions, technologies, and species. Nat.

322         Biotechnol. 36, 411–420 (2018).

323

324     7.  Haghverdi L, Lun ATL, Morgan MD, Marioni JC. Batch effects in single-cell RNA-

325         sequencing data are corrected by matching mutual nearest neighbors. Nat Biotechnol.

326         2018;36(5):421-427.

327

328     8.  Winkels H, Ehinger E, Vassallo M, et al.. Atlas of the immune cell repertoire in mouse

329         atherosclerosis defined by single-cell RNA-sequencing and mass cytometry.Circ Res.

330         2018;122:1675–1688. doi: 10.1161/CIRCRESAHA.117.3125

331

332     9.  Zheng GX, Terry JM, Belgrader P, et al. Massively parallel digital transcriptional

333         profiling of single cells. Nat Commun. 2017;8:14049. Published 2017 Jan 16.

334         doi:10.1038/ncomms14049

335

336     10. Xavier Glorot, Yoshua Bengio. Proceedings of the Thirteenth International Conference

337         on Artificial Intelligence and Statistics, PMLR 9:249-256, 2010.

338

339

340

**Figures and Tables**

**Fig 1. Neural network configuration.**

**Fig 2. Training and testing of the neural network on the Tabula Muris Atlas.** (a) Cell types obtained from the TMA. (b) Number of cells obtained for each cell type from each technique. (c) The same cell type tends to cluster separately by techniques. (d) Training and testing accuracy of the neural network when trained and tested using cells processed by the same technique. (e) Training and testing accuracy after each epoch when trained with 5,000 10X cells and tested with 5,000 SS2 cells. (f) Cost after each epoch when trained with 5,000 10X cells and tested with 5,000 SS2 cells.

**Fig 3. Neural network predicts cell types for human and mouse datasets.** (a) Cell types predicted by the neural network for the mouse leukocyte dataset. (b) Cell types identified by unsupervised clustering and canonical markers for the mouse leukocyte dataset. (c) Violin plots showing 3 genes' expression level in the NK and T cells from the mouse leukocytes. (d) Cell types predicted by the neural network for the human PBMC dataset. (e) TSNE plots showing 4 marker genes' expression for the human PBMC dataset.

**Fig 4. Neural network predicts sub cell types.** (a) TSNE plots showing 6 maker genes' expression for the reference T cell subtypes. (b) T cell subtypes obtained to train the neural network. (c) T cells from the human PBMC were grouped into 7 clusters by the unsupervised
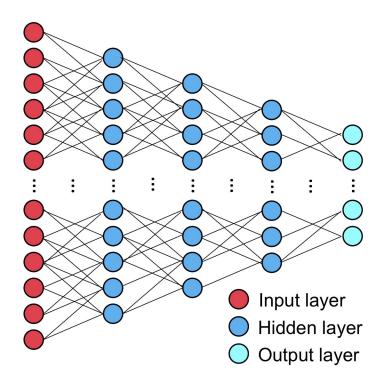
362    method. (d) Subtypes predicted for the T cells from the human PBMC. (e) Dot plot showing the

363    expression of 6 genes for the predicted subtypes, dot size represents the percentage of cells
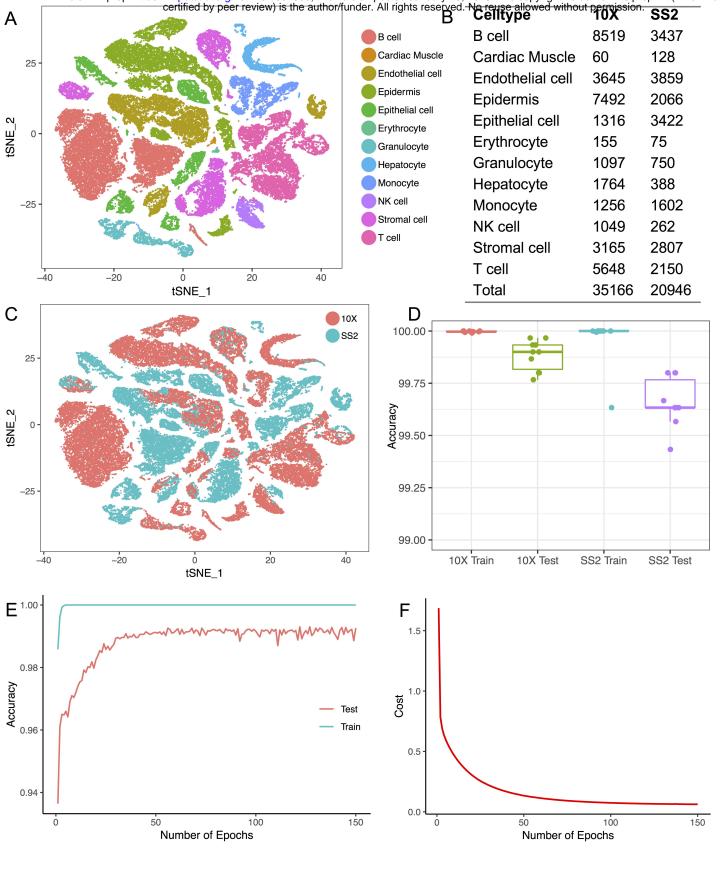
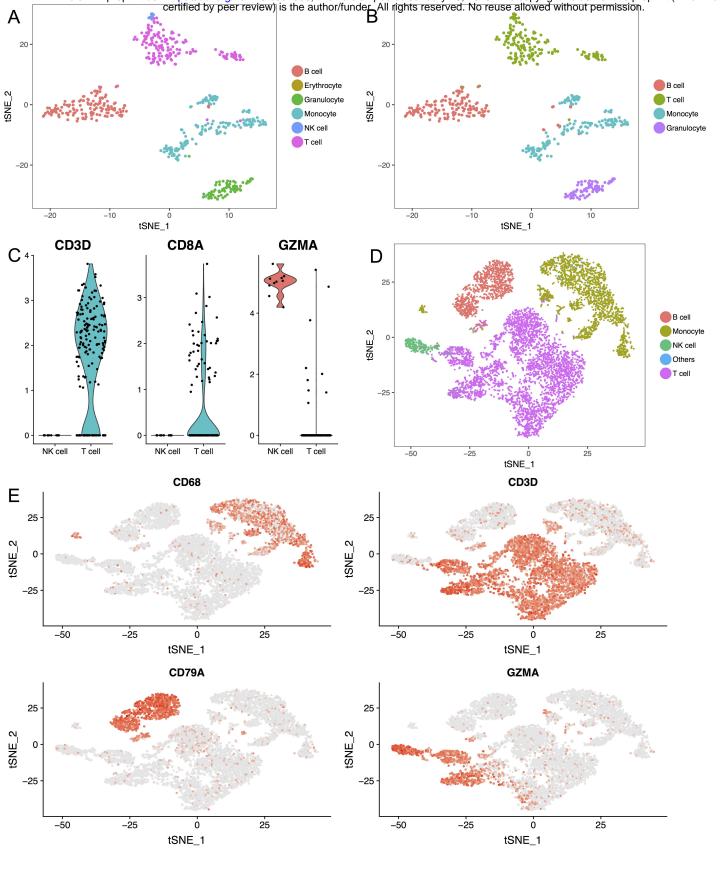364    expressing the gene, color scale represents the expression level of the gene.

365

## Supporting Information

367    **Supplementary File 1.** This file contains 4 tables that tells the number of accurately and

368    inaccurately predicted cells when the neural network was trained and test across the 10X and

369    SS2 datasets.

Input layer

Hidden layer

Output layer

A



B

| Celltype | 10X | SS2 |
|---|---|---|
| B cell | 8519 | 3437 |
| Cardiac Muscle | 60 | 128 |
| Endothelial cell | 3645 | 3859 |
| Epidermis | 7492 | 2066 |
| Epithelial cell | 1316 | 3422 |
| Erythrocyte | 155 | 75 |
| Granulocyte | 1097 | 750 |
| Hepatocyte | 1764 | 388 |
| Monocyte | 1256 | 1602 |
| NK cell | 1049 | 262 |
| Stromal cell | 3165 | 2807 |
| T cell | 5648 | 2150 |
| Total | 35166 | 20946 |

- B cell
- Cardiac Muscle
- Endothelial cell
- Epidermis
- Epithelial cell
- Erythrocyte
- Granulocyte
- Hepatocyte
- Monocyte
- NK cell
- Stromal cell
- T cell

C



D



E



F

A

CD3D
HLA-DQB1
CCR7
CD8A
S100A11
NKG7

B

Naive Cytotoxic T cell
Regulatory T cell
Naive T cell
Memory T cell
Cytotoxic T cell

Memory T cell
Naive Cytotoxic T cell
Naive T cell
Cytotoxic T cell
Regulatory T cell

C

0
1
2
3
4
5
6

D

Cytotoxic T cell
Memory T cell
Naive Cytotoxic T cell
Naive T cell
Regulatory T cell

E

pct.exp
25
50
75

avg.exp
1
0
-1