

1 **SSMD: A semi-supervised approach for a robust cell type identification and deconvolution**
2 **of mouse transcriptomics data**

3

4 Xiaoyu Lu^{1,2+}, Szu-Wei Tu^{1,2+}, Wennan Chang^{1,3}, Changlin Wan^{1,3}, Jiashi Wang¹, Yong Zang^{1,5},
5 Baskar Ramdas⁴, Reuben Kapur⁴, Xiongbin Lu^{1*}, Sha Cao^{1,2,5*}, Chi Zhang^{1,2*}

6

7 ¹Department of Medical and Molecular Genetics and Center for Computational Biology and
8 Bioinformatics, ⁴Department of Pediatrics, ⁵Department of Biostatistics, Indiana University
9 School of Medicine, Indianapolis, IN,46202, USA. ²Department of BioHealth Informatics,
10 Indiana University-Purdue University Indianapolis, Indianapolis, IN, 46202, USA. ³Department
11 of Electrical and Computer Engineering, Purdue University, Indianapolis, IN, 46202, USA
12

13 *To whom correspondence should be addressed: Chi Zhang, +1 317-278-9625; Email:
14 czhang87@iu.edu. Correspondence is also addressed to Xiongbin Lu, Email: xiolu@iu.edu; Sha
15 Cao, Email: shacao@iu.edu.

16 +These authors made equal contribution to this work.

17

18 **Authors' Biographical Note:**

19 Xiaoyu Lu is a PhD student in the Department of BioHealth Informatics, Indiana University-Purdue
20 University Indianapolis.

21 Szu-Wei Tu is a master student in the Department of BioHealth Informatics, Indiana University-Purdue
22 University Indianapolis.

23 Wennan Chang is a PhD student in the Department of Electrical and Computer Engineering, Purdue
24 University.

25 Changlin Wan is a PhD student in the Department of Electrical and Computer Engineering, Purdue
26 University.

27 Jiashi Wang is a research associate at the Biomedical Data Research Data (BDRD) Lab at Indiana
28 University School of Medicine.

29 Yong Zang is an Assistant Professor in the Department of Biostatistics and a member of the Center for
30 Computational Biology and Bioinformatics, Indiana University School of Medicine.

31 Baskar Ramdas is an Assistant Research Professor in the Department of Pediatrics, Indiana University
32 School of Medicine.

33 Reuben Kapur is Frieda and Albrecht Kipp Professor in the Department of Pediatrics, Indiana University
34 School of Medicine.

35 Xiongbin Lu is Vera Bradley Foundation Professor of Breast Cancer Innovation and Professor in the
36 Department of Medical and Molecular Genetics, Indiana University School of Medicine.

37 Sha Cao is an Assistant Professor in the Department of Biostatistics and a member of the Center for
38 Computational Biology and Bioinformatics, Indiana University School of Medicine.

39 Chi Zhang is an Assistant Professor in the Department of Medical and Molecular Genetics and a member
40 of the Center for Computational Biology and Bioinformatics, Indiana University School of Medicine.

41

42 **Key points:**

- 43 • We provide a novel tissue deconvolution method, namely SSMD, which is specifically designed for
44 mouse data to handle the variations caused by different mouse strain, genetic and phenotypic
45 background, and experimental platforms.
- 46 • SSMD is capable to detect data set and tissue microenvironment specific cell markers for more than
47 30 cell types in mouse blood, inflammatory tissue, cancer, and central nervous system.
- 48 • SSMD achieve much improved performance in estimating relative proportion of the cell types

49 compared with state-of-the-art methods.

- 50 • The semi-supervised setting enables the application of SSMD on transcriptomics, DNA methylation
51 and ATAC-seq data.
52 • A user friendly R package and a R shiny of SSMD based webserver are also developed.

53

54 **Keywords:**

55 Tissue Data Deconvolution, Cancer microenvironment, Semi-supervised Learning, Mouse Omics Data

56

57 **ABSTRACT**

58 Deconvolution of mouse transcriptomic data is challenged by the fact that mouse models
59 carry various genetic and physiological perturbations, making it questionable to assume fixed cell
60 types and cell type marker genes for different dataset scenarios. We developed a Semi-Supervised
61 **Mouse data Deconvolution (SSMD)** method to study the mouse tissue microenvironment (TME).
62 SSMD is featured by (i) a novel non-parametric method to discover data set specific cell type
63 signature genes; (ii) a community detection approach for fixing cell types and their marker genes;
64 (iii) a constrained matrix decomposition method to solve cell type relative proportions that is
65 robust to diverse experimental platforms. In summary, SSMD addressed several key challenges in
66 the deconvolution of mouse tissue data, including: (1) varied cell types and marker genes caused
67 by highly divergent genotypic and phenotypic conditions of mouse experiment, (2) diverse
68 experimental platforms of mouse transcriptomics data, (3) small sample size and limited training
69 data source, and (4) capable to estimate the proportion of 35 cell types in blood, inflammatory,
70 central nervous or hematopoietic systems. In silico and experimental validation of SSMD
71 demonstrated its high sensitivity and accuracy in identifying (sub) cell types and predicting cell
72 proportions comparing to state-of-the-arts methods. A user-friendly R package and a web server
73 of SSMD are released via <https://github.com/xiaoyulu95/SSMD>.

74

75 **INTRODUCTION**

76 The mouse has long served as the premier model organism for studying human biology
77 and disease, due to their striking genetic homologies and physiological similarity to humans, as
78 well as the relatively low cost of maintenance. Currently, thousands of unique inbred strains and
79 genetically engineered mutants have been made available for a wide array of specific disease types
80 [1]. Research on mouse models have provided added impetus and indispensable tool for studying
81 human disease, regarding its initiation, maintenance, progression and response to treatment, as
82 well as evaluating drug safety and efficacy [2] [3]. Amongst all, the ability to examine
83 physiological states and interactions between diseased cells and their microenvironment in vivo
84 represents the most important tool for studying disease dynamics. To this end, numerous omics
85 data have been collected from mouse that vary in terms of genetic perturbations, cell/tissue types,
86 and treatment conditions [4-7]. A strong computational capability is needed to study the
87 interactions of components within the mouse tissue microenvironment subject to different genetic
88 and physiological perturbations, the knowledge gained from which could be projected to human
89 disease scenarios and provide invaluable insight and guidance for effective human therapeutic
90 regimes.

91 Tissue transcriptomic data display convoluted signals from different cell types [8].
92 Deconvoluting cell components and identifying mouse strain-/tissue-/experimental condition-
93 specific cell types and gene expressions are crucial for understanding how experimentally
94 perturbed conditions are associated with cellular level characteristics and cell-cell interactions [9].
95 While multiple deconvolution methods have been developed for investigating the heterogeneous
96 cell types in human cancer or other tissues data [10-19], they may not be directly applicable to

97 mouse tissue data. First of all, the cell type specific genes for human cells differ from mouse cells;
98 secondly, compared with human, the variations among different mouse tissue samples may be
99 considerably higher, as they are collected from different strains with varied genetic background
100 and experimental conditions.

101 Currently, ImmuCC and its varied versions are the only method specifically focusing on
102 mouse data deconvolution [20]. The core computational algorithm, which was adapted from
103 CIBERSORT designed for human [13], assumes fixed cell type and signatures gene expressions
104 (subject to simple transformations) regardless of experimental conditions of the target data. This
105 assumption becomes problematic as mouse data, which are collected from different strains, have
106 varied genetic background, thus, it is expected the tissue compositions are highly adaptable
107 regarding the existent cell types and their expression profiles [21-23]. Aside from prominent
108 variability in the appearance of cell types and the expression levels of markers genes, mouse data
109 deconvolution also suffers from the following challenges: diverse experimental platforms,
110 prevalently small sample size of mouse experiments, and limited training data sets available for
111 deriving signature genes of cell types.

112 To address these challenges, we developed a novel semi-supervised deconvolution method,
113 namely Semi-Supervised Mouse data Deconvolution (SSMD), to infer data/tissue specific cell type
114 marker genes and their expression profiles and estimate their relative abundances from
115 transcriptomics data. SSMD is capable to infer the relative proportion of 35 cell types in the blood,
116 inflammatory, cancer, central nervous system and hematopoietic system. To the best of our
117 knowledge, SSMD is the only mouse data deconvolution method considering strain, tissue type
118 and data specificity of cell type specific gene markers. We demonstrated SSMD achieved a high
119 sensitivity in identifying the appearance of immune and stromal cell types in inflammatory tissue
120 and brain cell types in central nervous tissue, and with a high accuracy in estimating their relative
121 proportion on single cell RNA-seq simulated bulk tissue data sets. We also experimentally
122 validated that the cell populations inferred by SSMD accurately recapitulates the true cell
123 proportions measured by fluorescence-activated cell sorting (FACS) on a leukemia bone marrow
124 data. Applications of SSMD on a large collection of public mouse blood, brain, cancer, and other
125 inflammatory tissue data suggested that the method achieved a robust performance throughout
126 diverse types of experimental conditions and platforms including RNA-seq, microarray and
127 immuno-assay. In addition, the software of SSMD grants users to build in their own tissue/data
128 specific knowledge of cell type specific markers to reinforce the method. An R package of SSMD
129 is released through GitHub: <https://github.com/xiaoyulu95/SSMD> and a R Shiny based web server
130 of SSMD is available at <https://ssmd.cccb.iupui.edu/>.

131

132 RESULTS

133 *Mathematical consideration and problem formulation*

134 Denote $\tilde{X}_{M \times N}$ as a tissue data of M genes and N samples, a deconvolution analysis
135 assumes $\tilde{X}_{M \times N}$ as the following non-negative product form:

$$136 \quad \tilde{X}_{M_0 \times N} = \tilde{S}_{M_0 \times K_0} \cdot \tilde{P}_{K_0 \times N} + E, \tilde{S}_{M_0 \times K_0} \geq 0, \tilde{P}_{K_0 \times N} \geq 0 \quad (1)$$

137 Here, $\tilde{X}_{M_0 \times N}$ represents the observed gene expression matrix of M_0 selected genes (a subset in M)
 138 in N tissue samples, and columns in $\tilde{S}_{M_0 \times K_0}$, and rows in $\tilde{P}_{K_0 \times N}$, denote the expression signatures,
 139 and the relative proportions of the K_0 cell types respectively. In the conventional formulation of
 140 deconvolution analysis, with fixed M_0 and K_0 , $\tilde{S}_{M_0 \times K_0}$ and $\tilde{P}_{K_0 \times N}$ are solved to minimize the \mathcal{L}_2
 141 loss of the above linear equation. Because of the highly varied genetic and phenotypic background
 142 of mouse experiment, $\tilde{S}_{M_0 \times K_0}$, M_0 and K_0 are usually varied and unknown, i.e. for each $\tilde{X}_{M \times N}$
 143 collected from tissues of certain microenvironment, what cell types are present, what gene markers
 144 each cell type expresses and how much they were expressed, could vary drastically due to the
 145 genetic and physiological perturbations. Correctly specified cell types K_0 , and selected cell type
 146 marker genes M_0 can largely increase the prediction accuracy of $\tilde{P}_{K_0 \times N}$. **Table 1** lists the key
 147 mathematical definitions utilized in this study.

148 In this study, we define a cell type k is “transcriptomically identifiable” if its ground-truth
 149 proportion $P_{1 \times N}^k$ and estimated as $\tilde{P}_{1 \times N}^k$ have high correlation, i.e.. $cor(P_{1 \times N}^k, \tilde{P}_{1 \times N}^k) = 1 - \epsilon$ and
 150 ϵ is substantially small, where $\tilde{P}_{1 \times N}^k$ is the k th row of $\tilde{P}_{K_0 \times N}$, and K_0 as the number of “identifiable”
 151 cell types. A strong condition for a cell type to be identifiable is that it has uniquely expressed
 152 genes [24]. Here we provided a comprehensive mathematical derivation of the relationship
 153 between cell type unique expression and identifiability of cell proportion in the **Supplementary**
 154 **Notes**. We derived the identity of cell type uniquely expressed gene markers, denoted as the set
 155 G_k , is a necessary but non-sufficient condition for the identifiability of cell type k : – if k is
 156 “transcriptomically identifiable”, $\tilde{X}_{G_k \times T}$ must be a matrix of rank one, for $\forall T \subset \{1, \dots, N\}$. This
 157 condition forms the foundation of how SSMD discover cell type marker genes that are not fixed,
 158 but instead specific to each dataset. Fortunately, we do not need to scan for all the local rank-1
 159 matrices within $\tilde{X}_{M \times N}$, where M is usually to the tens of thousands. In fact, with an effective
 160 knowledge transfer of the gene labels derived from single or bulk cell training data, the genes that
 161 are more likely to be cell type specific markers of identifiable cell types can be detected, which
 162 forms the core algorithm of SSMD pipeline.

163

164 **Table1. Definition of mathematical terms**

Terminology	Mathematical Definition in this study
Rank-1 matrix	A matrix with rank = 1, i.e. the matrix is generated by the product of two vectors, $X = A \cdot B^T$. In this study, we consider all transcriptomics data are with error. Hence the rank-1 matrix is defined by $X = A \cdot B^T + E$, where the matrix rank of X is 1 can be computed by the bi-cross validation (BCV) algorithm detailed in Methods.
Local rank-1 matrix	A submatrix with rank = 1, i.e. denoting I and J as the indices of the submatrix, $X_{I \times J}$ is generated by the product of two vectors with error, $X_{I \times J} = A \cdot B^T + E$.
Transcriptomically identifiable cell type	The cell type with a high correlation between the true proportion $P_{1 \times N}^k$ and estimated $\tilde{P}_{1 \times N}^k$
Prediction accuracy	Pearson correlation between true proportion and predicted proportion of each cell type

Detection accuracy The number of true cell type signature genes were identified as signature genes of an identifiable cell type

Matrix total Rank The total rank of a data matrix that can be tested by the BCV algorithm

165

166 *SSMD Analysis pipeline*

167 SSMD is a semi-supervised method composed by (1) training a large candidate list of cell
168 type specific marker genes, (2) evaluating the identifiability of each cell type and confirming their
169 marker genes for each to-be-deconvolved data, and (3) estimating the proportion of each cell type.

170 The training step is to look for genes that are more likely to serve as cell type marker genes
171 through different tissue types and data sets, named as core marker lists. Specifically, we identified
172 the genes that are commonly over expressed in one cell type comparing to the others in bulk cell
173 data and commonly form rank-1 matrices in tissue data, by using a very extensive set of training
174 data sets collected from different mouse strains and tissue types (see details in Methods). **Fig 1A**
175 illustrates the procedure of SSMD to construct cell type core marker lists. On the bulk cell training
176 data, we adopted a random-walk based approach to detect genes that are significantly expressed in
177 higher quantities in one or a few cell types, than others (see details in Methods). As a result, a
178 labeling matrix that annotates cell type specifically expressed genes will be constructed, which
179 forms the first evidence of the potential marker genes for each cell type. Then on each bulk training
180 tissue dataset, we further identified marker genes that form rank-1 submatrices with a community
181 detection approach as detailed in methods. Only those modules, whose genes significantly and
182 consistently over-represent one and only one cell type across multiple training tissue datasets, are
183 selected to form the core marker list. Noted, variations caused by different experiment batches,
184 tissue types and mouse strains were handled by enabling certain errors in the random-walk based
185 cell type specific marker identification, i.e. identifying the genes overly expressed in the cell type
186 comparing to the others in a certain proportion of the collected bulk cell data. In addition, data
187 batch variation was also considered in the bulk data based training step, by identifying the genes
188 commonly serve as cell type specific marker in more than 50% of analyzed bulk tissue training
189 data. The goal of this training procedure is to summarize a relatively large list of commonly
190 observed cell type specific marker genes, which can be used to as semi-supervised information to
191 identify data set specific cell type marker for a further un-supervised deconvolution analysis.

192 Based on the cell type core markers, the deconvolution of any given bulk tissue dataset is
193 composed by the steps as illustrated in **Fig 1B**. SSMD first identifies all the rank-1 modules on the
194 target dataset by an iterative hierarchical clustering and bi-cross validation approach. Then SSMD
195 selects the rank-1 modules that are likely to be markers of a certain cell type for this data set, if
196 genes in the modules largely overlap with the core marker list of one and only one cell type.
197 Modules that are highly co-linear will be merged. Consequently, genes in each module is called
198 gene markers of one cell type, that satisfy the necessary condition for “transcriptomically
199 identifiable”. Notably, two modules may represent the same cell type, and they are treated as

200 marker genes of different subtypes of the cell type. Here, the total number of modules is an estimate
201 of the number of “identifiable” cell types, i.e., K_0 . Importantly, SSMD is an “semi-supervised”
202 approach, because the cell marker genes do not solely depend on the training data, but also the co-
203 expression patterns of the marker genes in the target dataset. In other words, SSMD addresses the
204 variability issue of signature genes from one dataset to another, and has the potential to discover
205 cell types not pre-defined. Algorithms of each computational step are detailed in Materials and
206 Methods. Complete flow chart of the SSMD pipeline is provided in **Supplementary Fig S1**.

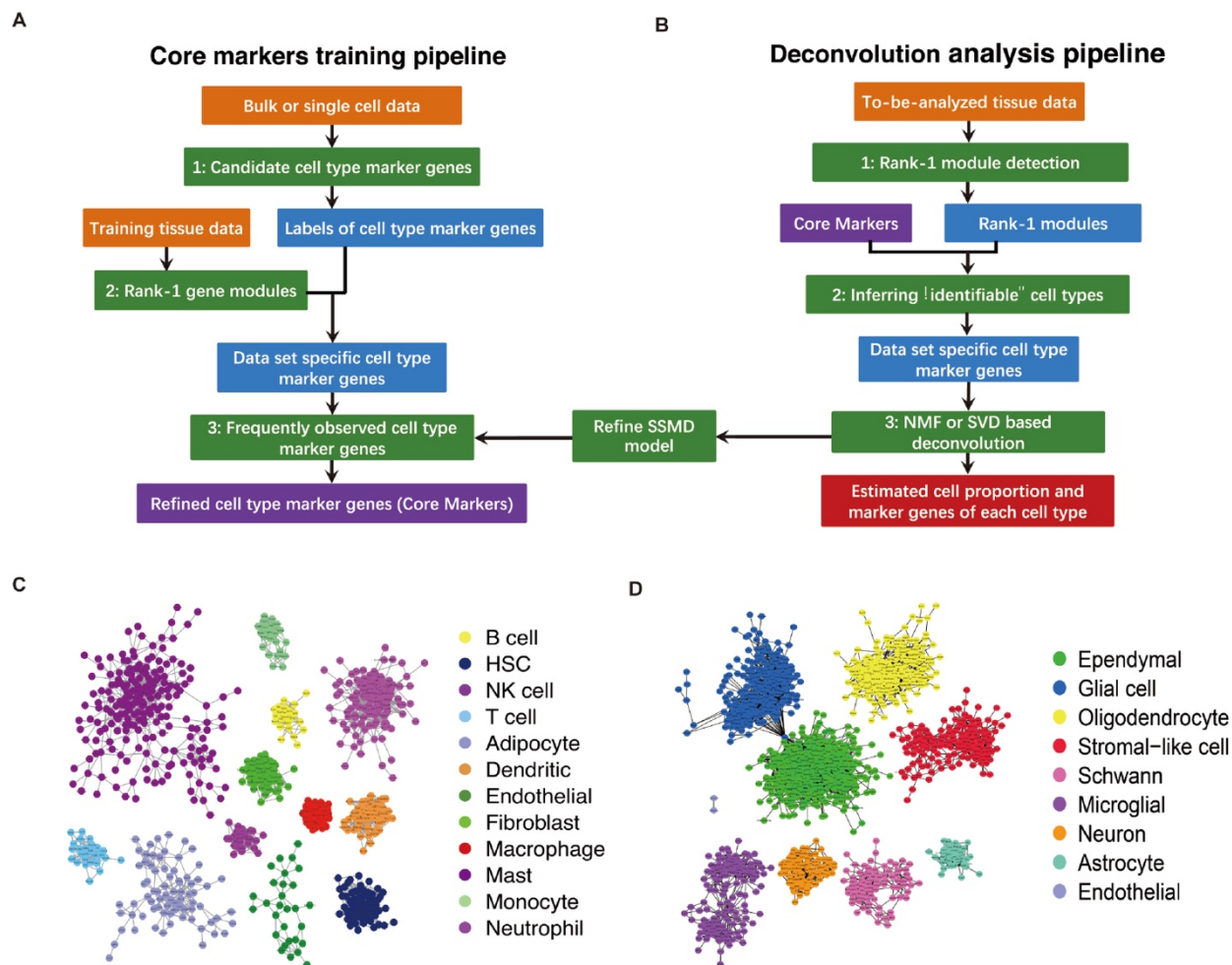
207 The prediction of the cell type proportions is conducted using a constrained Non-negative
208 Matrix Factorization (NMF) method by solving the following optimization problem:

$$209 \min_{\tilde{S}_{M_0 \times K_0}, \tilde{P}_{K_0 \times N}} \left(\|\tilde{X}_{M_0 \times N} - \tilde{S}_{M_0 \times K_0} \cdot \tilde{P}_{K_0 \times N}\|_F^2 + \lambda \cdot \text{trace} \left(\tilde{S}_{M_0 \times K_0}^T \cdot (\mathbf{1}_{M_0} \mathbf{1}_{K_0}^T - C_{M_0 \times K_0}) \right) \right) \quad (2)$$

210 , where $C_{M_0 \times K_0}[i, j] = 1$ if gene i is marker of the cell type j , and 0 otherwise. $\mathbf{1}_d$ denotes
211 an all-1 column vector of length d , λ is a hyperparameter selected by cross validation, and other
212 annotations follow equation (1). The constraint matrix $C_{M_0 \times K_0}$ is enforced upon the regular NMF
213 to guarantee similarity of the solved signature matrix $\tilde{S}_{M_0 \times K_0}$ and constraint $C_{M_0 \times K_0}$, namely, in
214 the k th column of $\tilde{S}_{M_0 \times K_0}$, it should have higher expressions for genes that are markers of cell type
215 k . The solution to (2) is by alternative update where each time one of $\tilde{S}_{M_0 \times K_0}, \tilde{P}_{K_0 \times N}$ is held fixed,
216 and the other is updated. λ can be tuned by using simulated tissue data with known cell proportion.
217 In this study, we tuned λ and empirically select λ as 10 when $\tilde{X}_{M_0 \times N}$ is log normalized microarray
218 data or log(X+1) normalized FPKM/CPM/TPM RNA-seq data.

219 Following these procedures, and on a large collection of mouse bulk cell and tissue training
220 data, we generated core marker gene lists for different tissue microenvironments: (1) for mouse
221 blood, solid cancer and inflammatory tissues, 980 genes of 12 cell types namely T cell, B cell, NK
222 cell, hematopoietic stem cell, monocyte, macrophage, neutrophil, mast cell, adipocytes, fibroblast,
223 dendritic cell, and endothelial cell were discovered (**Fig 1C**); (2) for mouse hematopoietic system,
224 2877 genes of 14 cell types namely hematopoietic stem cell, common lymphoid progenitor,
225 granulocyte-macrophage progenitors, megakaryocyte lineage-committed progenitor, erythroid cell,
226 megakaryocyte-erythrocyte progenitors, multipotent progenitors, early myeloid progenitor, mature
227 myeloid cell, pre colony forming unit erythroid, pre-megakaryocytic/erythroid progenitor, B cell,
228 CD4+ T and CD8+ T cell were discovered (**Supplementary Table S1**), and (3) for mouse central
229 nervous system tissue, 1570 genes of nine cell types namely ependymal cell, general glial cell,
230 oligodendrocyte, stromal-like cell, Schwann cell, microglial, neuron, and astrocyte were
231 discovered (Fig 1D). Complete lists of the core marker genes are given in **Supplementary Table**
232 **S1**. It is noteworthy that the size of core marker list ranges from 27 to 547 for different cell types.
233 However, our analysis suggested that more than 5 marker genes that form a rank-1 matrix is
234 sufficient for an accurate estimation of cell proportion. Note that, compared with conventional
235 regression based deconvolution analysis, SSMD only uses labels of the core markers as the semi-
236 supervised information and identifies data set specific cell type markers for a further unsupervised
237 estimation of cell types, which grants a flexibility and robustness to handle the variation of cell
238 type specific marker genes and their expression scale through different mouse strains, tissue types

239 and experimental platforms. In addition, the semi-supervised formulation of SSMD enables the
 240 inference of identifiability of each cell type and identification of rare or sub cell types.
 241
 242



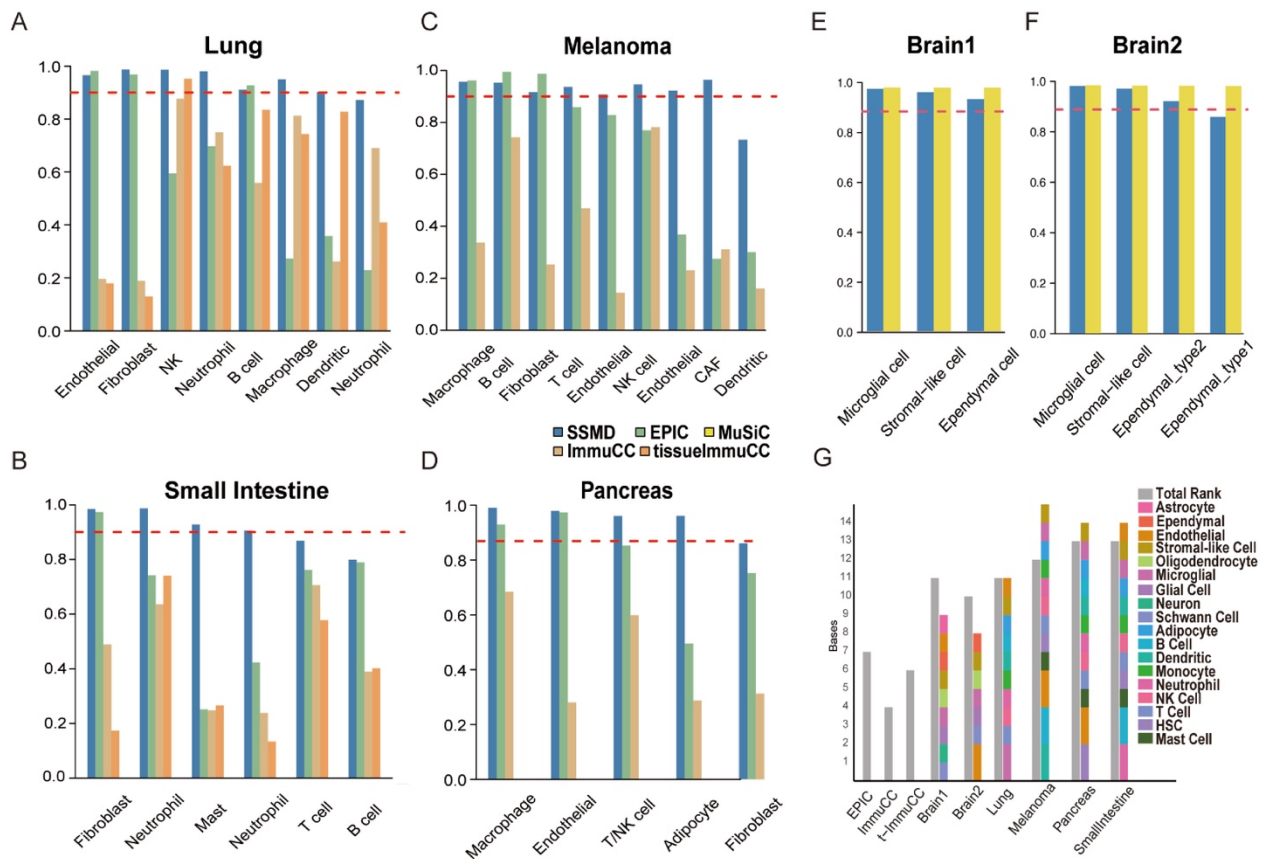
243
 244
 245
 246
 247
 248
 249
 250
 251
 252
 253
 254
 255
 256

Fig 1. Analysis pipeline of SSMD and core cell type specific markers. (A) Analysis pipeline of the core marker training procedure. (B) Analysis pipeline of the deconvolution procedure. In (A) and (B), input data including training and target data, computational procedure and key intermediate outputs were colored by orange, green and blue, respectively. (C) Core markers of 12 cell types in blood, solid cancer, and inflammatory tissue. An edge between two genes means the two genes are co-identified as markers of one cell type in more than 50% of the training data sets. (D) Core markers of 9 cell types in central nervous system. Noted, core markers for the endothelial cell in the inflammatory tissue and central nervous system were separately trained by comparing with other cell types in the same tissue system.

257 *Benchmarking based on artificial tissue data simulated by using single cell RNAseq data*

258 We first benchmarked SSMD on a set of artificial tissue data simulated from four single
 259 cell RNAseq (scRNA-seq) datasets of mouse lung, pancreas, small intestine and melanoma. For
 260 each data set, we simulated 100 tissue samples by randomly drawing and mixing cells of different
 261 types whose proportions follow random Dirichlet distributions. Prediction accuracy of each cell
 262 type was assessed by the Pearson correlation coefficients between its known mixing cell
 263 proportions and the predicted relative proportion. We compared SSMD with three state-of-arts
 264 deconvolution methods of mouse data, namely ImmuCC (ICC), tissue-ImmuCC (TICC) and EPIC
 265 [11]. Our analysis suggested that SSMD achieved 93.2% prediction accuracy on average in the
 266 four simulated data sets and 23 out of the 28 cell types (82.1%) are with higher than 0.9 prediction
 267 accuracy (**Fig 2A-D**). In contrast, EPIC, ICC and TICC achieved 69.7%, 45.2% and 48.5%
 268 averaged prediction accuracy on the cell types covered by these methods, and the proportion of
 269 cell types with higher than 0.9 prediction accuracy are 32.2% (9/28), 0% (0/28) and 7.2% (1/14),
 270 respectively. We also tested the popular human data deconvolution methods such as CIBERSORT
 271 (CIBERSORTx) and TIMER [9, 13], by using the known human and mouse homolog genes. Non-
 272 surprisingly, predictions made by CIBERSORT and TIMER on the mouse are less accurate than
 273 SSMD. TIMER and CIBERSORT achieved 49.25% and 47.5% averaged prediction accuracy, and
 274 the proportion of cell types with higher than 0.9 prediction accuracy are 17.9% (5/28) and 3.6%
 275 (1/28) (**Supplementary Table S4**).

276



277

278 **Fig 2. Method evaluation on scRNA-seq simulated tissue data.** (A-D) Correlation between true
279 and predicted cell proportions in the simulated Lung (A), Pancreas (B), Small Intestine (C), and
280 Mouse Melanoma (D) tissue data. The x-axis represents cell type and y axis represents prediction
281 accuracy. Predictions made by SSMD, EPIC, ImmuCC and tissue-ImmuCC were dark blue, green,
282 yellow and orange colored, respectively. The red dash line represents the 0.9 correlation cutoff.
283 (E-F) Correlation between true and predicted cell proportions in the two simulated brain tissue
284 data. (G) The total rank of the gene expression profile of selected marker genes in the six simulated
285 tissue data (grey), and the total number of cell types identified by SSMD in each data set or
286 assumed in other methods (left three grey bars).

287

288 It is noteworthy that the SSMD enables the detection of sub cell types defined as
289 transcriptomically identifiable. SSMD successfully identified two sub populations of fibroblast
290 cells in the melanoma data and different subtypes of neutrophils in lung and small intestine data.
291 In contrast, ICC, TICC and EPIC are not capable of providing cell subtype predictions due to their
292 fixed cell type assumption.

293 We also benchmarked SSMD on simulated brain tissue data using two scRNA-seq data of
294 central nervous systems. SSMD achieved more than 0.9 correlation in predicting the cell types
295 microglial, stromal-like, and ependymal subtypes in the simulated tissue data (**Fig 2E-F**). To the
296 best of our knowledge, SSMD is the first of its kind method to specifically target mouse central
297 nervous system decomposition. To benchmark SSMD, we selected MUSIC as the state-of-the-art
298 method, which requires an additional input of an scRNA-seq data to train context specific gene
299 signatures [25]. Here we first utilized the same scRNA-seq data for tissue data simulation and
300 signature training in MUSIC. Non-surprisingly, MUSIC achieved consistently good predictions
301 (averaged $\text{cor}=0.99$), and the predictions made by SSMD are very close to MUSIC with slightly
302 lower correlations compared with MUSIC under this ideal setup. In sight the possible disparity
303 caused by tissue, strain, and experimental platform variations between the target tissue data and
304 available scRNA-seq data for training cell markers, we also conducted a robustness test of MUSIC
305 and SSMD (see details in Supplementary Notes). Our analysis suggested that MUSIC highly
306 depends on the consistency of cell type specific marker genes and their expression scale between
307 the target tissue and the training scRNA-seq data. In contrast, the de novo data set specific marker
308 identification by SSMD enables a broader application to the tissue data without matched scRNA-
309 seq data. Because EPIC, ImmuCC and tissue-ImmuCC cannot analyze brain tissue data and the
310 melanoma and pancreas tissue were not covered by tissue-ImmuCC, we did not include the
311 comparison with these methods on the brain tissue data.

312 To further validate the specificity of SSMD, we tested the total rank of the identified marker
313 genes and compared with the number identified cell types (TIMER and CIBERSORT achieved
314 49.25% and 47.5% averaged prediction accuracy. and the proportion of cell types with higher than
315 0.9 prediction accuracy are 17.9% (5/28), and 3.6% (1/28).). We also compare the total matrix
316 rank of the marker genes used in other methods and the number of cell types assumed in those
317 methods. Comparing to the fixed number of cell types in other methods, the number of cell types
318 predicted by SSMD better matches the total rank of the expression profile of identified marker
319 genes. Our observation suggested SSMD can correctly estimate the number of cell types and select

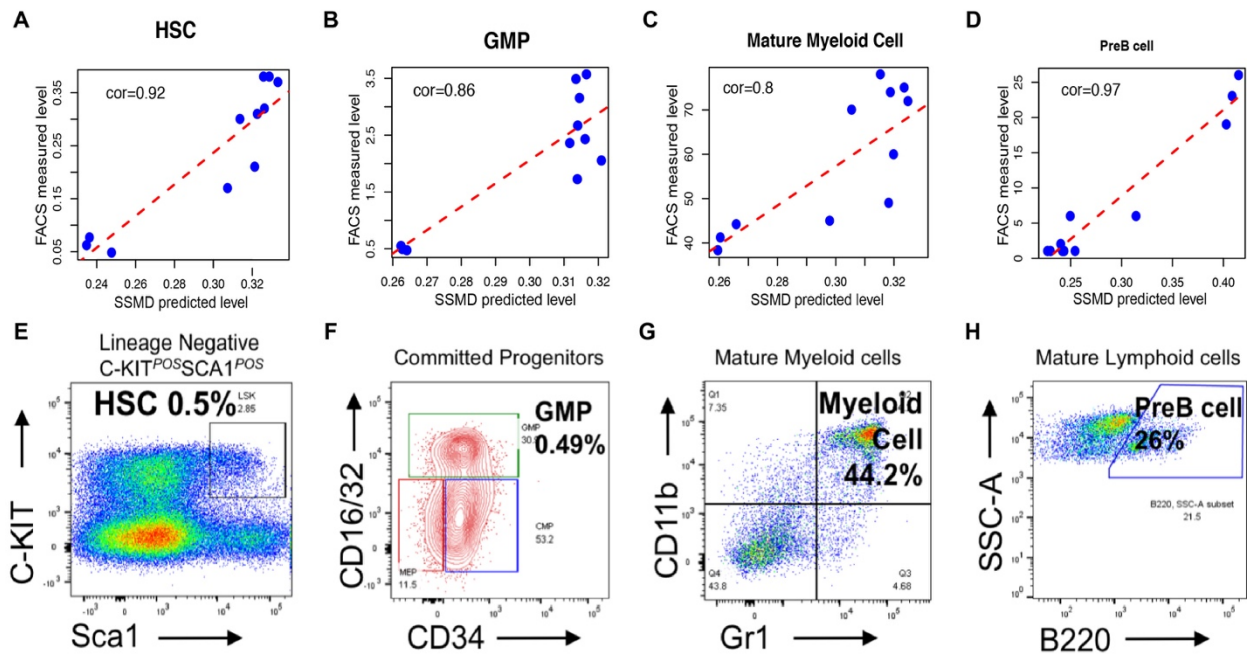
320 proper markers for cell type proportion estimation. It is noteworthy the predicted number of cell
 321 types may not exactly match the total rank of selected markers because possible co-linearity among
 322 the true proportion of the cell types.

323

324 *Experimental validation of SSMD by using matched RNA-seq and cell sorting data*

325 We generated a tissue RNA-seq data of 11 mouse bone marrow tissue samples with
 326 matched cell counting using Fluorescence activated cell sorting (FACS) (see details in Methods).
 327 Application of SSMD on the RNA-seq data identified hematopoietic stem cell (HSC), general
 328 myeloid progenitor (GMP), mature myeloid cell and Pre-B cells, and their cell type specific
 329 markers. We also observed that the correlation between SSMD predicted and FACS measured
 330 amount of HSC, GMP, mature myeloid cell and B cells are 0.92, 0.8, 0.86, and 0.97, respectively,
 331 suggesting a high prediction accuracy of SSMD. **Fig 3A-D** shows the correlation between the
 332 SSMD predicted cell proportion and the FACS measured cell proportion of the four cell types. **Fig**
 333 **3E-H** illustrate the FACS based cell counting of the four cell types. Complete cell type specific
 334 markers, cell proportions counted by FACS and predicted by SSMD were given in
 335 **Supplementary Table S2**. It is noteworthy that SSMD is not compared with other methods as
 336 none of the existing method is capable of predicting proportions of hematopoietic cell types.

337



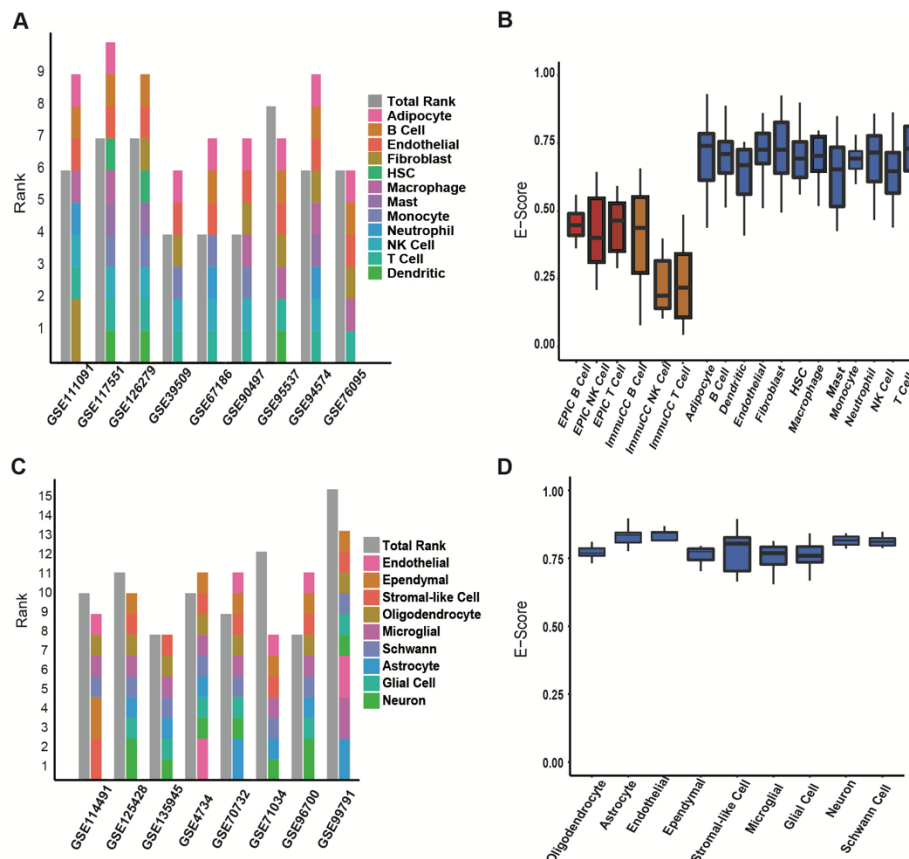
338

339 **Fig 3. Method evaluation on scRNA-seq simulated tissue data on hematopoietic tissue data.**

340 (A-D) Correlation between SSMD predicted (x-axis) and FACS identified (y-axis) cell proportions
 341 of HSC, GMP, mature myeloid cell and preB cell. (E-H) marker proteins utilized to identify the
 342 four cell types by using FACS. The x- and y- axis of the plots represent the level of cell type
 343 markers. The black block in (E), the green block in (F), the upper-right block in (G) and the block
 344 in (H) are the sorted HSC, GMP, Myeloid and Pre-B cell, respectively.

345 *Application of SSMD to real mouse tissue transcriptomics data*

346 We applied SSMD to nine cancer and eight central nervous system tissue data of four
 347 different experimental platforms, including one data set measured by immune-assay. On average,
 348 SSMD identified more than seven cell types in each of the cancer data, and the number of identified
 349 cell types is highly consistent with the total rank of the expression profile of the detected cell type
 350 specific marker genes (**Fig 4A**). This indicates that SSMD is capable of capturing the latent
 351 structure of the data. We further examined the explanation score (E-score), defined as the averaged
 352 absolute residual of the non-negative linear regression of each marker gene's expression on the
 353 predicted cell proportion, i.e. the average measure of how the predicted proportions could explain
 354 all the marker genes' expression levels. A high E-score is a necessary condition for an accurate
 355 cell proportion prediction. On average, the data set specific markers genes of each cell type
 356 identified by SSMD achieved 0.73 E-score while the average E-score of the marker genes used by
 357 EPIC and ImmuCC is 0.45 and 0.3 (**Fig 4B**). Similarly, application of SSMD on eight central
 358 nervous system tissue data identified more than seven cell types on average. The number of
 359 identified cell types is highly consistent with the total rank of the gene expression profile of the
 360 marker genes (**Fig 4C**). And the marker genes identified by SSMD achieved averaged 0.77 E-
 361 score for the cell types in central nervous system (**Fig 4D**). It is noteworthy that multiple marker
 362 sets of fibroblasts, myeloid or microglial cells that forming distinct rank-1 bases were identified in
 363 numerous data sets, suggesting the possible sub types of these cell types identified by SSMD.
 364



366
 367
 368
 369
 370
 371
 372
 373
 374
 375
 376
 377
 378
 379
 380
 381
 382
 383
 384
 385
 386
 387

Fig 4. Prediction of SSMD on real tissue data. (A, C) The total rank of the gene expression profile of selected marker genes (grey) in different (A) cancer tissue and (C) brain data, and the total number of cell types identified by SSMD in each data set (colored). (B, D) E-Score for different cell types identified by SSMD (blue) in (B) cancer and (D) brain data set or assumed in other methods (EPIC: red, ImmuCC: Yellow).

Robustness analysis

We first evaluated the variation of cell type specific markers through different mouse strains on one transcriptomic dataset of mouse liver tissue samples collected from 31 different mouse strains [26]. To the best of our knowledge, this is the only dataset in the public domain that systematically measured gene expression profiles of the same tissue type for different mouse strains by using the same experimental platform. SSMD was applied to the data of each mouse strain respectively. 9 cell and their sub types were commonly identified in the liver tissue of most strains. The identifiability of the cell types and the detected cell type markers among different strains were compared (Fig 5). We analyzed all the identified marker genes that form rank-1 modules, i.e. the necessary condition for gene markers of identifiable cell types, and noticed that only 9.1% of the identified marker genes are shared in more than 50% strains, while 58.4% of the identified marker genes only served as a cell type marker in less than 20% of the analyzed strains, suggesting a high variation of cell type specific markers among different mouse strains, and the necessity to consider strain or data set specificity in deconvolution analysis.

388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400

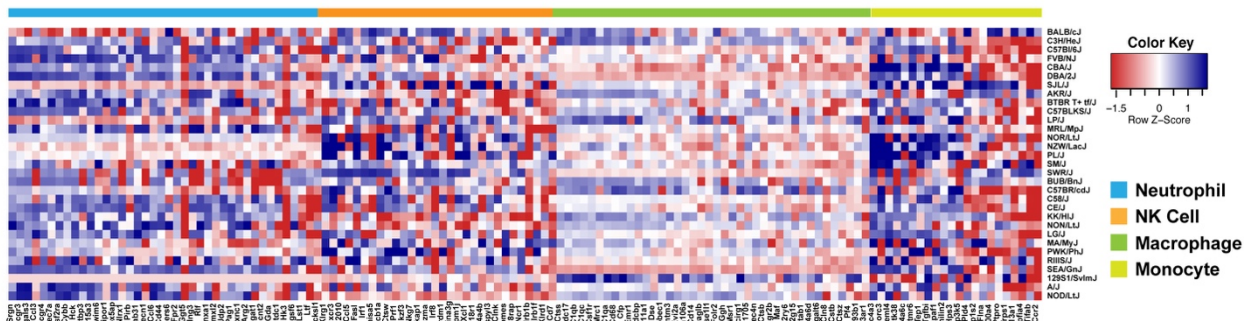
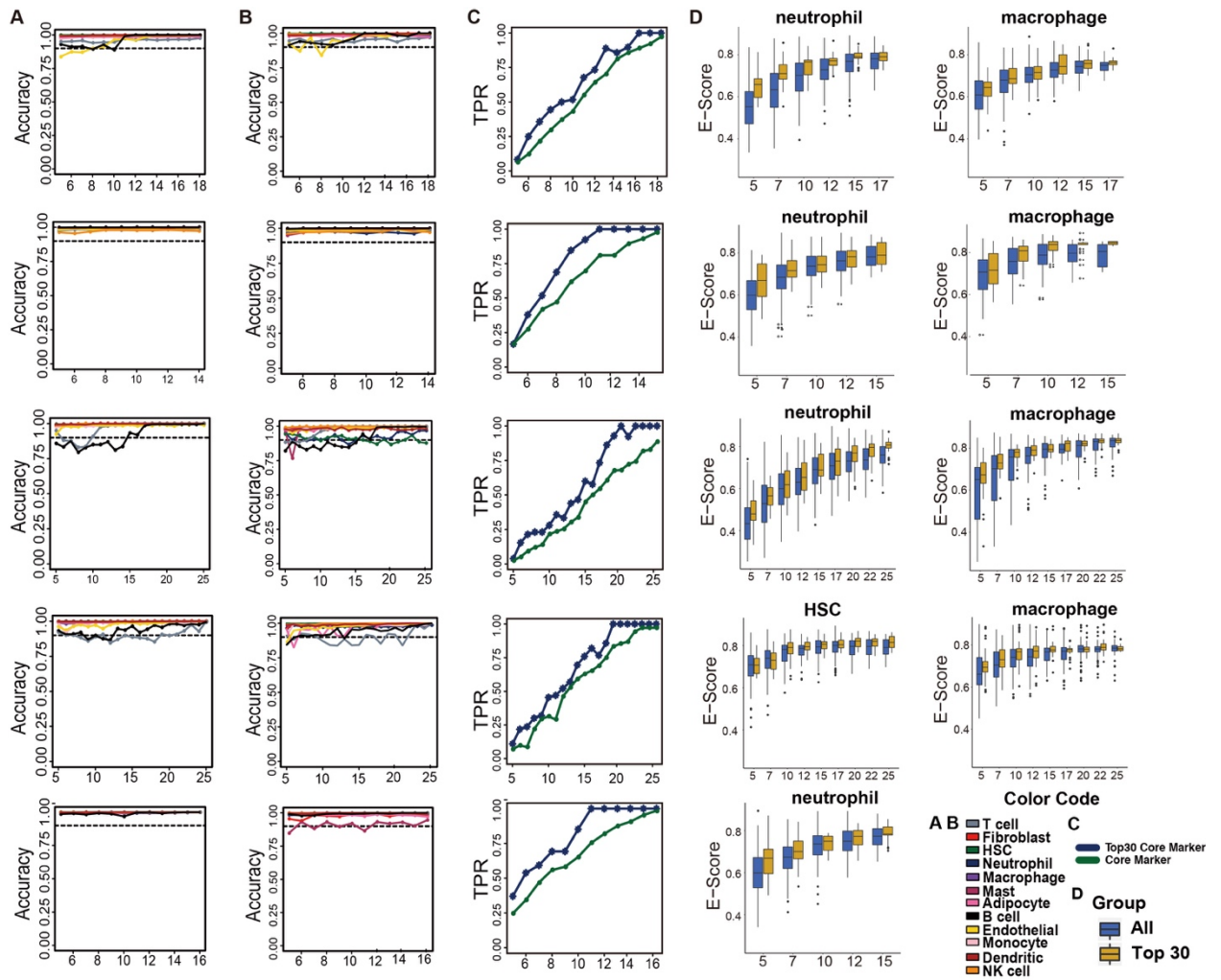


Fig 5. Correlation between expression level of strain specific cell type marker genes and predicted cell proportion. High correlation is a necessary but non-sufficient condition for the genes to serve as marker genes of the cell types in corresponding mouse strain. In the heatmap, x- and y-axis represent genes and mouse strains, respectively. Genes in the core marker list of four selected cell types, namely Neutrophil, Nature Kill (NK), Macrophage, and Monocyte, were colored on the column side bar.

We further examined the robustness of SSMD by evaluating its (1) sensitivity and (2) specificity in identifying cell types specific marker genes and its (3) accuracy in assessing of cell proportions, on the data of different sample sizes. Previous studies revealed that the robustness of the computation of co-expression correlation will decrease when the sample size is below 25. To comprehensively evaluate the method's robustness, we selected five data sets, namely GSE76095,

401 GSE67186, GSE90885, GSE94574, and GSE126279, with sample size ranging from 15 to 30 and
 402 randomly drew samples from each data set to build testing data sets of different sample size. We
 403 assumed the cell type markers and cell proportion inferred from whole data as “true” markers and
 404 proportions, and evaluated the consistency between the “true” ones and the ones predicted from
 405 small sub data sets. Accuracy in cell proportion prediction was assessed by the Pearson correlation
 406 between proportions predicted from small data and the “true” proportion on overlapped samples.
 407



408
 409
 410 **Fig 6. Performance evaluation of different sample size.** (A) Prediction accuracy (y-axis) in
 411 different sample size (x-axis) using all core markers. Accuracy is the Pearson correlation between
 412 predicted proportion using only selected small sample and using all samples. (B) Prediction
 413 accuracy (y-axis) in different sample size (x-axis) using selected stringent markers. (C) True
 414 positive rate (y-axis) of the cell type specific markers identified by using the stringent markers
 415 (blue) and core markers (green) with respect to different sample size (x-axis). (D) E-Score for
 416 using co-expression modules consisting of all core markers and only selected stringent markers.
 417 From top to bottom, the statistics were derived from GSE76095, GSE67186, GSE90885,
 418 GSE94574, and GSE126279.
 419

420 On average, all of the marker genes of the “true” cell types were also identified when
421 sample size is low (**Fig 6A**). In addition, the cell proportion of 92.3%, 94.6% and 98.9% of the
422 correctly identified cell types were with more than 0.9 correlation with their “true” proportions
423 when the sample size is 6, 12 and above 20 (**Fig 6A**). Our analysis suggested a high robustness of
424 the sensitivity and prediction accuracy of SSMD when sample size is as small as 6, i.e. the
425 commonly used sample size in two-condition-comparison experiment (3 samples vs 3 samples).
426 However, as a trade-off, there is a high false discovery rate of cell type specific modules when
427 sample size is small, due to the low specificity of gene co-express analysis. To control the false
428 discoveries on small data sets, we further derived a more “stringent” set of 341 cell type specific
429 marker genes among the core marker set (see details in Methods). Our method validation
430 demonstrated a slight drop of the sensitivity and prediction accuracy when using the stringent
431 marker set on small data set (**Fig 6B**), while the specificity of the identified cell type specific
432 markers increased to from 54.4% to 72.6% when sample size is above 12 (**Fig 6C**). **Fig 6D**
433 illustrates the E-score of the cell type specific marker genes identified by using the core and the
434 more stringent marker set with respect to different sample size. The E-score of the cell types
435 marker genes identified by using the more stringent marker set were significantly higher than the
436 ones identified by using the general core marker sets when sample size is below 10, also
437 demonstrating the stringent core marker sets can effectively increase the analysis specificity when
438 sample size is small.

439

440 **DISCUSSION**

441 Over the years, research using well-established mouse models to mimic human conditions
442 have provided extensive insight into the mechanisms underlying many human diseases. We
443 developed SSMD to study mouse tissue microenvironment of complex traits, to mine the
444 interactions of cell components in the microenvironment, which will feed back to studying human
445 microenvironment. In order to have a robust prediction of cell component abundance in mouse
446 tissue, SSMD detects a subset of the genes and identifiable cell types that are the most
447 representative to the tissues to be analyzed, instead of using fixed gene signatures and cell types
448 as in classic deconvolution schemes. The limitation in expression profiling and the intrinsic and
449 mysterious variability in microenvironments excludes the possibility to have a unified set of cell
450 type specific genes that have absolutely constant expression across all conditions. The way SSMD
451 flexibly defines cell type marker genes mitigates the impact of variable marker genes due to
452 experimental platforms and microenvironment alterations. This strategy allows our model to fully
453 recapitulate the disparity of cell types and their marker genes across different microenvironment
454 and data-generating platforms. In addition, the semi-supervised formulation enables the detection
455 of sub cell types, which has been validated on scRNA-seq data simulated tissue data. Hence, a
456 relatively coarse standard for categorizing the cell types was used in training the core marker list,
457 which enabled a high robustness of the core markers. The unsupervised constrained-NMF or SVD-
458 based deconvolution on the selected marker genes further excludes the adversarial batch effects.

459 It is noteworthy a successful identification of the rank-1 modules depends on a relatively
460 large samples (>25) sharing cell types and marker genes. Currently, SSMD cannot be applied to

461 the data with a single or small sample size. However, we consider such a tradeoff between sample
462 size and prediction robustness is highly worthwhile, especially considering using SSMD as an
463 exploratory tool in large scale publicly available mouse transcriptomics data. After all, the
464 predicted proportions are often to be associated with other biological and clinical features, which
465 will be severely underpowered with a small sample size.

466 We released a R package of SSMD via <https://github.com/xiaoyulu95/SSMD> and a web
467 server via <https://ssmd.ccbb.iupui.edu/>. As illustrated in **Supplementary Fig S2A**, the input data
468 is a mouse tissue transcriptomics data and user selected tissue specific cell type core marker sets.
469 Currently, SSMD offers general core and stringent marker sets of 6 cell types in blood system, 12
470 cell types in normal, inflammatory and cancer tissue, 9 cell types in central nerve systems, and 14
471 cell types in hematopoietic systems. **Supplementary Fig S2B** illustrates a practical guide for using
472 SSMD of different tissues and sample size. The input of SSMD is a mouse tissue expression data
473 set and user selected tissue environment category. The output of SSMD includes the identified
474 data set specific cell type markers and the estimated sample-wise relative proportion of each
475 identifiable cell type. We consider the currently included cell types are comprehensive enough to
476 cover major cell types in mouse. However, the tissue specific cell types (for example, liver cells
477 in liver tissue, colon cells in colon tissue, etc) were not included in our training scope. As forming
478 rank-1 pattern among marker genes is a necessary but non-sufficient condition of identifiable cell
479 types, SSMD R package can also output rank-1 modules that do not enrich the core markers of any
480 cell type, which could possibly be markers of rare cell types. The user could further investigate
481 whether the gene module corresponds to a real cell type or not. Another key feature of the
482 webserver is that users are welcome to contribute their data to reinforce the training of cell type
483 specific marker genes.

484 Potential future directions of SSMD include (1) enabling identification of cell type specific
485 varied functions, which is not generally available for tissue data analysis in the public domain, (2)
486 identifying data set specific cell type markers forming rank-1 submatrix in a subset of samples, i.e.
487 local rank-1 submatrix, which can benefit from state-of-the-arts subspace clustering methods [27-
488 29] and (3) extending and implementing the semi-supervised framework of SSMD with other state-
489 of-the-arts deconvolution methods by refining data set specific cell marker genes. We anticipate
490 that our computational concept, which is to identify data set specific and computationally
491 “identifiable” cell types and their marker genes, can provide high robustness in deconvolution
492 analysis, by which the predicted cell proportions can be reliably correlated with experimental
493 features to provide biologically meaningful interpretation of the roles of microenvironmental
494 changes in different disease tissues.

495

496

497 **MATERIALS AND METHODS**

498 *Random walk based identification of cell type specifically expressed genes from tissue data*

499 We applied a non-parametric random walk based approach to screen genes with higher
500 expression in certain cell types comparing to others, using bulk cell training data. On the combined

501 expression matrix containing M genes for N samples of K cell types, we first calculated the
502 expected frequency of each cell type, i.e. dividing the total number of samples for the cell type
503 ($N_k, k = 1, \dots, K$) by the total number of samples N, denoted as $E_k = N_k/N, k = 1, \dots, K$. For a
504 given gene g , denote \mathbf{x} and \mathbf{x}^k as vectors of expression profile for cells of all types and type k .
505 Denote O_{jk} as the percentage of values in \mathbf{x}^k that are no less than the j th largest value in vector \mathbf{x} .
506 A random walk vector $\mathbf{d}_{1 \times N}$ that describes the non-negative discrepancy between the observed
507 and expected cell type frequency of the gene was defined as $d_j = \sum_{k=1}^K (O_{jk} - E_k)^2, j = 1, \dots, N$,
508 which attains a minimum value of zero at N. A higher peak of the random walk $\mathbf{d}_{1 \times N}$ suggests
509 gene g is more enriched in certain cell types than the others. Denote m as the index of the
510 maximum of $\mathbf{d}_{1 \times N}$, i.e. $m = \text{argmax}(\mathbf{d}_{1 \times N})$, and the cell type frequency at m as $e_k^m = O_{mk} - E_k$.
511 Cell types were further ordered by e_k^m decreasingly, and a labeling matrix L was built such that
512 $L_{g,k} = 0, \text{ if } e_k^m \leq 0$; otherwise, $L_{g,k} = \frac{1}{p}, \text{ if } \mathbf{x}^k$ has the p th largest mean among $\mathbf{x}^1, \dots, \mathbf{x}^K$.

513 It is noteworthy the approach can be directly applied to scRNA-seq data for marker training.
514 In this study, due to the relatively limited availability of existing scRNA-seq data, especially the
515 mouse strain and tissue type coverage, we generate core marker list purely by using bulk cell data.

516

517 *Identification of rank-1 cell type uniquely expressed gene modules*

518 To screen genes that form tight rank-1 modules on various tissue training datasets, SSMD
519 performs a community detection method among the genes specifically expressed in each cell type
520 as stored the labeling matrix. A correlation matrix was first built among cell type specifically
521 expressed genes, and the significance cutoff of correlation was determined by random matrix
522 theory. Random matrix theory (RMT) has been widely used to understand the low rank structure
523 encoded in biological data. In this study, an RMT based approach developed by Luo et al was
524 used to determine the threshold of significant correlation for each dataset[30]. `rm.get.threshold`
525 functions in the `RMThreshold` R package was utilized. Specifically, RMT indicated that the nearest
526 neighbor spacing distribution of eigenvalues will have a characteristic change when the threshold
527 properly separates signal from noise. By removing all the below-threshold correlation elements,
528 the co-expression modules can be more robustly unraveled. Then, hierarchical clustering was
529 performed using the correlation matrix as similarity measure.

530 Specifically, SSMD gradually increases the height of the hierarchical clustering at which
531 the tree is cut. At each height, the number of genes, the average correlation among the genes, and
532 the rank of the matrix composed of the genes in each of the cluster, is calculated. Here, matrix rank
533 is determined by a modified bi-cross validation (BCV) algorithm. SSMD stops scanning the
534 hierarchical tree if all the clusters contain less than q_0 genes, or the three following criterion is met
535 for all the clusters: (1) with at least q_0 genes, (2) the average correlation among the genes is above
536 the threshold determined by RMT, and (3) the rank of the expression matrix profile of the genes
537 in the cluster is 1. In this study, $q_0=7$ is used. Such an iterative approach will eventually select the
538 clusters with at least q_0 genes, each of which is considered as possible cell specific marker genes
539 specific to this data set. SSMD merges modules until the canonical correlation between any pair

540 of module is lower than a cutoff cor_{cut} or the number of current modules is not larger than the
 541 total rank of the gene expression profile of the selected data set specific markers genes. In this
 542 study, we utilized $cor_{cut} = 0.9$.

543 *A modified Bi-cross validation rank test:* Bi-cross validation (BCV) has been developed to
 544 estimate the matrix rank for singular value decomposition (SVD) and Non-negative Matrix
 545 Factorization (NMF), which requires a prefixed low dimension K and two low rank matrices for
 546 the approximation $X_{M \times N} = W_{M \times K} \cdot H_{K \times N}$. The error distribution of gene expression data is usually
 547 non-identical/independent, mostly because a gene's expression can be affected by its major
 548 transcriptional regulators, other biological pathways and experimental bias. Hence undesired
 549 biological characteristics and experimental bias may form significant dimensions in a gene
 550 expression data [31]. In sight of this, we developed a modified BCV rank test (*Algorithm 1*) to
 551 minimize the effect of the non-i.i.d errors in assessing the matrix rank of a gene expression data.

552
 553
 554
 555
 556
 557

Algorithm 1: Modified Bi-cross validation matrix rank test

```

Input: Matrix  $X_{M \times N}$ , parameters  $M_0, N_0, R, msp$ .
For  $r = 1 \dots R$ 
    Sample row index set  $I_r = \{i_1, i_2, \dots, i_{M_0} | i_p \in \{1 \dots M\}\}$ ,  $\bar{I}_r = \{1 \dots M\} \setminus I_r$ 
    Sample column index set  $J_r = \{j_1, j_2, \dots, j_{N_0} | j_p \in \{1 \dots N\}\}$ ,  $\bar{J}_r = \{1 \dots N\} \setminus J_r$ 
    Split  $X$  into four submatrices  $\begin{bmatrix} A_r & B_r \\ C_r & D_r \end{bmatrix}$ , where  $A_r = X[I_r, J_r]$ ,  $B_r = X[I_r, \bar{J}_r]$ ,
     $C_r = X[\bar{I}_r, J_r]$ ,  $D_r = X[\bar{I}_r, \bar{J}_r]$ 
    For  $k = 1 \dots \min(M_0, N_0)$ 
        
$$BCV(k, r) = \sum_{i=1}^{M_0} \sum_{j=1}^{N_0} \left\| A_r - B_r \widehat{D}_r^{(k)+} C_r \right\|_F^2 (*)$$

    End
End
Rankx ← 0
For  $k = 1 \dots \min(M_0, N_0)$ 
    Do t test between  $\{BCV(k, r) | r = 1 \dots R\}$  and  $\{BCV(k+1, r) | r = 1 \dots R\}$ 
    if (p. value < 0.01 & mean(BCV(k+1, r)) - mean(BCV(k, r)) > msp)
        Rankx ← k
End
Return Rankx
(*) Denote the SVD of a matrix  $D$  as  $D = U \Sigma V'$ , and Moore-Penrose inverse of  $D$ 
as  $D^+, D^+ = V' \Sigma^+ U$ , where  $\Sigma^+$  is a diagonal matrix  $\text{diag}(\sigma_1^+, \sigma_2^+, \dots, \sigma_p^+)$  with  $\sigma_1^+ \geq$ 
 $\sigma_2^+ \geq \dots \geq \sigma_p^+ \geq 0$ . Define  $\widehat{D}^{(k)+} = \sum_{i=1}^k \sigma_i^+ v_i u_i$ 

```

558

559 After running the rank-1 module detection on all the training bulk tissue datasets, those
 560 genes commonly identified in the rank-1 modules in more than 40% (70%) data sets were selected
 561 as core (stringent) markers. The list of stringent marker sets was derived with more stringent
 562 criterion, which is particularly useful for the analysis of small sample sized target data. Core
 563 markers of cells in central nervous systems were identified by a similar approach on the brain
 564 training tissue datasets. Due to the limitation of hematopoietic system tissue training data, its core
 565 markers were selected as the genes specifically over expressed in each hematopoietic cell type, by
 566 using the criteria: the gene's expression level is above 10% quantile in one cell type and below

567 50% in the other cell types. Complete lists of selected core and stringent marker sets were given
568 in **Supplementary Table S1**.

569

570 *Estimation of cell proportion*

571 Two methods were utilized to estimate cell proportion: (1) SVD based computation. With
572 cell type specific markers derived, the first row base of the gene expression profile of the marker
573 genes is directly utilized as an estimation of the cell proportion, which can be directly computed
574 by SVD. (2) Constraint NMF based computation. With the number of identifiable cell types and
575 cell type specific markers identified, the signature matrix $\tilde{S}_{M_0 \times K_0}$ and proportion matrix $\tilde{P}_{K_0 \times N}$ can
576 be estimated by minimizing the following objective function:

$$577 \min_{\tilde{S}_{M_0 \times K_0}, \tilde{P}_{K_0 \times N}} \left(\|\tilde{X}_{M_0 \times N} - \tilde{S}_{M_0 \times K_0} \cdot \tilde{P}_{K_0 \times N}\|_F^2 + \lambda \cdot \text{trace} \left(\tilde{S}_{M_0 \times K_0}^T \cdot (\mathbf{1}_{M_0} \mathbf{1}_{K_0}^T - C_{M_0 \times K_0}) \right) \right)$$

578 , where $C_{M_0 \times K_0}[i, j] = 1$ if gene i is marker of the cell type j , and 0 otherwise. λ is the hyper
579 parameter. In this study, we tuned λ by using single cell data simulated tissue data. $\lambda=10$ is
580 empirically utilized in the analysis.

581

582 *Explanation score and Comparison with state-of-the-arts methods*

583 An explanation score (ES) was utilized to evaluate the goodness that each marker gene's
584 expression is fitted by the predicted cell proportions:

$$585 \text{EScore}(x) = 1 - \sum_{j=1}^N (x_j^* - \hat{x}_j)^2 / \sum_{j=1}^N (x_j^*)^2, \hat{x}_j = \sum_{k=1}^{k_x} \beta_k^x p_j^k, \beta_k^x \geq 0$$

586 where x_j^* is the observed expression of marker gene x in sample j , \hat{x}_j is the explainable expression
587 by cell proportions, obtained by a non-negative regression x on the predicted proportion p_j^k , $k =$
588 $1 \dots k_x$. Here, k_x represents the number of cell types that express x , and β_k^x are the non-negative
589 regression parameters. Intuitively, with correctly selected marker genes, the marker gene's
590 expression can be well explained by the predicted proportions of the cell types that express the
591 gene. Hence, a high ES score is a necessary but not sufficient condition for correctly selected
592 marker genes and predicted cell proportion.

593

594 *Data used in this study*

595 *Bulk cell training data sets:* for mouse blood, solid cancer and inflammatory tissue
596 microenvironment, we retrieved 116 datasets of sorted mouse cells of 12 selected cell types,
597 totaling 1106 samples from GEO database. For mouse brain tissue microenvironment, we collected
598 2130 bulk cell samples of the nine selected cell types in central nerve systems. For mouse
599 hematopoietic microenvironment, two datasets were available that cover 14 hematopoietic cell
600 types. All the bulk cell training data were generated by the Affymetrix GeneChip Mouse Genome
601 430 2.0 Array platform and normalized with MAS5 method [32]. Samples of the same cell type
602 were further merged together with batch effect removed using Combat [33].

603 *Single Cell RNA-sequencing data:* One mouse melanoma scRNAseq data set (6638, 9) was
604 acquired from the Human Cell Atlas database [34]. Three scRNA-seq datasets of lung (4485, 12),
605 pancreas (4405, 8), and small intestine (4764, 10) and two sets of brain tissue (3679, 7 and 1099,
606 6) were accessed from Mouse Cell Atlas (MCA) data portal [35]. The two numbers in the
607 parenthesis indicate the number of cell samples and cell types of each data set. We specifically

608 selected the cells with UMI more than 500 to exclude low quality cells. Cell labels were either
609 provided in the original data or curated using Seurat v3 with cell type specific genes [36, 37].

610 *Training tissue data from cancer and blood:* 33 cancer tissue datasets of 9 cancer types
611 generated by four popular experimental platforms were collected, namely Illumina HiSeq 2000
612 Mus musculus, Affymetrix Mouse Genome 430 2.0 Array, Illumina HiSeq 2500 Mus musculus
613 and Affymetrix Mouse Genome 430A 2.0 Array from GEO database. Each data set has at least 15
614 samples. We didn't consider datasets from immunodeficient mouse, mouse cell lines, and PDX
615 models, as only real cancer or blood micro-environment is considered. A data set of liver tissue
616 collected from 31 mouse strains (GSE55489) were utilized to evaluate the variation of cell type
617 specific markers through different mouse strains [26].

618 *Brain tissue data:* 14 datasets of mouse brain tissues generated by two experimental
619 platforms, namely Illumina HiSeq 2500 Mus musculus and Affymetrix Mouse Genome 430 2.0
620 Array were collected from Gene Expression Omnibus. Datasets were split into sub data sets of
621 different brain regions. Each data set has at least 40 samples. The complete training data
622 information are available in **Supplementary Table S3**.

623 *Hematopoietic System tissue and FACS data:* We generated a RNA-seq data set with
624 matched FACS data of bone marrow cells isolated from the hind limbs of C57BL/6, Tet2-/-
625 Flt3ITD, DNMT3A-/-Flt3ITD, and DNMT3A-/-Tet2-/-Flt3ITD mice (n=3 for each group). RNA
626 (600 ng/ sample) was used to prepare single indexed strand specific cDNA library using TruSeq
627 stranded mRNA library prep kit (Illumina). The library prep was assessed for quantity and size
628 distribution using Qubit and Agilent 2100 Bioanalyzer. The pooled libraries were sequenced with
629 75bp single-end configuration on NextSeq500 (Illumina) using NextSeq 500/550 high output kit.
630 The quality of sequencing was confirmed using a Phred quality score. The sequencing data was
631 next assessed using FastQC (Babraham Bioinformatics, Cambridge, UK) and then mapped to the
632 mouse genome (UCSC mm10) using STAR RNA-seq aligner [38], and uniquely mapped
633 sequencing reads were assigned by featureCounts. The data were normalized to RPKM. FACS
634 data were collected from same biological prep by IU School of Medicine Flowcytometry Core.
635 Hematopoietic stem cells were identified by lineage negative, C-Kit high and Sca1 high cells,
636 general myeloid progenitor cells were identified by Cd34 and Cd16/32 high cells, mature myeloid
637 cells were identified by Gr1 and Cd11b high cells, and PreB cells were identified by B220 and
638 SSC-A high cells.

639

640 *Generation of simulated bulk tissue data from scRNA-seq data*

641 Cell types in each scRNA-seq data were labeled by the cell clusters provided in the original
642 works or by using Seurat pipeline with default parameters. Detailed information of the scRNA-seq
643 data and cell type annotation is given in **Supplementary Table S3**. For each data set, we simulate
644 bulk tissue data by: (1) removing insignificantly expressed genes, (2) randomly generate the
645 proportion of each cell type, called true proportion in this paper, that follows a Dirichlet
646 distribution, and (3) draw cells randomly from the cell pool with replacement according to the cell
647 type proportion, and sum up the expression values of all cells to produce a pseudo bulk tissue data.
648 The insignificant expressed genes were identified by left truncated mixture Gaussian model [39,
649 40]. The Dirichlet distribution matrix was generated with R package "DirichletReg" [41].

650

651 SUPPLEMENTARY DATA

652 Supplemental Information can be found in supplementary files

653

654 ACKNOWLEDGEMENTS

655 C.Z thank Mr. Siyuan Qi from Indiana University School of Medicine for his help in the
656 early stage of this work. C.Z and S.C thank the Indiana University Center for Medical Genomics
657 for their support of this project.

658

659 FUNDING SUPPORTS

660 This work was supported by the National Science Foundation Div Of Information &
661 Intelligent Systems (No. 1850360). This work was also supported by an R01 award
662 #1R01GM131399-01 from the National Institute of General Medical Sciences. This work was also
663 supported by the Showalter Young Investigator Award from Indiana CTIS.

664

665 CONFLICT OF INTEREST

666 There are no conflicts to declare.

667

668 REFERENCES

- 669 1. Beck, J.A., et al., *Genealogies of mouse inbred strains*. Nature genetics, 2000. **24**(1): p. 23-25.
- 670 2. Rosenthal, N. and S. Brown, *The mouse ascending: perspectives for human-disease models*. Nature
671 cell biology, 2007. **9**(9): p. 993-999.
- 672 3. Van der Jeught, K., et al., *ST2 as checkpoint target for colorectal cancer immunotherapy*. JCI
673 insight, 2020. **5**(9).
- 674 4. Mund, J.A., et al., *Genetic disruption of the small GTPase RAC1 prevents plexiform neurofibroma
675 formation in mice with neurofibromatosis type 1*. Journal of Biological Chemistry, 2020: p. jbc.
676 RA119. 010981.
- 677 5. Huang, M., et al., *Sestrin 3 Protects Against Diet-Induced Nonalcoholic Steatohepatitis in Mice
678 Through Suppression of Transforming Growth Factor β Signal Transduction*. Hepatology, 2020.
679 **71**(1): p. 76-92.
- 680 6. Pandey, R., et al., *SHP2 inhibition reduces leukemogenesis in models of combined genetic and
681 epigenetic mutations*. Journal of Clinical Investigation, 2019. **129**(12): p. 5468-5473.
- 682 7. Zhang, C., S. Cao, and Y. Xu, *Population dynamics inside cancer biomass driven by repeated
683 hypoxia-reoxygenation cycles*. Quantitative Biology, 2014. **2**(3): p. 85-99.
- 684 8. Hackl, H., et al., *Computational genomics tools for dissecting tumour-immune cell interactions*.
685 Nature Reviews Genetics, 2016. **17**(8): p. 441.
- 686 9. Li, B., et al., *Comprehensive analyses of tumor immunity: implications for cancer immunotherapy*.
687 Genome biology, 2016. **17**(1): p. 174.
- 688 10. Wang, X., et al., *Bulk tissue cell type deconvolution with multi-subject single-cell expression
689 reference*. 2019. **10**(1): p. 380.
- 690 11. Racle, J., et al., *Simultaneous enumeration of cancer and immune cell types from bulk tumor gene
691 expression data*. Elife, 2017. **6**.
- 692 12. Newman, A.M., et al., *Determining cell type abundance and expression from bulk tissues with
693 digital cytometry*. Nat Biotechnol, 2019. **37**(7): p. 773-782.
- 694 13. Newman, A.M., et al., *Robust enumeration of cell subsets from tissue expression profiles*. Nature
695 methods, 2015. **12**(5): p. 453.
- 696 14. Li, B., et al., *Comprehensive analyses of tumor immunity: implications for cancer immunotherapy*.
697 Genome Biol, 2016. **17**(1): p. 174.

- 698 15. Gaujoux, R. and C.J.B. Seoighe, *CellMix: a comprehensive toolbox for gene expression*
699 *deconvolution*. 2013. **29**(17): p. 2211-2212.
- 700 16. Frishberg, A., et al., *Cell composition analysis of bulk genomics using single-cell data*. Nat
701 Methods, 2019. **16**(4): p. 327-332.
- 702 17. Finotello, F. and Z.J.C.I. Trajanoski, Immunotherapy, *Quantifying tumor-infiltrating immune cells*
703 *from transcriptomics data*. 2018. **67**(7): p. 1031-1040.
- 704 18. Abbas, A.R., et al., *Deconvolution of blood microarray data identifies cellular activation patterns*
705 *in systemic lupus erythematosus*. 2009. **4**(7): p. e6098.
- 706 19. Abbas, A., et al., *Immune response in silico (IRIS): immune-specific genes identified from a*
707 *compendium of microarray expression data*. 2005. **6**(4): p. 319.
- 708 20. Chen, Z., et al., *Inference of immune cell composition on the expression profiles of mouse tissue*.
709 Scientific reports, 2017. **7**: p. 40508.
- 710 21. Marques, S., et al., *Oligodendrocyte heterogeneity in the mouse juvenile and adult central nervous*
711 *system*. Science, 2016. **352**(6291): p. 1326-1329.
- 712 22. La Manno, G., et al., *Molecular Diversity of Midbrain Development in Mouse, Human, and Stem*
713 *Cells*. Cell, 2016. **167**(2): p. 566-580 e19.
- 714 23. Codeluppi, S., et al., *Spatial organization of the somatosensory cortex revealed by osmFISH*. Nat
715 Methods, 2018. **15**(11): p. 932-935.
- 716 24. Chang, W., et al., *ICTD: A semi-supervised cell type identification and deconvolution method for*
717 *multi-omics data*. bioRxiv, 2019: p. 426593.
- 718 25. Wang, X., et al., *Bulk tissue cell type deconvolution with multi-subject single-cell expression*
719 *reference*. Nature communications, 2019. **10**(1): p. 1-9.
- 720 26. Church, R.J., et al., *A systems biology approach utilizing a mouse diversity panel identifies genetic*
721 *differences influencing isoniazid-induced microvesicular steatosis*. Toxicological Sciences, 2014.
722 **140**(2): p. 481-492.
- 723 27. Wan, C., et al., *Denoising individual bias for a fairer binary submatrix detection*. arXiv preprint
724 arXiv:2007.15816, 2020.
- 725 28. Wan, C., et al., *Fast And Efficient Boolean Matrix Factorization By Geometric Segmentation*. arXiv,
726 2019: p. arXiv: 1909.03991.
- 727 29. Chang, W., et al., *Supervised clustering of high dimensional data using regularized mixture*
728 *modeling*. arXiv preprint arXiv:2007.09720, 2020.
- 729 30. Luo, F., et al., *Constructing gene co-expression networks and predicting functions of unknown*
730 *genes by random matrix theory*. BMC bioinformatics, 2007. **8**(1): p. 299.
- 731 31. Lopez, R., et al., *Deep generative modeling for single-cell transcriptomics*. 2018. **15**(12): p. 1053.
- 732 32. Pepper, S.D., et al., *The utility of MAS5 expression summary and detection call algorithms*. BMC
733 bioinformatics, 2007. **8**(1): p. 273.
- 734 33. Johnson, W.E., C. Li, and A. Rabinovic, *Adjusting batch effects in microarray expression data*
735 *using empirical Bayes methods*. Biostatistics, 2007. **8**(1): p. 118-127.
- 736 34. Regev, A., et al., *Science forum: the human cell atlas*. Elife, 2017. **6**: p. e27041.
- 737 35. Han, X., et al., *Mapping the mouse cell atlas by microwell-seq*. Cell, 2018. **172**(5): p. 1091-1107.
738 e17.
- 739 36. Stuart, T., et al., *Comprehensive Integration of Single-Cell Data*. Cell, 2019.
- 740 37. Butler, A., et al., *Integrating single-cell transcriptomic data across different conditions,*
741 *technologies, and species*. Nature biotechnology, 2018. **36**(5): p. 411.
- 742 38. Dobin, A., et al., *STAR: ultrafast universal RNA-seq aligner*. Bioinformatics, 2013. **29**(1): p. 15-
743 21.
- 744 39. Wan, C., et al., *LTMG: a novel statistical modeling of transcriptional expression states in single-*
745 *cell RNA-Seq data*. Nucleic acids research, 2019. **47**(18): p. e111-e111.
- 746 40. Zhang, Y., et al., *M3S: A comprehensive model selection for multi-modal single-cell RNA*
747 *sequencing data*. BMC bioinformatics, 2019. **20**(24): p. 1-5.
- 748 41. Maier, M.J., *DirichletReg: Dirichlet regression for compositional data in R*. 2014.
749

750 **FIGURE LEGENDS**

751 **Fig 1. Analysis pipeline of SSMD and core cell type specific markers.** (A) Analysis pipeline of
752 the core marker training procedure. (B) Analysis pipeline of the deconvolution procedure. In (A)
753 and (B), input data including training and target data, computational procedure and key
754 intermediate outputs were colored by orange, green and blue, respectively. (C) Core markers of 12
755 cell types in blood, solid cancer, and inflammatory tissue. An edge between two genes means the
756 two genes are co-identified as markers of one cell type in more than 50% of the training data sets.
757 (D) Core markers of 9 cell types in central nervous system. Noted, core markers for the endothelial
758 cell in the inflammatory tissue and central nervous system were separately trained by comparing
759 with other cell types in the same tissue system.

760
761 **Fig 2. Method evaluation on scRNA-seq simulated tissue data.** (A-D) Correlation between true
762 and predicted cell proportions in the simulated Lung (A), Pancreas (B), Small Intestine (C), and
763 Mouse Melanoma (D) tissue data. The x-axis represents cell type and y axis represents prediction
764 accuracy. Predictions made by SSMD, EPIC, ImmuCC and tissue-ImmuCC were dark blue, green,
765 yellow and orange colored, respectively. The red dash line represents the 0.9 correlation cutoff.
766 (E-F) Correlation between true and predicted cell proportions in the two simulated brain tissue
767 data. (G) The total rank of the gene expression profile of selected marker genes in the six simulated
768 tissue data (grey), and the total number of cell types identified by SSMD in each data set or
769 assumed in other methods (left three grey bars).

770
771 **Fig 3. Method evaluation on scRNA-seq simulated tissue data on hematopoietic tissue data.**
772 (A-D) Correlation between SSMD predicted (x-axis) and FACS identified (y-axis) cell proportions
773 of HSC, GMP, mature myeloid cell and preB cell. (E-H) marker proteins utilized to identify the
774 four cell types by using FACS. The x- and y- axis of the plots represent the level of cell type
775 markers. The black block in (E), the green block in (F), the upper-right block in (G) and the block
776 in (H) are the sorted HSC, GMP, Myeloid and Pre-B cell, respectively.

777
778 **Fig 4. Prediction of SSMD on real tissue data.** (A, C) The total rank of the gene expression
779 profile of selected marker genes (grey) in different (A) cancer tissue and (C) brain data, and the
780 total number of cell types identified by SSMD in each data set (colored). (B, D) E-Score for
781 different cell types identified by SSMD (blue) in (B) cancer and (D) brain data set or assumed in
782 other methods (EPIC: red, ImmuCC: Yellow).

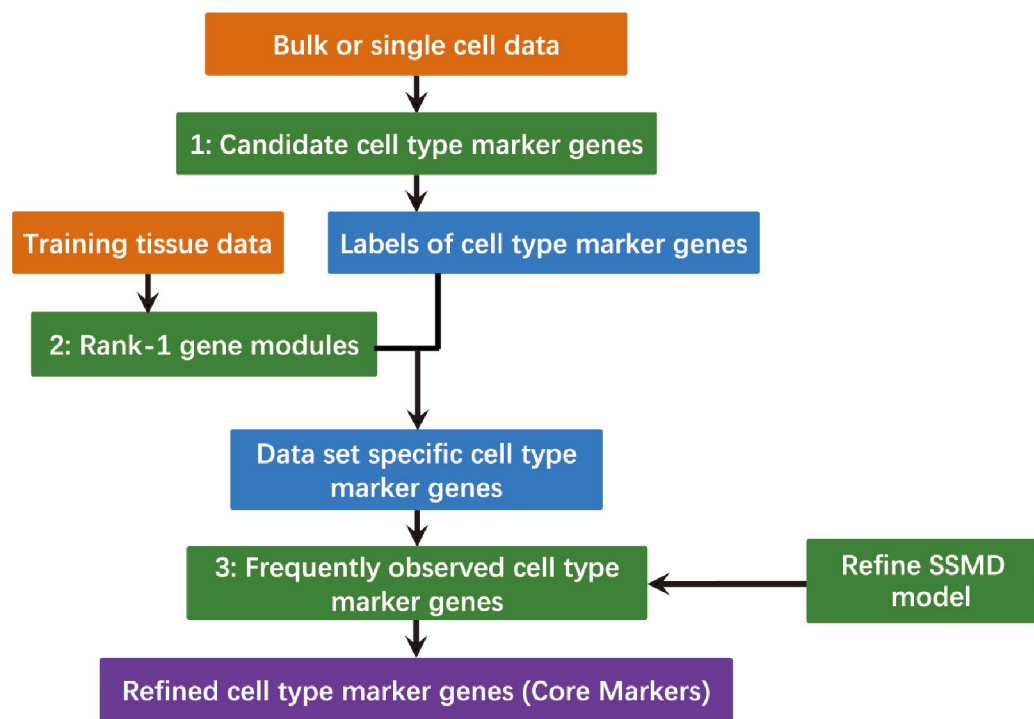
783
784 **Fig 5. Correlation between expression level of strain specific cell type marker genes and**
785 **predicted cell proportion.** High correlation is a necessary but non-sufficient condition for the
786 genes to serve as marker genes of the cell types in corresponding mouse strain. In the heatmap, x-
787 and y-axis represent genes and mouse strains, respectively. Genes in the core marker list of four
788 selected cell types, namely Neutrophil, Nature Kill (NK), Macrophage, and Monocyte, were
789 colored on the column side bar.

790
791 **Fig 6. Performance evaluation of different sample size.** (A) Prediction accuracy (y-axis) in
792 different sample size (x-axis) using all core markers. Accuracy is the Pearson correlation between
793 predicted proportion using only selected small sample and using all samples. (B) Prediction

794 accuracy (y-axis) in different sample size (x-axis) using selected stringent markers. (C) True
795 positive rate (y-axis) of the cell type specific markers identified by using the stringent markers
796 (blue) and core markers (green) with respect to different sample size (x-axis). (D) E-Score for
797 using co-expression modules consisting of all core markers and only selected stringent markers.
798 From top to bottom, the statistics were derived from GSE76095, GSE67186, GSE90885,
799 GSE94574, and GSE126279.
800

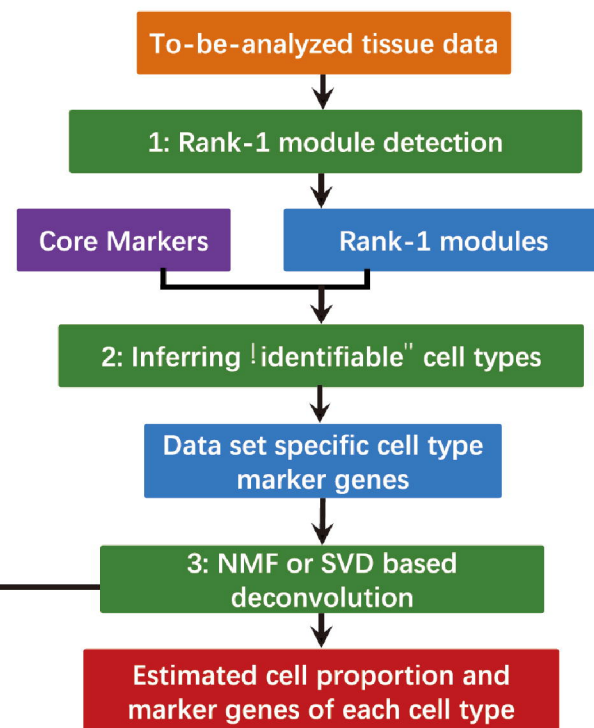
A

Core markers training pipeline

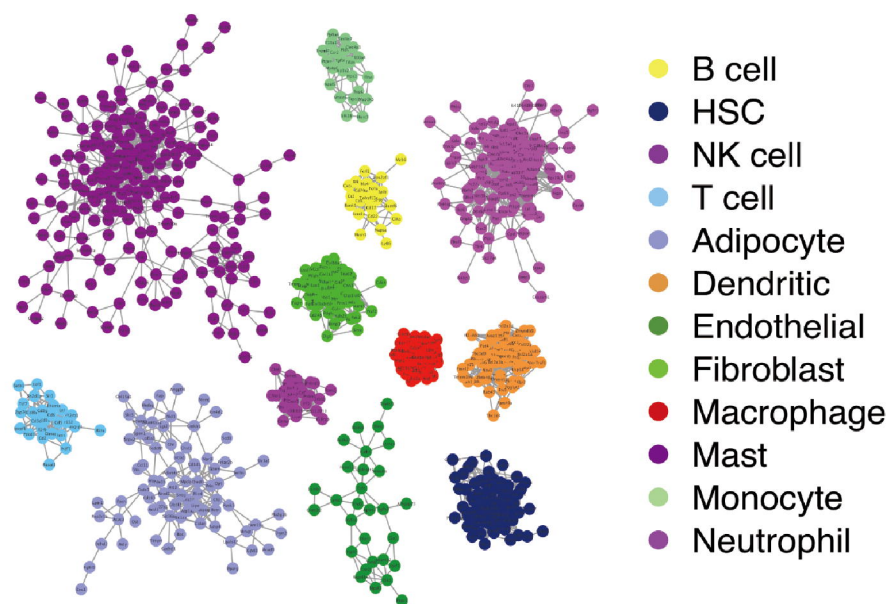


B

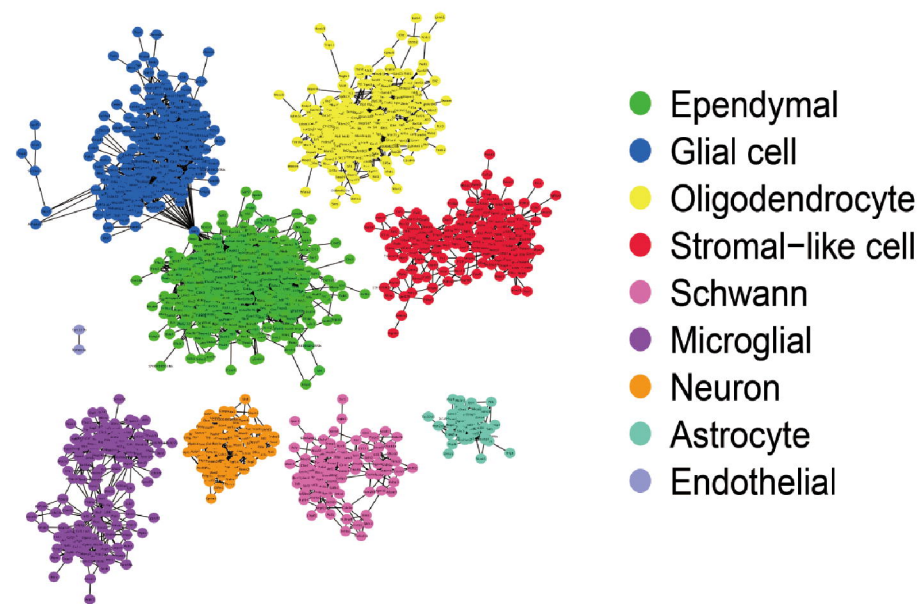
Deconvolution analysis pipeline



C

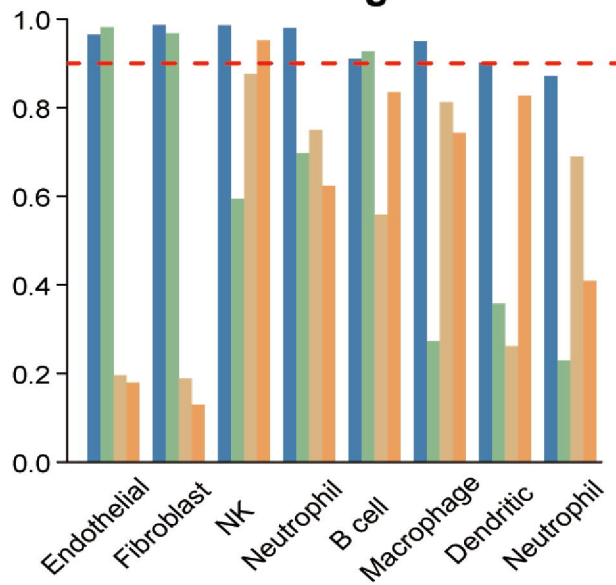


D



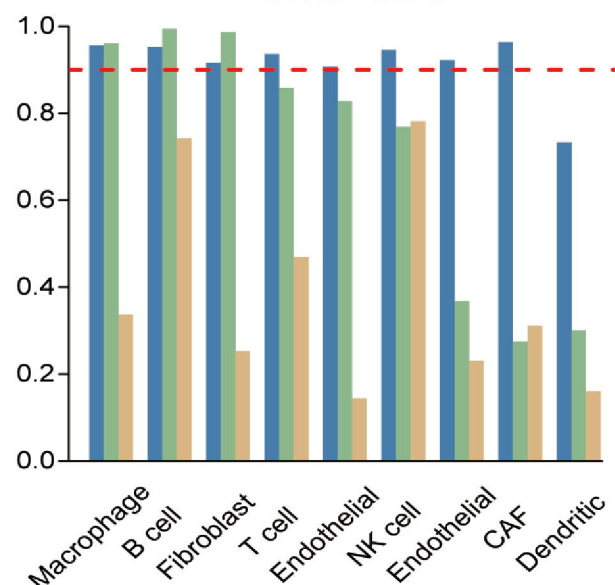
A

Lung



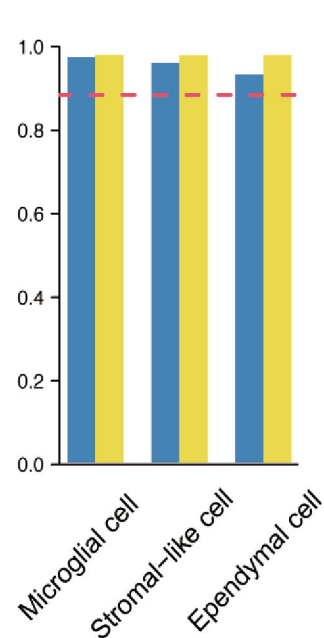
C

Melanoma



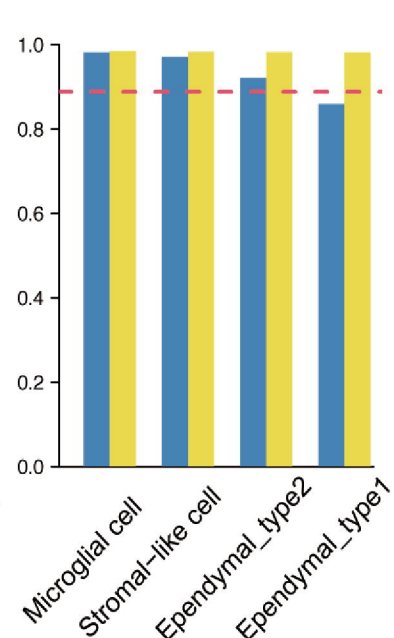
E

Brain1



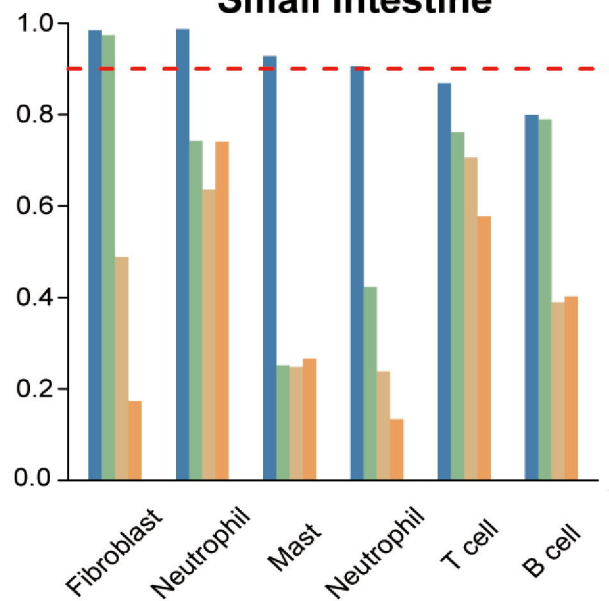
F

Brain2



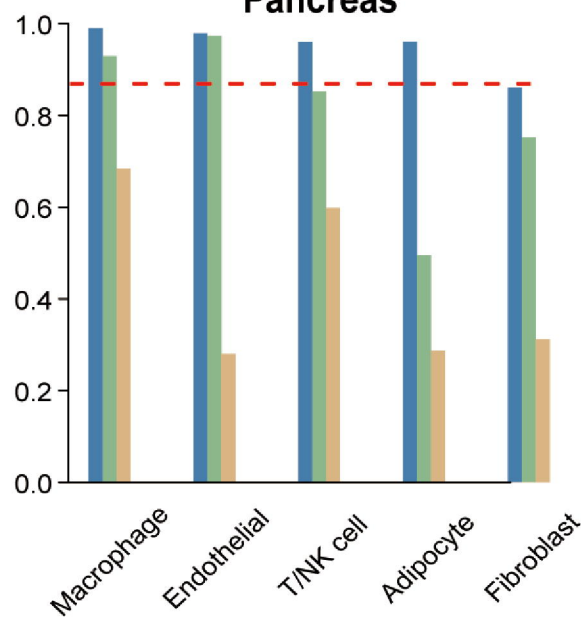
B

Small Intestine

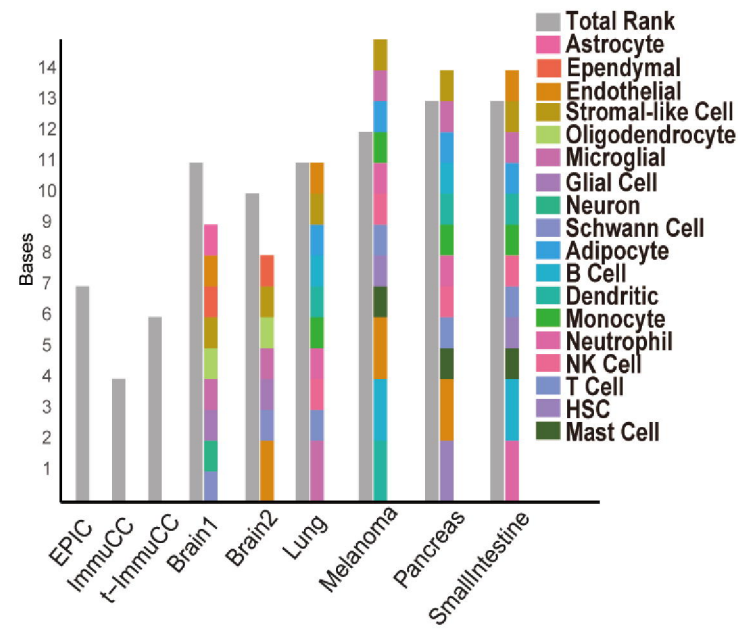


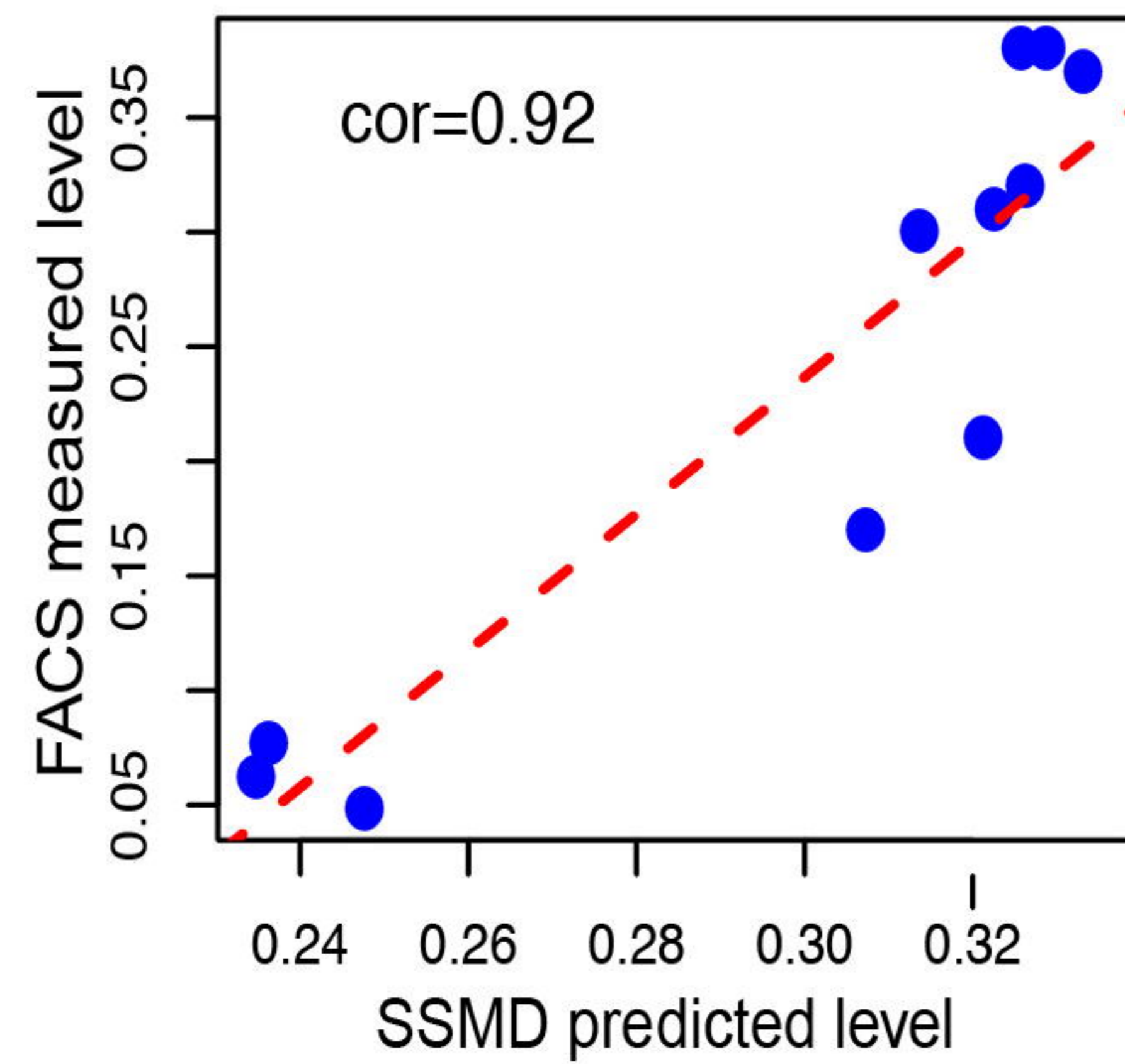
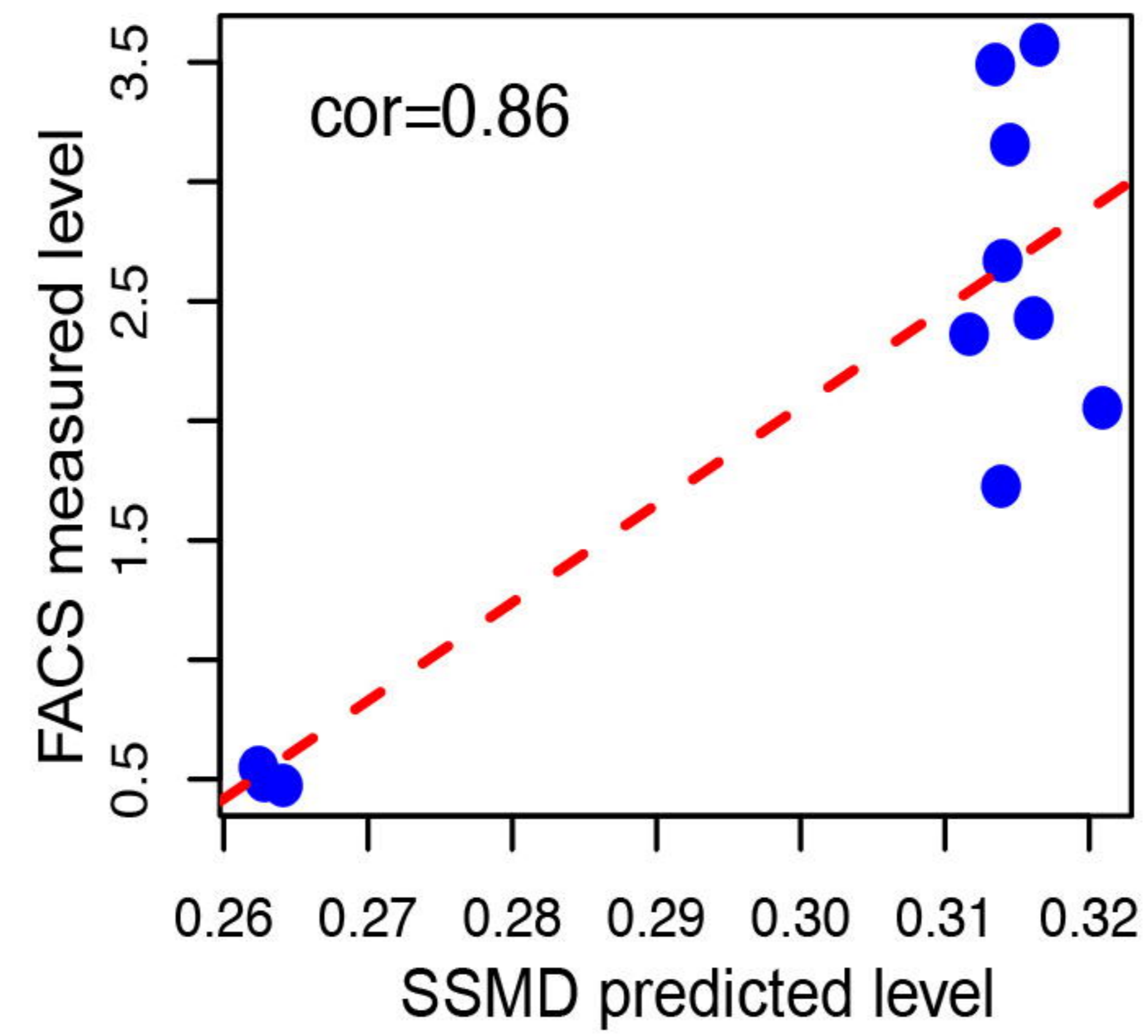
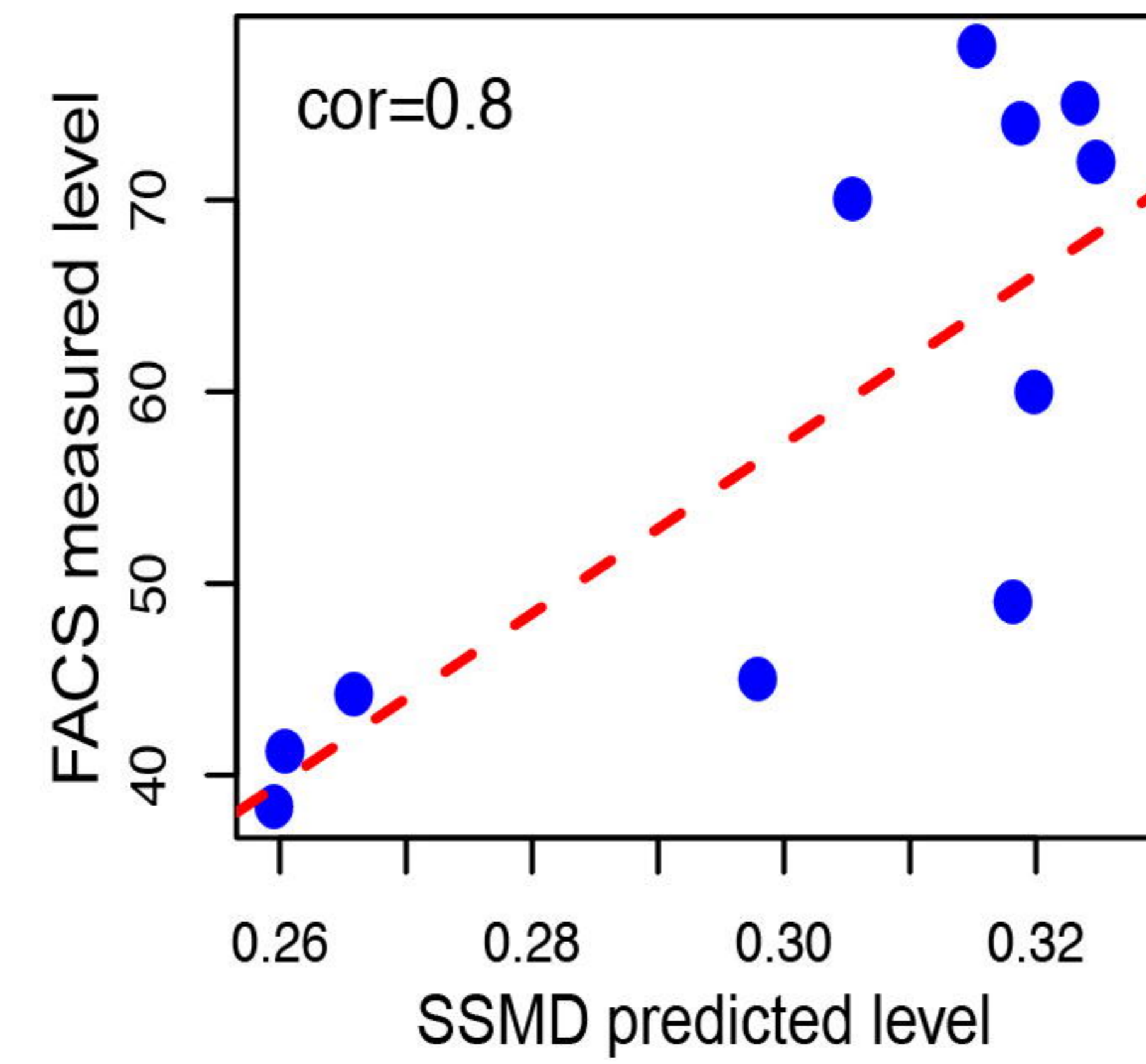
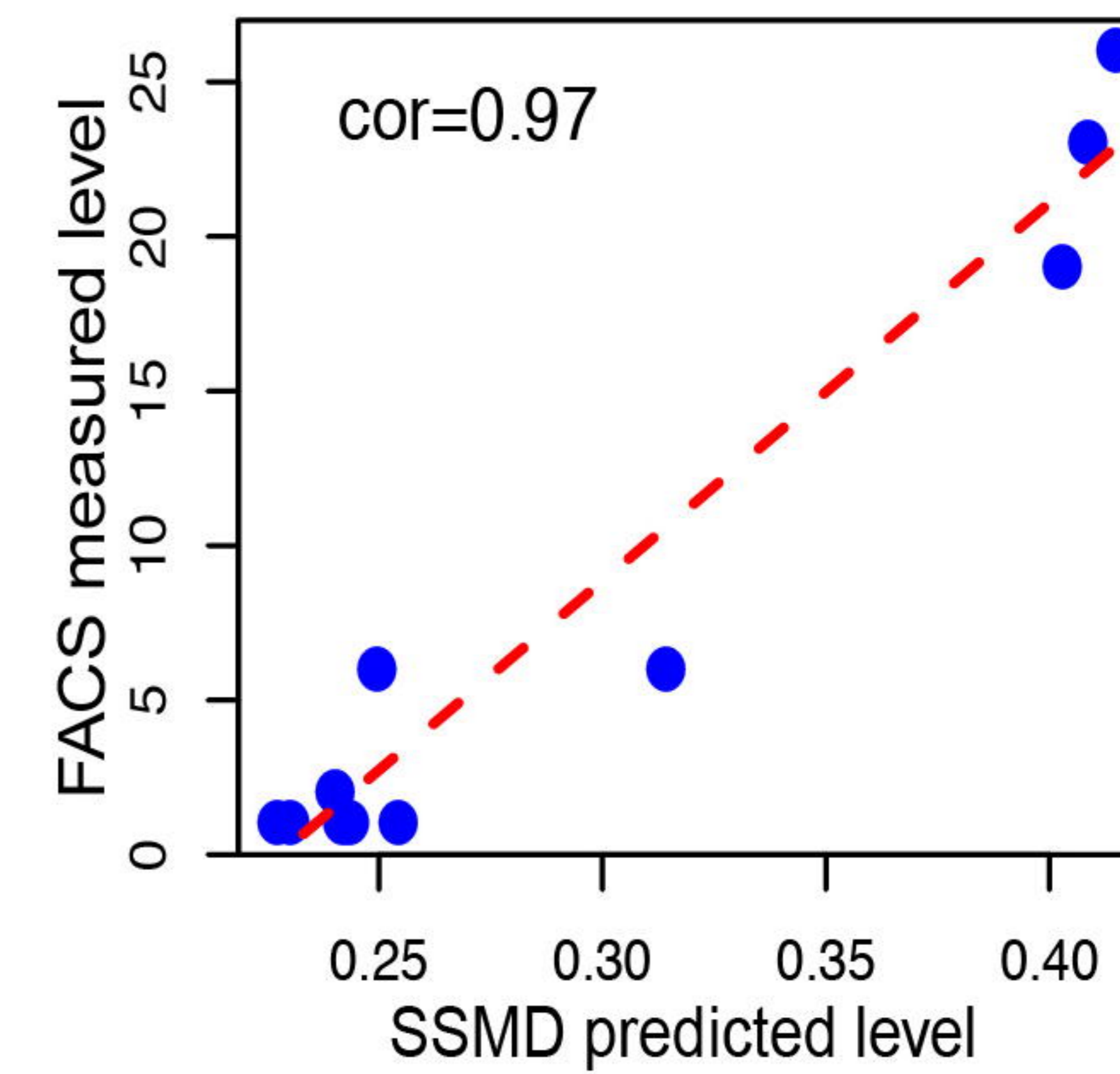
D

Pancreas

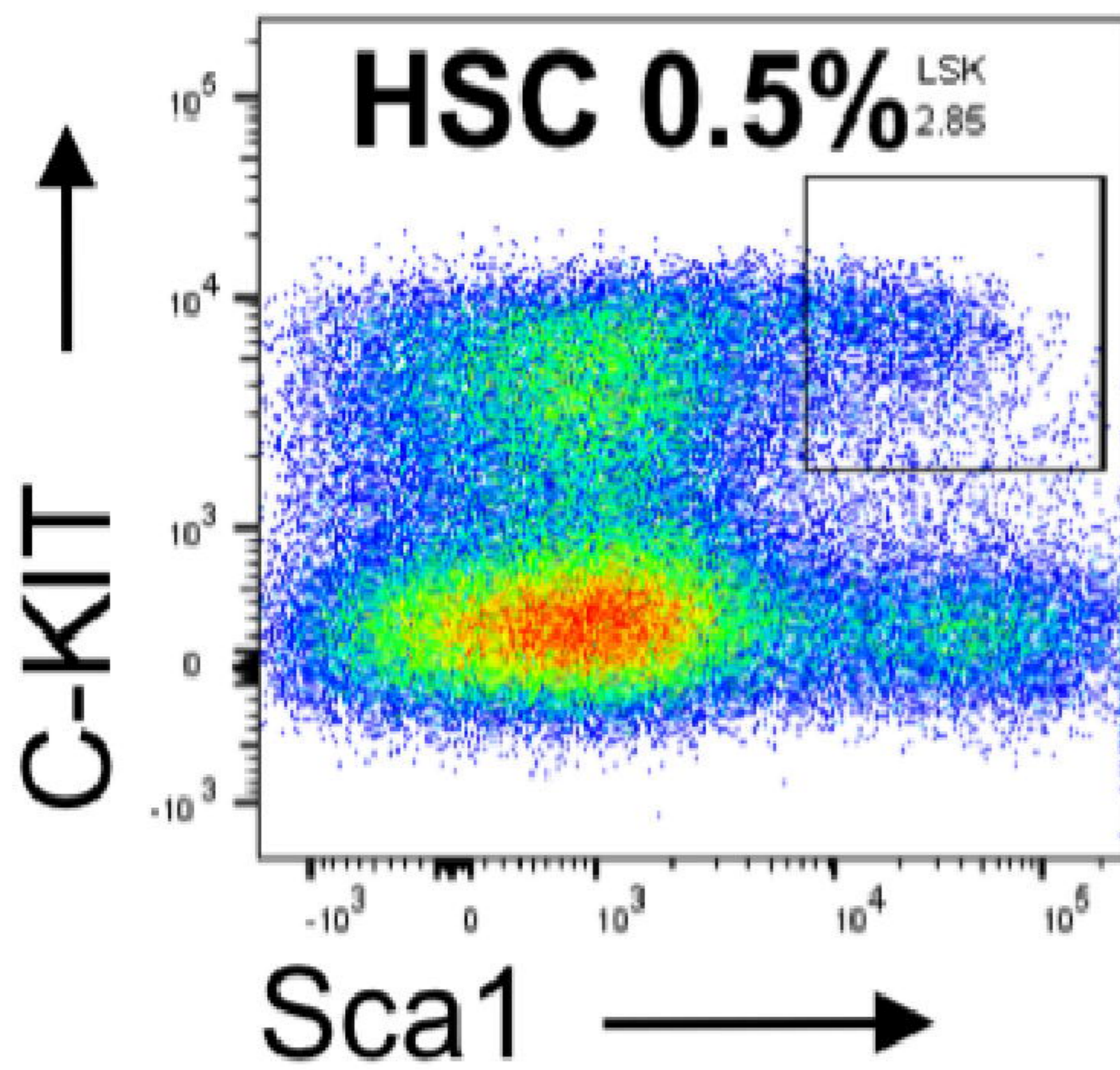
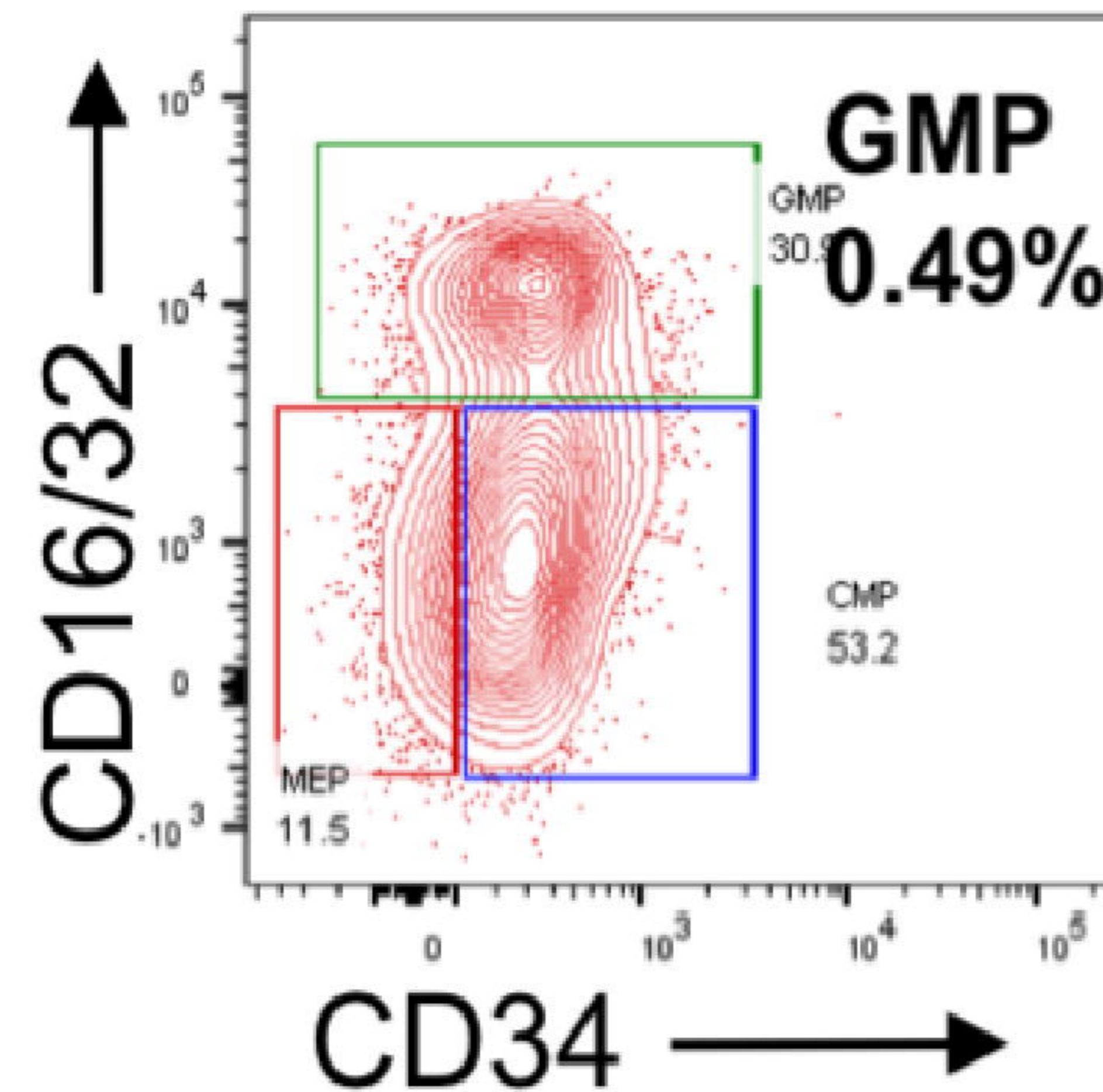
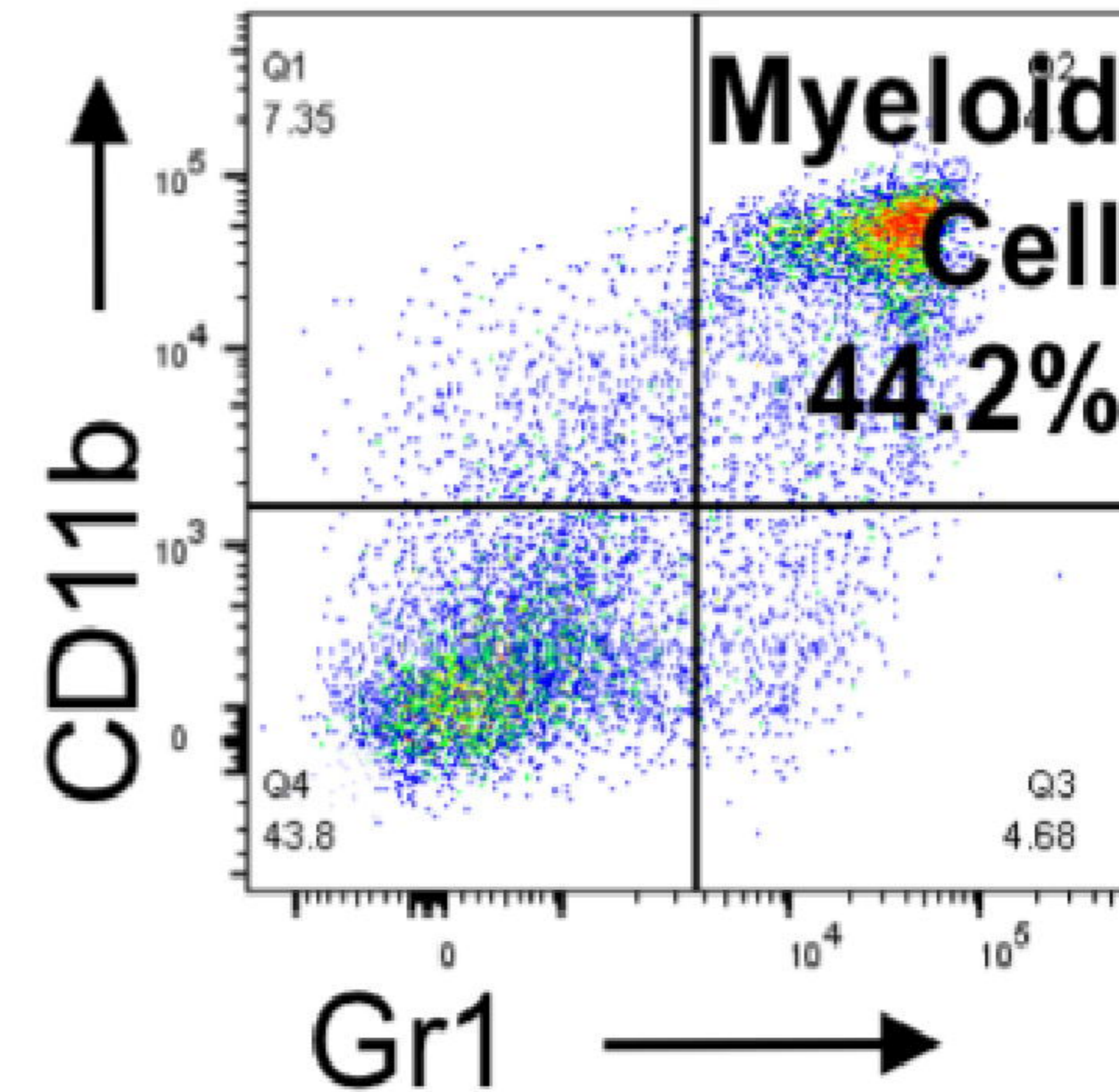
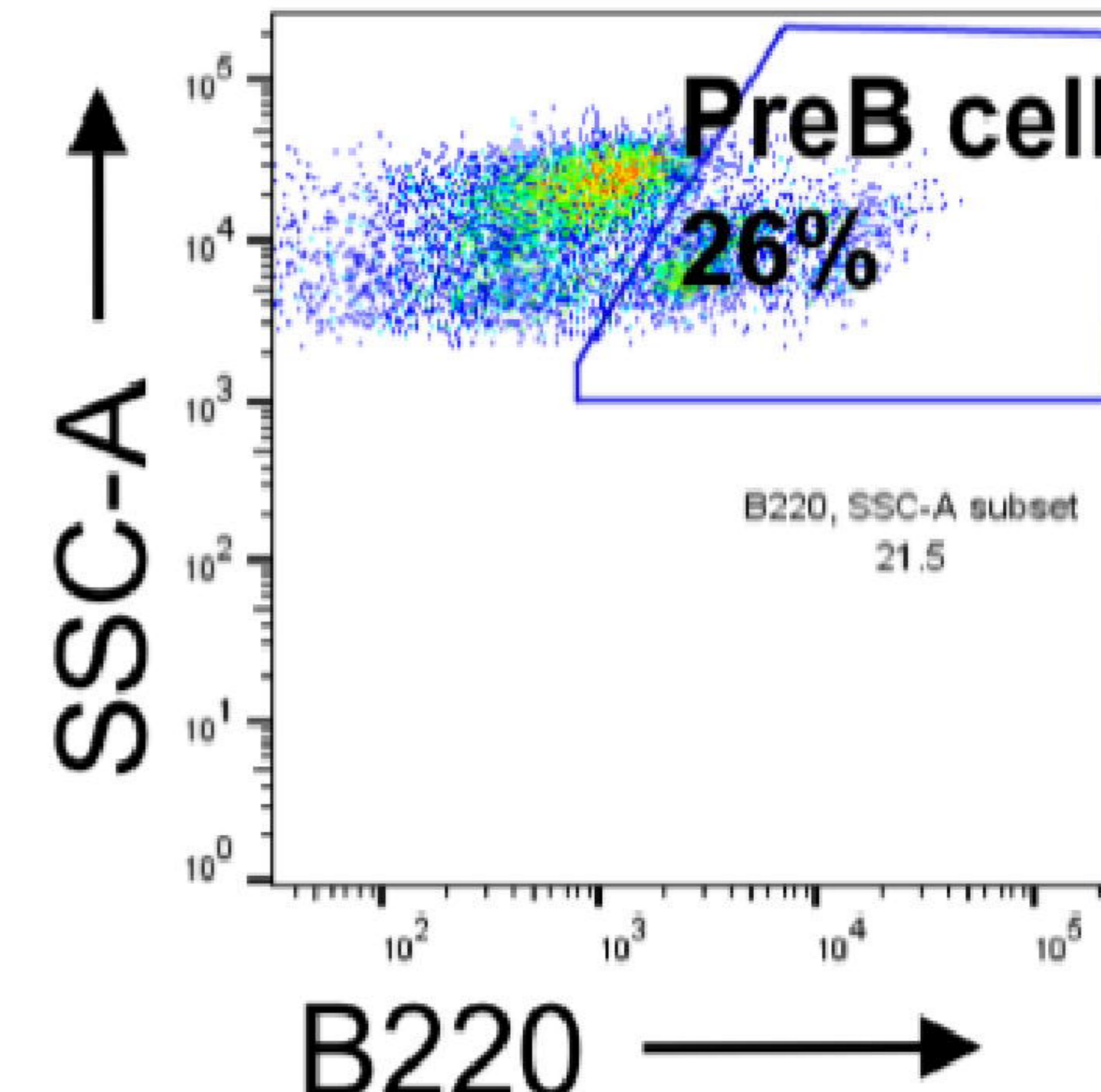


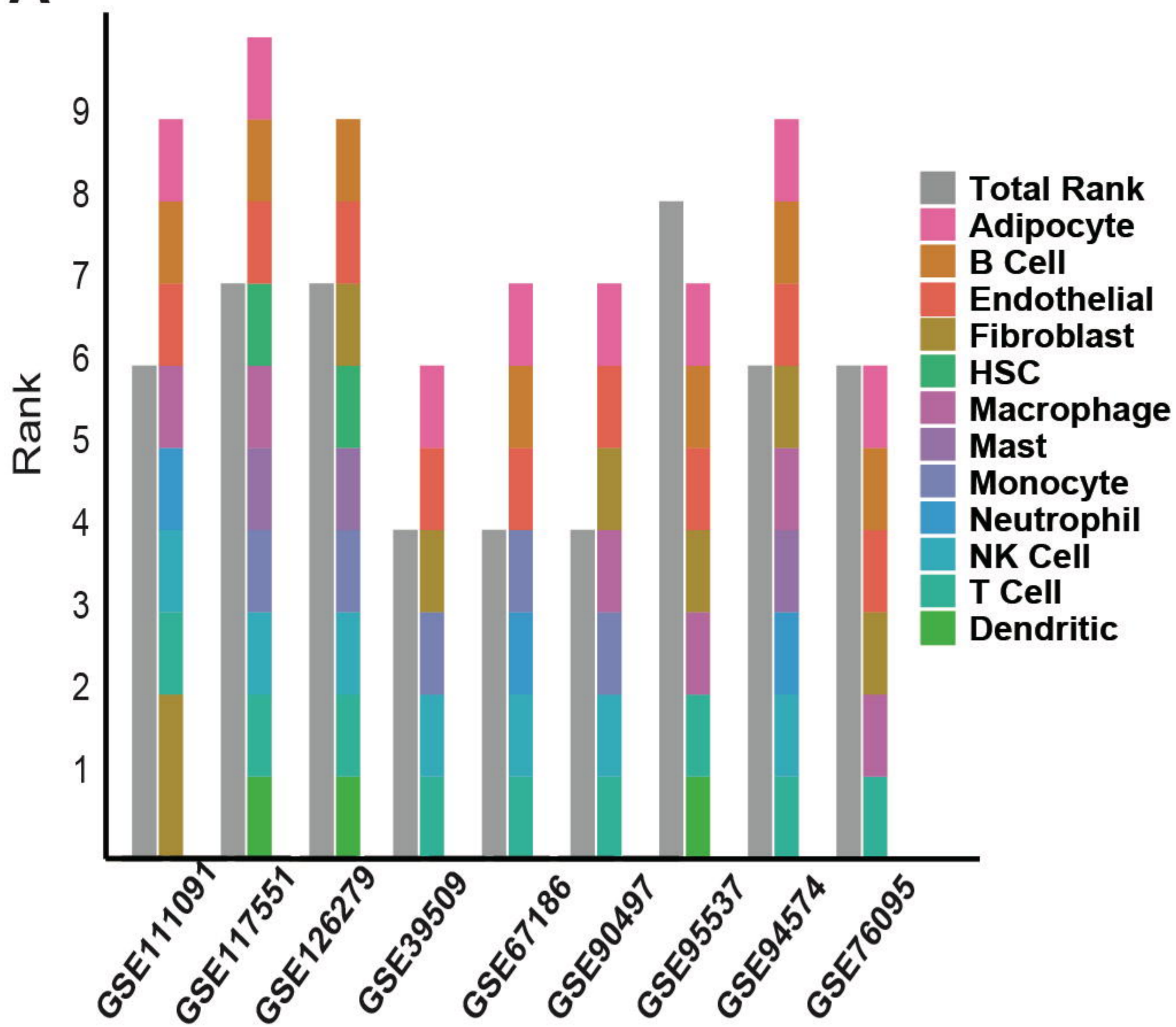
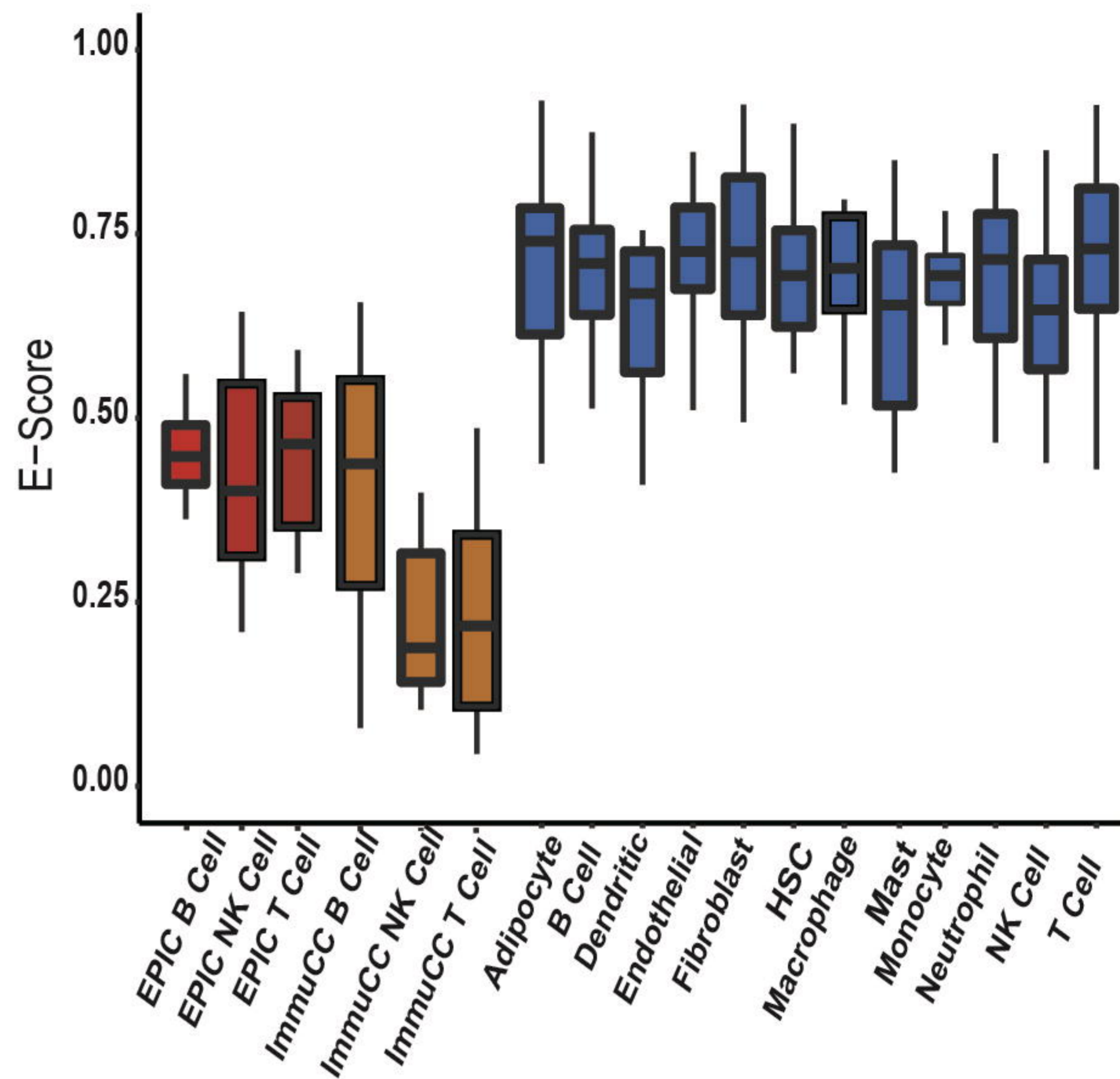
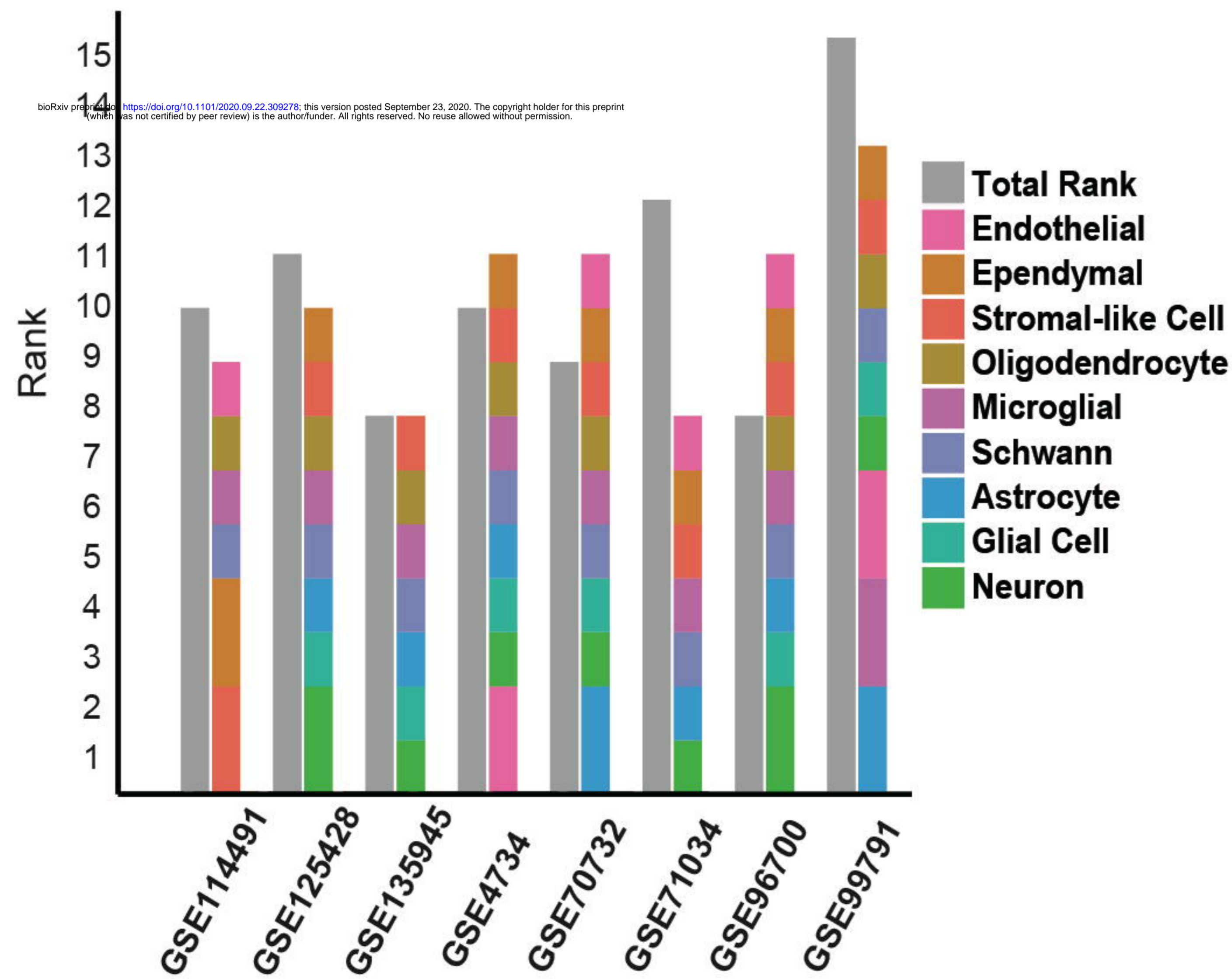
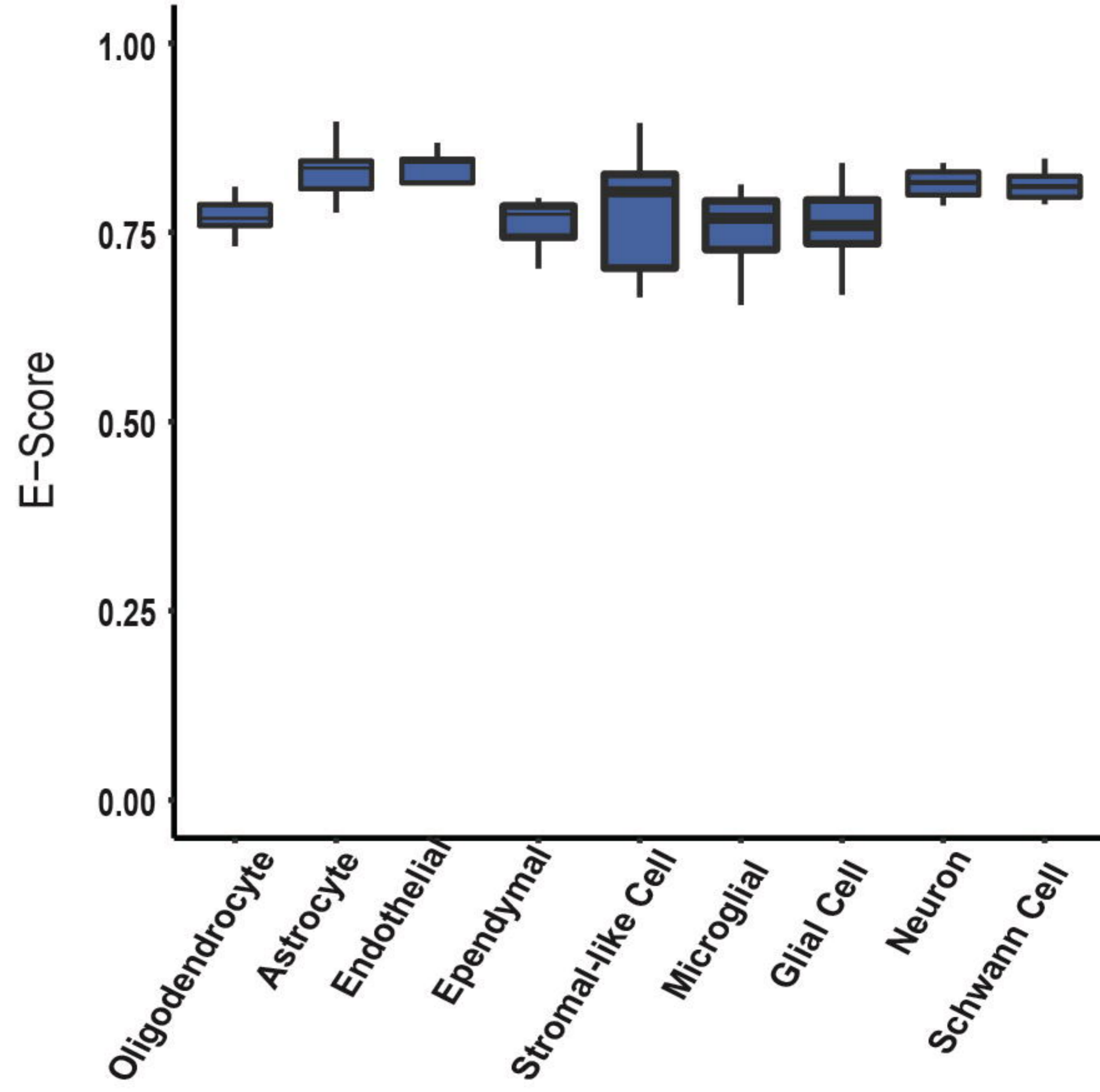
G

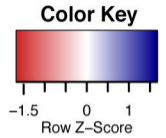
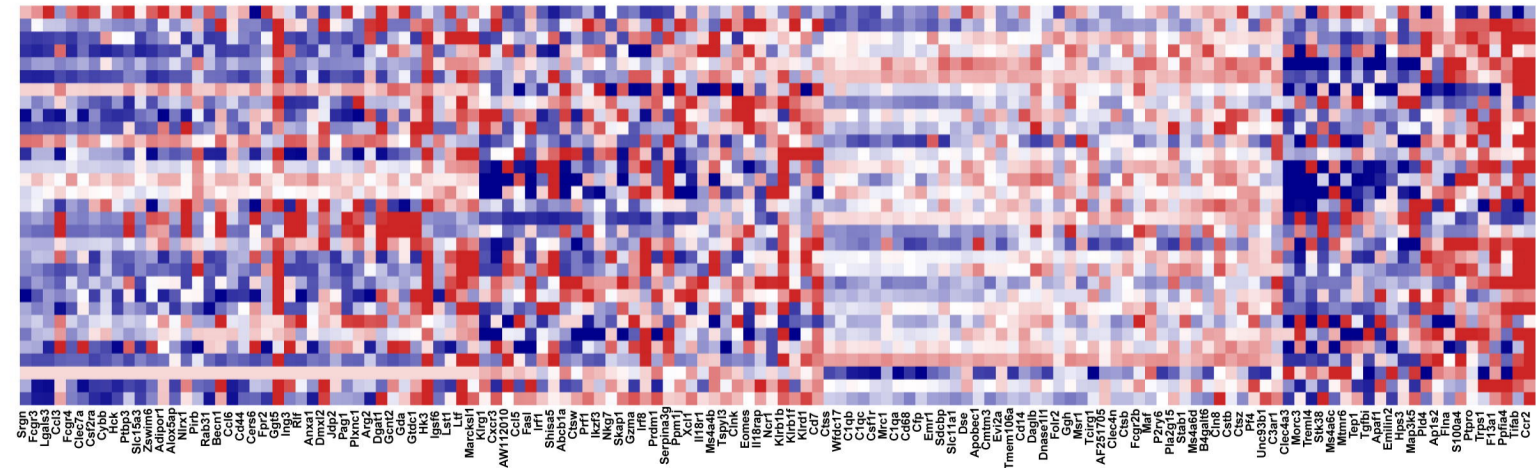


A**HSC****B****GMP****C****Mature Myeloid Cell****D****PreB cell****E**

Lineage Negative
C-KIT^{POS}SCA1^{POS}

**F****Committed Progenitors****G****Mature Myeloid cells****H****Mature Lymphoid cells**

A**B****C****D**



- Neutrophil
- NK Cell
- Macrophage
- Monocyte

