

# SSMD: a semi-supervised approach for a robust cell type identification and deconvolution of mouse transcriptomics data

Xiaoyu Lu<sup>†</sup>, Szu-Wei Tu<sup>†</sup>, Wennan Chang, Changlin Wan, Jiashi Wang, Yong Zang, Baskar Ramdas, Reuben Kapur, Xiongbin Lu, Sha Cao and Chi Zhang

Corresponding authors: Chi Zhang, Department of Medical and Molecular Genetics, Department of Bio Health Informatics, Center for Computational Biology and Bioinformatics, Indiana University School of Medicine, 410 W. 10th Street, Indianapolis, IN 46202. Tel: +1 317-2789625; E-mail: czhang87@iu.edu; Sha Cao, Department of Biostatistics, Center for Computational Biology and Bioinformatics, Indiana University School of Medicine, 410 W. 10th Street, Indianapolis, IN 46202. Tel: +1-3172742602; E-mail: shacao@iu.edu; Xiongbin Lu, Department of Medical and Molecular Genetics, Center for Computational Biology and Bioinformatics, Indiana University School of Medicine, 410 W. 10th Street, Indianapolis, IN 46202. Tel: +1-317-2744398; E-mail: xiolu@iu.edu

<sup>†</sup>These authors made equal contribution to this work.

## Abstract

Deconvolution of mouse transcriptomic data is challenged by the fact that mouse models carry various genetic and physiological perturbations, making it questionable to assume fixed cell types and cell type marker genes for different data set scenarios. We developed a Semi-Supervised Mouse data Deconvolution (SSMD) method to study the mouse tissue microenvironment. SSMD is featured by (i) a novel nonparametric method to discover data set-specific cell type signature genes; (ii) a community detection approach for fixing cell types and their marker genes; (iii) a constrained matrix decomposition method to solve cell type relative proportions that is robust to diverse experimental platforms. In summary, SSMD addressed several key challenges in the deconvolution of mouse tissue data, including: (i) varied cell types and marker genes caused by highly divergent genotypic and phenotypic conditions of mouse experiment; (ii) diverse experimental platforms of mouse transcriptomics data; (iii) small sample size and limited training data source and (iv) capable to estimate the proportion of 35 cell types in blood, inflammatory, central nervous or hematopoietic systems. *In silico* and experimental validation of SSMD demonstrated its high sensitivity and accuracy in identifying (sub) cell types and

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

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

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

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

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

Yong Zang is an assistant professor in the Department of Biostatistics and a member of the Center for Computational Biology and Bioinformatics, Indiana University School of Medicine.

Baskar Ramdas is an assistant research professor in the Department of Pediatrics, Indiana University School of Medicine.

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

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

Sha Cao is an assistant professor in the Department of Biostatistics and a member of the Center for Computational Biology and Bioinformatics, Indiana University School of Medicine.

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

Submitted: 21 August 2020; Received (in revised form): 18 September 2020

© The Author(s) 2020. Published by Oxford University Press. All rights reserved. For Permissions, please email: journals.permissions@oup.com

predicting cell proportions comparing with state-of-the-arts methods. A user-friendly R package and a web server of SSMD are released via <https://github.com/xiaoyulu95/SSMD>.

**Key words:** tissue data deconvolution; cancer microenvironment; semi-supervised learning; mouse omics data

## Introduction

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

Tissue transcriptomic data display convoluted signals from different cell types [8]. Deconvoluting cell components and identifying mouse strain-/tissue-/experimental condition-specific cell types and gene expressions are crucial for understanding how experimentally perturbed conditions are associated with cellular level characteristics and cell–cell interactions [9]. While multiple deconvolution methods have been developed for investigating the heterogeneous cell types in human cancer or other tissues data [10–19], they may not be directly applicable to mouse tissue data. First of all, the cell type-specific genes for human cells differ from mouse cells; secondly, compared with human, the variations among different mouse tissue samples may be considerably higher, as they are collected from different strains with varied genetic background and experimental conditions.

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

To address these challenges, we developed a novel semi-supervised deconvolution method, namely Semi-Supervised

Mouse data Deconvolution (SSMD), to infer data-/tissue-specific cell type marker genes and their expression profiles and estimate their relative abundances from transcriptomics data. SSMD is capable to infer the relative proportion of 35 cell types in the blood, inflammatory, cancer, central nervous system and hematopoietic system. To the best of our knowledge, SSMD is the only mouse data deconvolution method considering strain, tissue type and data specificity of cell type-specific gene markers. We demonstrated SSMD achieved a high sensitivity in identifying the appearance of immune and stromal cell types in inflammatory tissue and brain cell types in central nervous tissue, and with a high accuracy in estimating their relative proportion on single-cell RNA-sequencing (scRNA-seq) simulated bulk tissue data sets. We also experimentally validated that the cell populations inferred by SSMD accurately recapitulates the true cell proportions measured by fluorescence-activated cell sorting (FACS) on a leukemia bone marrow data. Applications of SSMD on a large collection of public mouse blood, brain, cancer and other inflammatory tissue data suggested that the method achieved a robust performance throughout diverse types of experimental conditions and platforms including RNA-seq, microarray and immunoassay. In addition, the software of SSMD grants users to build in their own tissue-/data-specific knowledge of cell type-specific markers to reinforce the method. An R package of SSMD is released through GitHub: <https://github.com/xiaoyulu95/SSMD> and an R Shiny-based web server of SSMD is available at <https://ssmd.cccb.iupui.edu/>.

## Results

### Mathematical consideration and problem formulation

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

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

Here,  $\tilde{X}_{M_0 \times N}$  represents the observed gene expression matrix of  $M_0$  selected genes (a subset in  $M$ ) 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, and the relative proportions of the  $K_0$  cell types, respectively. In the conventional formulation of 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  $\mathcal{L}_2$  loss of the above linear equation. Because of the highly varied genetic and phenotypic background of mouse experiment,  $\tilde{S}_{M_0 \times K_0}$ ,  $M_0$  and  $K_0$  are usually varied and unknown, i.e. for each  $\tilde{X}_{M \times N}$  collected from tissues of certain microenvironment, what cell types are present, what gene markers each cell type expresses and how much they were expressed, could vary drastically due to the genetic and physiological perturbations. Correctly specified cell types  $K_0$  and selected cell type marker genes  $M_0$  can largely increase the prediction accuracy of  $\tilde{P}_{K_0 \times N}$ .

Table 1. 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, $X = A \bullet B^T$ . In this study, we consider all transcriptomics data are with error. Hence, the rank-1 matrix is defined by $X = A \bullet B^T + E$ , where the matrix rank of $X$ is 1 can be computed by the BCV algorithm detailed in Materials and 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, $X_{I \times J} = A \bullet B^T + E$ .
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

Table 1 lists the key mathematical definitions utilized in this study.

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

### SSMD analysis pipeline

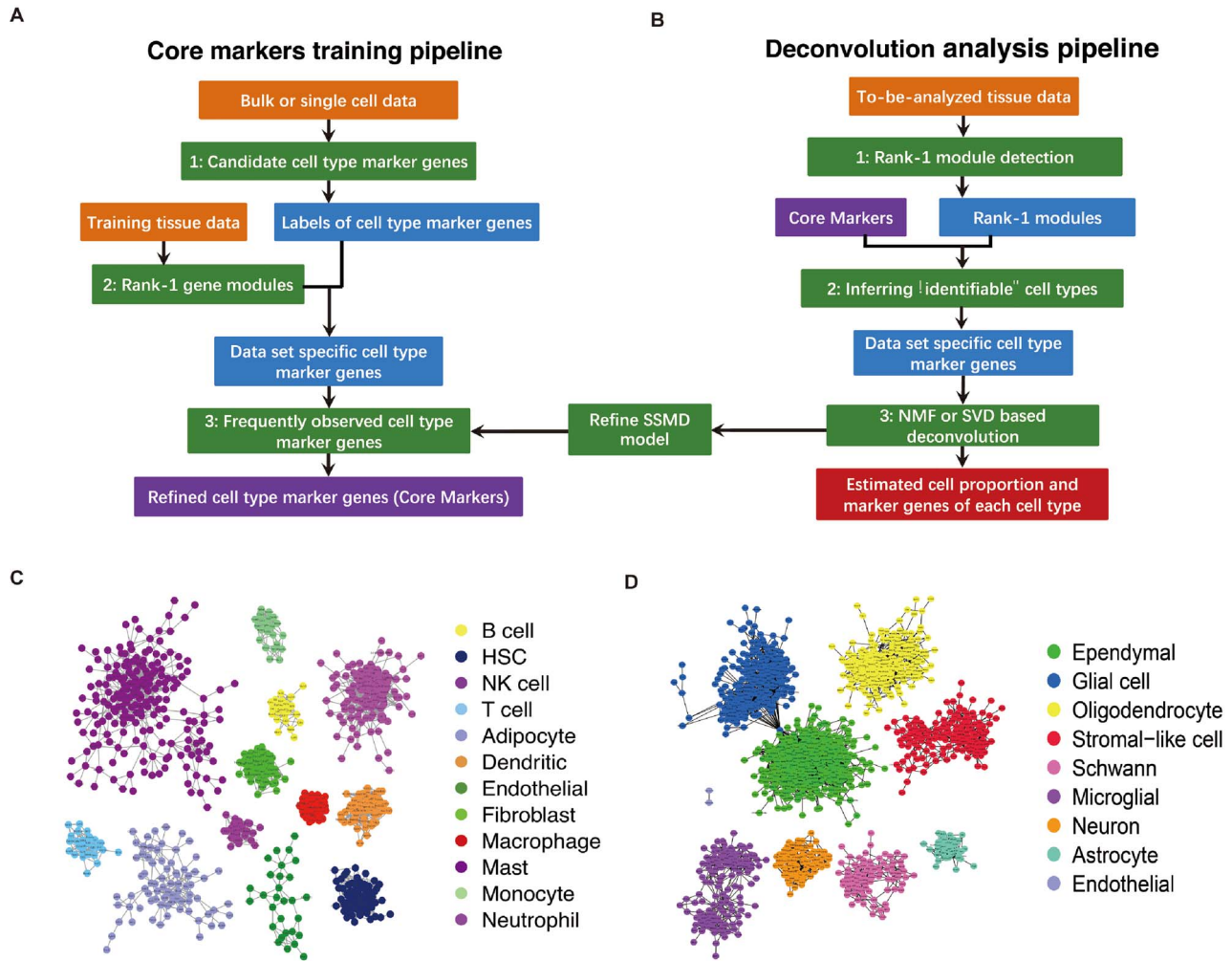
SSMD is a semi-supervised method composed by (i) training a large candidate list of cell type-specific marker genes, (ii) evaluating the identifiability of each cell type and confirming their marker genes for each to-be-deconvolved data and (iii) estimating the proportion of each cell type.

The training step is to look for genes that are more likely to serve as cell type marker genes through different tissue types and data sets, named as core marker lists. Specifically, we identified the genes that are commonly overexpressed in one cell type comparing with the others in bulk cell data and commonly form rank-1 matrices in tissue data, by using a very extensive set of training data sets collected from different mouse strains and tissue types (see details in Materials and methods). Figure 1A illustrates the procedure of SSMD to construct cell type core marker lists. On the bulk cell training data, we adopted a random walk-based approach to detect genes that are significantly expressed in higher quantities in one or a few cell types, than others (see details in Materials and methods). As

a result, a labeling matrix that annotates cell type specifically expressed genes will be constructed, which forms the first evidence of the potential marker genes for each cell type. Then, on each bulk training tissue data set, we further identified marker genes that form rank-1 submatrices with a community detection approach as detailed in Materials and methods. Only those modules, whose genes significantly and consistently overrepresent one and only one cell type across multiple training tissue data sets, are selected to form the core marker list. Notably, variations caused by different experiment batches, tissue types and mouse strains were handled by enabling certain errors in the random walk-based cell type-specific marker identification, i.e. identifying the genes overly expressed in the cell type comparing with the others in a certain proportion of the collected bulk cell data. In addition, data batch variation was also considered in the bulk data-based training step, by identifying the genes commonly serve as cell type-specific marker in more than 50% of analyzed bulk tissue training data. The goal of this training procedure is to summarize a relatively large list of commonly observed cell type-specific marker genes, which can be used to as semi-supervised information to identify data set-specific cell type marker for a further unsupervised deconvolution analysis.

Based on the cell type core markers, the deconvolution of any given bulk tissue data set is composed by the steps as illustrated in Figure 1B. SSMD first identifies all the rank-1 modules on the target data set by an iterative hierarchical clustering and bi-cross-validation (BCV) approach. Then, SSMD selects the rank-1 modules that are likely to be markers of a certain cell type for this data set, if genes in the modules largely overlap with the core marker list of one and only one cell type. Modules that are highly colinear will be merged. Consequently, genes in each module are called gene markers of one cell type, which satisfy the necessary condition for 'transcriptomically identifiable'. Notably, two modules may represent the same cell type, and they are treated as marker genes of different subtypes of the cell type. Here, the total number of modules is an estimate of the number of 'identifiable' cell types, i.e.  $K_0$ . Importantly, SSMD is an 'semi-supervised' approach, because the cell marker genes do not solely depend on the training data, but also the coexpression patterns of the marker genes in the target data set. In other words, SSMD addresses the variability issue of signature genes from one data set to another and has the potential to discover cell types not predefined. Algorithms of each computational step are detailed in Materials and methods. Complete flowchart of the SSMD pipeline is provided in Supplementary Figure S1.

The prediction of the cell type proportions is conducted using a constrained non-negative matrix factorization (NMF) method



**Figure 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 coincident as markers of one cell type in more than 50% of the training data sets. (D) Core markers of nine cell types in central nervous system. Notably, 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.

by solving the following optimization problem:

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

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 an all-1 column vector of length  $d$ ,  $\lambda$  is a hyperparameter selected by cross-validation, and other annotations follow equation (1). The constraint matrix  $C_{M_0 \times K_0}$  is enforced upon the regular NMF 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 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  $k$ . The solution to (equation 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, and the other is updated.  $\lambda$  can be tuned by using simulated tissue data with known cell proportion. In this study,

we tuned  $\lambda$  and empirically select  $\lambda$  as 10 when  $\tilde{X}_{M_0 \times N}$  is log-normalized microarray data or  $\log(X + 1)$  normalized FPKM/CPM/TPM RNA-seq data.

Following these procedures, and on a large collection of mouse bulk cell and tissue training data, we generated core marker gene lists for different TMEs: (i) for mouse blood, solid cancer and inflammatory tissues, 980 genes of 12 cell types, namely T cell, B cell, nature kill (NK) cell, hematopoietic stem cell (HSC), monocyte, macrophage, neutrophil, mast cell, adipocytes, fibroblast, dendritic cell and endothelial cell were discovered (Figure 1C); (ii) for mouse hematopoietic system, 2877 genes of 14 cell types namely HSC, common lymphoid progenitor, granulocyte-macrophage progenitors, megakaryocyte lineage-committed progenitor, erythroid cell, megakaryocyte-erythrocyte progenitors, multipotent progenitors, early myeloid progenitor, mature myeloid cell, precolony-forming unit erythroid, premegakaryocytic/erythroid progenitor, B cell, CD4+ T and CD8+ T cell were discovered (Supplementary Table S1) and (iii) for mouse central nervous system tissue, 1570 genes of nine cell types namely ependymal cell, general glial cell, oligodendrocyte, stromal-like cell, Schwann cell, microglial,



neuron and astrocyte were discovered (Figure 1D). Complete lists of the core marker genes are given in Supplementary Table S1. It is noteworthy that the size of core marker list ranges from 27 to 547 for different cell types. However, our analysis suggested that more than five marker genes that form a rank-1 matrix is sufficient for an accurate estimation of cell proportion. Note that, compared with conventional regression-based deconvolution analysis, SSMD only uses labels of the core markers as the semi-supervised information and identifies data set-specific cell type markers for a further unsupervised estimation of cell types, which grants a flexibility and robustness to handle the variation of cell type-specific marker genes and their expression scale through different mouse strains, tissue types and experimental platforms. In addition, the semi-supervised formulation of SSMD enables the inference of identifiability of each cell type and identification of rare or sub cell types.

### Benchmarking based on artificial tissue data simulated by using scRNA-seq data

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

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

We also benchmarked SSMD on simulated brain tissue data using two scRNA-seq data of central nervous systems. SSMD achieved more than 0.9 correlation in predicting the cell types microglial, stromal-like and ependymal subtypes in the simulated tissue data (Figure 2E and F). To the best of our knowledge, SSMD is the first of its kind method to specifically target mouse central nervous system decomposition. To benchmark SSMD, we selected MUSIC as the state-of-the-art method, which requires an additional input of scRNA-seq data to train context-specific gene signatures [25]. Here, we first utilized the same scRNA-seq data for tissue data simulation and signature training in MUSIC.

Nonsurprisingly, MUSIC achieved consistently good predictions (averaged  $\text{cor}=0.99$ ), and the predictions made by SSMD are very close to MUSIC with slightly lower correlations compared with MUSIC under this ideal setup. In sight the possible disparity caused by tissue, strain and experimental platform variations between the target tissue data and available scRNA-seq data for training cell markers, we also conducted a robustness test of MUSIC and SSMD (see details in Supplementary Notes). Our analysis suggested that MUSIC highly depends on the consistency of cell type-specific marker genes and their expression scale between the target tissue and the training scRNA-seq data. In contrast, the *de novo* data set-specific marker identification by SSMD enables a broader application to the tissue data without matched scRNA-seq data. Because EPIC, ICC and TICC cannot analyze brain tissue data and the melanoma and pancreas tissue were not covered by TICC, we did not include the comparison with these methods on the brain tissue data.

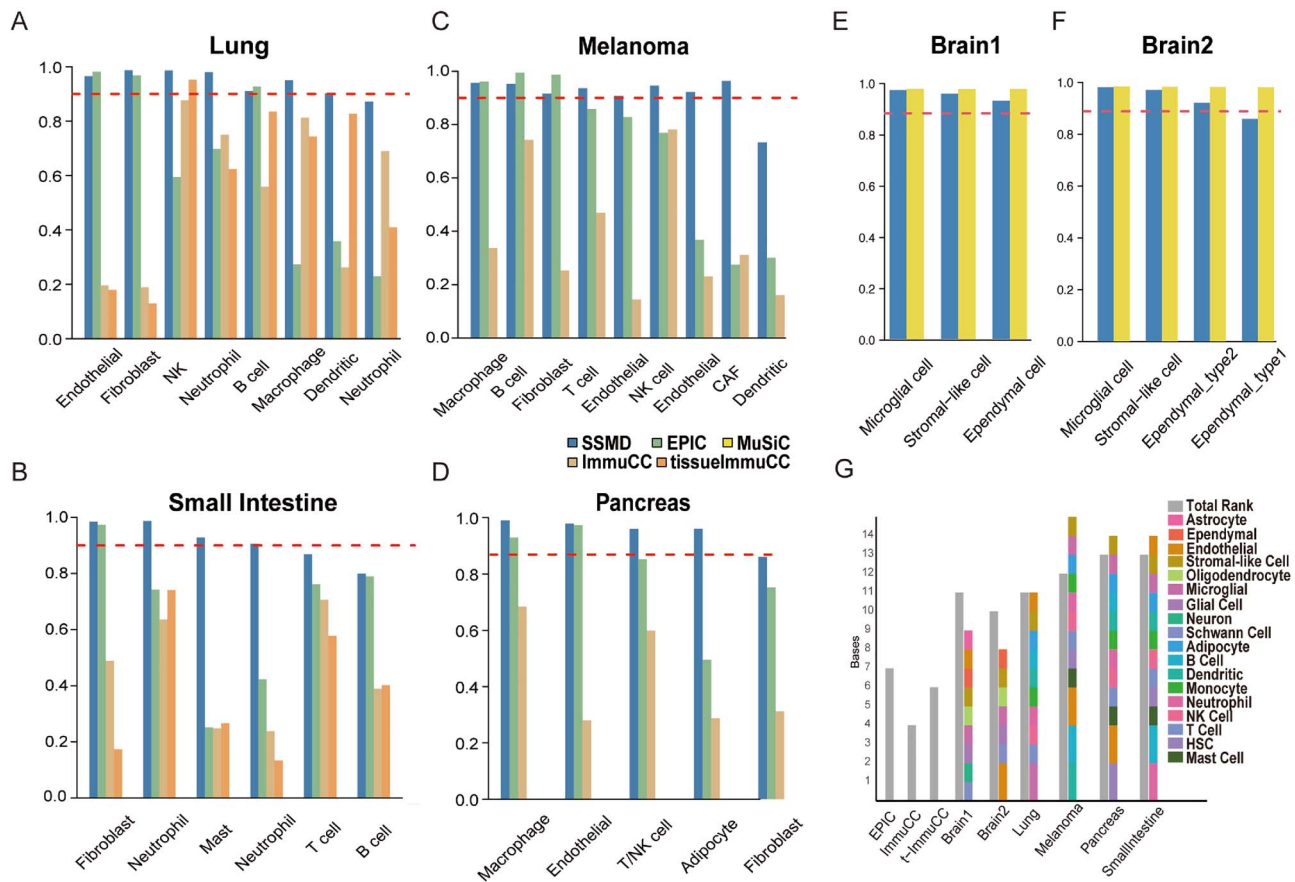
To validate the specificity of SSMD, we tested the total rank of the identified marker genes and compared with the number identified cell types (TIMER and CIBERSORT achieved 49.25% and 47.5% averaged prediction accuracy, and the proportion of cell types with higher than 0.9 prediction accuracy are 17.9% (5/28), and 3.6% (1/28)). We also compare the total matrix rank of the marker genes used in other methods and the number of cell types assumed in those methods. Comparing with the fixed number of cell types in other methods, the number of cell types predicted by SSMD better matches the total rank of the expression profile of identified marker genes. Our observation suggested SSMD can correctly estimate the number of cell types and select proper markers for cell type proportion estimation. It is noteworthy the predicted number of cell types may not exactly match the total rank of selected markers because possible colinearity among the true proportion of the cell types.

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

We generated tissue RNA-seq data of 11 mouse bone marrow tissue samples with matched cell counting using FACS (see details in Materials and methods). Application of SSMD on the RNA-seq data identified HSC, general myeloid progenitor (GMP), mature myeloid cell and pre-B cells and their cell type-specific markers. We also observed that the correlation between SSMD predicted and FACS measured amount of HSC, GMP, mature myeloid cell and B cells are 0.92, 0.8, 0.86 and 0.97, respectively, suggesting a high prediction accuracy of SSMD. Figure 3A–D shows the correlation between the SSMD predicted cell proportion and the FACS measured cell proportion of the four cell types. Figure 3E–H illustrates the FACS-based cell counting of the four cell types. Complete cell type-specific markers, cell proportions counted by FACS and predicted by SSMD were given in Supplementary Table S2. It is noteworthy that SSMD is not compared with other methods as none of the existing method is capable of predicting proportions of hematopoietic cell types.

### Application of SSMD to real mouse tissue transcriptomics data

We applied SSMD to nine cancer and eight central nervous system tissue data of four different experimental platforms, including one data set measured by immunoassay. On average, SSMD identified more than seven cell types in each of the cancer data, and the number of identified cell types is highly consistent with the total rank of the expression profile of the detected cell type-specific marker genes (Figure 4A). This indicates that



**Figure 2.** Method evaluation on scRNA-seq simulated tissue data. (A–D) Correlation between true and predicted cell proportions in the simulated lung (A), pancreas (B), small intestine (C) and mouse melanoma (D) tissue data. The x-axis represents cell type and y-axis represents prediction accuracy. Predictions made by SSMD, EPIC, ICC and TICC were dark blue, green, yellow and orange colored, respectively. The red dash line represents the 0.9 correlation cutoff. (E, F) Correlation between true and predicted cell proportions in the two simulated brain tissue data. (G) The total rank of the gene expression profile of selected marker genes in the six simulated tissue data (gray), and the total number of cell types identified by SSMD in each data set or assumed in other methods (left three gray bars).

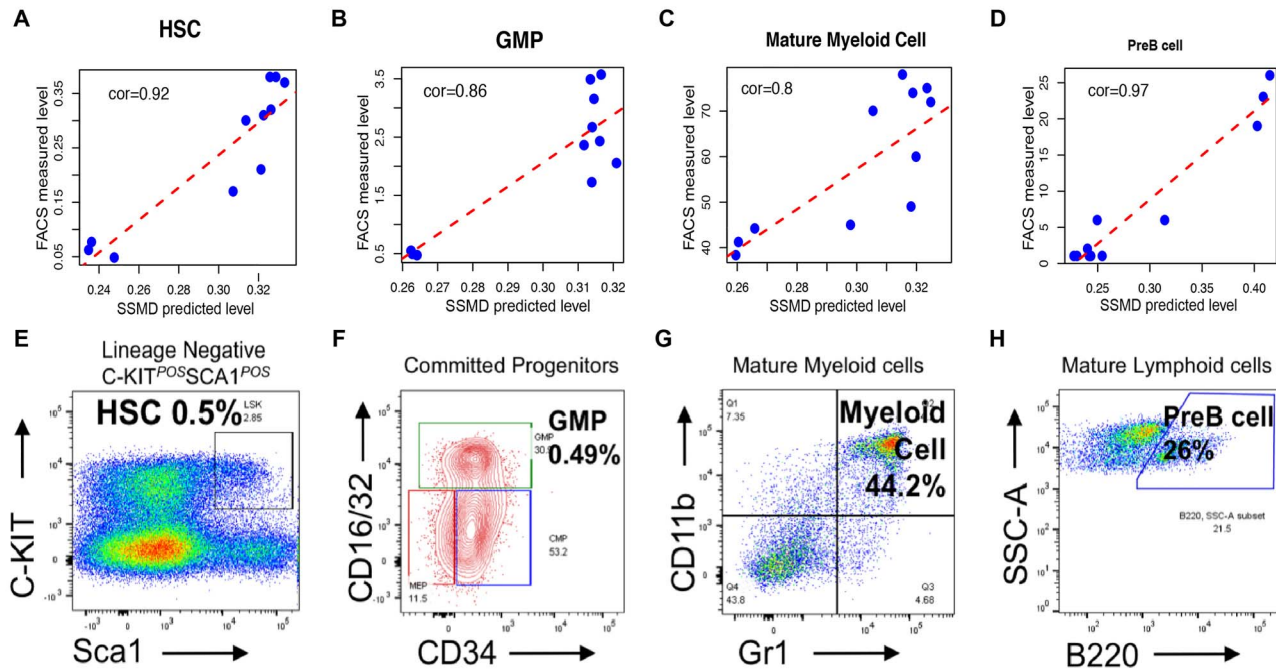
SSMD is capable of capturing the latent structure of the data. We further examined the explanation score (*E*-score), defined as the averaged absolute residual of the non-negative linear regression of each marker gene's expression on the predicted cell proportion, i.e. the average measure of how the predicted proportions could explain all the marker genes' expression levels. A high *E*-score is a necessary condition for an accurate cell proportion prediction. On average, the data set-specific markers genes of each cell type identified by SSMD achieved 0.73 *E*-score, whereas the average *E*-score of the marker genes used by EPIC and ICC is 0.45 and 0.3 (Figure 4B). Similarly, application of SSMD on eight central nervous system tissue data identified more than seven cell types on average. The number of identified cell types is highly consistent with the total rank of the gene expression profile of the marker genes (Figure 4C). In addition, the marker genes identified by SSMD achieved averaged 0.77 *E*-score for the cell types in central nervous system (Figure 4D). It is noteworthy that multiple marker sets of fibroblasts, myeloid or microglial cells that forming distinct rank-1 bases were identified in numerous data sets, suggesting the possible subtypes of these cell types identified by SSMD.

### Robustness analysis

We first evaluated the variation of cell type-specific markers through different mouse strains on one transcriptomic data set

of mouse liver tissue samples collected from 31 different mouse strains [26]. To the best of our knowledge, this is the only data set 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. Nine cell and their subtypes 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 (Figure 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, whereas 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.

We further examined the robustness of SSMD by evaluating its (i) sensitivity and (ii) specificity in identifying cell type-specific marker genes and its (iii) accuracy in assessing of cell proportions on the data of different sample sizes. Previous studies revealed that the robustness of the computation of coexpression correlation will decrease when the sample size is below 25. To comprehensively evaluate the method's robustness, we selected five data sets, namely GSE76095, GSE67186, GSE90885,



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

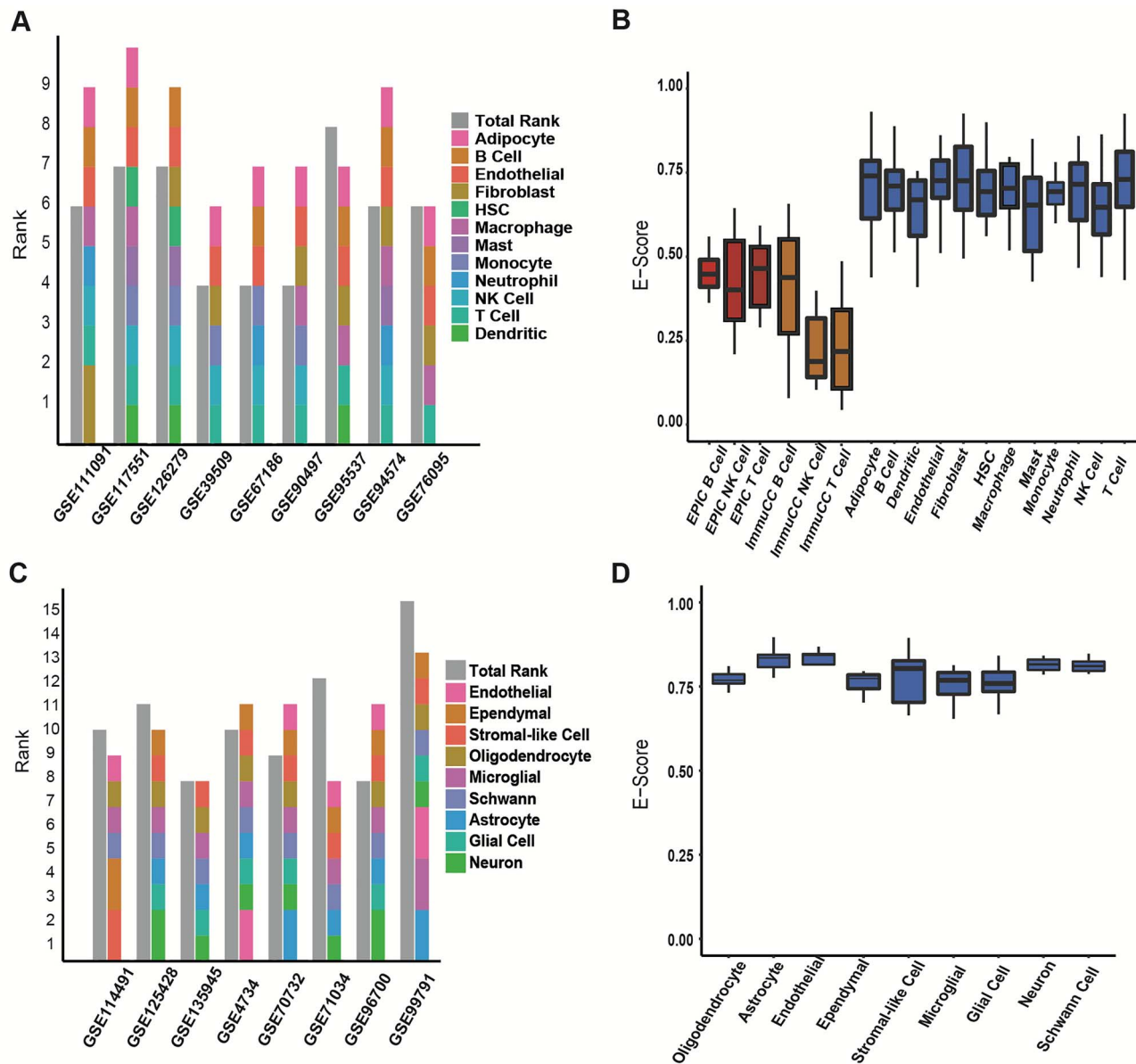
GSE94574 and GSE126279, with sample size ranging from 15 to 30 and randomly drew samples from each data set to build testing data sets of different sample size. We assumed the cell type markers and cell proportion inferred from whole data as ‘true’ markers and proportions and evaluated the consistency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

On average, all of the marker genes of the ‘true’ cell types were also identified when sample size is low (Figure 6A). In addition, the cell proportion of 92.3%, 94.6% and 98.9% of the correctly identified cell types were with more than 0.9 correlation with their ‘true’ proportions when the sample size is 6, 12 and above 20 (Figure 6A). Our analysis suggested a high robustness of the sensitivity and prediction accuracy of SSMD when sample size is as small as six, i.e. the commonly used sample size in two-condition-comparison experiment (three samples versus three samples). However, as a trade-off, there is a high false discovery rate of cell type-specific modules when sample size is small, due to the low specificity of gene coexpress analysis. To control the false discoveries on small data sets, we further derived a more ‘stringent’ set of 341 cell type-specific marker genes among the core marker set (see details in Materials and methods). Our method validation demonstrated a slight drop of the sensitivity and prediction accuracy when using the stringent marker set on small data set (Figure 6B), whereas the specificity of the identified cell type-specific markers increased to from 54.4% to 72.6% when sample size is above 12 (Figure 6C). Figure 6D illustrates the E-score of the cell type-specific marker genes identified by using the core and the more stringent marker set with respect to different sample size. The E-score of the cell types marker genes identified by using the more stringent

marker set were significantly higher than the ones identified by using the general core marker sets when sample size is below 10, also demonstrating the stringent core marker sets can effectively increase the analysis specificity when sample size is small.

## Discussion

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



**Figure 4.** Prediction of SSMD on real tissue data. (A, C) The total rank of the gene expression profile of selected marker genes (gray) 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, ICC: yellow).

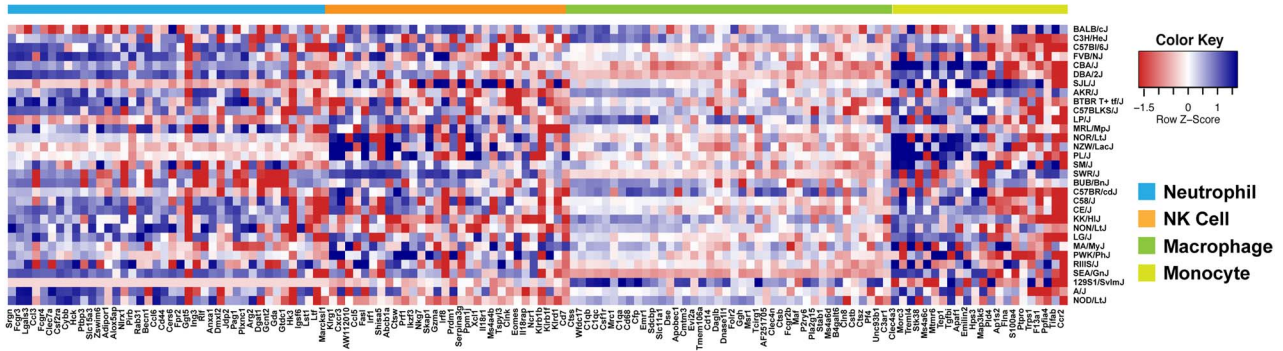
unsupervised constrained-NMF or singular value decomposition (SVD)-based deconvolution on the selected marker genes further excludes the adversarial batch effects.

It is noteworthy a successful identification of the rank-1 modules depends on a relatively large samples (>25) sharing cell types and marker genes. Currently, SSMD cannot be applied to the data with a single or small sample size. However, we consider such a trade-off between sample size and prediction robustness is highly worthwhile, especially considering using SSMD as an exploratory tool in large scale publicly available mouse transcriptomics data. After all, the predicted proportions are often to be associated with other biological and clinical features, which will be severely underpowered with a small sample size.

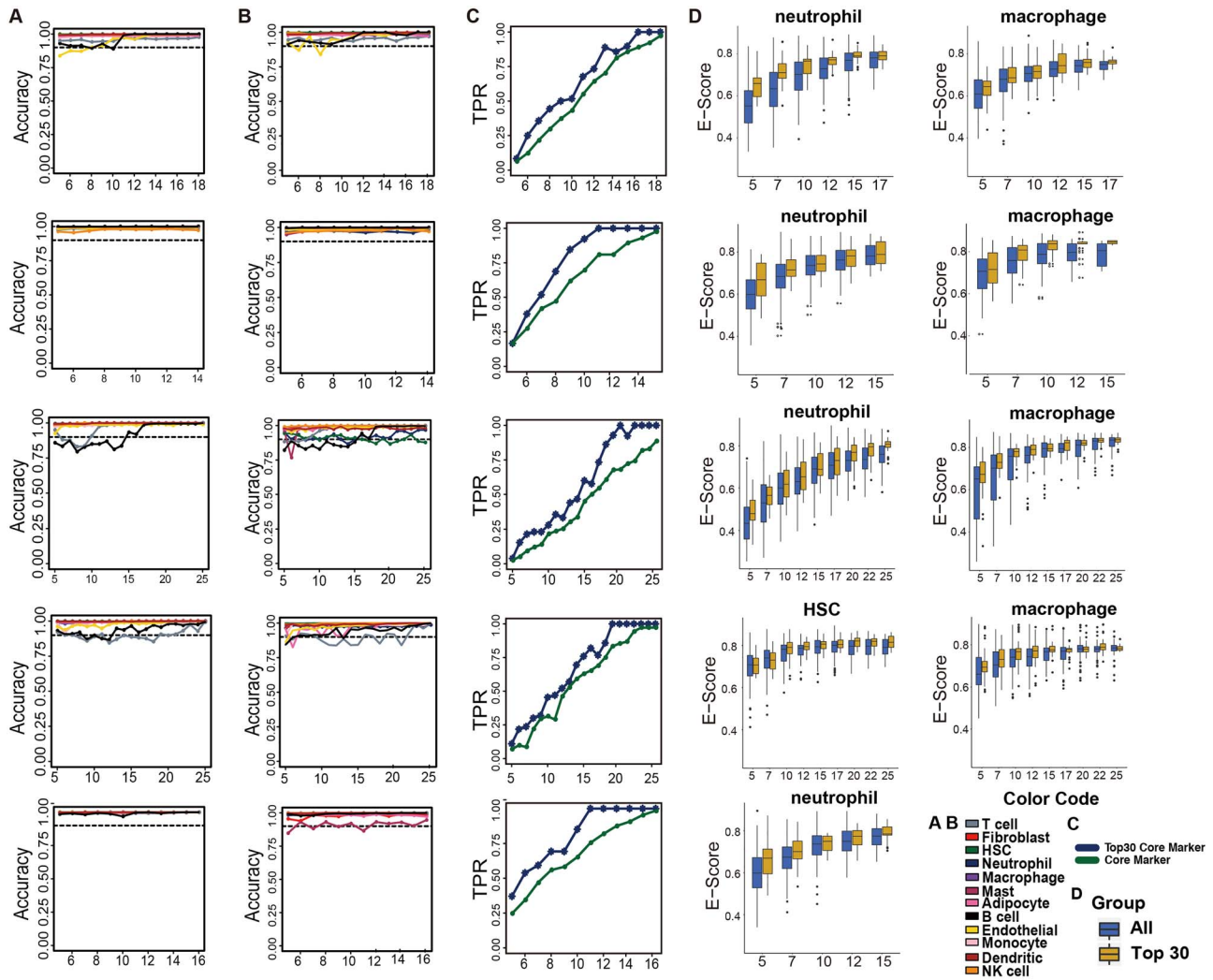
We released an R package of SSMD via <https://github.com/xiaoyulu95/SSMD> and a web server via <https://ssmd.cccb.iu>

[pui.edu/](http://pui.edu/). As illustrated in [Supplementary Figure S2A](#), the input data are a mouse tissue transcriptomics data and user-selected tissue-specific cell type core marker sets. Currently, SSMD offers general core and stringent marker sets of 6 cell types in blood system, 12 cell types in normal, inflammatory and cancer tissue, 9 cell types in central nerve systems and 14 cell types in hematopoietic systems. [Supplementary Figure S2B](#) illustrates a practical guide for using SSMD of different tissues and sample size. The input of SSMD is a mouse tissue expression data set and user-selected tissue environment category. The output of SSMD includes the identified data set-specific cell type markers and the estimated sample-wise relative proportion of each identifiable cell type. We consider the currently included cell types are comprehensive enough to cover major cell types in mouse. However, the tissue-specific cell types (e.g. liver cells in liver tissue, colon cells in colon tissue, etc.) were not included





**Figure 5.** Correlation between expression level of strain-specific cell type marker genes and predicted cell proportion. High correlation is a necessary but nonsufficient 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, NK, macrophage and monocyte, were colored on the column side bar.



**Figure 6.** Performance evaluation of different sample size. (A) Prediction accuracy (y-axis) in different sample size (x-axis) using all core markers. Accuracy is the Pearson correlation between predicted proportion using only selected small sample and using all samples. (B) Prediction accuracy (y-axis) in different sample size (x-axis) using selected stringent markers. (C) True positive rate (y-axis) of the cell type-specific markers identified by using the stringent markers (blue) and core markers (green) with respect to different sample size (x-axis). (D) E-score for using coexpression modules consisting of all core markers and only selected stringent markers. From top to bottom, the statistics were derived from GSE76095, GSE67186, GSE90885, GSE94574 and GSE126279.

in our training scope. As forming rank-1 pattern among marker genes is a necessary but non-sufficient condition of identifiable cell types, SSMD R package can also output rank-1 modules that do not enrich the core markers of any cell type, which could possibly be markers of rare cell types. The user could further investigate whether the gene module corresponds to a real cell type or not. Another key feature of the web server is that users are welcome to contribute their data to reinforce the training of cell type-specific marker genes.

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

## Materials and methods

### Random walk-based identification of cell type specifically expressed genes from tissue data

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

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

### Identification of rank-1 cell type uniquely expressed gene modules

To screen genes that form tight rank-1 modules on various tissue training data sets, SSMD performs a community detection method among the genes specifically expressed in each

cell type as stored the labeling matrix. A correlation matrix was first built among cell type specifically expressed genes, and the significance cutoff of correlation was determined by random matrix theory (RMT). RMT has been widely used to understand the low-rank structure encoded in biological data. In this study, an RMT-based approach developed by Luo et al. [30] was used to determine the threshold of significant correlation for each data set. `rm.get.threshold` functions in the `RMThreshold` R package was utilized. Specifically, RMT indicated that the nearest neighbor spacing distribution of eigenvalues will have a characteristic change when the threshold properly separates signal from noise. By removing all the below-threshold correlation elements, the coexpression modules can be more robustly unraveled. Then, hierarchical clustering was performed using the correlation matrix as similarity measure.

Specifically, SSMD gradually increases the height of the hierarchical clustering at which the tree is cut. At each height, the number of genes, the average correlation among the genes and the rank of the matrix composed of the genes in each of the cluster is calculated. Here, matrix rank is determined by a modified BCV algorithm. SSMD stops scanning the hierarchical tree if all the clusters contain less than  $q_0$  genes or the three following criteria are met for all the clusters: (i) with at least  $q_0$  genes; (ii) the average correlation among the genes is above the threshold determined by RMT and (iii) the rank of the expression matrix profile of the genes in the cluster is 1. In this study,  $q_0=7$  is used. Such an iterative approach will eventually select the clusters with at least  $q_0$  genes, each of which is considered as possible cell-specific marker genes specific to this data set. SSMD merges modules until the canonical correlation between any pair of module is lower than a cutoff  $\text{cor}_{\text{cut}}$  or the number of current modules is not larger than the total rank of the gene expression profile of the selected data set-specific markers genes. In this study, we utilized  $\text{cor}_{\text{cut}} = 0.9$ .

### A modified BCV rank test

BCV has been developed to estimate the matrix rank for SVD and NMF, which requires a prefixed low-dimension  $K$  and two low-rank matrices for the approximation  $X_{M \times N} = W_{M \times K} \bullet H_{K \times N}$ . The error distribution of gene expression data is usually non-identical/independent, mostly because a gene's expression can be affected by its major transcriptional regulators, other biological pathways and experimental bias. Hence, undesired biological characteristics and experimental bias may form significant dimensions in a gene expression data [31]. In sight of this, we developed a modified BCV rank test (Algorithm 1) to minimize the effect of the non-i.i.d errors in assessing the matrix rank of a gene expression data.

After running the rank-1 module detection on all the training bulk tissue data sets, those genes commonly identified in the rank-1 modules in more than 40% (70%) data sets were selected as core (stringent) markers. The list of stringent marker sets was derived with more stringent criterion, which is particularly useful for the analysis of small sample-sized target data. Core markers of cells in central nervous systems were identified by a similar approach on the brain training tissue data sets. Due to the limitation of hematopoietic system tissue training data, its core markers were selected as the genes specifically overexpressed in each hematopoietic cell type, by using the criteria: the gene's expression level is above 10% quantile in one cell type and below 50% in the other cell types. Complete lists of selected core and stringent marker sets were given in [Supplementary Table S1](#).

**Algorithm 1.** Modified BCV 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  $k = 1 \dots \min(M_0, N_0)$ 
         $BCV(k, r) = \sum_{i=1}^{M_0} \sum_{j=1}^{N_0} \left\| A_r - B_r D_r^{\hat{\Sigma}^{(k)+}} C_r \right\|_F^2 (*)$ 
    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  $\sigma_1^+ \geq \sigma_2^+ \geq \dots \geq \sigma_p^+ > 0$ . Define  $\hat{D}^{(k)+} = \Sigma^{k+}$ 

```

**Estimation of cell proportion**

Two methods were utilized to estimate cell proportion: (i) SVD-based computation. With cell type-specific markers derived, the first row base of the gene expression profile of the marker genes is directly utilized as an estimation of the cell proportion, which can be directly computed by SVD. (ii) Constraint NMF-based computation. With the number of identifiable cell types and 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 be estimated by minimizing the following objective function:

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

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

**E-score and comparison with state-of-the-arts methods**

An E-score was utilized to evaluate the goodness that each marker gene's expression is fitted by the predicted cell proportions:

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

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

**Data used in this study****Bulk cell training data sets**

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

**Single-cell RNA-sequencing data**

One mouse melanoma scRNA-seq data set (6638, 9) was acquired from the Human Cell Atlas database [34]. Three scRNA-seq data sets of lung (4485, 12), pancreas (4405, 8) and small intestine (4764, 10) and two sets of brain tissue (3679, 7 and 1099, 6) were accessed from Mouse Cell Atlas data portal [35]. The two numbers in the parenthesis indicate the number of cell samples



and cell types of each data set. We specifically selected the cells with UMI more than 500 to exclude low-quality cells. Cell labels were either provided in the original data or curated using Seurat v3 with cell type-specific genes [36, 37].

#### Training tissue data from cancer and blood

In total, 33 cancer tissue data sets of nine cancer types generated by four popular experimental platforms were collected, namely Illumina HiSeq 2000 *Mus musculus*, Affymetrix Mouse Genome 430 2.0 Array, Illumina HiSeq 2500 *M. musculus* and Affymetrix Mouse Genome 430A 2.0 Array from GEO database. Each data set has at least 15 samples. We did not consider data sets from immunodeficient mouse, mouse cell lines and PDX models, as only real cancer or blood microenvironment is considered. A data set of liver tissue collected from 31 mouse strains (GSE55489) were utilized to evaluate the variation of cell type-specific markers through different mouse strains [26].

#### Brain tissue data

Fourteen data sets of mouse brain tissues generated by two experimental platforms, namely Illumina HiSeq 2500 *M. musculus* and Affymetrix Mouse Genome 430 2.0 Array were collected from Gene Expression Omnibus. Data sets were split into sub-data sets of different brain regions. Each data set has at least 40 samples. The complete training data information is available in [Supplementary Table S3](#).

#### Hematopoietic system tissue and FACS data

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

#### Generation of simulated bulk tissue data from scRNA-seq data

Cell types in each scRNA-seq data were labeled by the cell clusters provided in the original works or by using Seurat pipeline with default parameters. Detailed information of the scRNA-seq data and cell type annotation is given in [Supplementary Table S3](#). For each data set, we simulate bulk tissue data by: (i) removing insignificantly expressed genes, (ii) randomly generate the proportion of each cell type, called true proportion in this paper, which follows a Dirichlet distribution

and (iii) draw cells randomly from the cell pool with replacement according to the cell type proportion, and sum up the expression values of all cells to produce a pseudo bulk tissue data. The insignificant expressed genes were identified by left truncated mixture Gaussian model [39, 40]. The Dirichlet distribution matrix was generated with R package 'DirichletReg' [41].

#### Key Points

- We provide a novel tissue deconvolution method, namely SSMD, which is specifically designed for mouse data to handle the variations caused by different mouse strain, genetic and phenotypic background and experimental platforms.
- SSMD is capable to detect data set- and tissue microenvironment-specific cell markers for more than 30 cell types in mouse blood, inflammatory tissue, cancer and central nervous system.
- SSMD achieve much improved performance in estimating relative proportion of the cell types compared with state-of-the-art methods.
- The semi-supervised setting enables the application of SSMD on transcriptomics, DNA methylation and ATAC-seq data.
- A user-friendly R package and an R shiny of SSMD-based web server are also developed.

#### Supplementary Data

[Supplementary data](#) are available online at [https://academic.oup.com/bib](https://academic.oup.com/bib/article/22/4/bbaa307/5998844).

#### Acknowledgements

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

#### Funding

National Science Foundation Div Of Information & Intelligent Systems (No. 1850360); National Institute of General Medical Sciences (R01 award #1R01GM131399-01); Showalter Young Investigator Award from Indiana CTSI.

#### References

1. Beck JA, Lloyd S, Hafezparast M, et al. Genealogies of mouse inbred strains. *Nat Genet* 2000;24(1):23–5.
2. Rosenthal N, Brown S. The mouse ascending: perspectives for human-disease models. *Nat Cell Biol* 2007;9(9):993–9.
3. Van der Jeught K, Sun Y, Fang Y, et al. ST2 as checkpoint target for colorectal cancer immunotherapy. *JCI Insight* 2020;5(9).
4. Mund JA, Park S-J, Smith AE, et al. Genetic disruption of the small GTPase RAC1 prevents plexiform neurofibroma formation in mice with neurofibromatosis type 1. *J Biol Chem* 2020; p. jbc. RA119. 010981.
5. Huang M, Kim HG, Zhong X, et al. Sestrin 3 protects against diet-induced nonalcoholic steatohepatitis in mice through suppression of transforming growth factor  $\beta$  signal transduction. *Hepatology* 2020;71(1):76–92.



6. Pandey R, Ramdas B, Wan C, et al. SHP2 inhibition reduces leukemogenesis in models of combined genetic and epigenetic mutations. *J Clin Invest* 2019;129(12):5468–73.
7. Zhang C, Cao S, Xu Y. Population dynamics inside cancer biomass driven by repeated hypoxia-reoxygenation cycles. *Quant Biol* 2014;2(3):85–99.
8. Hackl H, Charoentong P, Finotello F, et al. Computational genomics tools for dissecting tumour-immune cell interactions. *Nat Rev Genet* 2016;17(8):441.
9. Li B, Severson E, Pignon J-C, et al. Comprehensive analyses of tumor immunity: implications for cancer immunotherapy. *Genome Biol* 2016;17(1):174.
10. Wang X, Park J, Susztak K, et al. Bulk tissue cell type deconvolution with multi-subject single-cell expression reference. 2019;10(1):380.
11. Racle J, de Jonge K, Baumgaertner P, et al. Simultaneous enumeration of cancer and immune cell types from bulk tumor gene expression data. *Elife* 2017;6:e26476.
12. Newman AM, Steen CB, Liu CL, et al. Determining cell type abundance and expression from bulk tissues with digital cytometry. *Nat Biotechnol* 2019;37(7):773–82.
13. Newman AM, Liu CL, Green MR, et al. Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods* 2015;12(5):453.
14. Li B, Severson E, Pignon JC, et al. Comprehensive analyses of tumor immunity: implications for cancer immunotherapy. *Genome Biol* 2016;17(1):174.
15. Gaujoux R, Seoighe CJB. CellMix: a comprehensive toolbox for gene expression deconvolution. 2013;29(17):2211–12.
16. Frishberg A, Peshes-Yaloz N, Cohn O, et al. Cell composition analysis of bulk genomics using single-cell data. *Nat Methods* 2019;16(4):327–32.
17. Finotello F, Trajanoski ZJCI. Immunotherapy, Quantifying tumor-infiltrating immune cells from transcriptomics data. 2018;67(7):1031–40.
18. Abbas AR, Wolslegel K, Seshasayee D, et al. Deconvolution of blood microarray data identifies cellular activation patterns in systemic lupus erythematosus. 2009;4(7):e6098.
19. Abbas A, Baldwin D, Ma Y, et al. Immune response in silico (IRIS): immune-specific genes identified from a compendium of microarray expression data. 2005;6(4):319.
20. Chen Z, Huang A, Sun J, et al. Inference of immune cell composition on the expression profiles of mouse tissue. *Sci Rep* 2017;7:40508.
21. Marques S, Zeisel A, Codeluppi S, et al. Oligodendrocyte heterogeneity in the mouse juvenile and adult central nervous system. *Science* 2016;352(6291):1326–29.
22. La Manno G, Gyllborg D, Codeluppi S, et al. Molecular diversity of midbrain development in mouse, human, and stem cells. *Cell* 2016;167(2):566–80 e19.
23. Codeluppi S, Borm LE, Zeisel A, et al. Spatial organization of the somatosensory cortex revealed by osmFISH. *Nat Methods* 2018;15(11):932–35.
24. Chang W, Wan C, Lu X, et al. ICTD: A semi-supervised cell type identification and deconvolution method for multi-omics data. *bioRxiv* 2019;426593.
25. Wang X, Park J, Susztak K, et al. Bulk tissue cell type deconvolution with multi-subject single-cell expression reference. *Nat Commun* 2019;10(1):1–9.
26. Church RJ, Wu H, Mosedale M, et al. A systems biology approach utilizing a mouse diversity panel identifies genetic differences influencing isoniazid-induced microvesicular steatosis. *Toxicol Sci* 2014;140(2):481–92.
27. Wan C, Chang W, Zhao T, et al. Denoising individual bias for a fairer binary submatrix detection. 2020; arXiv preprint arXiv:2007.15816.
28. Wan C, Chang W, Zhao T, et al. Fast and efficient boolean matrix factorization by geometric segmentation. *arXiv* 2019; arXiv: 1909.03991.
29. Chang W, Wan C, Zang Y, et al. Supervised clustering of high dimensional data using regularized mixture modeling. 2020; arXiv preprint arXiv: 2007.09720.
30. Luo F, Yang Y, Zhong J, et al. Constructing gene co-expression networks and predicting functions of unknown genes by random matrix theory. *BMC Bioinformatics* 2007;8(1):299.
31. Lopez R, Regier J, Cole MB, et al. Deep generative modeling for single-cell transcriptomics. 2018;15(12):1053.
32. Pepper SD, Saunders EK, Edwards LE, et al. The utility of MAS5 expression summary and detection call algorithms. *BMC Bioinformatics* 2007;8(1):273.
33. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics* 2007;8(1):118–27.
34. Regev A, Teichmann SA, Lander ES, et al. Science forum: the human cell atlas. *Elife* 2017;6:e27041.
35. Han X, Wang R, Zhou Y, et al. Mapping the mouse cell atlas by microwell-seq. *Cell* 2018;172(5):1091–107 e17.
36. Stuart T, Butler A, Hoffman P, et al. Comprehensive integration of single-cell data. *Cell* 2019.
37. Butler A, Hoffman P, Smibert P, et al. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat Biotechnol* 2018;36(5):411.
38. Dobin A, Davis CA, Schlesinger F, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 2013;29(1):15–21.
39. Wan C, Chang W, Zhang Y, et al. LTMG: a novel statistical modeling of transcriptional expression states in single-cell RNA-Seq data. *Nucleic Acids Res* 2019;47(18):e111–1.
40. Zhang Y, Wan C, Wang P, et al. M3S: A comprehensive model selection for multi-modal single-cell RNA sequencing data. *BMC Bioinformatics* 2019;20(24):1–5.
41. Maier MJ. DirichletReg: dirichlet regression for compositional data in R. Research report series/departement of statistics and mathematics, 125. WU Vienna University of Economics and Business, Vienna. <http://epub.wu.ac.at/4077/>. 2014.