# Bulk Tissue Cell Type Deconvolution with Multi-Subject Single-Cell Expression Reference

Xuran Wang[1], Jihwan Park[2], Katalin Susztak[2], Nancy R. Zhang[3]*, and Mingyao Li[4]*

1) Graduate Group in Applied Mathematics and Computational Science, University of Pennsylvania, Philadelphia, PA

2) Departments of Medicine and Genetics, University of Pennsylvania, Philadelphia, PA

3) Department of Statistics, The Wharton School, University of Pennsylvania, Philadelphia, PA

4) Department of Biostatistics, Epidemiology & Informatics, University of Pennsylvania, Philadelphia, PA


* Correspondence:
 Nancy R. Zhang
 nzh@wharton.upenn.edu
 (215) 898-8007

 Mingyao Li
 mingyao@pennmedicine.upenn.edu
 (215) 746-3916

**Abstract**

**We present MuSiC, a method that utilizes cell-type specific gene expression from single-cell RNA sequencing (RNA-seq) data to characterize cell type compositions from bulk RNA-seq data in complex tissues. When applied to pancreatic islet and whole kidney expression data in human, mouse, and rats, MuSiC outperformed existing methods, especially for tissues with closely related cell types. MuSiC enables characterization of cellular heterogeneity of complex tissues for identification of disease mechanisms.**

Bulk tissue RNA-seq is a widely adopted method to understand genome-wide transcriptomic variations in different conditions such as disease states. Bulk RNA-seq measures the average expression of genes, which is the sum of cell type-specific gene expression weighted by cell type proportions. Knowledge of cell type composition and their proportions in intact tissues is important, because certain cell types are more vulnerable for disease than others. Characterizing the variation of cell type composition across subjects can identify cellular targets of disease, and adjusting for these variations can clarify downstream analysis.

The rapid development of single-cell RNA-seq (scRNA-seq) technologies have enabled cell type-specific transcriptome profiling. Although cell type composition and proportions are obtainable from scRNA-seq, scRNA-seq is still costly, prohibiting its application in clinical studies that involve a large number of subjects. Furthermore, scRNA-seq is not well suited to characterizing cell type proportions in a solid tissue, because the cell dissociation step is biased towards certain cell types[1].

48  Computational methods have been developed to deconvolve cell type proportions using
49  cell type-specific gene expression references[2]. CIBERSORT[3], based on support vector
50  regression, is a widely used method designed for microarray data. More recently,
51  BSEQ-sc[4] extended CIBERSORT to allow the use of scRNA-seq gene expression as a
52  reference. TIMER[5], developed for cancer data, focuses on the quantification of immune
53  cell infiltration. These methods rely on pre-selected cell type-specific marker genes, and
54  thus are sensitive to the choice of significance threshold. More importantly, these
55  methods ignore cross-subject heterogeneity in cell type-specific gene expression as
56  well as within-cell type stochasticity of single-cell gene expression, both of which cannot
57  be ignored based on our analysis of multiple scRNA-seq datasets (**Supplementary**
58  **Figure 1a**).
59
60  Here we introduce a new MUlti-Subject SIngle Cell deconvolution (MuSiC) method
61  (https://github.com/xuranw/MuSiC) that utilizes cross-subject scRNA-seq to estimate
62  cell type proportions in bulk RNA-seq data (**Figure 1**). A key concept in MuSiC is
63  "marker gene stability". We show that, when using scRNA-seq data as a reference for
64  cell type deconvolution, two fundamental types of stability must be considered: cross-
65  subject and cross-cell, in which the first is to guard against bias in subject selection, and
66  the second is to guard against bias in cell capture in scRNA-seq. By incorporating both
67  types of stability, MuSiC allows for scRNA-seq datasets to serve as effective references
68  for independent bulk RNA-seq datasets involving different individuals.
69
70  Rather than pre-selecting marker genes from scRNA-seq based only on mean
71  expression, MuSiC gives weight to each gene, allowing for the use of a larger set of
72  genes in deconvolution. The weighting scheme prioritizes stable genes across subjects:
73  up-weighing genes with low cross-subject variance (informative genes) and down-
74  weighing genes with high cross-subject variance (non-informative genes). This
75  requirement on cross-subject stability is critical for transferring cell type-specific gene
76  expression information from one dataset to another.
77
78  Solid tissues often contain closely related cell types, and correlation of gene expression
79  between these cell types leads to collinearity, making it difficult to resolve their relative
80  proportions in bulk data. To deal with collinearity, MuSiC employs a tree-guided
81  procedure that recursively zooms in on closely related cell types. Briefly, we first group
82  similar cell types into the same cluster and estimate cluster proportions, then recursively
83  repeat this procedure within each cluster (**Figure 1**). At each recursion stage, we only
84  use genes that have low within-cluster variance, a.k.a. the cross-cell stable genes. This
85  is critical as the mean expression estimates of genes with high variance are affected by
86  the pervasive bias in cell capture of scRNA-seq experiments, and thus cannot serve as
87  reliable reference. See online methods for details.
88
89  To demonstrate and evaluate MuSiC, we started with a well-studied tissue, the islets of
90  Langerhans, which are clusters of endocrine cells within the pancreas that are essential
91  for blood glucose homeostasis. Pancreatic islets contain five endocrine cell types
92  (α,β,δ,ϵ, and γ), of which β cells, which secrete insulin, are gradually lost during type 2
93  diabetes (T2D). We applied MuSiC to bulk pancreatic islet RNA-seq samples from 89

94  donors from Fadista et al.[6], to estimate cell type proportions and to characterize their
95  associations with hemoglobin A1c (HbA1c) level, an important biomarker for T2D. We
96  were motivated to re-analyze this data because, as shown in **Figure 2** and in Baron et
97  al.[4], existing methods failed to recover the correct β cell proportions, which should be
98  around 50-60%[7], and also failed to recover their expected negative relationship with
99  HbA1c level. As reference, we experimented with scRNA-seq data from two sources: 6
100  healthy and 4 T2D adult donors from Segerstolpe et al.[8], and 12 healthy and 6 T2D
101  adult donors from Xin et al.[9]. All bulk and single-cell datasets in this analysis are
102  summarized in **Supplementary Table 1**.
103
104  First, to systematically benchmark, we applied MuSiC and three other methods
105  (Nonnegative least squares (NNLS), CIBERSORT, and BSEQ-sc) to artificial bulk RNA-
106  seq data constructed by simply summing the scRNA-seq read counts across cells for
107  each single-cell sequenced subject. In this case, true cell type proportions are known,
108  which allows the evaluation of accuracy. More details on artificial bulk construction are
109  described in the **Supplementary Note**. **Figure 2a**, **Supplementary Figure 1c** and
110  **Supplementary Figure 2b** show the estimation results when the artificial bulk and the
111  single-cell reference data are from the same study, either both from Segerstolpe et al.[8]
112  or both from Xin et al.[9]. MuSiC achieves improved accuracy over existing procedures.
113  **Figure 2b** and **Supplementary Figure 2a** show the estimation results when the
114  artificial bulk and the single-cell reference data are from different studies. This is a more
115  challenging but more realistic scenario, since library preparation protocols vary across
116  labs and bulk deconvolution analyses are often performed using single-cell reference
117  generated by others.  MuSiC still maintains high accuracy, while other methods perform
118  substantially worse.  Further comparisons show that, unlike existing methods that rely
119  on pre-selected marker genes, MuSiC gives accurate results when the cell type
120  composition in the bulk data is substantially different from that of the single cell
121  reference (**Supplementary Figure 2c** and **Supplementary Note 2**), and when the bulk
122  tissue contains minority cell types that are missing in the reference (**Supplementary
123  Figure 3** and **Supplementary Note 3**). MuSiC's ability to transfer knowledge across
124  data sources is derived from its consideration of marker gene stability.
125
126  We now turn to the deconvolution of bulk RNA-seq data from Fadista et al.[6]. We used
127  the scRNA-seq data from Segerstolpe et al. as reference for all methods. MuSiC
128  recovers the expected ~50-60% β cell proportion for the healthy subjects[7], whereas
129  other methods grossly overestimate the proportion of α cells and underestimate the
130  proportion of β cells. Furthermore, MuSiC detects a significant association of β cell
131  proportion with HbA1c level (p-value 0.00126, **Figure 2d**). Based on clinical standard,
132  HbA1c level <6.0% is classified as normal, and >6.5% is classified as diabetic. After
133  adjusting for age, gender and body mass index, MuSiC estimates suggest that 0.5%
134  increase in HbA1c level, representing the magnitude of increase from normal to the
135  diabetes cutoff, corresponds to a drop of 6.14% ± 4.98% in β cell proportion.
136
137  As a second tissue example, we used the kidney, a complex organ consisting of several
138  anatomically distinct segments each playing critical roles in the filtration and
139  reabsorption of electrolytes and small molecules of the blood. Chronic kidney disease

3

140  (CKD), the gradual loss of kidney function, is increasingly recognized as a major health
141  problem, affecting 10-16% of the global adult population. We aim to characterize how
142  kidney cell type composition changes during CKD. Fibrosis is the histologic hallmark
143  common to all CKD models, and hence, we analyzed the bulk RNA-seq data from three
144  mouse models for renal fibrosis: unilateral ureteric obstruction induced by surgical
145  ligation of the ureter (UUO, Arvaniti et al.[10]), toxic precipitation in the tubules induced by
146  high dose folic acid injection (FA, Craciun et al.[11]), or genetic alteration by transgenic
147  expression of genetic risk variant APOL1 in podocytes (APOL1 transgenic mice[12]). As
148  reference, we used the mouse kidney specific scRNA-seq data from Park et al.[1]. Details
149  of all datasets are summarized in **Supplementary Table 2**. We systematically
150  benchmarked all methods on artificial bulk experiments performed using the Park et al.
151  scRNA-seq data, finding similar trends as those in **Figure 2a-b** (**Supplementary Figure
152  4a-b**).
153
154  Hierarchical clustering of the cell types in the single cell reference reveals that, apart
155  from neutrophils and podocytes, kidney cells fall into two large groups: Immune cell
156  types (macrophages, fibroblasts, T lymphocytes, B lymphocytes, and natural killer cells)
157  and kidney-specific cell types (proximal tubule, distal convoluted tubule, loop of Henle,
158  two cell types forming the collecting ducts, and endothelial cells). Of these, proximal
159  tubule (PT) is the dominant cell type in kidney, and the proportion of PT cells is known
160  to decrease with CKD progression. MuSiC finds this decrease in all three mouse
161  models (**Figure 3b-d**). Other methods also detect this association for the APOL1 and
162  UUO mouse models, but showed ambiguous results for the FA model.
163
164  Distal convoluted tubule cells (DCT) are known to be the second most numerous cell
165  type in kidney, with an expected proportion of ~10-20%[1]. Yet, CIBERSORT did not
166  detect DCT in any of the three bulk datasets; BSEQ-sc missed it in two datasets and
167  grossly over-estimated its proportion in the third dataset at the cost of a grossly
168  underestimated PT proportion. This is due to the high similarity between DCT and PT,
169  observable in **Figure 3a**. Through its tree-guided recursive algorithm, MuSiC first
170  estimates the combined proportion of kidney cell types versus immune cell types using
171  stable genes for these two large groups, and then zooms in and deconvolves the kidney
172  cell types using genes re-selected for each kidney cell type. This allows MuSiC to
173  successfully separate PT and DCT cells in all three bulk datasets, recovering a
174  consistent DCT proportion between 8-20%, matching expectations. Interestingly, unlike
175  for PT, the proportion of DCT cells show a consistent increase with disease progression
176  across all three mouse models. This may seem counterintuitive given that loss of kidney
177  function is expected to be associated with the loss of kidney cell types. But given the
178  substantial drop of the dominant PT cell type, the proportion of DCT cells relative to the
179  whole may increase, even if its absolute count drops.
180
181  Next, consider immune cells, known to play a central role in the pathogenesis of CKD.
182  MuSiC found the largest immune sub-type to be macrophage, and all methods detected
183  the expected increase of macrophage proportion with disease progression. Apart from
184  this, MuSiC also found fibroblasts, B-, and T-lymphocytes to increase in proportion with
185  disease progression, giving a consistent immune signature that is reproduced across

4

186 mouse models. These findings are consistent with clinical and histological observations,
187 indicating tissue inflammation is a consistent feature of kidney fibrosis. Such
188 reproducible signatures were not found by other methods, which show much less
189 agreement across mouse models.
190
191 Finally, to illustrate MuSiC's cross-species applicability, we used the mouse kidney
192 scRNA-seq reference from Park et al.[1] to deconvolve the bulk rat RNA-seq data from
193 Lee et al.[13], which contains 105 samples obtained from 14 segments spaced along the
194 renal tubule. We mapped samples to their physical locations, and computed correlations
195 between their cell type proportions (**Figure 3e**). Reassuringly, cell types recovered by
196 MuSiC for each segment agree with knowledge about the dominant cell type at its
197 mapped position, e.g. DCT cells come from the DCT region. Correlation between
198 samples is also high within anatomically distinct segments.
199
200 Knowledge of cell type composition in disease relevant tissues is an important step
201 towards the identification of cellular targets in disease. Although most scRNA-seq data
202 do not reflect true cell type proportions in intact tissues, they do provide valuable
203 information on cell type-specific gene expression. Harnessing multi-subject scRNA-seq
204 reference data, MuSiC reliably estimates cell type proportions from bulk RNA-seq. As
205 bulk tissue data are more easily accessible than scRNA-seq, MuSiC allows the
206 utilization of the vast amounts of disease relevant bulk tissue RNA-seq data for
207 elucidating cell type contributions in disease.
208

215 **Author Contributions**
216 This study was conceived of and led by N.R.Z. and M.L. Jointly with N.R.Z. and M.L.,
217 X.W. designed the model and estimation algorithm, implemented the MuSiC software,
218 designed the *in silico* experiments, and led the data analysis. J.P. and K.S. performed
219 the mouse scRNA-seq experiment and provided scientific insight on chronic kidney
220 disease and data interpretation. X.W., N.R.Z. and M.L. wrote the paper with feedback
221 from J.P. and K.S.
222
223 **Competing Financial Interests Statement**
224 The authors declare no competing interests.
225
226
227 **Reference**
228
229 1.    Park, J. et al. Single-cell transcriptomics of the mouse kidney reveals potential
230       cellular targets of kidney disease. *Science*, eaar2131 (2018).

5

2.    Avila Cobos, F., Vandesompele, J., Mestdagh, P. & De Preter, K. Computational deconvolution of transcriptomics data from mixed cell populations. *Bioinformatics* (2018).

3.    Newman, A.M. et al. Robust enumeration of cell subsets from tissue expression profiles. *Nature methods* **12**, 453 (2015).

4.    Baron, M. et al. A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas Reveals Inter- and Intra-cell Population Structure. *Cell Syst* **3**, 346-360 e344 (2016).

5.    Li, B. et al. Comprehensive analyses of tumor immunity: implications for cancer immunotherapy. *Genome biology* **17**, 174 (2016).

6.    Fadista, J. et al. Global genomic and transcriptomic analysis of human pancreatic islets reveals novel genes influencing glucose metabolism. *Proceedings of the National Academy of Sciences* **111**, 13924-13929 (2014).

7.    Cabrera, O. et al. The unique cytoarchitecture of human pancreatic islets has implications for islet cell function. *Proc Natl Acad Sci U S A* **103**, 2334-2339 (2006).

8.    Segerstolpe, Å. et al. Single-cell transcriptome profiling of human pancreatic islets in health and type 2 diabetes. *Cell metabolism* **24**, 593-607 (2016).

9.    Xin, Y. et al. RNA sequencing of single human islet cells reveals type 2 diabetes genes. *Cell metabolism* **24**, 608-615 (2016).

10.   Arvaniti, E. et al. Whole-transcriptome analysis of UUO mouse model of renal fibrosis reveals new molecular players in kidney diseases. *Scientific reports* **6**, 26235 (2016).

11.   Craciun, F.L. et al. RNA sequencing identifies novel translational biomarkers of kidney fibrosis. *Journal of the American Society of Nephrology*, ASN. 2015020225 (2015).

12.   Beckerman, P. et al. Transgenic expression of human APOL1 risk variants in podocytes induces kidney disease in mice. *Nat Med* **23**, 429-438 (2017).

13.   Lee, J.W., Chou, C.-L. & Knepper, M.A. Deep sequencing in microdissected renal tubules identifies nephron segment–specific transcriptomes. *Journal of the American Society of Nephrology*, ASN. 2014111067 (2015).

**Figure Legends**

**Figure 1**: Overview of MuSiC framework.
MuSiC starts from scRNA-seq data from multiple subjects, classified into cell types (shown in different colors), and constructs a hierarchical clustering tree reflecting the similarity between cell types. Based on this tree, the user can determine the stages of recursive estimation and which cell types to group together at each stage. MuSiC then determines the group-stable genes and calculates cross-subject mean (red to blue) and cross-subject variance (black to white) for these genes in each cell type. MuSiC up-weighs genes with low cross-subject variance and down-weighs genes with high cross-subject variance. In the example shown, deconvolution is performed in two stages, only cluster proportions are estimated for the first stage. Constrained by these cluster proportions, the second stage estimates cell type proportions, illustrated by the length of

277  the bar with different colors. The deconvolved cell type proportions can then be
278  compared across disease cohorts.
279
280  **Figure 2**: Pancreatic islet cell type composition in healthy and T2D human samples.
281  **a** and **b** Benchmarking of deconvolution accuracy on bulk data constructed by
282  combining together scRNA-seq samples. **a.** The bulk data is constructed for 10 subjects
283  from Segerstolpe et al. while the single cell reference is taken from the same dataset.
284  The cell type proportions of healthy subjects are estimated by leave-one-out single cell
285  reference. The subject names are relabeled; the table shows average root mean square
286  error (RMSD), mean absolute deviation (mAD), and Pearson correlation (R) across all
287  samples and cell types. **b.** The bulk data is constructed for 18 subjects from Xin et al.
288  while the single cell reference is 6 healthy subjects from Segerstolpe et al.  **c.** Jitter plots
289  of estimated cell type proportions for Fadista et al subjects, color-coded by
290  deconvolution method. Of the 89 subjects from Fadista et al., only the 77 that have
291  recorded HbA1c level are plotted, and T2D subjects are denoted as triangles. **d.** HbA1c
292  vs beta cell type proportions estimated by each of 4 methods. The reported p-values are
293  from single variable regression β cell proportion ~ HbA1c. Multivariable regression
294  results are reported in **Supplementary Table 3**.
295
296  **Figure 3**: Cell type composition in kidney of mouse CKD models and rat.
297  **a.** Cluster dendrogram showing similarity between 13 cell types that were confidently
298  characterized in Park et al. Abbreviations: Neutro: neutrophils, Podo: podocytes, Endo:
299  endothelials, LOH: loop of Henle, DCT: distal convoluted tubule, PT: proximal tubule,
300  CD-PT: collecting duct principal cell, CD-IC: CD intercalated cell, Macro: macrophages,
301  Fib: fibroblasts, NK: natural killers. **b, c and d.** Average estimated proportions for 6 cell
302  types in bulk RNA-seq samples taken from 3 different studies, each study based on a
303  different mouse model for chronic kidney disease.  Results from three different
304  deconvolution methods (MuSiC, BSEQ-sc and CIBERSORT) are shown by different
305  colors. **Supplementary Figure 5a-c** show complete estimation results of all 13 cell
306  types. **b.** Bulk samples are from Beckerman et al., who sequenced 6 control and 4
307  APOL1 mice. **c.** Bulk data are from Craciun et al.[9], where samples are taken before (C)
308  and at 1, 2, 3, 7, 14 days after administering folic acid.  Line plot shows cell type
309  proportion changes over time (days), averaged over 3 replicates at each time point. **d.**
310  Bulk data are from Arvaniti et al.[10], where samples are taken from mice after Sham
311  operation (C), 2 days after UUO operation (D2), and 8 days after UUO operation (D8).
312  The average proportions at each time point are plotted. **e.** MuSiC estimated cell type
313  proportions of rat renal tubule segments.  The estimated cell type proportions (left) and
314  the proportions correlations between samples (right) are shown as heatmap. Segment
315  names are color coded and aligned according to their physical positions along the renal
316  tubule. **Supplementary Figure 6a-c** show NNLS, BSEQ-sc and CIBERSORT results.
317  Segment name abbreviation: S1: S1 proximal tubule; S2: S2 proximal tubule; S3: S3
318  proximal tubule; SDL: Short descending limb; LDLOM: Long descending limb, outer
319  medulla; LDLIM: Long descending limb, inner medulla; tAL: Thin ascending limb; mTAL:
320  Medullary thick ascending limb; cTAL: Cortical thick ascending limb; DCT: Distal
321  convoluted tubule; CNT: connecting tubule; CCD: Cortical collecting duct; OMCD: Outer
322  medullary collecting duct; IMCD: Inner medullar collecting duct.

7

323
324
325
326 **Online Methods**
327
328 <u>MuSiC model set-up</u>
329 In this section, we derive the relationship between gene expression in bulk tissue and
330 cell type-specific gene expression in single cells. This relationship forms the basis of our
331 regression-based deconvolution. For gene $g$, let $X_{jg}$ be the total number of mRNA
332 molecules in subject $j$ of the given tissue, which is composed of $K$ cell types.
333 Then, $X_{jg} = \sum_{k=1}^{K} \sum_{c \in C_j^k} X_{jgc}$, where $X_{jgc}$ is the number of mRNA molecules of gene $g$ in
334 cell $c$ of subject $j$, and $C_j^k$ is the set of cell index for cell type $k$ in subject $j$ with $m_j^k = $
335 $|C_j^k|$ being the total number of cells in this set. The relative abundance of gene $g$ in
336 subject $j$ for cell type $k$ is

$$\theta_{jg}^k = \frac{\sum_{c \in C_j^k} X_{jgc}}{\sum_{c \in C_j^k} \sum_{g'=1}^{G} X_{jg'c}} \ . \tag{1}$$

337 We can show that

$$X_{jg} = \sum_{k=1}^{K} m_j^k S_k^j \theta_{jg}^k = m_j \sum_{k=1}^{K} p_j^k S_j^k \theta_{jg}^k, \tag{2}$$

339 where, for subject $j$, $S_j^k = \frac{\sum_{c \in C_j^k} \sum_{g'=1}^{G} X_{jg'c}}{m_j^k}$ is the average number of total mRNA
340 molecules for cells of cell type $k$ (also referred to as "cell size" below), $m_j = \sum_{k=1}^{K} m_j^k$ is
341 the total number of cells in the bulk tissue, and $p_j^k = \frac{m_j^k}{m_j}$ is the proportion of cells from
342 cell type $k$. Let $Y_{jg} = \frac{X_{jg}}{\sum_{g'=1}^{G} X_{ig'}}$ be the relative abundance of gene $g$ in the bulk tissue of
343 subject $j$. Equation (2) implies
344

$$Y_{jg} \propto \sum_{k=1}^{K} p_j^k S_j^k \theta_{jg}^k. \tag{3}$$

345
346 Thus, across $G$ genes in subject $j$, we have
347

$$\begin{bmatrix} Y_{j1} \\ \vdots \\ Y_{jG} \end{bmatrix} \propto \begin{bmatrix} \theta_{j1}^1 & \cdots & \theta_{j1}^K \\ \vdots & \ddots & \vdots \\ \theta_{jG}^1 & \cdots & \theta_{jG}^K \end{bmatrix} \cdot \begin{bmatrix} S_j^1 & & \\ & \ddots & \\ & & S_j^K \end{bmatrix} \cdot \begin{bmatrix} p_j^1 \\ \vdots \\ p_j^K \end{bmatrix}. \tag{4}$$

348
349 The goal of MuSiC is to estimate $p_j^k$ using data from scRNA-seq and bulk RNA-seq.
350
351 <u>Model assumptions</u>
352 If scRNA-seq data were available for subject $j$, we would be able to obtain the cell size
353 factor $S_j^k$ and cell type-specific relative abundance $\theta_{jg}^k$. With bulk RNA-seq data in

8

354 subject $j$, we get the bulk tissue relative abundance $Y_{jg}$, and, if $\theta_{jg}^k$ and $S_j^k$ were known,

355 we would be able to perform a regression to estimate $p_j^k$. However, since scRNA-seq is

356 still costly, most studies cannot afford the sequencing of a large number of individuals

357 using scRNA-seq. To make deconvolution possible for a broader range of studies, it is

358 desirable to utilize cell type-specific gene expression from other studies or from a

359 smaller set of individuals in the same study. This is feasible under the following two

360 assumptions: (A1) Individuals with scRNA-seq and bulk RNA-seq are from the same

361 population, with their cell-type specific relative abundances $\theta_{jg}^k$ in equation (1) following

362 the same distribution with means $\theta_g^k$ and variances $\sigma_{gk}^2$,

363

$$\theta_{jg}^k \sim F\big(\theta_g^k, \sigma_{gk}^2\big). \tag{5}$$

364

365 Under this assumption, deconvolution can use available single cell data from other

366 subjects or even subjects from other studies as reference for cell type proportion

367 estimation. (A2) The ratio of average cell size $S_k^j$ across cell types are the same

368 regardless of subjects and studies

369

$$\frac{S_j^k}{S_j^{k'}} = \frac{S_{j'}^k}{S_{j'}^{k'}} \quad \text{for all } j, j' \in \{1, \dots, N\} \text{ and } k, k' \in \{1, \dots, K\}. \tag{6}$$

370

371 The second assumption allows us to replace $S_j^k$ by a common value $S^k$ across subjects.

372 In MuSiC, we use the average cell size and relative abundance across all subjects from

373 the scRNA-seq data to estimate $S_j^k$ and $\theta_g^k$.

374

375 <u>Cell type proportion estimation</u>

376 To estimate cell type proportions $\boldsymbol{p}_j = \{p_j^k, k = 1, \dots, K\}$, we need to consider two

377 constraints: (C1) Non-negativity: $p_j^k \geq 0$ for all $j, k$; (C2) Sum-to-one: $\sum_{k=1}^K p_j^k = 1$ for

378 all $j$. Because the bulk tissue and single-cell relationship derived in equation (5) is a

379 "proportional to" relationship, to satisfy the (C2) constraint, we need a normalizing

380 constant $C$ so that

381

$$Y_{jg} = C \cdot \sum_{k=1}^K p_{jk} S_k \theta_{jg}^k + \epsilon_{jg}, \tag{8}$$

382 where $\epsilon_{jg} \sim N(0, \delta_{jg}^2)$ represents bulk tissue RNA-seq gene expression measurement

383 noise. When cell type proportions $\boldsymbol{p}_j = \{p_j^k, k = 1, \dots, K\}$ and subject-specific relative

384 abundances $\boldsymbol{\theta}_{jg} = \{\theta_{jg}^k, k = 1, \dots, K\}$ are known, the variance of bulk tissue gene

385 expression measurement is

386

$$Var\big[Y_{jg} \mid \boldsymbol{p}_j, \boldsymbol{\theta}_{jg}\big] = \delta_{jg}^2. \tag{9}$$

387 Given only cell type proportions, the variance is

9

388

$$
\begin{aligned}
Var[Y_{jg}|\,\boldsymbol{p}_j] &= E[Var[Y_{jg}|\boldsymbol{p}_j,\boldsymbol{\theta}_{jg}]] + Var[E[Y_{jg}|\boldsymbol{p}_j,\boldsymbol{\theta}_{jg}]] \\
&= \delta_{jg}^2 + Var\left[C \cdot \sum_{k=1}^{K} p_{jk}\, S_k\, \theta_{jg}^k\right] \\
&= \delta_{jg}^2 + C^2 \cdot \sum_{k=1}^{K} p_{jk}^2\, S_k^2\, Var[\theta_{jg}^k] = \delta_{jg}^2 + C^2 \sum_{k=1}^{K} p_{jk}^2 S_k^2\, \sigma_{gk}^2 \\
&= \frac{1}{w_{jg}}
\end{aligned}
\tag{10}
$$

389
390  Because of the heteroscedasticity of gene expression over genes, including the weight
391  $w_{jg}$ can improve estimates. Since $\delta_{jg}^2$ is unknown, we will estimate the weight $w_{jg}$
392  iteratively, initialized by NNLS.
393
394  MuSiC is a weighted non-negative least squares regression (W-NNLS), which does not
395  require pre-selected marker genes. Indeed, the iterative estimation procedure
396  automatically imposes more weight on informative genes and less weight on non-
397  informative genes. Because it is a linear regression-based method, genes showing less
398  cross cell type variations will have low leverage, thus having less influence on the
399  regression, whereas the most influential genes are those with high weight and high
400  leverage. To illustrate this point, we also performed benchmarking experiments to show
401  that applying MuSiC using all genes gives more accurate results than applying MuSiC
402  using pre-selected marker genes, thus demonstrating that MuSiC's weighting scheme
403  makes marker gene pre-selection unnecessary (**Supplementary Figure 1c**,
404  **Supplementary Figure 2**).
405
406  Recursive tree-guided deconvolution for closely related cell types
407  Complex solid tissues often include closely related cell types with similar gene
408  expression levels. Correlation in gene expression can lead to collinearity, making it
409  difficult to reliably estimate cell type proportions, especially for less frequent and rare
410  cell types. Although the collinearity problem can be improved by selecting marker genes
411  through support vector regression, as is done in CIBERSORT[3] and BSEQ-sc[4], these
412  approaches still have limited power to resolve similar cell types. In MuSiC, we introduce
413  a recursive tree-guided deconvolution procedure based on a cell type similarity tree,
414  which can be easily obtained through hierarchical clustering. In stage 1 of this
415  procedure, cell types in the design matrix are divided into high-level clusters by
416  hierarchical clustering with closely related cell types clustered together. Proportion for
417  these cell type clusters are estimated using genes with small intra-cluster variance
418  (cluster-stable genes) using the above described W-NNLS. In stage 2, for cell types in
419  each cluster, the cell type proportions are estimated using W-NNLS with genes
420  displaying small intra-cell type variance, subject to the constraint on the pre-estimated
421  cluster proportions.  If necessary, more than 2 stages of recursion can be applied, with
422  each stage separating the cell types within each large cluster into finer clusters, and
423  using cluster-stable genes to do W-NNLS subject to the constraint that fixes higher-level
424  cluster proportions.

10

425

426 To illustrate this recursive tree-guided deconvolution procedure, we start with a simple
427 case with four cell types and $G$ genes. Let $X_1, X_2, X_3, X_4$ represent cell type-specific
428 expression in the design matrix, obtained from scRNA-seq, and let $Y$ be the gene
429 expression vector in the bulk RNA-seq data. The relationship of bulk and single-cell
430 data can be written as

431

$$\begin{pmatrix} Y^{(1)} \\ Y^{(2)} \end{pmatrix} = \begin{pmatrix} X_1^{(1)} & X_2^{(1)} & X_3^{(1)} & X_4^{(1)} \\ X_1^{(2)} & X_2^{(2)} & X_3^{(2)} & X_4^{(2)} \end{pmatrix} \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix} + \begin{pmatrix} \epsilon^{(1)} \\ \epsilon^{(2)} \end{pmatrix}, \qquad (11)$$

432

433 where the superscripts (1) and (2) indicate two sets of genes. Suppose the four cell
434 types are grouped into two clusters, $(X_1, X_2)$ and $(X_3, X_4)$. The first set of genes are those
435 showing small intra-cluster variance in gene expression, that is, $X_1^{(1)} \approx X_2^{(1)}$ and $X_3^{(1)} \approx$
436 $X_4^{(1)}$, whereas the second set of genes are the remaining genes.

437

438 *Stage 1*: Estimate cluster proportions $\pi_1 = p_1 + p_2$ and $\pi_2 = p_3 + p_4$,
439

$$Y^{(1)} = X_1^{(1)} \pi_1 + X_3^{(1)} \pi_2 + \epsilon^{(1)}. \qquad (12)$$

440 The cluster proportions, $\hat{\pi}_1$ and $\hat{\pi}_2$, are estimated by W-NNLS using intra-cluster
441 homogenous genes.

442

443 *Stage 2*: Estimate cell type proportions $(p_1, p_2, p_3, p_4)$,
444

$$Y^{(2)} = X_1^{(2)} p_1 + X_2^{(2)} p_2 + X_3^{(2)} p_3 + X_4^{(2)} p_4 + \epsilon^{(2)}. \qquad (13)$$

445

446 The cell type proportions are estimated by W-NNLS using the remaining genes subject
447 to the constraint that

448

$$\hat{p}_1 + \hat{p}_2 = \hat{\pi}_1, \text{ and } \hat{p}_3 + \hat{p}_4 = \hat{\pi}_2. \qquad (14)$$

449

450 Construction of benchmark datasets and evaluation metrics
451 To evaluate MuSiC and compare with other deconvolution methods, we need bulk RNA-
452 seq data with known cell type proportions. Therefore, we construct artificial bulk tissue
453 data from a scRNA-seq dataset in which the bulk data is obtained by summing up gene
454 counts from all cells in the same subject. Relative abundance is calculated by equation
455 (1). The true cell type proportions in the artificial bulk data can be directly obtained from
456 the scRNA-seq data and this allows us to use this artificially constructed bulk data as a
457 benchmark dataset to evaluate the performance of different deconvolution methods.
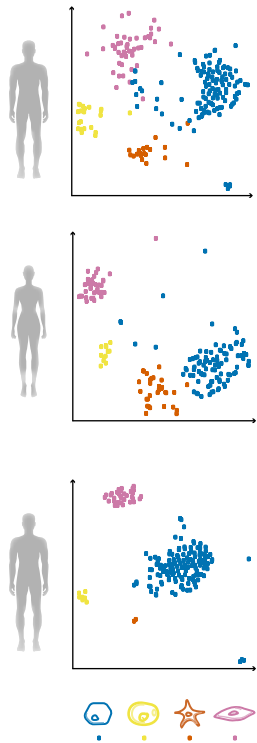458 Denote the true cell type proportions by $\boldsymbol{p}$ and the estimated proportions by $\hat{\boldsymbol{p}}$.
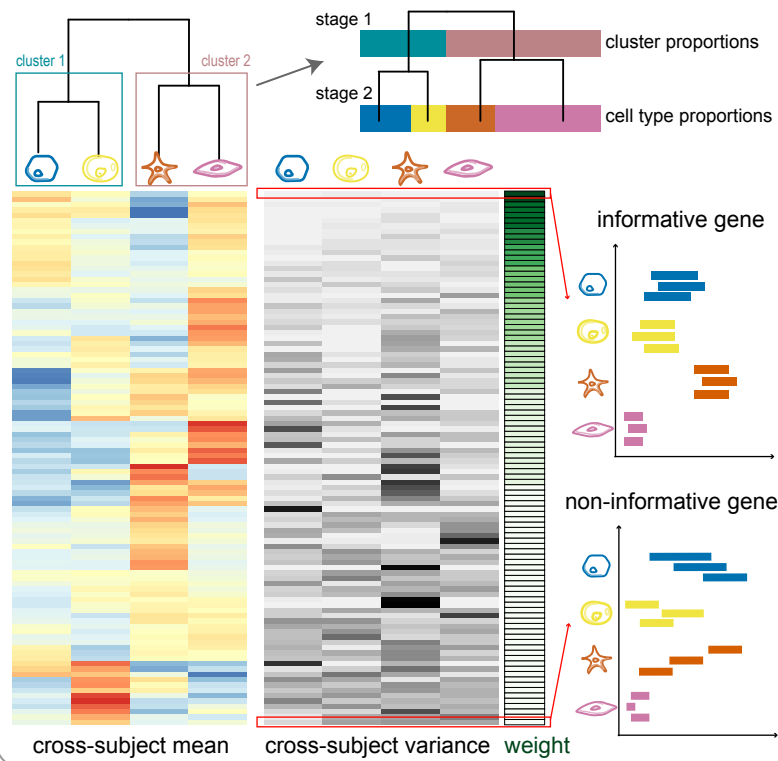459 Deconvolution methods are evaluated by the following metrics.
460     (i)        Pearson correlation, $R = Cor(\boldsymbol{p}, \hat{\boldsymbol{p}})$.
461     (ii)      Root mean squared deviation, RMSD $= \sqrt{avg(\boldsymbol{p} - \hat{\boldsymbol{p}})^2}$;

462      (iii)      Mean absolute deviation, mAD = $avg(|\boldsymbol{p} - \widehat{\boldsymbol{p}}|)$.
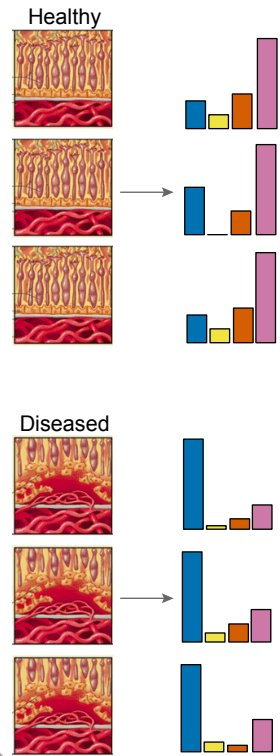
| Multi-subject scRNA-seq | Cell type specific gene expression reference from scRNA-seq | Bulk tissue deconvolution |

stage 1 — cluster proportions

stage 2 — cell type proportions

cluster 1   cluster 2

informative gene

non-informative gene

cross-subject mean   cross-subject variance   weight

Healthy

Diseased

a

| Method | MuSiC | NNLS | BSEQ-sc | CIBERSORT |
|--------|-------|------|---------|-----------|
| RMSD | 0.040 | 0.098 | 0.099 | 0.085 |
| mAD | 0.029 | 0.064 | 0.068 | 0.061 |
| R | 0.97 | 0.85 | 0.86 | 0.89 |

b

| Method | MuSiC | NNLS | BSEQ-sc | CIBERSORT |
|--------|-------|------|---------|-----------|
| RMSD | 0.10 | 0.17 | 0.21 | 0.21 |
| mAD | 0.06 | 0.12 | 0.15 | 0.15 |
| R | 0.94 | 0.82 | 0.79 | 0.76 |

**a** Cluster log(Design Matrix)

**e** MuSiC estimated cell type proportions and correlations