1  **ICTD: A semi-supervised cell type identification and deconvolution method for multi-omics data**

2  Wennan Chang[1, 2+], Changlin Wan[1, 2+], Xiaoyu Lu[1], Szu-wei Tu[1], Yifan Sun[1], Xinna Zhang[1], Yong Zang[3], Anru
3  Zhang[4], Kun Huang[5], Yunlong Liu[1], Xiongbin Lu[1*], Sha Cao[4*], Chi Zhang[1, 2*]

4  [1]Department of Medical and Molecular Genetics and Center for Computational Biology and Bioinformatics,
5  [4]Department of Medicine, [3]Department of Biostatistics, [5]Department of Medicine, Indiana University School of
6  Medicine, Indianapolis, IN,46202, USA.

7  [2]Department of Electrical and Computer Engineering, Purdue University, Indianapolis, IN, 46202, USA

8  [4]Department of Statistics, University of Wisconsin–Madison, Madison, WI 53706, USA

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

11  +These authors have equal contribution to this work.

12  **Abstract**

13  We developed a novel deconvolution method, namely **I**nference of **C**ell **T**ypes and **D**econvolution (ICTD) that
14  addresses the fundamental issue of identifiability and robustness in current tissue data deconvolution problem.
15  ICTD provides substantially new capabilities for omics data based characterization of a tissue microenvironment,
16  including (1) maximizing the resolution in identifying resident cell and sub types that truly exists in a tissue, (2)
17  identifying the most reliable marker genes for each cell type, which are tissue and data set specific, (3) handling
18  the stability problem with co-linear cell types, (4) co-deconvoluting with available matched multi-omics data, and
19  (5) inferring functional variations specific to one or several cell types. ICTD is empowered by (i) rigorously derived
20  mathematical conditions of identifiable cell type and cell type specific functions in tissue transcriptomics data and
21  (ii) a semi supervised approach to maximize the knowledge transfer of cell type and functional marker genes
22  identified in single cell or bulk cell data in the analysis of tissue data, and (iii) a novel unsupervised approach to
23  minimize the bias brought by training data. Application of ICTD on real and single cell simulated tissue data
24  validated that the method has consistently good performance for tissue data coming from different species, tissue
25  microenvironments, and experimental platforms. Other than the new capabilities, ICTD outperformed other state-
26  of-the-art devolution methods on prediction accuracy, the resolution of identifiable cell, detection of unknown sub
27  cell types, and assessment of cell type specific functions. The premise of ICTD also lies in characterizing cell-
28  cell interactions and discovering cell types and prognostic markers that are predictive of clinical outcomes.

29  **Introduction**

30  Tissue deconvolution aims to disentangle the cell composition in terms of their relative quantities, based on
31  which, the cell type specific functions and their cross-talks in the tissue microenvironment could be studied [1][2][3]
32  [4] . Existing deconvolution algorithms usually assume the observed expression matrix as a product of a cell type
33  signature matrix S and proportion matrix P [2][3][4] . Independent training data is usually needed to impose prior on
34  S via certain information transfer [2][5][6][7] . The recent emergence of single cell RNA-seq (scRNA-seq) allows
35  researchers to uncover new biological traits in cell populations of bulk tissue[8]. Regardless, the knowledge
36  transfer from training single/bulk cell data to target bulk tissue should be carefully handled, as the gene
37  expression distribution of the two domains could be highly variable, which tend to be oversimplified in current
38  deconvolution methods [9] . Novel or rare cell subtypes are of great interest to researchers [10]. However, current
39  deconvolution algorithms usually assume a fixed pool of cell types, which clearly is incapable of identifying novel
40  sub cell types [2][3][4] . Moreover, certain cell types such as immune cells tend to co-infiltrate in a real tissue,
41  suggesting that the proportions of these cell populations are highly co-linear [11] . As a result, estimating
42  proportions with plain linear regression model or non-negative factorization would suffer from multi-collinearity,
43  leading to highly unstable predictions [12][13] . Recent methods such as Cell Population Mapping (CPM) and
44  CIBERSORTx have been developed to predict cell type specific functions [14][9] . However, they rely on precisely
45  predicted cell proportions, and matched scRNA-seq profiles of similar tissues, which limited their applications to
46  a wider extent. It is also noteworthy that none of the existing deconvolution methods is designed to handle highly
47  varied tissue microenvironments or multi-omics data. Here, we summarize the key challenges of deconvolution
48  methods as (i) detect the resident (sub) cell types and their true marker genes dependent on the tissue [15] (ii)
49  handle systematic expression variations from training to target data domain; (iii) deal with the prevalent co-

50  linearity in the cell type specific expression signatures and cell proportions; (iv) define expression patterns that
51  represent varied cell type specific functions; (v) enable application to a variety of tissue microenvironment and
52  multi-omics data types. More detailed discussions and comparisons of the formulations of existing methods are
53  provided in the **Supplementary Notes**.

54  Based on a preliminary evaluation of the variations of known cell type signature genes in a large set of single
55  and bulk cell data, we first derived mathematical conditions for a cell type to be "identifiable" in a tissue omics
56  data: (1) the cell type has uniquely expressed genes, the expression values of which over any subset of samples
57  form a rank-1 matrix (a matrix with matrix rank equals to one), or (2) there are genes expressed by the cell type
58  and  other cell types satisfy (1), and the expression values contributed by the cell type over any subset of samples
59  form a rank-1 matrix. And a cell type-specific function is "identifiable" if there are marker genes of the function
60  forming a rank-1 submatrix in a subset of samples with significant presence of the cell type. These "identifiability"
61  conditions grant the potential to detect novel cell subtypes or cell functions via the detection of low rank matrices.
62  Detailed mathematical considerations and derivations were given in **Online Methods and Supplementary**
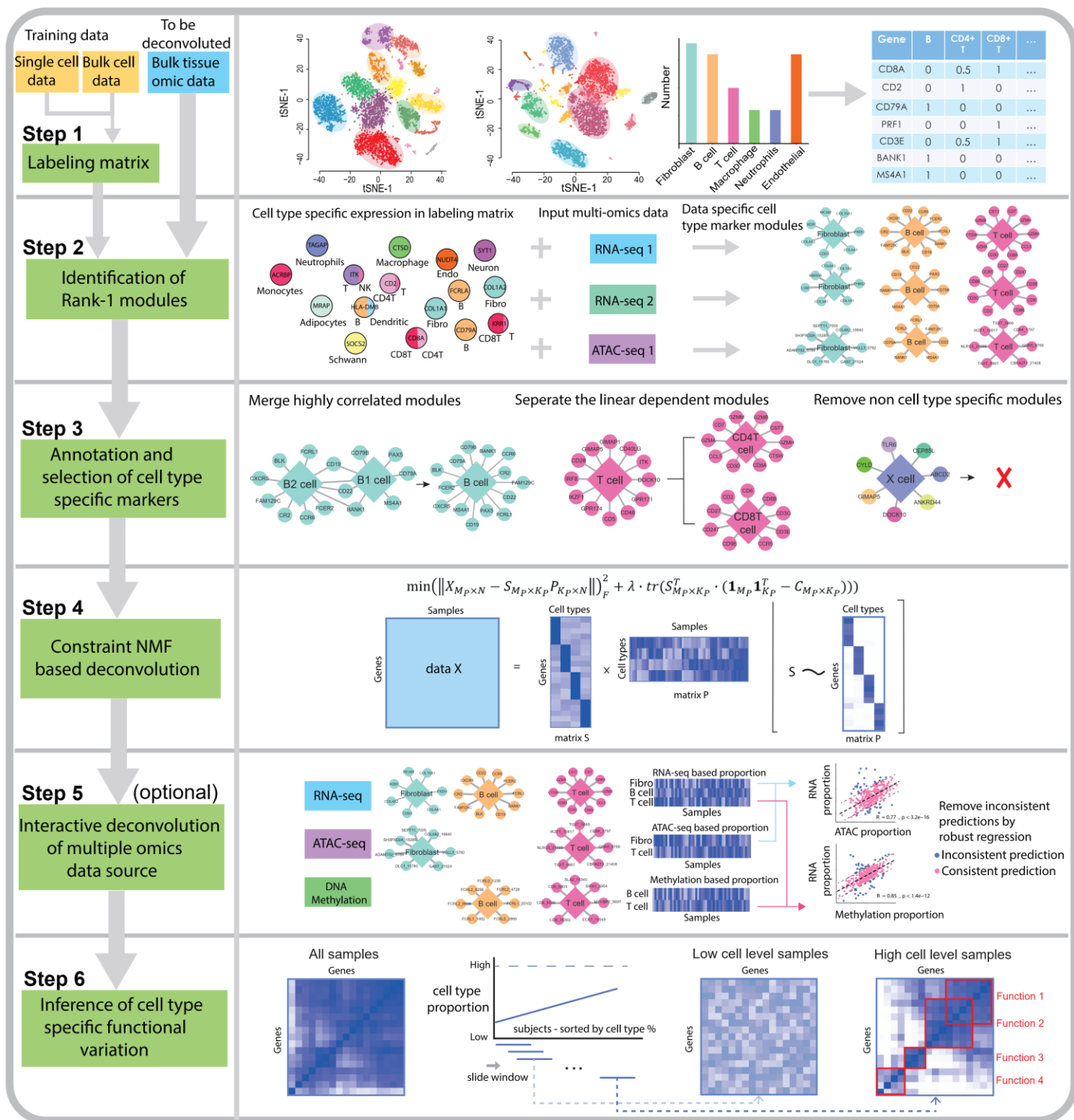63  **Notes**.

64  Based on the rigorously derived mathematical conditions, we developed a semi-supervised method, namely
65  inference of cell types and deconvolution (ICTD), featured by: (1) a semi-supervised detection of "identifiable"
66  cell types and marker genes specific to each omics dataset and tissue micro-environment; (2) a novel
67  nonparametric detection and annotation of cell type signature genes, which is used as information basis to
68  annotate the identified cell types; (3) a novel constrained non-negative matrix factorization (NMF) method to
69  decrease the bias caused by knowledge transfer from training data, as well as to effectively handle the co-
70  occurring cells; (4) a robust regression based approach to interactively deconvolute multi-omics data of matched
71  samples, and (5) a local-low-rank screening approach to identify cell type specific functions, which altogether
72  offers a systematic solution of the five key challenges.

73  **Results**

74  Our core algorithm ICTD consists of six steps (**Fig 1**): <u>(1) Compute the relative specificity of all genes for all cell</u>
75  <u>types in a given microenvironment.</u> A labeling matrix $L_{M \times K}$ of M genes and K selected cell types is first
76  constructed based on training single or bulk cell transcriptomics data, where $L_{i,j} = \frac{1}{l}, l = 1, ..., K - 1$, if gene $i$ is
77  significantly expressed in in cell type $j$ and its expression is significantly lower than in $l - 1$ other cell types, and
78  $L_{i,j} = 0$ otherwise (Supplementary Table S1). Without loss of generality, we assume that all the M genes are
79  specific to one or a few cell types, namely,  $\Sigma_j L_{i,j} > 0, \forall\, i = 1, ..., M$. (2) <u>Detect all gene modules within which the</u>
80  <u>gene expression vectors are linearly dependent and form rank-1 matrix, and the modules present evidence of</u>
81  <u>"identifiable" cell types.</u> For each gene module detected on the target tissue expression matrix among the M
82  genes, if its member genes are all highly expressed in one or several cell types according to the labeling matrix,
83  the module will be considered as evidence of potential cell type(s) by ICTD. In this step, ICTD can exclude
84  undesired cell types, such as cancer or other disease cells, from further analysis, by a non-negative projection
85  of the input data to the complementary of the space spanned by the marker genes of undesired cell types. (3)
86  <u>Infer the "identifiable" cell types and signature genes.</u> Non-negative linear dependency among the selected
87  modules is evaluated and each module is annotated by the genes' significant enrichment to a cell type based on
88  the labeling matrix. Modules are merged with high inter-dependency, and further filtered such that modules
89  enriching none of the cell types are removed. The total number of "identifiable" cell types is computed as the
90  total rank of the expression matrix composed by all genes in the remaining rank-1 modules, the genes in each
91  module will be considered as markers of the corresponding cell type. <u>(4) Predict cell proportions using</u>
92  <u>constrained NMF.</u> With the "identifiable" cell types and their marker genes, a constraint matrix $C_{M \times K}$ can be
93  constructed. Specifically, for cell type $k, k = 1..K$ with $M_k$ marker genes, $C_{(\sum_{k=1,..,K} M_k) \times K}[i, j] = 1$, if gene $i$ is
94  marker of the cell type $j$, and 0 otherwise. The constraint matrix is then enforced upon the regular NMF
95  formulation to guarantee similarity of the signature matrix with the constraint matrix, namely, we solve
96  $\min_{S,P} \left( \|X - S \cdot P\|_F^2 + \lambda \cdot \text{trace}\left(S^T \cdot (\mathbf{1}_M \mathbf{1}_k^T - C)\right) \right)$, where $\mathbf{1}_d$ denotes an all-1 column vector of length $d$. (5) Co-
97  deconvolution of matched multi-omics data. The semi-supervised property of ICTD enable its application to multi-
98  omics data. A robust regression approach is further applied to identify the cell types and samples, in which the
99  cell proportions inferred from different omics-data are highly consistent. <u>(6) Estimate cell type specific functions.</u>
100  For each cell type detected, ICTD screens the rank of the expression matrix containing a group of samples which

101 are stratified by their cell abundance levels, and pins down marker genes of a varied cell type specific function if
102 they form at least one distinct dimension.

103



104

**Figure 1. Analysis pipeline of ICTD.** ICTD first constructs labeling matrix to store genes' relative specificity to different cell types using bulk or single cell training data (Step 1). Rank-1 modules were detected among the cell type marker genes in each input omics dataset (Step 2). Similar modules were merged, modules that do not (non-negatively) depend on other modules are kept, and modules that do not overrepresent any cell type markers are removed. The number of cell types of the target deconvolution is determined as the total rank of the expression matrix of genes in the remaining modules (Step 3). A constrained NMF is conducted to regularize the signature matrix $S$, such that values in $S$ are shrunken towards 0 if the corresponding entries in the constraint matrix is 0. (Step 4). If matched multi-omics data are available, robust

112 regression among cell proportions inferred from different omics data set is performed to remove outlier samples (Step 5,
113 optional). Marker genes of cell type specific functions are further identified by looking for local low rank submatrices in
114 sample groups stratified by different level of the cell proportion (Step 6).

115 The core algorithms for each step are described in the **Online Methods**. Detailed algorithms, data used for
116 method validation, and model comparisons with other methods, are provided in the **Supplementary Notes and**
117 **Methods**. Below we present the application of ICTD on simulated bulk data using single cell RNA-seq data (**Fig.**
118 **2**) and real tissue data (**Fig. 3**). We demonstrated (1) the ability of ICTD to identify both known and novel (sub)
119 cell types with high accuracy, (2) the overall competitive performance of ICTD in analyzing data of different tissue
120 microenvironment and experimental platforms, (3) the robustness of ICTD in cases where cell types have highly
121 co-linear proportions, (4) ICTD's capability in interactive deconvolution of matched multi-omics data, (5) inference
122 of cell type specific functions, and (6) explorative findings derived by correlating ICTD predicted cell and
123 functional levels with other omics, imaging and clinical data.

124 ### *Validation on single cell simulated bulk tissue data*

125 We benchmarked ICTD on predicting the types of resident cells and their relative proportions against three state-
126 of-art deconvolution methods, namely CIBERSORT, TIMER, and EPIC (**Online Methods**), using single cell
127 simulated bulk tissue datasets. The bulk tissue datasets were simulated by RNA-seq data of single cells or single
128 nucleus from different tissue microenvironments, including five from human solid cancer (namely, breast, colon,
129 head and neck, lung, and melanoma), five from human central nervous system (namely glioblastoma,
130 oligodendroglioma, astrocytoma and two normal brain), three from human immune system (monocyte and
131 dendritic cell, lymphoid, and myeloid progenitor cells), and one from mouse melanoma. On all five human solid
132 cancer microenvironment, all mixing cell types were detected as "identifiable" by ICTD. In addition, ICTD
133 achieved significantly higher accuracy in predicting total B-, T-, mast, fibroblast, endothelial cells and
134 macrophage proportions comparing to other methods. On 23 out of the 25 cells type in the simulated bulk cancer
135 datasets, ICTD predicted relative proportions achieved higher than 0.95 Pearson correlation coefficient (PCC)
136 with true proportions, while the average PCC is 0.86, 0.63 and 0.52 for EPIC, TIMER and CIBERSORT,
137 respectively (**Fig 2a**). On the five human brain microenvironments, ICTD successfully detected astrocyte,
138 oligodendrocyte and progenitors, exhibitory and inhibitory neuron, microglial and Schwann cells as identifiable
139 cell types, all with at least 0.9 PCC with true proportions (**Fig 2b**). Similarly, ICTD also accurately identified sub
140 ell types from the mixture of multiple classes of monocyte and dendritic cells, human lymphoid and myeloid
141 progenitors, and the immune and stromal cells in mouse melanoma microenvironment, with reliable prediction
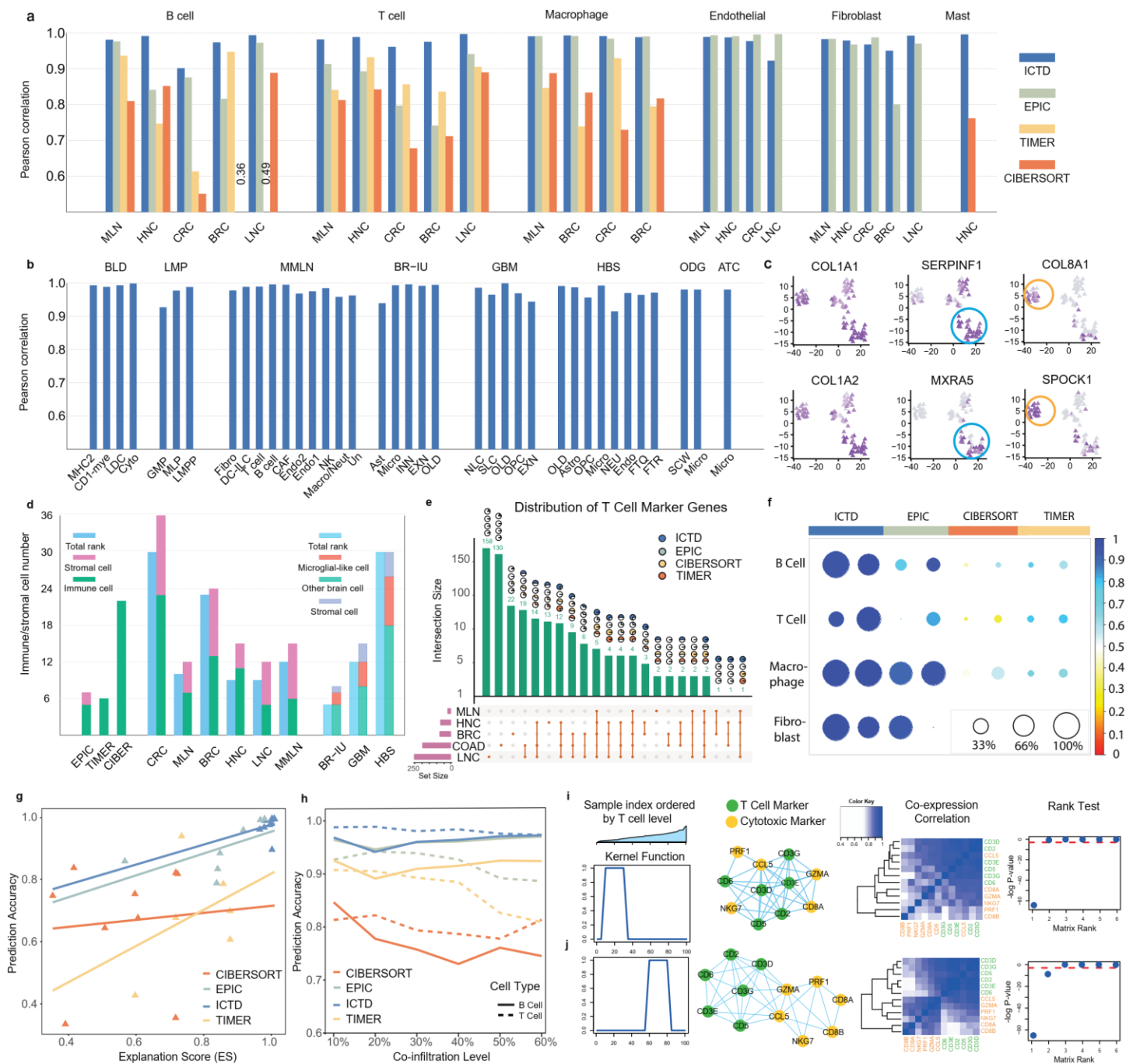142 of proportions (**Fig 2b**).

143 *Novel cell types.* A unique feature of ICTD is its capability to automatically detect cell and sub cell types along
144 with cell marker genes for effective cell (sub)type annotation. Our analysis on simulated cancer tissue data
145 suggested that each of the rank-1 module corresponds to one cell or sub-cell type (Supplementary Fig S1). On
146 the simulated human solid cancer datasets, ICTD was able to identify subtypes of immune/stromal cells,
147 including CD4+ and CD8+ T cells, novel subtypes of fibroblast and myeloid cells. The sub cell type markers were
148 further validated by the tSNE visualization, where the expression level of each marker set turns out to be specific
149 to the associated cell types or subtypes (Supplementary Fig S2). As illustrated in **Fig 2c**, among the three
150 fibroblast rank-1 modules identified by ICTD, one clearly corresponds to the general fibroblast (with COL1A1
151 expression) type, and the other two correspond to two fibroblast subtypes (with COL8A1 or SERPINF1
152 expression) in the simulated human melanoma data. We confirmed all the rank-1 modules identified by ICTD
153 from all the single cell simulated tissue data are specifically expressed in only one cell type, suggesting the high
154 specificity of ICTD in identifying true cell types.

155 *Variability of cell types and their marker genes.* It is noteworthy that the number of identifiable cell types could
156 vary through disease contexts and data sets. Comparing to the fixed cell types assumed in most of the
157 deconvolution methods, the number of cell types identified by ICTD highly matches the number of mixing cell
158 types in each single cell simulated tissue data set (**Fig 2d**). We further investigated the level of variation for cell
159 type markers through different disease contexts and data set. As shown in **Fig 2e**, there is a strong disease
160 context specificity of T cell markers: only four T cell markers were shared by all the five cancer data sets, and
161 19 T cell markers were shared in four out of the five data sets. We observed on average 93.75%, 90.36% and
162 83.33% of the T cell markers utilized in CIBERSORT, TIMER and EPIC are specific to only three or less cancer
163 types and only 65.21%, 69.57% and 13.04% of the common cell type marker genes were included in their
164 signature matrix. Similar patterns are also seen for B and fibroblast cells (Supplementary Fig S3). In contrast,

165 ICTD considers the variations in both the identifiable cell types and cell type markers in different tissues and
166 datasets, resulting in a better prediction accuracy throughout different scenarios (**Fig 2e-f**). An explanation score
167 (ES) is defined for each marker gene to evaluate the goodness of fitting of the gene's expression by the predicted
168 proportions of the cell types expressing the gene (**Online Methods**). High ES scores of the marker genes for
169 one cell type is a necessary condition for the high prediction accuracy and specificity of the marker genes. We
170 observed strong positive correlations between the ES scores and prediction accuracy using ICTD and EPIC, as
171 these two methods rely on cell type uniquely expressed genes. Similarly, for CIBERSORT and TIMER, positive
172 associations were also observed (**Fig 2g**). Analysis of six major immune and stromal cell types in five simulated
173 bulk cancer data sets suggested that in general, when ES is below 0.8, the prediction accuracy is lower than 0.8;
174 on the other hand, when ES is above 0.9, the prediction accuracy tends to higher than 0.9 (**Fig 2g**). We observed
175 the ES of all the cell type specific markers identified by ICTD on the simulated cancer tissue data are all above
176 0.95. It is noteworthy ES can partially evaluate the performance of a deconvolution method without knowing true
177 cell proportions.



178

179 **Figure 2. Validation of ICTD by using single cell simulated bulk tissue data.** (a) PCC between true and predicted proportion
180 of six cell types by ICTD, EPIC, TIMER and CIBERSORT, in the bulk tissue data simulated using scRNA-seq data collected
181 from Melanoma (MLN), Head and Neck Cancer (HNC), Colorectal Cancer (CRC), Breast Cancer (BRC), and Lung Cancer (LNC).
182 (b) PCC between true and predicted proportion of cell types and subtypes identified by ICTD in the bulk tissue data
183 simulated by scRNA-seq data of myeloid and dendritic cell mixture (BLD), lymphoid and myeloid progenitor mixture (LMP),
184 mouse melanoma (MMLN), normal brain cells nucleic sequencing generated in this study (BR-IU), glioblastoma (GBM),
185 human normal brain (HBS), oligodendroglioma (ODG), and astrocytoma (ATC). Detailed cell type codes are given in
186 **Supplementary Note**. (c) *t*-SNE plot of the marker genes of fibroblast subtypes in MLN scRNA-seq data, which were
187 identified by ICTD from simulated human melanoma tissue data. In each panel, darker color denotes higher expression of
188 the gene in a cell. (d) Consistency of the number of ICTD identified cell types and the matrix rank of the expression profile
189 of the marker genes of identified cell types, i.e. the number of identifiable cell types, in each simulated tissue data. (e)
190 Distribution of the true T cell marker genes identified in the five cancer data and their overlap with the actually used T cell
191 signature genes in CIBERSORT, TIMER and EPIC. Each bar and number represent the number of genes specifically expressed
192 by T cells in each of five cancer types, which is labeled in the dot plot on the bottom. The pie charts illustrate the proportion
193 of the T cell marker genes used by ICTD (data adaptive) and CIBERSORT, TIMER and EPIC (held fixed). (f) Re-evaluation of
194 robustness of cell type specific markers used by each method. The circle size represents the ratio of true marker genes
195 among all genes used as marker genes for each cell type (row) for each method (column).  The color represents the E-
196 score level. The two columns of each method show the results of simulated MLN (left) and HNC (right) tissue data. The
197 plots of the other three cancer types were shown in Supplementary Fig S6. (g) Dependency between explanation score (x-
198 axis) and prediction accuracy (y-axis) of the cell type proportions given by the four methods. (h) PCC (y-axis) between true
199 and predicted T and B cell proportions on simulated data with different level of T and B cell co-infiltration (x-axis). (i-j)
200 prediction of varied T cell cytotoxicity level in simulated HNC data. From left to right, the four plots illustrate the kernel
201 function used for local low rank screening, co-expression network of T cell and cytotoxic marker genes, heatmap of
202 correlations between T cell and cytotoxic marker genes, and p values of the expression matrix rank of the T cell and
203 cytotoxic marker genes, in the samples of low T cell infiltration (i) and high T cell infiltration level (j).
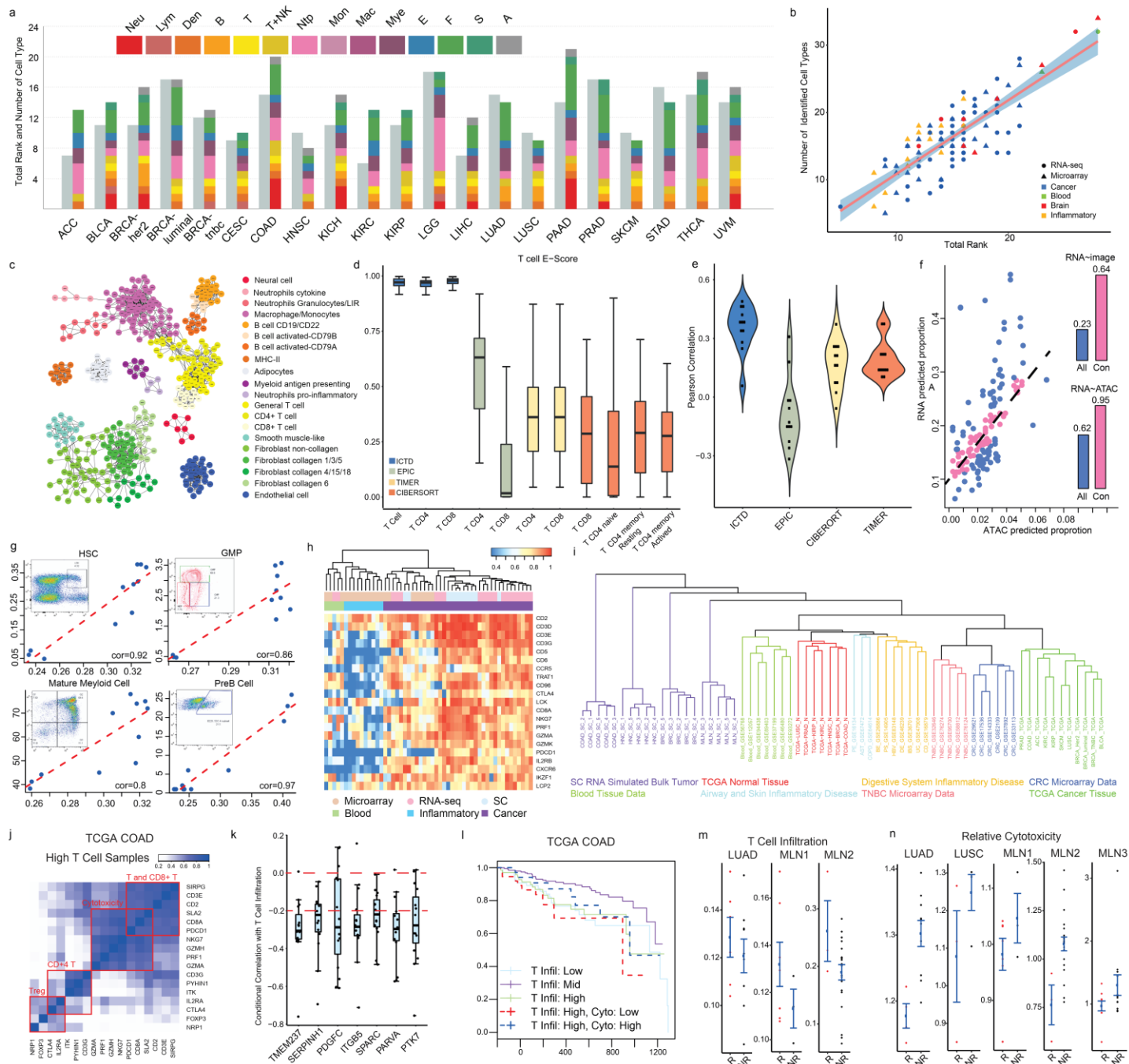
204 *Cell type co-linearity.* ICTD also demonstrated its superiority in handling co-linearity of cell proportions, caused
205 by cells' functional dependencies. Our preliminary analysis on TCGA data suggested correlation among the
206 immune and stromal cells to be as high as 0.94 (**Supplementary Notes**). We simulated batches of bulk tissue
207 samples in each of which the cell proportions are intentionally set to have different levels of correlations to mimic
208 the dependencies of different cell types in cancer microenvironment (**Online Methods**). Not surprisingly, while
209 performance of regression based methods dropped significantly when co-linearity level was high, ICTD achieved
210 high robustness and prediction accuracy at different levels of co-linearity. This owes to the data-adaptive
211 selection of cell type specific markers and constrained NMF formulation adopted by ICTD. The four methods'
212 prediction accuracy of B and T cells across different co-linearity levels in simulated human melanoma tissue data
213 is shown in **Fig 2h**. In addition, significant correlations among ES, prediction accuracy, and co-linearity of cell
214 proportions were identified (Supplementary Fig S4).

215 *Cell type specific functions.* ICTD can identify varied function of a certain cell type using a local low rank
216 identification approach [16] . In the human head and neck cancer data, we identified the expression level of
217 cytotoxic gene in the CD8+ T cells vary considerably in patient stratifications of different T cell abundances,
218 suggesting mixed T cell exhaustion levels (**Supplementary Notes**). To evaluate the capability of ICTD in
219 identifying varied T cell cytotoxicity level, we simulated bulk tissue data with different proportion and cytotoxicity
220 level of T cells (**Online Methods**). ICTD conducted a local low rank screening with a kernel function along
221 samples ordered by predicted T cell proportions. Our analysis clearly identified the linear space spanned by the
222 T cell and cytotoxicity marker genes switches from rank-1 to rank-2 throughout the samples with low to high T
223 cell levels, suggesting the identifiability of the varied cytotoxic level in the samples of high T cell infiltrations (**Fig
224 2i-j**). On average, correlation level of 0.86 between the true cytotoxicity level per unit T cell and the prediction
225 made by ICTD was observed (Supplementary Fig S5).

226 ### *Implications from real tissue data*

227 We then applied ICTD on a collection of human cancer, normal, blood and inflammatory tissue (CNBI) data,
228 including 28 cancer and 11 normal tissue types from TCGA, 17 colorectal cancer, 7 triple negative breast cancer,
229 7 blood tissue, and 11 human inflammatory disease data sets from GEO (Supplementary Table S3). We
230 identified rank-1 markers of B, T, dendritic, general myeloid, macrophage, monocytes, neutrophil, fibroblast,

endothelial and adipocyte cell and their sub cell types in each dataset (**Fig 3a** and Supplementary Fig S7). A strong association between the number of identified cell types and the total rank of the matrix of marker genes was observed (**Fig 3b**). It is noteworthy that the types of resident cells most variable across different cancer types are the subtypes of adipocytes, fibroblast, and myeloid cells, which seem to be most prevalent in breast, colorectal, lung, pancreatic and stomach cancers, commonly known to have considerable stromal components. The complete set of cell types and their marker genes identified in each data set were summarized in Supplementary Table S4. In the TCGA datasets, 21 commonly "identifiable" cell and subtype types have been observed in more than 10 cancer types, including CD19/CD22 expressing regulatory-like B cell and CD79A/CD79B expressing activated B cell; total, CD8+, and CD4+ T cell; Neurexin and Caytaxin expressing Neuron cell; myofibroblast-like cell; Collagen 1/3/5, Collagen 4/15/18, Collagen 6, and Non-collagen expressing Fibroblast; Endothelial cell; MHC class II antigen presenting cell; MHC class I, pro-inflammatory cytokine releasing, chemokine and cytokine releasing Myeloid cells; complement pathway activated Macrophage and Monocytes; granulocytes; and adipocytes (**Fig 3c**).

245 **Figure 3. Application of ICTD on real bulk tissue transcriptomics data.** (a) The number of identifiable cell types (colored
246 bar) and the matrix rank (grey bar) of their marker genes identified by ICTD through all the TCGA cancer data; (b) Scatter
247 plot of the number of identifiable cell types (y-axis) and the matrix rank (x-axis) of their marker genes identified by ICTD
248 through all the analyzed data. (c) Network of the marker genes of the commonly identified cell (sub) types in the TCGA
249 data. An edge between two genes means the two genes are both identified as markers of one cell type in more than 10
250 analyzed TCGA data. (d) E-score of the T cell marker genes identified by ICTD and those used by EPIC, TIMER and
251 CIBERSORT in TCGA data. E-score of other cell types are given in Supplementary Fig S11. (e) Correlation (y-axis) between
252 the imaging data derived tumor infiltrated lymphocyte level and T cell proportion predicted by the four methods (x-axis)
253 in 11 TCGA cancer. (f) Scatter plot of the T cell proportions predicted by TCGA BRCA RNA-seq and ATAC-seq data. Samples
254 with highest consistency identified by the robust regression were pink colored. The bar plots represent the correlations
255 of the proportions inferred by the RNA-Seq vs ATAC-Seq (or RNA-Seq vs imaging) in all the samples (PCC=0.62, or 0.23)
256 and the most consistent samples (PCC=0.96, or 0.64). (g) Consistency of ICTD predicted (x-axis) and FACS measured (y-
257 axis) cell proportions of four hematopoietic cell types. (h) Evaluation of T cell markers identified in CNBI data. In the
258 heatmap, each row is the commonly identified T cell markers and each column is one data set. Color in the heatmap
259 represents the E-score of each gene in each data set. Statistics of other cell types are given in Supplementary Fig S12. (i)
260 Clustering of datasets from different microenvironment from different platforms based on a distance measure of the
261 marker gene expression profiles of identifiable cell types (see **Online Methods)**. This is to show the relative impact of
262 technological platforms and tissue microenvironment on the variability of gene markers expressions. (j) Co-expression
263 between T cell, CD8+ T cell, cytotoxic function, CD4+ T cell and T-reg marker genes in the samples with high T cell
264 infiltration in TCGA COAD data. (k) Correlation between fibroblast cell expressing genes and T cell infiltration level
265 conditional on the fibroblast cell level in 15 cancer types. (l) Survival curves of the TCGA COAD patients with low, medium
266 and high T cell infiltration, and the high T cell infiltration patients with low and high cytotoxicity functions predicted by
267 ICTD. (m) Variation of T cell infiltration level in response (R) and non-response (NR) patients in three independent
268 checkpoint inhibitor treated clinical data. (n) Variation of T cell relative cytotoxic level in response (R) and non-response
269 (NR) patients in five independent checkpoint inhibitor treated clinical data. LUAD, LUSC and MLN* represents different
270 sets of lung adenocarcinoma, lung squamous cell carcinoma and melanoma.

271 We confirmed markers of each commonly identifiable cell types in cancer microenvironment do have significant
272 overlaps with the immune and stromal cell markers identified in normal microenvironment, suggesting these
273 marker genes truly belong to immune and stromal cells rather than cancer cells (Supplementary Table S5). On
274 average, the ICTD marker genes of each cell type have ES higher than 0.9, while the ES scores of the signature
275 genes used by CIBERSORT, TIMER and EPIC are 0.22, 0.39 and 0.26, respectively. **Fig 3d** illustrate the ES of
276 T cell (sub)type markers of the four methods. The level of tumor infiltrated lymphocytes (TIL) in 12 TCGA cancer
277 types have been previously assessed by imaging data [17]. On average, the correlation between imaging predicted
278 TIL and ICTD predicted T cell level is 0.4, comparing to 0.14, 0.2, and -0.11 with CIBERORT, TIMER, and EPIC
279 predicted T cell level (**Fig 3e**). For other cell types, with a lack of ground truth, we rely on evaluating the ES
280 scores of 3,552 known immune and stromal cells marker genes. It turns out that ICTD-predicted cell proportions
281 achieved on average 0.56 $R^2$ value in explaining the 3,552 known immune and stromal cells marker genes, while
282 the $R^2$ is 0.2, 0.24, and 0.18 for CIBERSORT, TIMER, and EPIC (Supplementary Table S6).

283 ICTD enables interactive deconvolution of matched multi-omics data. We co-deconvoluted the RNA-seq, ATAC-
284 seq and DNA methylation data of five TCGA cancer types with available data (Online Methods). On average,
285 more than 70% of the cell types identified from RNA-seq data were also identified in ATAC-seq or methylation
286 data, including adipocytes, B cell, CD4+ and CD8+ T cell, macrophage, fibroblast, endothelial and dendritic cells
287 (Supplementary Table S7). The correlations between cell proportions inferred from different data types are higher
288 than 0.6. Fig 3f illustrated the strong consistency between the T cell proportion inferred from TCGA BRCA RNA-
289 seq and ATAC-seq data. It is noteworthy the samples used in multi-omics experiments were from different parts
290 of a tumor tissue, and some are less representative of the whole tumor tissue. ICTD utilizes a robust regression
291 approach to remove such samples with inconsistent cell proportions inferred from the multiple data sources. As
292 a result, the correlation between RNA-seq and imaging inferred T cell proportion was increased from 0.23 to
293 0.64, wherein the imaging based proportion is deemed as a reliable reference here. This suggests the interactive
294 co-deconvolution of multi-omics data has the potential to increase the robustness of the prediction.

295  Application of ICTD on 7 human normal brain, 5 neuro-degenerative disease and 4 brain cancer data sets
296  identified 23 common cell types in central nervous system, including two astrocyte, three general glial,  two
297  oligodendrocyte, oligodendrocyte progenitor, exhibitory and inhibitory neuron, MHC class I and II antigen
298  presenting cells, general myeloid, macrophage, neutrophil and stromal like microglial cells, one endothelial, one
299  epithelial, three ependymal, and one collagen expressing stromal like cell types (Supplementary Fig S8).

300  To experimentally validate ICTD in identifying rare sub cell types and predicting cell proportions in complex tissue
301  system, we generated an RNA-seq data set of 12 mouse bone marrow tissue samples each with flow cytometry
302  (FACS) measured cell numbers (see details in **Supplementary Notes**). ICTD successfully identified all the four
303  hematopoietic cell types measured by FACS, namely hematopoietic stem cell, general myeloid progenitor,
304  mature myeloid cell and pre-B cell, and achieved correlations of 0.92, 0.86, 0.8 and 0.97 between predicted and
305  FACS measured cell proportions. Complete statistics including labeling matrix of mouse hematopoietic cell types,
306  cell type specific markers identified by ICTD, cell proportions predicted by ICTD and measured by FACS were
307  given in **Fig 3g**, Supplementary Table S8 and Supplementary Fig S9.

308  ICTD considers the variability of resident cell types and their marker genes across tissue microenvironments and
309  technology platforms. **Fig 3h** illustrate the ES of T cell expressing genes in different CNBI data sets, suggesting
310  a significant variation of the T cell markers in the microenvironment of different cancer, inflammatory disease
311  and blood tissue, as well as under different experimental platforms [18] [19]. To further investigate how the data set
312  specific makers vary by disease/tissue micro-environments or experimental platforms, we further computed the
313  averaged Jaccard distance between the marker genes of same cell types identified in any two CNBI or single
314  cell simulated bulk datasets (**Supplementary Methods**). As illustrated in **Fig 3i**, the cell type marker genes vary
315  drastically between cancer, normal inflammatory and blood tissues. Three distinct clusters were observed (1)
316  TCGA cancer and other cancer, (2) single cell simulated cancer, and (3) TCGA normal and other inflammatory
317  disease, and blood tissue. Among the cancer data, TCGA and other RNA-seq based cancer data sets is well
318  separated from scRNA-seq simulated cancer data and the Microarray cancer data sets, and the later one is
319  further divided into two sub-clusters containing independent CRC and TNBC data sets. Similarly, the TCGA
320  RNA-seq and microarray data of normal, inflammatory conditions, and blood tissue form three distinct sub-
321  clusters. Among the microarray data of chronic inflammatory conditions, the disease of digestive system and
322  airway and skin tissues from two sub-clusters.

323  ICTD detected general T cell, fibroblast, and myeloid cells in all 28 analyzed TCGA cancer types, while the CD8+
324  T, non-collagen extracellular component expressing fibroblast, and oxidative stress producing myeloid cells were
325  identified as distinct cell types in only 10, 12, and 15 cancer types, respectively. We found that the markers of
326  these functional sub cell types are detected as cell type specific functions instead of a cell type in some cancer
327  types by the local low rank screening function. For the 19 cancer types where CD8+ T cell is not identified as a
328  cell type, CD8+ T cell markers were treated as one T cell specific function in 15 cancer types, while in 4 cancer
329  types, high concordance is observed between total T cell and CD8+T cell markers in all the samples, making the
330  CD8+ T subtype not differentiable from the general T cell. **Fig 3j** illustrated the marker genes of general T, CD8+
331  T, CD4+ T and T-reg cells form a distinct rank-4 submatrix in samples with high T cell infiltration, while the genes
332  were less distinguishable in the complete TCGA COAD data (Supplementary Fig 10). This suggests the "locality"
333  of finding identifiable cell types and functions, and hence it is necessary to implement a local low rank module
334  detection approach. Similar locality was also observed for the marker genes of non-collagen expressing
335  fibroblast and NADPH oxidase expressing myeloid cells in certain TCGA cancer types and other analyzed CRC
336  and TNBC data sets (Supplementary Fig S10). We also conducted comprehensive screening to identify unknown
337  immune/stromal cell type specific functional genes (**Online Methods**). 84 major functional modules were
338  identified as common cell type specific functions in TCGA data (**Supplementary Notes**).

339  *Cell-cell interaction*. The prediction of cell proportions and functions by ICTD makes it possible to computationally
340  characterize cell-cell interactions. We observed co-infiltrations among immune and stromal cell types with PCC
341  in the range of -0.2-0.94 in all the analyzed TCGA cancer data (Supplementary Table S9). More importantly, the
342  functional promotion or inhibition of cell type A to cell type B could now be examined by the correlations between
343  the abundance level of A and the activity level of the function in B, conditional on the predicted proportion of B.
344  We found seven genes expressed by fibroblast cells with significant negative conditional correlation with T cell
345  infiltration in at least 10 out of 15 cancer types with high level of stromal cells (p<0.01) (**Fig 3k**). The seven genes
346  execute functions related to the modification and synthesis of collagen and extracellular polysaccharide,
347  suggesting a possible role of the dysregulated extracellular matrix composition in directing T cell infiltration.
348  Similarly, the interactions of functions in two cell types can be computed by the correlation of the activity levels

349  of the two functions conditional to their proportions. We identified a low conditional correlation among CD8 T cell
350  markers such as CD8A/CD8B and cytotoxic genes, and a high conditional correlation among general T, CD8+
351  T, and cytotoxic genes in 4 cancer types, suggesting possibly perturbed cytotoxicity of T cells in the first 19
352  cancer types, namely T cell exhaustion. We also observed a significant negative correlation (p <0.01) between
353  the NADPH oxidase and T cell cytotoxicity levels conditional to the total myeloid and T cell in 11 out of the 25
354  TCGA cancer types (Supplementary Table S10). This is consistent with previous observation that NADPH
355  oxidases produce reactive oxygen species (ROS) on the surface of myeloid-derived suppressor cells that
356  suppress the cytotoxic function of T cells [20].

357  *Clinical implications.* ICTD enables investigation of the impact on clinical prognosis by microenvironment. We
358  conducted association analysis between the predicted cell proportions and varied functions with patient's overall
359  survival in TCGA data, as well as patients' response in five clinical trial data with immune checkpoint inhibitor
360  treatment (**Supplementary Methods**). We identified significant associations of patients' overall survival with T
361  cell infiltration and relative cytotoxicity levels in 12 and 7 TCGA cancer types, respectively. More interestingly, in
362  colorectal and ovarian cancer, we observed that patients with moderate level of T cell infiltration have the best
363  overall survival comparing to the patients with high and low T cell levels (**Fig 3l**). We define the T cell's relative
364  cytotoxicity (RC) level as the predicted cytotoxic function. level divided by the predicted total T cell level in each
365  sample and observed patients with higher RC have significantly better overall survival. This clearly suggests the
366  existence of T cell exhaustion and its association with poor prognosis. On the five clinical trial data, we noticed
367  that patients with high T cell infiltration have better response to the treatment (**Fig 3m**), which is consistent with
368  previously reported [21]. Moreover, the level of T cell cytotoxicity was observed to vary significantly in four datasets
369  of melanoma, lung adenocarcinoma and lung squamous carcinoma. We observed the patients with lower RC
370  tend to have better clinical response (**Fig 3n**), possibly due to more PD-1/PD-L1-mediated immuno-suppression
371  in these tumors. It is noteworthy that association between T cell infiltration and patients' clinical outcome, and
372  the identifiability of varied cytotoxic function show a high consistency between TCGA and the clinical trial data
373  (Supplementary Table S10).

## Discussion

375  Our semi-supervised deconvolution method ICTD brought up a novel notion called "identifiability" of a cell type
376  and cell type specific function, which was mathematically rigorously defined. By adaptively defining detectable
377  cell types and selecting cell type markers based on the input data resolution, ICTD highly reduces the estimation
378  bias, and also enables detection of novel cell (sub) types, and cell type functional activities. These features are
379  particularly favorable when the goal is to computationally characterize the cell-cell interactions in large-scale
380  tissue transcriptomic profiles. It is noteworthy that the "transcriptionally identifiable" cell types differ from those
381  defined by cell differentiation lineage: some cell types on the lineage map may not be identifiable, while an
382  "identifiable" cell type can be a certain cell or cell subtype, or the total of several cell types on the lineage map
383  that express same gene markers. We believe the liberty of ICTD in its deconvoluted cell types makes it entirely
384  data-driven, less biased to the training data, and it thus grants more sensible findings for downstream correlation
385  analysis with other clinical and biological features.

386  ICTD is flexible in utilizing different types of training data to construct the labeling matrix, and we noticed using
387  scRNA-seq profiles of cells from the real microenvironment of a certain cancer type, we are able to derive more
388  tissue specific cell type markers than using microarray expression profiles of primary cells collected from healthy
389  donors (**Supplementary Notes**). It is also worthy of mention that since ICTD is not fully supervised, we suggest
390  at least 10 samples is needed for the method to work. While the method has increased type II error when the
391  sample size is small, the identified rank-1 gene modules can be informative in guiding the flexible selection of
392  cell type signature genes. Based on this, our ICTD R package was integrated with a regression based approach
393  specifically for small sized samples with data-guided gene markers. When multi-omics data is available, we
394  showed that co-deconvolution of matched multi-omics data could improve the prediction robustness, by
395  excluding certain "outlier" samples with unstably predicted proportions using robust regression, and this function
396  is available in the ICTD R package. The R package and web server version of ICTD are available at
397  https://github.com/changwn/ICTD and https://shiny.ph.iu.edu/ICTD/.

398  Application of ICTD on TCGA pan-cancer data identified variations of T cell marker, cytotoxic marker and T cell
399  exhaustion level, association between fibroblast expressing genes and T cell infiltration level, and association
400  between ROS produced by myeloid cell and T cell cytotoxic level in different cancer types, suggesting the

401 capacity of ICTD in providing a comprehensive evaluation of tissue specific cell types, cell type specific function,
402 and cell-cell interactions. Nevertheless, the sensitivity of detecting cell type varied function can be largely
403 improved if more prior knowledge of functional marker genes is available. And additionally, more novel cell type
404 functions can be predicted if the rank-1 module detection approach could be optimized such that certain modules
405 may exist with respect to only a subset of samples, considering the prevalence of disease heterogeneity and
406 subtype specificity. In other words, co-expression modules local to subset of samples may be desirable in
407 revealing more cell type functions.

## Online Methods

*Single cell, bulk cell and tissue transcriptomics data sets used in this study*

We collected bulk cell data of 11 types in human blood, inflammatory and cancer tissue microenvironment, 8 types in human central nervous system, all generated by Affymetrix UA133 plus 2.0 Array; and 13 types in mouse inflammatory and tissue microenvironment, generated by Affymetrix Mouse Genome 430 2.0 Array. Detailed cell types include: *human stromal and immune cells:* fibroblast (34, 387), adipocytes (3, 26), endothelial cell (29, 606), B cell (20, 404), CD4+ T cell (23, 443), CD8+ T cell (9, 130), natural killer cell (9, 141), dendritic cell (32, 410), monocytes (22, 477), macrophages (21, 277), and neutrophil (10, 257); *human central nervous system:* neuron (16, 243), Schwann cell (2, 14), astrocyte (10, 57), ependymal cell (1, 39), oligodendrocyte (4, 30), and microglial cells (43,754), endothelial (29, 606), and stromal-like cell (34, 387); *mouse stromal and immune cells:* fibroblast (28, 277), adipocytes (3, 63), myocytes (myocyte), endothelial cell (10, 56), B cell (6, 31), CD4+ T cell (6, 80), CD8+ T cell (3, 34), natural killer cell (7, 35), dendritic cell (12, 84), monocytes (10, 46), macrophages (8, 102), neutrophils (11, 36), and mast cell (3, 31). The two numbers in the parenthesis indicate the number of datasets and samples of each cell type. We believe these cell types, together with tissue primary cells can cover major cell populations in the microenvironment of solid cancer, inflammatory disease, central nervous and hematopoietic system. 2854 samples of cancer cell line, human and mouse tissue index, and other cancer and normal tissue data were utilized as background to exclude the genes expressed by cancer or tissue primary cells.

The method was validated on single cell simulated bulk data. 13 single cell RNA-seq data sets generated by either C1/SMART-seq2 or 10x Genomics pipelines are used, and the cells are collected from (1) the TME of human solid cancer melanoma (8, 4486), breast (7, 535), colorectal (8, 375), head and neck (9, 5902), and lung cancer (8, 6630), (2) human glioma (5, 751), oligodendroglioma (7, 2728), and astrocytoma (7, 5171), (3) one public (8, 420) and one in-house (5, 1239) human normal brain sets, (4) human myeloid cell lineage and lymphoid cell lineage (3, 318) and monocyte/dendritic cell populations (4, 700), and (5) the TME of mouse melanoma (9, 2903). The two numbers indicate the number of cell types and cells of each data set.

We applied ICTD on real bulk tissue transcriptomic data of (1) 28 TCGA cancer types, (2) 11 TCGA normal tissue data, (3) 17 independent microarray data sets of colorectal cancer measured by different platforms; (4) metabric and 6 other triple negative breast cancer data sets; (5) 7 blood tissue RNA-seq and microarray data; (6) 11 human inflammatory disease data sets generated by Affymetrix UA133 plus 2.0 Array, and (7) 7 human normal brain, 5 neuro-degenerative disease and 4 brain cancer types. Detailed information of the bulk cell, scRNA-seq and bulk tissue data were provided in Supplementary Table S3. The sample information and selection, downloading and processing procedures of the public data, and sample and sequencing information of the inhouse generated data were given in **Supplementary Notes**.

*Preliminary derivation of the mathematical conditions of "Identifiable" cell types and cell type specific functions*

As detailed in **Supplementary Notes**, we analyzed the following characteristics of the cell type signature genes in the scRNA-seq and bulk tissue data of different disease context, experimental platforms and batches: (1) the consistency of cell type uniquely expressed genes were evaluated by their averaged expression level in different cell types of different scRNA-seq data sets; (2) inter- and intra- sample variations of cell type signature genes were characterized by the "drop-out" rates and multimodality of each gene's expression profile in the scRNA-seq data of different samples; (3) matrix rank and expression scale of cell type uniquely expressed genes in bulk tissue data were evaluated by using BCV based rank test and Kolmogorov Smirnov (KS) test, and (4) immune and stromal cell co-infiltrations in cancer and inflammatory tissues were further assessed by using the averaged co-expression correlations among a small number of known cell type uniquely expressed genes.

Our evaluation suggested that NMF solution may not be unique if the used marker gene set are expressed by more than one cell type due to the prevalent co-linearity of cell proportions (**Supplementary Notes**). Hence only the cell type with uniquely expressed genes are transcriptomically "identifiable", and the markers genes should also be stably expressed through cells of the same type so that its tissue level expression can reflect the cell's population in the tissue. Specifically, if gene $i$ is uniquely and stably expressed in cell type $k$, its gene expression can be expressed as $X_{i,\cdot} = S_i^k \cdot P_{k,\cdot} + e$, where $S_i^k$ is the unit expression of $i$ in $k$, and $P_{k,\cdot}$ is the relative proportion of cell type $k$ across all the samples. This shows that genes uniquely expressed by a cell type forms a (matrix) rank-1 submatrix, which form a necessary condition of "transcriptomically identifiable" cell type. On the other hand, a significant rank-1 structure of the expression profile of multiple genes $X_{i,\cdot}, i = 1 \dots m$ suggests that these

459 genes are highly possibly expressed by a dominating cell type in the current tissue microenvironment or the
460 genes are with similar expression pattern in several cell types.

461 Noting cell type specific functional activities, such as the T cell cytotoxicity, are highly varied through different
462 patients, it is not feasible to use constant gene expressions level to characterize their activities. Denote the
463 averaged level of a functional gene $i$ in cell type k in the sample j as $S_{i,j}^k$, our evaluation suggested that the
464 function is identifiable only if there exists a group of marker genes $i = 1 \dots K$ satisfy $S_{i,j}^k \cdot P_{k,j}, j = 1 \dots N$ form a
465 rank-1 matrix. Specifically, the cell type specific functional genes should share the same rank-1 space with the
466 cell type markers if there is no variation while the functional genes can be identified as the markers of a cell type
467 if $S_{i,j}^k$ varied in all samples. If only a subset of samples has the functional variation, the low rank structure of the
468 functional genes will be absorbed by the cell type markers and diminish on the co-expression network of all the
469 samples. For such a case, the linear base of the varied function can be distinguished when the computation was
470 limited to the samples with the functional variation, i.e. a local low rank identification method is needed (See
471 more discussions in **Supplementary Notes**).

472 *A modified Bi-cross validation (BCV) based test of matrix rank*

473 Bi-cross validation (BCV) has been developed to estimate the matrix rank for singular value decomposition (SVD)
474 and Non-negative Matrix Factorization (NMF), which requires a prefixed low dimension $K$ and two low rank
475 matrices for the approximation $X_{M \times N} = W_{M \times K} \cdot H_{K \times N}$.. The error distribution of gene expression data is usually non-
476 identical/independent, majorly because a gene's expression can be affected by its major transcriptional
477 regulators, other biological pathways and experimental bias. Hence undesired biological characteristics and
478 experimental bias may form significant dimensions in a gene expression data [22]. In sight of this, we developed a
479 modified BCV rank test (**Algorithm 1**) to minimize the effect of the non-i.i.d errors in assessing the matrix rank
480 of a gene expression data.

---

481 **Algorithm 1: Modified Bi-cross validation matrix rank test**

---

482 $For \ r = 1 \dots R$

483 $Sample$ row index set $I_r = \left\{ i_1, i_2, \dots, i_{\left[\frac{M}{c}\right]} | i_p \in \{1 \dots M\} \right\}, \overline{I_r} = \{1 \dots M\} \backslash I$

484 $Sample$ column index set $J_r = \left\{ j_1, j_2, \dots, j_{\left[\frac{N}{c}\right]} | j_p \in \{1 \dots N\} \right\}, \overline{J_r} = \{1 \dots N\} \backslash J$

485 $Split \ X$ into for submatrix $\begin{vmatrix} A_r & B_r \\ C_r & D_r \end{vmatrix}, where \ A_r = X[I_r, J_r], B_r = X[I_r, \overline{J_r}],$

486 $C_r = X[\overline{I_r}, J_r], D_r = X[\overline{I_r}, \overline{J_r}]$

487 $For \ k = 1 \dots \min\left( \left[\frac{M}{c}\right], \left[\frac{N}{c}\right] \right)$

488 $\text{BCV}(k, r) = \sum_{i=1}^{\left[\frac{M}{c}\right]} \sum_{j=1}^{\left[\frac{N}{c}\right]} \left\| A_r - B_r \widehat{D}^{(k)^+} C \right\|_F^2 \ (*)$

489 $\text{Rank}_x \leftarrow 0$

490 $For \ k = 1 \dots \min\left( \left[\frac{M}{c}\right], \left[\frac{N}{c}\right] \right)$

491 $Do$ t test between $\{\text{BCV}(k, r) | r = 1 \dots R\}$ and $\{BCV(k + 1, r) | r = 1 \dots R\}$

492 $if$ (p. value $< 0.01$ & mean $(\text{BCV}(k + 1, r))$ − mean $(\text{BCV}(k, r)) < \text{msp})$

493 $\text{Rank}_x \leftarrow k + 1$

494 $Return \ \text{Rank}_x$

495 $(*)$ Denote the SVD of a matrix $D$ as $D = U\Sigma V'$, and Moore– Penrose inverse of $D$

496 as $D^+, D^+ = V'\Sigma^+ U$, where $\Sigma^+$ is a diganol matrix $\text{diag}(\sigma_1^+, \sigma_2^+, \dots \sigma_p^+)$ with $\sigma_1^+ \geq$

497 $\sigma_2^+ \geq \cdots \geq \sigma_p^+ \geq 0$. Define $\widehat{D}^{(k)^+} = \sum_{i=1}^k \sigma_i^+ v_i u_i$

---

498  *ICTD Step 1: Construction of labeling matrix to represent TME specific cell type marker genes*

499  A labeling matrix $L_{M \times K}$ was first constructed to represent the genes that are overly expressed in a certain cell
500  type, where $M$ is the number of genes and $K$ is the number of cell types, $L_{i,j} = \frac{1}{R}$ stands for the gene $G_i's$
501  expression in cell type $C_j$ is the $R$th highest among its expression in all the cells, and $L_{i,j} = 0$ stands for $G_i$ is not
502  a significant signature of cell type $C_j$. Two different approaches were developed to construct the labeling matrix
503  by using scRNA-seq or bulk cell data:

504  (1) scRNA-seq data:

505  For a scRNA-seq data set with annotated cell labels of $K$ cell types and a given gene $g$, denote the expression
506  profile of $g$ in cell type k as $x_{g,\cdot}^k$, its mean as $x_g^k = mean(x_{g,\cdot}^k)$, and the Z score of $x_g^k$ as $z_g^k$. The cell type order
507  vector $o$ was further computed, where $o_j = k$, if the $j$th largest value of $x_g^k$ happens to be of cell type $k$. Then for
508  cell type $o_1$ to $o_K$, the labeling matrix was built by

509
$$L_{g,z_k} = \begin{cases} 0, & if\ z_g^k < -1.96 \\ \frac{1}{k}, if\ x_{g,\cdot}^{o_k} < x_{g,\cdot}^{o_{k-1}}, z_g^k \geq -1.96 \\ \frac{1}{p},\ if\ x_{g,\cdot}^{o_k} < x_{g,\cdot}^{o_{p-1}}\ and\ x_{g,\cdot}^{o_k} \not< x_{g,\cdot}^{o_p}, z_g^k \geq -1.96, 1 \leq p \leq k-1 \end{cases}$$

510  , where $x_{g,\cdot}^{o_i} < x_{g,\cdot}^{o_j}$ denotes $g$ is significant over expressed in cell type $o_j$ compare to cell type $o_i$, which is tested
511  by using MAST [23].

512  (2) bulk cell data:

513  We applied a non-parametric random walk based approach to identify if a gene has higher expression in certain
514  cell types comparing to others, i.e. a signature gene of the cell types, by using the training data set composed
515  by a large independent data sets of the cell types. ICTD enables the user to select the cell types specific to a
516  tissue microenvironment. For examples, bulk cell data of normal breast cell, breast cancer cell lines and breast
517  cancer tissue samples were selected as background to train the marker genes of immune and stromal cells for
518  analyzing breast cancer tissue data. The labeling matrix used in this paper were computed by using human
519  CCLE cell line, human body index and more than 20 human cancer tissue data as the background data. Batch
520  effect of the training data of each cell type were first removed by using COMBAT [24] and the expression profile of
521  each sample was further normalized by its mean.

522  Denote the combined expression matrix containing $M$ genes for $N$ samples of $K$ cell types, and for each cell type,
523  we first calculated the expected frequency of the cell type, i.e. dividing the total number of samples for the cell
524  type ($N_k, k = 1, ..., K$) by the total number of samples $N$, denoted by $E_i = N_k/N$, $i = 1, ..., K$. For a given gene $g$,
525  denote $x_{g,\cdot}$ **and** $x_{g,\cdot}^k$ as its expression profile of all cell types and cell type $k$. We order the corresponding cell type
526  labels of these samples based on the expression value from large to small, denoted by vector $z$, where $z_j = k$,
527  if the $j$th largest expression value in $x_{g,\cdot}$ happens to be of cell type $k$. Denote $O_k$ as the cumulative frequency of
528  cell type $k$ over the expression order of $x_{g,\cdot}$, which is calculated as:

529
$$O_{jk} = \frac{\sum_{m=1}^{j} \delta_{z_m=k}}{j}, j = 1, ..., N$$

530  , where $\delta_{z_m=k}$ is the indicating function for $z_m = k$. A discrepancy score vector $d$ between the observed and
531  expected cell type frequency was further defined as

532
$$d_j = \sum_{k=1}^{K} (O_{jk} - E_k)^2, j = 1, ..., N$$

533  , where $d$ is a non-negative vector of length $N$, and it attains a minimum value of zero at $N$. The larger the
534  maximum value $d$ suggests the expression values are more enriched in certain cell types than the others. Denote

535    $m$ as the index of the maximum of $d$, i.e. $d_m = \max(d_j)$, and the cell type frequency at the best discrepancy as

536    $e_k^m = O_{mk} - E_k$, the cell types were further ordered by $e_k^m$ from large to small and denoted as $\boldsymbol{o}$, where $\boldsymbol{o}_j = k$ if

537    the $j$th largest value of $e_k^m$ happens to be of cell type $k$. Then for cell type $\boldsymbol{o}_1$ to $\boldsymbol{o}_K$, the labeling matrix was built

538    by $L_{g,z_k} = \begin{cases} 0, & if\ e_k^m \leq 0 \\ \frac{1}{k}, & if\ \boldsymbol{x}_{g,\cdot}^{\boldsymbol{o}_k} < \boldsymbol{x}_{g,\cdot}^{\boldsymbol{o}_{k-1}}, e_k^m > 0 \\ \frac{1}{p}, & if\ \boldsymbol{x}_{g,\cdot}^{\boldsymbol{o}_k} < \boldsymbol{x}_{g,\cdot}^{\boldsymbol{o}_{p-1}}\ and\ \boldsymbol{x}_{g,\cdot}^{\boldsymbol{o}_k} \nless \boldsymbol{x}_{g,\cdot}^{\boldsymbol{o}_p}, e_k^m > 0, 1 \leq p \leq k-1 \end{cases}$ , where $\boldsymbol{x}_{g,\cdot}^{\boldsymbol{o}_i} < \boldsymbol{x}_{g,\cdot}^{\boldsymbol{o}_j}$ denotes $g$ is significant

539    over expressed in cell type $\boldsymbol{o}_j$ compare to cell type $\boldsymbol{o}_i$, which is tested by Mann Whitney test.

540    *Exclusion of the expression of undesired cells*

541    ICTD can eliminate the expression signal from undesired cell types to excluder those cells from further analysis.
542    To do this, ICTD first identifies gene co-expression modules from the decentralized expression matrix of their
543    marker genes by using WGCNA and computes the first row base of each module by using SVD [25]. Then for each
544    gene that is positively co-expressed with one or several module(s) of the undesired cell type, its expression are
545    further projected to the complementary space spanned by the first row base of each of such modules (s). Denote
546    a decentralized tissue data as X, the data of pseudo-code of exclusion of the expression of undesired cells are
547    given below:

---

548    **Algorithm 2: Remove the low rank space of undesired cell types**

---

549    $Modules_C \leftarrow WGCNA(X_C)$

550    $for\ i\ in\ Modules_C$

551        $U_i \Sigma_i V_i^T = SVD(X_i)$

552        $RB_c[i,] \leftarrow V_i^T[,1]$

553    $for\ each\ gene\ in\ X$

554        $for\ k\ in\ 1{:}K$

555           $if\,(\max(cor(RB_c, X_{genes})) > 0)$

556           $i \leftarrow \underset{i}{\operatorname{argmax}}(cor(RB_c[i,], X_{genes}))$

557           $X_{gene} \leftarrow X_{gene} - X_{gene}\dfrac{RB_c[i,]RB_c[i,]^t}{||RB_c[i,]||^2}$

558    $return(X)$

---

559    In this paper, we first identified 1089 cancer cell genes, as evidenced by their consistent up-regulation in 11
560    cancer types of TCGA data and significant expression in CCLE cell line data (Supplementary Table S10).
561    Differential gene expression analysis was conduct by using Mann-Whitney test with FDR<0.05 as the significant
562    cutoff and significant expression in cancer cell line data is determined by log(FPKM)>2. In the analysis of one
563    specific cancer type, gene co-expression modules of the cancer genes were first identified. The linear space
564    spanned by the modules were further excluded by the complementary space projection. Our analysis on single
565    cell simulated and real bulk tissue data validated that such an elimination procedure can largely remove the
566    expression of the genes stably expressed in cancer cells while retaining the low rank structure of the gene
567    expressions from other cells (See **Supplementary Notes**).

568    *ICTD Step 2: Identification of rank-1 modules*

569    Highly co-expressed modules were identified using our in-house method, namely MRHCA [26] [27]. More details
570    about the MRHCA based module identification and its rationality in our case are given in **Supplementary**
571    **Methods.**

572 The BCV test described in **Algorithm 1** is further applied to find the modules of rank-1, which possibly
573 correspond to marker genes of identifiable cell types. The matrix rank of a module centered by a cell type
574 uniquely expressed genes always increases with the module size, due to the genes less co-expressed with the
575 hub may be expressed by other cell types. In this paper, we selected the modules of with hub significance p<1e-
576 3, average co-expression correlation>0.8, rank=1 (p<1e-3) and with at least seven genes, as possible markers
577 of identifiable cell types.

578 *ICTD Step 3: Determine the number and select Rank-1 modules of "identifiable" cell types*

579 After identifying all sets of rank-1 marker genes, ICTD further determines the number of identifiable cell types,
580 eliminates redundant and insignificant cell type marker genes, annotates each set of marker genes with a most
581 likely cell type by using the labeling matrix, and build a marker gene – cell type representing matrix for the
582 downstream deconvolution analysis.

583 Denote a rank-1 marker set $G_i = \{g_1, \ldots, g_{n_i}\}$ and labeling matrix $L_{M \times K}$, we first compute $S_i = \{s_{i,1}, \ldots, s_{i,K}\}$, where
584 $s_{i,k} = \sum_{j=1}^{n_i} L_{g_j,k}$ representing the enrichment level of $G_i$ to the genes top expressed in cell type k. The
585 significance level of $s_{i,k}$, $p_{s_{i,k}}$, is assessed by a permutation test, and $G_i$ is annotated as cell type with the minimal
586 $p_{s_{i,k}}$ if $\min(FDR(p_{s_{i,k}})) < Cutoff_{CES}$. In this study, $Cutoff_{CES}$ is selected as 0.01. The rank-1 markers annotated
587 without a significant cell type annotation are excluded from further analysis. It is noteworthy that a larger
588 $Cutoff_{CES}$ can be selected for identification of possible unknown cell types.

589 Rather than predefining the cell types, ICTD determines the cell types that are "identifiable". In some
590 circumstance, the proportion of the cell type with a lower resolution is a non-negative linear sum of the proportion
591 of several cell types with higher resolutions, such as the myeloid cell proportion equals to the sum of macrophage
592 and neutrophils when these two cell types dominate the myeloid cell populations in the tissue [28]. This linear
593 dependency may correspond to a linear dependency between the row base of marker genes of cell types of
594 different resolutions, which may result in number of identifiable cell types exceeding the rank of the linear space
595 generated by the identified rank-1 markers.

596 To determine the number of identifiable cell types covered by the rank-1 marker genes, ICTD first construct a
597 tree structure to represent the linear dependency among the identified rank-1 marker sets. A rank-1 marker set
598 is considered as a root node if its row base can be non-negatively fitted by the row bases of other nodes with
599 $R^2 > Cutoff_{R^2}$. In this study, $Cutoff_{R^2} = 0.9$ is selected. The rank-1 marker sets fitting each other with $R^2 >$
600 $Cutoff_{R^2}$ are merged together. All the root rank-1 marker sets are considered as markers of "identifiable" cell
601 types and excluded from the further analysis. ICTD further computes the rank of the expression matrix of all the
602 non-root rank-1 maker genes. Denoting the number of non-root rank-1 maker sets and their total rank as $P$ and
603 $\hat{P}$. The total number of "identifiable" cell types among the non-root rank-1 marker sets is determined as $\hat{P}$.

604 A marker gene – cell type representation matrix is further computed for the downstream NMF analysis. Denote
605 a selected rank-1 marker set as $G_i = \{g_1, \ldots, g_{n_i}\}, i = 1 \ldots P$, its gene expression profile as $X_{G_i}$, and ot SVD as
606 $X_{G_i} = U_i \Sigma_i V_i^t$, $G_i$'s self-explanation score is defined as $\frac{\sum_{g \in G_i} cor(X_g, V_i[,1])^2}{|G_i|}$, i.e. the averaged R square of the genes'
607 expression fitted by their first row base. The marker gene – cell type representation matrix C is constructed by
608 **Algorithm 3**:

609 **Algorithm 3: Construction of representation matrix**

610 $for\ i\ in\ 1 \ldots P$

611     $Compute\ the\ SVD\ of\ X_{G_i}\ as\ U_i \Sigma_i V_i^t$

612     $Conduct\ a\ hierachical\ clustering\ of\ G_i\ in\ to\ \hat{P}\ clusters\ C_j, i =$
613 $1 \ldots \hat{P}, by\ using\ eucliean\ distance\ between\ V_i[,1]$

614 $for\ j\ in\ 1 \ldots \hat{P}$

615 $$Select\ rank\ 1\ marker\ set\ G_{k_j}\ by\ \underset{j_k}{argmax}(\frac{\sum_{g\in G_{j_k}} cor(X_g, V_{j_k}[,1])^2}{|G_{j_k}|}\ |G_{j_k} \in C_j)$$

616 $$C_{\sum_{j=1\ldots\hat{P}} n_{j_k}\times\hat{P}}[i,j] = \begin{cases} 0, if\ gene\ i \notin G_{k_j} \\ 1, if\ gene\ i \in G_{k_j} \end{cases}$$

617 $$return(C_{\sum_{j=1\ldots\hat{P}} n_{j_k}\times\hat{P}})$$

618 This step assigns marker genes of identifiable cell types that highly determines the prediction accuracy of the
619 deconvolution analysis. ICTD also includes three other options in constructing marker genes and C matrix of
620 identifiable cell types. The computational details and performance comparison of these methods were given in
621 **Supplementary Methods.**

622 *ICTD Step 4: Constrained Non-negative Matrix Factorization*

623 With the NMF constraint matrix $CS_{X\times K}^{NMF}$, each of the K cell type is assigned with at least one cell type uniquely
624 expressed gene (see derivations in method), hence the constraint NMF problem $X_{M\times N} = S_{M\times K} \cdot P_{K\times N}, S[I, k] \geq$
625 $0, P[k, j] \geq 0, S[I, k] = 0\ if\ CS^{NMF}[I, k] = 0$ does have a unique solution [29]. The rationale here is that the analysis
626 only focuses on cell types with uniquely expressed markers that form rank-1 structure, and the analysis is robust
627 to collinearity of cell proportions due to the uniqueness of solution. Specifically, for the $p$th disconnected
628 subgraph with $M_p$ genes, rank= $K_p$, and constraint matrix $C_{M\times K}$, the NMF of $X_{M\times N} = S_{M\times K} \cdot P_{K\times N}$ is solved by
629 $\underset{S,P}{min}(\|X_{M\times N} - S_{M\times K} \cdot P_{K\times N}\|_F^2 + \lambda \cdot tr(S_{M\times K}^T \cdot (1 - C_{M\times K})))$, where $S_{M\times K}$ and $P_{M\times K}$ are the predicted signature and
630 proportion of K cell types. Variables with fitted S that are highly varied from C are further removed. Detailed
631 solution of the constrained NMF problem was given in **Supplementary Methods**. It is noteworthy when $\lambda \to \infty$,
632 $P_{i,j}$ is the first row base of the SVD of $Diag(C_{.,j}) \cdot X$, where $Diag(C_{.,j})$ is the diagonal matrix generated by $C_{.,j}$. In
633 this study, $\lambda$ was selected based the best prediction accuracy trained on single cell simulated bulk data.

634 *ICTD Step 5: Co-deconvolution of matched multi-omics data*

635 Multi-omics, including epigenetic and chromatin profiles, provide equally important characterization of tissue
636 compositions as transcriptomic profiles. When multiple omics data are available for the same tissues, it is
637 reasonable to assume that cell relative proportions deduced from each of the omics profile should be strongly
638 associated. Based on this, co-deconvolution of matched multi-omics data could be used to cross-validate and
639 robustify the proportion predictions, as detailed in **Algorithm 4**:

---

640 **Algorithm 4: Co-deconvolution of matched multi-omics data**

---

641 **Input:**

642 $U^{(0)} = \emptyset$, denoting the set of outlier samples.

643 For $i = 1 \ldots N_{iter}$

644       Run deconvolution on each of omic profile $l$ where only samples not in $U^{(i-1)}$ are used, denoted the
645 predicted proportion matrix as $P^{(i),l}$, $l = 1, \ldots, L$, of dimension $K \times N_i$, where $K$ is the total number of cell types, and
646 $N_i$ the total number of tissues of the current run;

647       Perform robust mixture regression using robust trimmed likelihood estimation (TLE) approach, between
648 the $r_1$th row of $P^{(i),l_1}$ and $r_2$th row of $P^{(i),l_2}$

649       Collect all the outlier samples based on the robust TLE approach for all the runs, and denote the union
650 set of outlier samples as $U^{(i)}$

651       Repeat 1-3, and stop if $U^{(i)} = \emptyset$

---

652 *ICTD Step 6: Conditional local low rank test of cell type varied function*

653 Identifiable cell type specific function is defined by a group of genes that form a local rank-1 structure conditional
654 on the estimated proportion of the cell type. A kernal function based local low rank structure screening method
655 is developed for identification of such local rank-1 structures. Denote $P_k = \{p_1^k, p_2^k, ..., p_n^k\}$ as predicted proportion
656 of cell type k through the n samples and $P_{(k)} = \{p_{k(1)}, ..., p_{k(n)}\}$ as sorted $P_k$ with an increasing order, $G_{I_k}$ as the
657 rank-1 marker genes of cell type k, and $G_{F_k}$ is a gene set containing possible marker genes of a varied function
658 of k, the level of functional activity and its associated marker genes can be identified by **Algorithm 5**:

---

**Algorithm 5: BCV screening of a local low rank structure**

---

660 For a given data $X$ and cell proportion $P_k = \{p_1^k, p_2^k, ..., p_n^k\}$

661 $Sort\ P_{(k)}$ by increasing order: $P_{(k)} = \{p_{k(1)}, ..., p_{k(n)}\}$

662 Reorder the samples in $X$ into $X^0$ by the order of $P_{(k)}$

663 For $i = 1 ... N$

664 Do BCV test of $X_i \triangleq X^0[(G_{I_k}, G_{F_k}),] \cdot diag(K_i)$ $(*)$

665 $p_{ij} =$ FDR correted p value of the rank j of $X_i$

666 If $\exists\ i^*$ and $j > 1$,

667 s.t. $p_{ij} < 0.05$ for all $i \geq i^*$ and $p_{ij} \geq 0.05$ for all $i < i^*$

668 $\rightarrow G_{F_k}$ contains marker genes of a varied function

669 Identify gene froming the rank 1 matrices in $X[G_{F_k}, (i^* ... N)]$

670 $(*)$ $K_i$ is a nonnegative kernal function centered at $i$:

671 $K_i(z) = \begin{cases} 0 & , if\ |z - i| \geq C_1 \\ \frac{|z-i|-C_2}{C_1-C_2}, & if\ C_1 < |z - i| < C_2, C_2 < C_1, z = 1..N \\ 1 & , if\ |z - i| \leq C_2 \end{cases}$

---

672 The idea of this algorithm is that the genes of a cell type specific function may form additional ranks in the
673 samples with high proportion of the cells, which can be identified by the BCV test when only looking at those
674 samples. The kernel function is to smooth the inter-sample variation in cell proportions (see more details in
675 **Supplementary Methods**).

676 In this paper, $G_{F_k}$ is selected for each cell type k by the genes annotated as top expressed by cell type $k$ in the
677 labeling matrix and with more than 0.8 co-expression correlation with the cell type $k$'s proportion. ICTD enables
678 users to predefine $G_{F_k}$ and select proportion of cell type $k$ for a specified analysis, such as using known markers
679 for prediction of T cell cytotoxicity [30]. The functional activity level of each set of gene markers are then predicted
680 by its first row base in the samples i ≥ i* by SVD. Averaged activity level per cell is further estimated by dividing
681 the predicted functional activity level by the predicted cell type proportion.

682 *Single cell simulated Bulk Tissue data*

683 Cell types in each scRNA-seq data were labeled by the cell clusters provided in the original works or by using
684 Seurat pipeline with default parameters. Detailed information of the scRNA-seq data and cell type annotation is
685 given in **Supplementary Table S3 and Notes**. For each data set, we simulate bulk tissue data with three steps:
686 (1) randomly generate the proportion of each cell type, called true proportion in this paper, that follows a Dirichlet
687 distribution, (2) enforce a certain co-infiltration level of two selected cell types, and (3) draw cells randomly from
688 the cell pool with replacement according to the cell type proportion, and sum up the expression values of all cells
689 to produce a pseudo bulk tissue data. More details are provided in **Algorithm 6**:

---

**Algorithm 6: simulate bulk data using single cell**

---

691 Input: single cell gene expression matrix $S^{m \times n}$; cell type label vector $\boldsymbol{l}$; patient number $p$; total cell number $N$.

692    (optional) $CoF \in \{0,1\}$; $Corr \in (0,1]$; $row1, row2$.

693        1. Find the cell type number $k$ from $\boldsymbol{l}$.

694        2. Generate $D^{k \times p}$, s.t. $\mathbf{d}_{.,i} \sim \mathbf{Dirichlet}(\boldsymbol{\alpha}), i = 1, \dots, p$; $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_k) \sim \boldsymbol{U}(0,1)$.

695        3. If $CoF$ is TURE go to step 4, else go to step 7.

696        4. $v1, v2 \leftarrow D_{ij}, i \in \{row1, row2\}, j \in \{1, \dots, p\}$.

697        5. Generate $u1, u2$, s.t. $|cor(u1, u2) - x| \leq 0.05$.

698        6. $v1 \leftarrow u1, v2 \leftarrow u2$.

699        7. For $i = 1, \dots, p; j = 1, \dots, k$:

700            i) $\boldsymbol{sc}_{i,j} \sim Sample\ D_{j,i} \cdot N$ cells from the pool of cell type $j$ with replacement ;

701            ii) $B_{.,i} = \dfrac{\sum_{t \in \boldsymbol{sc}_{i,j}} S_{.,t}}{length(\boldsymbol{sc}_{i,j})}$

702        8. Return $B^{m \times p}$, $D^{k \times p}$.

703    , in which:

704    $Corr$ is the coinfiltration level parameter if $CoF$ is TRUE;

705    $CoF$ is the coinfiltration flag to indicate whether adding dependency to two cell types or not;

706    $row1, row2$ are the cell type location that indicate two selected cell types adding $Corr$ dependency;

707    $\boldsymbol{sc}_{i,j}\ i \in \{1, \dots, p\}, j \in \{1, \dots, k\}$ is the selected single cells sampling randomly from the cell pool with replacement;

708    $B^{m \times p}$ is the simulated bulk tissue expression value matrix;

709    $D^{k \times p}$ is the true proportion matrix.

710

711    The Dirichlet distribution matrix was generated with R package "DirichletReg" (version 3.5.3). In order to evaluate
712    the robustness of the deconvolution method while co-infiltration exists, we add different levels of co-infiltrations
713    in our simulated bulk data to four pairs of cells that are commonly known to co-infiltrate in cancer tissue, namely,
714    B/T cell, T/NK cell, Fibroblast/Endothelial cell, and B/Dendritic cell. (Supplementary Figure S13). For a robust
715    method evaluation, five replicates were generated in the simulation of each data set and at each co-infiltration
716    parameter.

717

718    *Explanation Score to evaluate the performance of our deconvolution method*

719    We assessed the methods' performance by the correlation between predicted and known proportion of each cell
720    type in simulated data, which is inapplicable in the real tissue data. Thus, we developed

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

723
$$EScore(x) = 1 - \sum_{j=1}^{N}(x_j^* - \hat{x}_j)^2 \Big/ \sum_{j=1}^{N}(x_j^*)^2$$

724
$$\hat{x}_j = \sum_{k=1}^{k_x} \beta_k^x p_j^k, \beta_k^x \geq 0$$

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

730

# References

731

732 1 Hackl, H., Charoentong, P., Finotello, F. & Trajanoski, Z. Computational genomics tools for dissecting tumour-
733 immune cell interactions. *Nat Rev Genet* **17**, 441-458, doi:10.1038/nrg.2016.67 (2016).
734 2 Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods* **12**, 453-
735 457, doi:10.1038/nmeth.3337 (2015).
736 3 Li, B. *et al.* Comprehensive analyses of tumor immunity: implications for cancer immunotherapy. *Genome Biol*
737 **17**, 174, doi:10.1186/s13059-016-1028-7 (2016).
738 4 Racle, J., de Jonge, K., Baumgaertner, P., Speiser, D. E. & Gfeller, D. Simultaneous enumeration of cancer and
739 immune cell types from bulk tumor gene expression data. *Elife* **6**, doi:10.7554/eLife.26476 (2017).
740 5 Vallania, F. *et al.* Leveraging heterogeneity across multiple datasets increases cell-mixture deconvolution
741 accuracy and reduces biological and technical biases. **9**, 4735 (2018).
742 6 Abbas, A. *et al.* Immune response in silico (IRIS): immune-specific genes identified from a compendium of
743 microarray expression data. **6**, 319 (2005).
744 7 Abbas, A. R., Wolslegel, K., Seshasayee, D., Modrusan, Z. & Clark, H. F. J. P. o. Deconvolution of blood microarray
745 data identifies cellular activation patterns in systemic lupus erythematosus. **4**, e6098 (2009).
746 8 Wang, X., Park, J., Susztak, K., Zhang, N. R. & Li, M. J. N. c. Bulk tissue cell type deconvolution with multi-subject
747 single-cell expression reference. **10**, 380 (2019).
748 9 Newman, A. M. *et al.* Determining cell type abundance and expression from bulk tissues with digital cytometry.
749 *Nat Biotechnol* **37**, 773-782, doi:10.1038/s41587-019-0114-2 (2019).
750 10 Finotello, F. & Trajanoski, Z. J. C. I., Immunotherapy. Quantifying tumor-infiltrating immune cells from
751 transcriptomics data. **67**, 1031-1040 (2018).
752 11 Varn, F. S., Wang, Y., Mullins, D. W., Fiering, S. & Cheng, C. Systematic Pan-Cancer Analysis Reveals Immune Cell
753 Interactions in the Tumor Microenvironment. *Cancer Res* **77**, 1271-1282, doi:10.1158/0008-5472.CAN-16-2490
754 (2017).
755 12 Li, B., Liu, J. S. & Liu, X. S. Revisit linear regression-based deconvolution methods for tumor gene expression
756 data. *Genome Biol* **18**, 127, doi:10.1186/s13059-017-1256-5 (2017).
757 13 Dormann, C. F. *et al.* Collinearity: a review of methods to deal with it and a simulation study evaluating their
758 performance. **36**, 27-46 (2013).
759 14 Frishberg, A. *et al.* Cell composition analysis of bulk genomics using single-cell data. *Nat Methods* **16**, 327-332,
760 doi:10.1038/s41592-019-0355-5 (2019).
761 15 Li, Z. & Wu, H. J. G. b. TOAST: improving reference-free cell composition estimation by cross-cell type differential
762 analysis. **20**, 190 (2019).
763 16 Lee, J., Kim, S., Lebanon, G., Singer, Y. & Bengio, S. J. T. J. o. M. L. R. LLORMA: Local low-rank matrix
764 approximation. **17**, 442-465 (2016).
765 17 Saltz, J. *et al.* Spatial organization and molecular correlation of tumor-infiltrating lymphocytes using deep
766 learning on pathology images. **23**, 181-193. e187 (2018).
767 18 Venteicher, A. S. *et al.* Decoupling genetics, lineages, and microenvironment in IDH-mutant gliomas by single-cell
768 RNA-seq. **355**, eaai8478 (2017).
769 19 Tirosh, I. *et al.* Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-seq. **352**, 189-
770 196 (2016).
771 20 Gabrilovich, D. I. & Nagaraj, S. J. N. r. i. Myeloid-derived suppressor cells as regulators of the immune system. **9**,
772 162 (2009).
773 21 Jiang, P. *et al.* Signatures of T cell dysfunction and exclusion predict cancer immunotherapy response. 1 (2018).
774 22 Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. J. N. m. Deep generative modeling for single-cell
775 transcriptomics. **15**, 1053 (2018).
776 23 Finak, G. *et al.* MAST: a flexible statistical framework for assessing transcriptional changes and characterizing
777 heterogeneity in single-cell RNA sequencing data. **16**, 278 (2015).
778 24 Johnson, W. E., Li, C. & Rabinovic, A. J. B. Adjusting batch effects in microarray expression data using empirical
779 Bayes methods. **8**, 118-127 (2007).
780 25 Langfelder, P. & Horvath, S. J. B. b. WGCNA: an R package for weighted correlation network analysis. **9**, 559
781 (2008).
782 26 Zhang, Y. *et al.* MRHCA: a nonparametric statistics based method for hub and co-expression module
783 identification in large gene co-expression network. **6**, 40-55 (2018).

784   27   Zhang, C., Liu, C., Cao, S. & Xu, Y. J. J. o. m. c. b. Elucidation of drivers of high-level production of lactates
785        throughout a cancer development.  **7**, 267-279 (2015).
786   28   Chen, C.-Z., Li, L., Lodish, H. F. & Bartel, D. P. J. s. MicroRNAs modulate hematopoietic lineage differentiation.
787        **303**, 83-86 (2004).
788   29   Huang, K., Sidiropoulos, N. D. & Swami, A. J. I. T. o. S. P. Non-negative matrix factorization revisited: Uniqueness
789        and algorithm for symmetric decomposition.  **62**, 211-224 (2013).
790   30   Van Acker, H. H., Capsomidis, A., Smits, E. L. & Van Tendeloo, V. F. J. F. i. i. CD56 in the immune system: more
791        than a marker for cytotoxicity?  **8**, 892 (2017).

792