


METHOD

Open Access



# GTM-decon: guided-topic modeling of single-cell transcriptomes enables sub-cell-type and disease-subtype deconvolution of bulk transcriptomes

Lakshmiapuram Seshadri Swapna<sup>1</sup>, Michael Huang<sup>1</sup> and Yue Li<sup>1\*</sup> 

\*Correspondence:  
yueli@cs.mcgill.ca

<sup>1</sup> School of Computer Science,  
McGill University, Montreal, QC,  
Canada

## Abstract

Cell-type composition is an important indicator of health. We present Guided Topic Model for deconvolution (GTM-decon) to automatically infer cell-type-specific gene topic distributions from single-cell RNA-seq data for deconvolving bulk transcriptomes. GTM-decon performs competitively on deconvolving simulated and real bulk data compared with the state-of-the-art methods. Moreover, as demonstrated in deconvolving disease transcriptomes, GTM-decon can infer multiple cell-type-specific gene topic distributions per cell type, which captures sub-cell-type variations. GTM-decon can also use phenotype labels from single-cell or bulk data to infer phenotype-specific gene distributions. In a nested-guided design, GTM-decon identified cell-type-specific differentially expressed genes from bulk breast cancer transcriptomes.

**Keywords:** Cell-type deconvolution, Single-cell transcriptome, Topic models, Variational inference, Bayesian modeling, Type 2 diabetes, Cancer transcriptome, Disease biomarkers

## Background

Cell-type composition and its relative proportions in a tissue is an indicator of health. For example, several studies have shown that type 2 diabetes is characterized by reduced beta cell mass and number in pancreatic tissue [1, 2]. In acute myeloid leukemia (AML), cell-type abundance variation between patients was indicative of degree of malignancy [3]. Experimental approaches such as fluorescence-activated cell sorting (FACS) and immunohistochemistry (IHC) are used to elucidate cell-type composition of biological samples. Single-cell RNA-sequencing (scRNA-seq) technology enables high-resolution cell-type-specific (CTS) transcriptome analysis, providing molecular insights into the cell-type composition, cell-state behavior, and cell-type heterogeneity [4–7]. In the context of cancer research, scRNA-seq has led to



© The Author(s) 2023. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

identification of distinct cancer cell states. These are shown to occur across a range of cancer types and are implicated in tumor progression, with high cell-type heterogeneity associated with poor prognostic outcomes [8, 9]. However, challenges such as high cost, lower throughput, and difficulties in dissociation of cell types in solid samples make it hard to apply these experimental approaches at the patient population scale.

On the other hand, bulk RNA-seq has been the workhorse behind transcriptome research over the past decades. Its falling costs and ease of experimental setup make it an attractive tool to work with any organism [10]. Several databases host tremendous amounts of bulk RNA-seq data such as Gene Expression Omnibus (GEO) [11, 12], Genotype-Tissue Expression (GTEx) [13], and The Cancer Genome Atlas (TCGA), a repository of bulk RNA-seq data for more than 11,000 primary cancer samples [14]. However, bulk RNA-seq data are mixture of gene expression profiles in the tissue. Computational approaches have been developed to deconvolve these bulk RNA-seq profiles into their constituent cell types in the form of cell-type proportions, since these are substantially cheaper and easier to obtain than conducting scRNA-seq experiments. Moreover, deconvolving the bulk data into constituent cell-type components can not only yield the cell-type proportions facilitating clinical investigation but also enable high-resolution differential analysis of gene expression. Studying the differentially expressed genes at the cell-type or cell-state level can help uncover gene regulatory programs that drive different tissue states. Many deconvolution methods were developed to this end.

Most of the deconvolution approaches require a set of gene markers for each cell type of interest. These marker genes are derived from expert knowledge or differential expression analysis of purified samples of specific cell types. Early methods that leverage these marker genes could achieve good performance in deconvolving mixtures with highly distinct cell types such as blood [15]. CIBERSORT improved on these approaches by incorporating a feature selection step, where genes are adaptively selected from the signature matrix based on the input bulk RNA-seq data [16]. It uses a linear support vector regression (SVR) to delineate closely related cell types such as leukocytes. BSEQ-sc combined with CIBERSORT was an early method that used scRNA-seq data as a reference for bulk deconvolution [17]. CIBERSORTx also uses scRNA-seq as reference profiles, along with improved normalization schemes to suppress cross-platform variation, and an adaptive noise filter to eliminate unreliably estimated genes [18]. MuSiC adopts a weighted non-negative least-squares regression approach and addresses the issue of cross-subject heterogeneity as well as within-cell-type variation of gene expression [19]. EPIC accommodates user-defined reference profiles to account for the presence of uncharacterized cell types in the target bulk samples [20]. Bisque learns gene-specific transformations of the bulk data based on the single-cell reference profiles and the corresponding cell proportions to account for their differences [21]. BayesPrism uses Bayesian inference to model scRNA-seq data jointly with bulk RNA-seq data to infer cell-type composition and their proportions [22]. The joint modeling overcomes biases that may arise due to technical and biological differences. Beyond cell-type deconvolution, some recently developed methods can estimate CTS gene expression from the bulk samples [23, 24]. However, these methods rely on an external cell-type deconvolution method like those aforementioned ones and do not utilize or model the distribution of

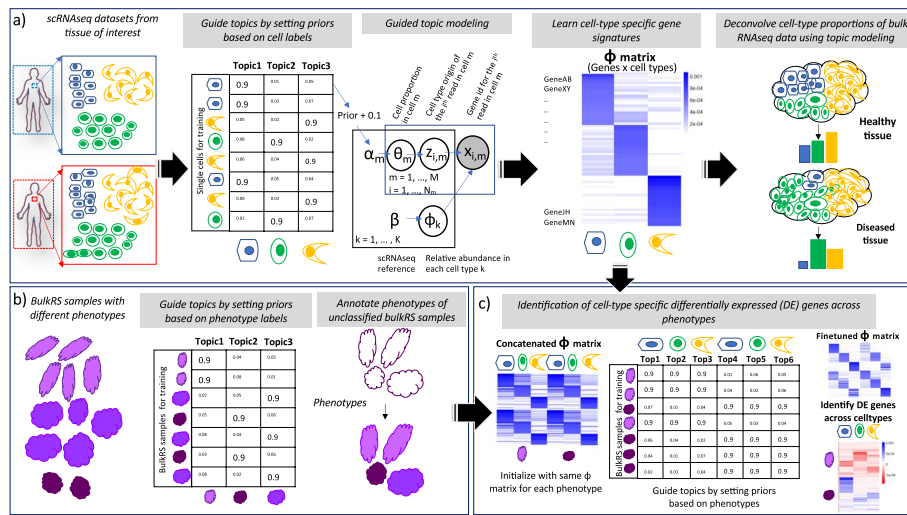
the scRNA-seq reference data to properly express statistical uncertainty while leveraging their information richness.

In this study, we present a guided topic model for cell-type deconvolution (GTM-decon). As an overview of our analysis, we first benchmark GTM-decon on deconvolving simulated and real bulk data in comparison with the state-of-the-art deconvolution methods. We then train GTM-decon on pancreatic and breast scRNA-seq datasets to deconvolve bulk RNA-seq datasets from pancreatic and breast tissues, respectively. When applied to deconvolving cancer bulk transcriptomes, GTM-decon successfully identifies the cell type of origin for pancreatic and breast cancer datasets. Interestingly, the results for human pancreatic cancer are recapitulated using the CTS topics inferred from the mouse pancreas scRNA-seq data, postulating cross-species deconvolution as an option where it is difficult to obtain scRNA-seq data due to technical or ethical challenges. Furthermore, our GTM framework also enables the inference of phenotype-specific topic distributions from bulk RNA-seq data by using the phenotypes (e.g., breast cancer subtypes) from single cell or bulk RNA-seq data as a guide for the topic inference. We leverage this capability to distinguish basal from estrogen receptor (ER) positive (ER+) breast tumor samples, not only achieving high classification accuracy but also identifying the genes and pathways that segregate the cancer subtypes. By fine-tuning the inferred CTS topic distributions guided by the breast cancer subtypes, we deconvolve the average differential gene expression into CTS expression changes, thereby enabling discovery of the subtype-specific aberrance of the gene regulatory programs.

## Results

### GTM-decon model overview

In GTM-decon, we have made *three methodological contributions*. As our *first and the main contribution*, GTM-decon is a marker-free method and automatically infers the contribution of each gene for each cell type in the form of CTS categorical distributions, which we define as “topics” [25], without using marker gene information. Each CTS topic distribution is related to the transcriptional rate of each gene for each cell type. For instance, B cells have higher transcription rate for *CD19* compared to alpha cells, which have relatively high rate for *FXYS5*. Conceptually, we consider genes as vocabulary and cells as documents whose word tokens (i.e., scRNA-seq reads) are sampled from the vocabulary with the CTS topic probabilities. We incorporate the observed cell-type labels for each cell in the form of topic prior to guide the inference of CTS topic mixture, which reflects the uncertainty of the noisy cell-type label. Specifically, the cell-type mixture for cell  $m$  follows a K-dimensional asymmetric Dirichlet distribution,  $\theta_m \sim \text{Dir}(\alpha_m + 0.1)$ , with the hyperparameter  $\alpha_{m,k}$  set to a relatively high value (i.e., 0.9 by default) given the cell-type label  $y_m = k$ ; the rest of the K-1  $\alpha_{m,k'}$  values, where  $y_m \neq k'$ , are set to a relatively low values (i.e., randomly sampled from a range between 0.1 and 0.01). As a result, each topic is automatically identified with exactly one cell type. This differs from the standard topic model, where topics are not directly associated with any known concept and require post hoc manual inspection based on their top scoring words to interpret them. Given the CTS gene distributions, we can infer the CTS topic mixtures from the bulk transcriptomes, which are the desired cell-type mixing proportion in the context of deconvolution (Fig. 1a).



**Fig. 1** GTM-decon overview. **a** Inferring cell-type-specific (CTS) topics from scRNA-seq reference data. In brief, GTM-decon infers CTS topics from scRNA-seq data by using a guided topic modeling approach utilizing cell-type labels from the reference. High prior values are assigned to the topic corresponding to the cell type, and lower prior values are assigned to the other topics, enabling it to learn a genes-by-CTS topics matrix, with each topic anchored to a specific cell type. This matrix is used to infer cell-type proportions in bulk RNA-seq data using standard topic modeling, capturing variations in cell type proportions between healthy and diseased tissue. The probabilistic graphical model (PGM) diagram depicts the data generative process assumed by the proposed guided topic model. Suppose there are  $K$  cell types in the scRNA-seq data. For each cell indexed by  $m \in \{1, \dots, M\}$ , we use  $K$ -dimensional Dirichlet-distributed cell-type topic mixture  $\theta_m \sim \text{Dir}(\alpha_m)$  to represent the statistical uncertainty of the noisy cell-type label  $y_m \in \{1, \dots, K\}$ . Specifically, we clamp the Dirichlet hyperparameter  $\alpha_{m,y_m}$  of the Dirichlet variable to a relatively high value while setting the rest of the values of  $\alpha_{m,k}$  ( $y_m \neq k$ ) to relatively low values (i.e., 0.9 and [0.01, 0.1], respectively in the cartoon illustration of  $M=8$  cells and  $K=3$  cell types). The non-zero prior values for the  $K-1$  unobserved cell types allow the cell-type mixture variable  $\theta_m$  to have non-zero density over those cell types as dictated by the scRNA-seq data likelihood and therefore account for potentially mislabeled cell types. Suppose there are in total  $N_m$  reads in cell  $m$ . Each scRNA-seq read  $i \in \{1, \dots, N_m\}$  is assumed to be originated from one of the  $K$  CTS topics with the categorical rates fixed to the cell-type mixture, i.e.,  $z_{i,m} \sim \text{Cat}(\theta_m)$ . Given cell-type topic assignment  $z_{i,m} \in \{1, \dots, K\}$ , the  $i$ th read is then mapped to one of the  $G$  genes as indexed by  $x_{i,m}$  with categorical rates set to be  $\phi_{z_{i,m}}$ , which itself is a  $G$ -dimensional Dirichlet variable of flat hyperparameter  $\beta$ , i.e.,  $x_{i,m} \sim \text{Cat}(\phi_{z_{i,m}})$ . To infer the latent variables, namely cell-type mixture proportion  $\theta_m \sim \text{Dir}(\alpha)$ , CTS topic assignments for each read  $z_{i,m}$ , and CTS topic distributions  $\Phi$ , we employ an efficient collapsed variational Bayes algorithm as detailed in the “Methods” section. The genes-by-CTS-topic  $\hat{\Phi}$  matrix estimated from the scRNA-seq reference then serves as a template when it comes to infer the cell-type mixing proportions  $\theta_j$  of a bulk RNA-seq sample  $j$  using essentially the same inference algorithm as in the scRNA-seq data modeling except for having a flat hyperparameter for the prior (e.g.,  $\alpha_k = 1 \forall k$  by default) while fixing  $\hat{\Phi}$  and only inferring the expected total reads allocated for each CTS topics (i.e.,  $E_q[n_{j,k}] = E_q[\sum_i [z_{ij} = k]]$ ). **b** Phenotype-guided modeling of bulk RNA-seq data. GTM-decon can also use phenotype labels as a guide for topic inference to model sparsified bulk transcriptomes in a disease study. In this design, instead of having each row as a cell and each column as a cell type, each row corresponds to a bulk sample and each column to a phenotype class. For each subject  $j$ , we set the topic hyperparameter  $\alpha_{j,y}$  based on the noisy phenotype label  $y_j$  of the subject. The inference algorithm is the same as in modeling the scRNA-seq reference data. Given a test subject  $j$ , the inferred topic mixture  $\theta_j$ , represents the phenotypic probabilities of the subject. **c** Nested-guided topic model for detecting cell-type-specific differentially expressed genes between phenotypes. In this nested design, we treat the phenotype as level 1 and the cell types as level 2. The pretrained genes-by-CTS-topic distribution  $\hat{\Phi}$  learned from panel **a** are used to initialize the topic distributions for each phenotype in a sparsified bulk transcriptome disease study. As illustrated in the cartoon, for example, for 2 phenotypes and 3 cell types, there are 6 topics. GTM-decon then fine-tunes the combined cell-type-specific topic distribution by running the same algorithm described in panel **b**. The resulting topic distributions reflect the phenotypic influences on CTS gene distributions, which are the statistics for conducting differential expression analysis in a case–control study design

As our *second contribution*, we extend GTM-decon to infer multiple topics per cell type. The rationale is that cells of the same cell type can manifest in different cell states due to the changes of environments or stimuli. As a result, these cells may exhibit expression patterns that are different from the canonical CTS expression pattern. While there are sophisticated hierarchical topic models involving Dirichlet Processes [26], we took a simple and elegant design. Specifically, we extend the basic GTM-decon model to infer sub-cell-type topics by dedicating multiple topics per cell type (Additional file 1: Fig. S1). As our *third contribution*, we extend GTM-decon to infer phenotype-specific (PTS) topic distributions using the phenotype label (e.g., cancer subtypes or cancer stages) available in the single-cell or bulk transcriptome data as a guide to detect PTS gene expression (Fig. 1b). We then further extend it to a nested-guided topic model to conduct CTS differential expression analysis in the single-cell or bulk patient cohort data (Fig. 1c). To that end, we use the phenotype labels as the level-1 guide and the cell-type labels as the level 2 guide. Through the same guided topic mechanism, GTM-decon updates the CTS topic distributions under each phenotype by fitting the data likelihood of the transcriptomes from either the single-cell or bulk data. The algorithmic details for the 3 contributions were described in “[Methods](#).”

#### **Experimentation of data preprocessing and GTM-decon model configurations**

We experimented gene selection, data normalization, hyperparameter settings, and number of topics per cell type. We find that GTM-decon works the best with raw read count data using all genes (Additional file 1: Figs. S2 and S3), and it is fairly robust to different hyperparameter values for the topic mixture prior (Figs. S4–S6) and the CTS topic prior values (Fig. S7). In general, GTM-decon confers better deconvolution performance using multiple topics per cell type than the baseline GTM-decon with one topic per cell type (Fig. S8). Please refer to Additional file 1 Section S1–S5 for more details.

#### **Evaluation of deconvolution of simulated bulk from scRNA-seq data**

To quantitatively benchmark GTM-decon, we compared it against five existing deconvolution methods, namely Bisque [21], Bseq-sc [17], CIBERSORTx [18], MuSiC [19], and BayesPrism [22] on artificially simulated bulk RNA-seq datasets from the scRNA-seq data (Additional file 1: Table S1). To simulate bulk data, three human scRNA-seq datasets (Pancreas—E-MTAB-5061, Breast Tissue with GEO accession number GSE113197, and Rheumatoid Arthritis (RA) Synovium—SDY998), generated using different technologies (Smart-seq2, 10 × Genomics Chromium, CEL-Seq2) were used. The artificial bulk data for each individual was constructed by summing up the counts for each gene from all cells in that individual [19]. This allowed us to use the cell-type proportions from the single-cell data as the ground-truth proportions. Artificially constructed bulk data from scRNA-seq data appear to be a good surrogate of the real bulk data, as observed from the excellent correlation of the log-transformed artificial counts with the log-transformed counts from real bulk data for each gene (Additional file 1: Fig. S9). We performed leave-one-out cross-validation (LOOCV), to avoid any leakage from training data, and used the single-cell RNA-seq of the left-out individual to simulate the bulk RNA-seq (i.e., the total read counts of each gene for that sample) as the validation data and its cell-type proportions as the ground-truth mixing proportions. GTM-decon performs better than

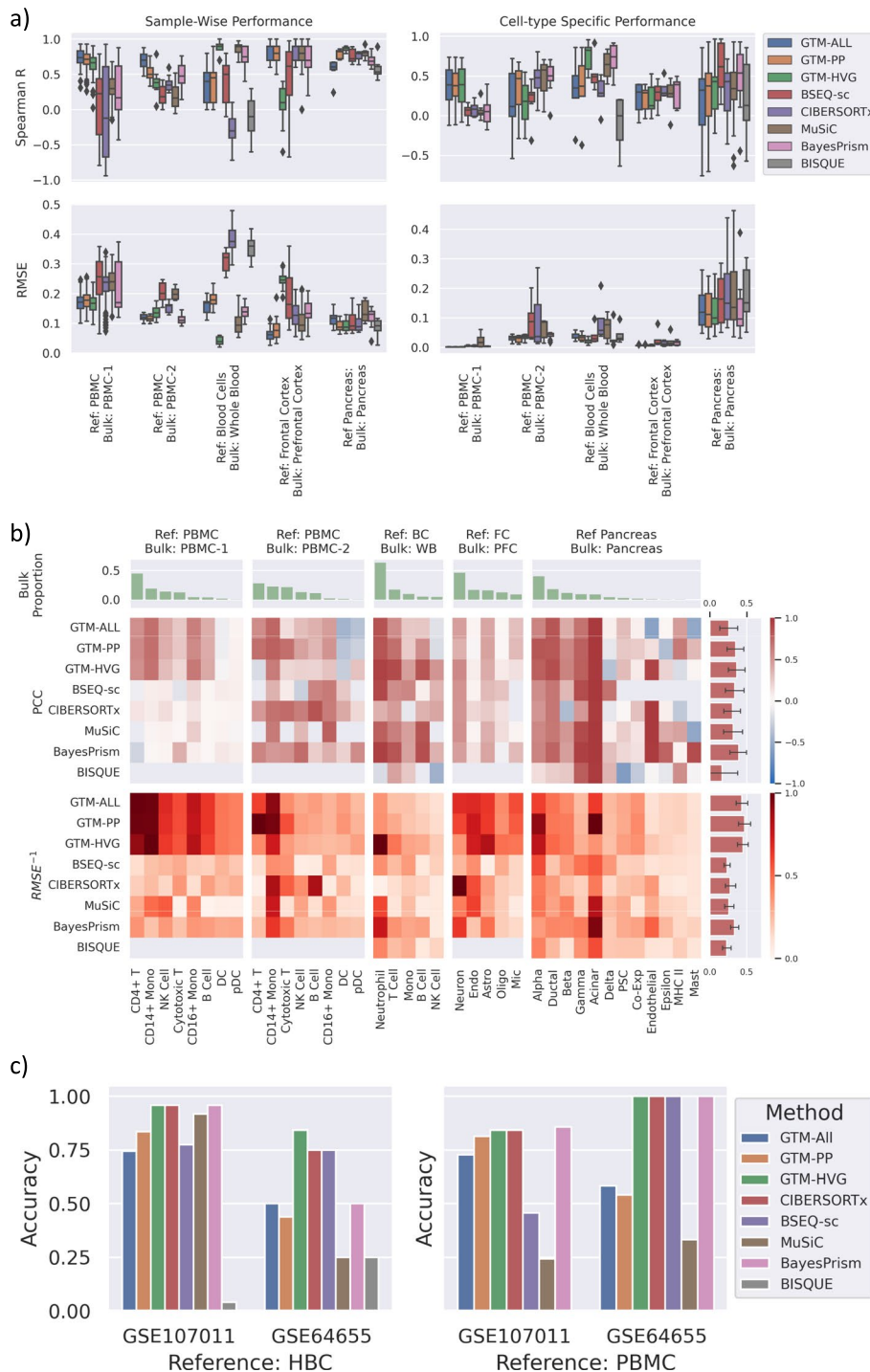
other models on the Pancreas and Breast Tissue datasets in terms of both Spearman Rank-based correlation (Spearman R) and Cross Entropy and conferred comparable performance on the RA Synovium dataset (Additional file 1: Fig. S10). Notably, GTM-decon achieves smaller variance for both the pancreatic and breast tissue datasets. We further ascertained the qualities of the predicted cell-type proportions of each method against the ground-truth cell-type proportions (Fig. S10a) and observed that GTM-decon recapitulates cell-type proportions well for pancreatic (Fig. S10b), breast tissues (Fig. S10c), and RA synovium dataset (Fig. S10d).

### Evaluation of deconvolution of real bulk with ground-truth cell-type proportions

We also benchmarked GTM-decon on 5 real bulk RNA-seq data with known ground-truth cell-type proportions from 3 different tissue types (Additional file 1: Table S1 and S2; “Methods”). We evaluated deconvolution performance using Spearman R and Root Mean Square Error (RMSE), by comparing the inferred proportions for the matching cell types against the ground-truth proportions for each sample. GTM-decon conferred on-par or superior performance compared to the existing methods for all the datasets (Fig. 2a left). In particular, GTM-ALL performed the best for deconvolving PBMC-1 and PBMC-2. For deconvolving whole blood (WB), GTM-HVG performed the best in terms of both metrics and MuSiC is a close second. While deconvolving bulk RNA-seq from the prefrontal brain region from ROSMAP dataset, all methods except Bseq-sc and GTM-HVG performed reasonably well. It is possible that the HVG are the genes having high variance within the cell types, which caused the poor performance of GTM-HVG on this dataset. For the pancreatic dataset with the paired single-cell and bulk RNA-seq data collected from the same individuals, GTM-HVG performs the best with GTM-PP as the runner-up in terms of Spearman

(See figure on next page.)

**Fig. 2** Evaluation of deconvolution performance on real bulk data. **a** Evaluation of sample deconvolution. We evaluated the deconvolution performance of GTM-decon using all genes (GTM-ALL), preprocessed genes (GTM-PP), and highly variable genes (GTM-HVG) with five SOTA methods—CIBERSORTx, MUSIC, BSEQ-sc, BISQUE, BayesPrism. The 3 immune bulk data and the brain data were deconvolved using independent references of a similar tissue, while the pancreas bulk data is deconvolved using single-cell reference from the same individuals in a leave-one-out cross-validation (LOOCV) manner. The bulk labeled PBMC-1 corresponds to SDY67 dataset, PBMC-2 corresponds to S13 cohort, whole blood to whole blood dataset, and prefrontal cortex to ROSMAP dataset (Additional file 1: Table S2). For each test bulk sample, Spearman correlation and root mean square error (RMSE) were computed between its ground truth and predicted cell-type proportions by each method. The box and whiskers in each boxplot indicate the 25–75% quartile and min–max of the evaluation scores over all samples in a dataset, respectively. The boxplot on the left displays the evaluation across cell types per sample, and the boxplot on the right displays the evaluation across samples per cell type. **b** Heatmaps comparing the cell-type-specific deconvolution performance of GTM-decon against existing methods on 5 different real bulk datasets with known ground truth mixing proportions. The cell types are ordered from most to least prevalent in the bulk data (green barplots in first row indicate average proportion for each cell type in the bulk data). The middle row shows the Pearson correlation coefficient between the predicted and known cell-type proportions. The lower row shows the inverse RMSE (higher is better, scaled between 0 and 1), per cell type per dataset. The barplots on the right show the average performance over all cell types for each method. For each cell type, Pearson correlation and RMSE were computed between its ground truth and predicted cell-type proportion for each dataset by each method. **c** Cell-type prediction accuracy of the purified immune bulk RNA-seq samples. The two panels indicate the use of different independent immune references, for the deconvolution of two purified bulk immune datasets (Accession Numbers: GSE107011, GSE64655). For each purified bulk sample, the cell type corresponding to the highest inferred cell-type proportion by each method was used as the predicted cell type. The barplots show the prediction accuracy as the percentage of the correctly predicted samples



**Fig. 2** (See legend on previous page.)

R, and RMSE are similar among all methods except MuSiC with notably higher error. In summary, GTM-ALL performed the best in 3 datasets; GTM-HVG performed the best in the other two datasets, where the CTS gene expression might exhibit more distinct inter-cell-type variability. Furthermore, we also evaluated the deconvolution

performance for separate cell types in terms of the correlation between ground-truth proportions and predicted proportions for each cell type across samples. Overall, GTM-decon conferred competitive performance with the runner-up method being BayesPrism (Fig. 2a right and b; Additional File 1: Figs. S11-S15). Therefore, the results suggest the general effectiveness of topic modeling in cell-type deconvolution and the additional benefits conferred by GTM-decon due to its algorithmic innovations.

Additionally, we evaluated the deconvolution accuracy for purified bulk RNA-seq data of immune cells (GEO accession numbers: GSE107011 and GSE64655; Table S2) using two different independent references (HBC and PBMC2; Table S1). Overlapping cell types were used to evaluate the performance on the purified bulk samples, whereby the highest deconvolution proportion was used as the predicted cell-type label for computing the prediction accuracy. GTM-HVG achieved the highest accuracy across all four experiments (Fig. 2c). Moreover, since some cell types present in the bulk are missing in the scRNA-seq reference, a robust model should either find the closest-matched cell types or properly express statistical uncertainty in this situation. To this end, we examined the inferred cell proportions of the purified bulk for those missing cell types. Plasmablast samples are present among the purified PBMC samples (GSE107011) but absent in the HBC reference (Additional file 1: Fig. S16a). GTM-decon assigned Plasmablast samples with high cell-type proportions for B-cell, a cell type that shares Plasmablast cell lineage. The inferred cell-type proportions for Basophils purified samples were split between HSPCs (immune progenitor cells) and neutrophils, which is also classified as granulocyte. Granulocytes are the most common white blood cells, consisting of 3 specific cell types—neutrophils, eosinophils, and basophils. Using PBMC2 as a reference to deconvolve the same purified bulk immune samples led to similar deconvolution patterns for the missing cell types of plasmablast and basophils (Additional file 1: Fig. S16b). Interestingly, neutrophils were absent in the PMBC2 reference and inferred to be monocytes, which are related to the granulocyte family—a class of immune cells that include basophils, eosinophils, and neutrophils [27]. Another missing cell type HSPC (hematopoietic stem and progenitor cells) have their signal spread across all the cell types. These results suggest that when some cell types are missing in the reference, GTM-decon either finds the closest match or appropriately expresses uncertainty.

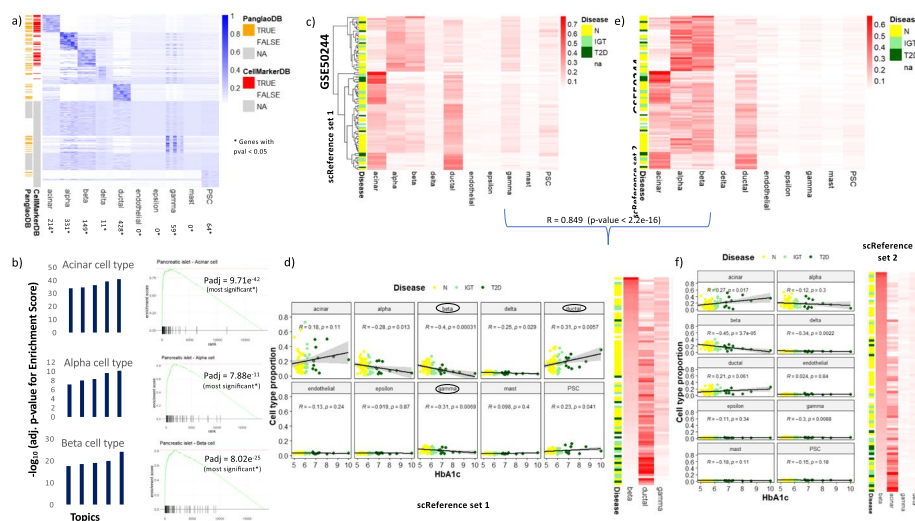
Finally, we performed benchmarking on the time and memory usage of our GTM-decon software. GTM-decon scales linearly with both the number of topics per cell type and the number of cells (Additional file 1: Section S6; Fig. S17), which is what we expected since its time and space complexity are both  $O(N \times G \times K)$  for  $N$  cells,  $G$  genes, and  $K$  topics. For large number of cells, we can also perform stochastic variational inference [28] to rapidly update model parameters based on mini-batches of cells with much lower memory overhead. It also compares favorably with BISQUE, BSEQ-sc, and MuSiC in terms of running time and memory usage.

#### **GTM-decon automatically learns CTS gene signatures from scRNA-seq reference**

We evaluated the performance of GTM-decon in recapitulating cell-type-specific information as well as deconvolution using pancreatic tissue as a reference. The pancreas consists of several cell types including exocrine and endocrine. While the former aids digestion by secreting several enzymes, the latter regulates glucose uptake and



processing by secreting hormones. With a vast literature documenting the biological roles of several cell types and their behavior in healthy and diseased conditions, such as diabetes and cancer, the pancreatic tissue serves as a good benchmark to assess GTM-decon. We trained GTM-decon on an scRNA-seq reference dataset of pancreatic tissue from Segerstolpe et al. [29], consisting of 2209 cells, corresponding to 14 cell types from 10 individuals. GTM-decon captured distinct sets of CTS gene signatures, as shown by the gene-by-topic probability distributions (i.e., the matrix  $\Phi$ ) for the top 20 genes in each topic (Fig. 3a). Indeed, each topic recovers a large number of marker genes for the corresponding cell types based on two databases, namely CellMarker database [30], a manually curated resource of cell markers in human and mouse and the PanglaoDB



**Fig. 3** Cell-type-specific topic inference and deconvolution of pancreatic tissue. **a** Gene signatures of cell-type-specific topics in pancreas. We trained GTM-decon on the PI-Segerstolpe scRNA-seq dataset of pancreas tissue. We used 5 topics per cell type, allowing sub-cell-type inference. For the inferred genes-by-cell-type matrix  $\Phi$ , we took the top 20 genes under each topic and visualized their topic distributions in heatmap. Whenever available from CellMarkerDB and PanglaoDB, cell-type marker genes are indicated on the left. For the cell types where marker genes are not available, “NA” were indicated on the left. The number of statistically significantly different genes in each cell type based on their topic scores ( $p$ -value < 0.05; permutation test) is shown below. **b** Gene set enrichment analysis (GSEA) of inferred topics based on known marker genes. Cell-type-specific topics for acinar, alpha, and beta were evaluated based on whether the top genes are enriched for the known marker genes under that cell type. The bar plots show the  $-\log_{10}(p.\text{adj})$  values of enrichment score for each of the 5 topics. The leading-edge plot for the topic with the best adj.  $p$ -value for that cell type is shown on the right. In each of the leading-edge plots, genes were ordered in decreasing order from left to right. The green curves indicate the running scores of enrichments. The barcode bars indicate cell-type marker genes. Adjusted  $p$ -values based on the GSEA enrichment scores are indicated in each panel. The three large panels display the most significantly enriched topic of among the five topics for each cell type and the 12 small panels display the remaining topics. **c, e** Deconvolution of bulk RNA-seq samples of 89 human pancreatic islet donors. The GTM-decon models separately trained on the Segerstolpe pancreas islet dataset (i.e., panel **c**) and Baron pancreas islet data (i.e., panel **e**) reference datasets were used to deconvolve the 89 bulk transcriptomes. As indicated by the legend, the 89 subjects consist of 51 normal, 15 impaired glucose tolerance, and 12 T2D individuals. In the heatmap, the rows represent subjects, and the columns represent cell types; the color intensity are proportional to the inferred cell-type proportions. **d, f** Deconvolved cell-type proportions as a function of Hemoglobin A1c (HbA1c) level. GTM-decon were trained on Segerstolpe (i.e., panel **d**) and Baron scRNA-seq (i.e., panel **f**) reference datasets. Each of the 10 panels displays a scatter plot of inferred cell-type proportion ( $y$ -axis) and HbA1c level ( $x$ -axis). The color legend indicates the 3 phenotypes. The heatmap on the right shows the deconvolved proportion of 3 most indicative cell types with subjects (rows) ordered on the basis of inferred beta cell-type proportions

[31], a database of marker genes generated from scRNA-seq datasets (Fig. 3a). We further ascertained the cell-type coherence of each topic by Gene Set Enrichment Analysis (GSEA), while using the probabilities learnt for each cell type against the CellMarkerDB. For the three cell types with abundant marker genes—acinar, alpha, and beta, each of the 5 topics recovers the exact cell type as the top-most hit in the analysis, with the adjusted  $p$ -value  $\leq 1 \times 10^{-15}$  (permutation test) in most cases (Fig. 3b). Furthermore, the enrichment of known marker genes for the main cell types suggested that GTM-decon with 5 topics per cell type best captures the cell-type-specific signatures (Additional file 1: Fig. S18). We also evaluated the effect of different number of cells per cell type. As expected, the topic confidence scores as measured by the average probabilities over the CTS gene distributions increase with the increasing number of cells for that cell type (i.e., evidence) (Additional file 1: Fig. S19).

#### **GTM-decon delineates the variations of cell-type proportions in pancreatic tissues of healthy and T2D subjects**

Based on the inferred CTS topic distributions  $\phi$ , we used GTM-decon with 5 topics per cell type to infer the cell-type proportions of bulk RNA-seq data from a cohort of 89 human pancreatic islet donors with and without type 2 diabetes (GEO accession number: GSE50244) [32] (Fig. 3c). The dataset consists of 51 individuals with normal glucose tolerance (N), 15 with impaired glucose tolerance (IGT) and 12 with type 2 diabetes (T2D); it also has a good segregation of males ( $N_1 = 54$ ) and females ( $N_2 = 35$ ). As expected, the inferred proportions of most types of cells for these two sets of individuals are similar since they came from the same tissue. However, GTM-decon predicted a significant reduction in beta cells in T2D individuals (Pearson correlation coefficient (PCC) with HbA1c =  $-0.4$ ,  $p$ -value =  $0.00031$ ; t-distribution with  $n-2$  degrees of freedom) (Fig. 3d). IGT and T2D individuals exhibit low Beta cell-type proportions (Fig. 3d), which was supported by the literature [1, 2]. The increase in ductal cells is possibly caused by their regulation of glucose uptake (Fig. 3d). These results are consistent when using a different pancreatic dataset as reference (i.e., the PI-Baron reference dataset, generated via Drop-seq instead of PI-Segerstolpe, generated via Smart-seq) (Fig. 3e, f).

#### **Deconvolving human pancreatic data from mouse pancreas scRNA-seq reference**

For cases where scRNA-seq data were not available for the organism of interest, due to either ethical, technical, or financial challenges, there is a need to leverage scRNA-seq collected from a model organism. To this end, we investigated the possibility of deconvolving bulk RNA-seq data by training GTM-decon on a mouse scRNA-seq data. We separately trained two GTM-decon models on the human and mouse pancreatic datasets. Specifically, the human datasets include the Segerstolpe and Baron datasets (Additional file 1: Table S1), which consist of 2209 cells from 10 individuals and 8569 cells from 4 individuals, respectively; the mouse dataset consists of 1886 cells collected from 2 mice. For comparative analysis, we focused on only the common set of high-confidence orthologous genes between the two species mapped by the Ensembl database [33]. We visualized the cell-type proportions comparing against the ground-truth values (Additional file 1: Fig. S20a). As expected, the cell-type proportions deconvolved using the two human datasets accurately recapitulate the ground-truth proportions

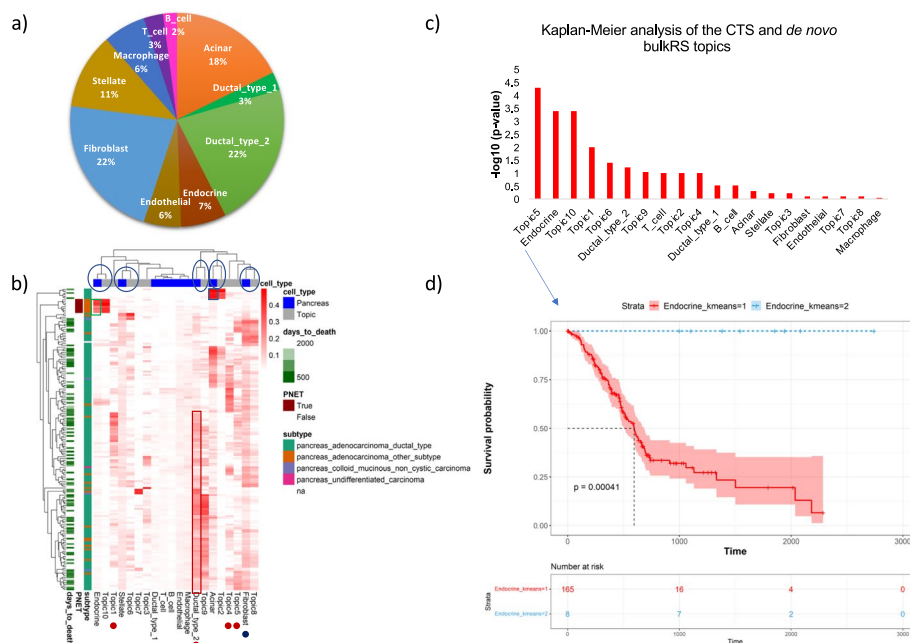
(median PCC of 0.94 and 0.97). Interestingly, GTM-decon trained on the non-reference human dataset performed better than the one trained on the reference-matched dataset, which was probably due to the tenfold higher number of cells in the former scRNA-seq dataset. Moreover, GTM-decon trained on the mouse reference dataset also performed quite well in terms of the concordant proportions of the shared cell types (PCC of 0.94).

### Deconvolving pancreatic cancer transcriptomes identified tumor cell-type origin

We next turned to deconvolving pancreatic adenocarcinoma (PAAD) (also known as pancreatic ductal adenocarcinoma or PDAC) tumor bulk samples from TCGA. For this application, we also used GTM-decon with 5 topics per cell type. Since the tumor microenvironment is known to be infiltrated with immune cells [34, 35], we sought to train GTM-decon on a single-cell reference dataset derived from individuals with pancreatic cancer, in order to capture the cell types of both the tissue of interest and the immune cells in its tumor microenvironment. To this end, we trained GTM-decon on an scRNA-seq dataset comprised of the transcriptomic profiles of about 57,000 cells from 24 primary PAAD tumors and 11 healthy control pancreas samples [36] in order to deconvolve the 174 bulk RNA-seq profiles from the TCGA-PAAD tumor samples. Additionally, we also sought to identify possible novel cell types or pathways present in the bulk RNA-seq, which are not represented in the reference profiles. This is achieved by running an unguided topic model (i.e., a standard LDA) on the sparsified bulk RNA-seq data to detect de novo bulk RNA-seq (bulkRS) topics (“Methods”). We empirically chose the number of de novo bulkRS topics based on how well they could explain the variation observed in the clinical phenotypes.

We observe that the most prevalent cell types are 4 main pancreatic cell types, namely ductal (type 2), acinar, endocrine (alpha and beta cell types), and fibroblasts (Fig. 4a). Notably, the cell type of tumor origin is correctly predicted for the samples: Ductal cells have the highest proportion among the PAAD samples (Fig. 4b; brown rectangle), and acinar for a subset of the PAAD samples (Fig. 4b; blue rectangle). This recapitulates the literature remarkably well, as ductal cells are known to be the site of the tumor origin for most cases of PAAD; however, a subset arises from acinar cells [37]. More significantly, an unknown subtype is predicted to originate from endocrine cells (Fig. 4b; green rectangle). This is supported by the recent literature, which reported these samples as the derivatives of the pancreatic neuroendocrine tumor (PNET) were in fact misclassified as PAAD [38]. PNET are supposed to originate from alpha or beta cells (endocrine cells) [39]. In addition, most of the samples are predicted to have a high proportion of fibroblasts, which are known to be prevalent in pancreatic cancer (Fig. 4b; blue dot) [40–42]. Interestingly, deconvolving the pancreatic cancer PAAD dataset using mouse reference dataset also conferred high-quality patient clustering comparable to that of human dataset (Additional file 1: Fig. S20b). Notably, the PAAD “other subtype” was predicted to originate from alpha cells (an endocrine cell), which mirrors the results from human reference data (Fig. S20b; shown in blue rectangles).

Interestingly, the de novo bulkRS topics cluster with the most abundant reference topics inferred from the scRNA-seq reference data (Fig. 4b, blue circles). For example, bulkRS Topic 2 corresponds to the CTS topic for acinar cell type, Topic 10 to the



**Fig. 4** Deconvolution of bulk RNA-seq samples for pancreatic cancer from TCGA-PAAD. GTM-decon was first trained on an scRNA-seq dataset from individuals with pancreatic cancer. The trained GTM-decon model was then used to infer the cell-type proportions of the 174 TCGA-PAAD bulk RNA-seq profiles. **a** The average inferred cell-type proportion across the TCGA-PAAD tumor samples. We summed up inferred cell-type proportions over all samples followed by normalization. The pie chart displays the resulting percentage of cell-type proportions. **b** Inferred cell-type proportions of individual TCGA-PAAD tumor samples. To complement the inferred proportions of known cell types, we also ran unguided topic model (i.e., LDA) on the TCGA-PAAD bulk RNA-seq profiles directly to detect de novo topics that are not present in scRNA-seq reference. The heatmap visualizes the combined deconvolution results based on the 10 pancreatic cell types, and 10 de novo topics (i.e., columns). Each of the 174 rows represents a subject. Three types of demographic or clinical phenotypes were shown in the legend to aid result interpretation. These include “days to death,” cancer subtype, and whether the cancer type is PNET or not. The regions in the highlighted boxes were discussed in more details in the main text. **c** Survival analysis of the CTS and de novo topics using inferred cell-type proportions. The 174 subjects were divided into two groups based on K\*-means clustering with K\* set to 2 (not to be confused with the K cell types or topics). Kaplan–Meier curves were generated for these groups and compared using log-rank test. The plot shows the  $-\log_{10}(p\text{-value})$  from the log-rank test for all the CTS and de novo topics, in decreasing order of significance. **d** Kaplan–Meier curve for endocrine cell-type proportion. The curve and shaded area represent the mean and standard deviation of the cell-type proportions in the two groups, respectively. The number of subjects for each cluster was indicated in the bottom panel

CTS topic for endocrine cell type, Topic 9 to Ductal\_type\_2, Topic 8 to Fibroblast, and Topic 6 to stellate cells. However, there are a few de novo bulkRS topics, such as bulkRS Topic 1, Topic 4, and Topic 5, capturing distinct distributions for specific subsets of samples (Fig. 4b; red dots). These topics could correspond to either novel cell types or gene pathways in the bulk but not implicated in the scRNA-seq reference.

We next estimated whether variation in cell-type proportions or bulkRS topics is indicative of survival time. To this end, we performed Cox Regression to regress the number of days patients lived since their cancer diagnosis on their inferred cell-type proportions as well as the de novo bulkRS topics. Overall, the Cox Regression model is statistically significant compared to the bias term (adjusted  $p\text{-value} = 9 \times 10^{-5}$  based on likelihood ratio test). To explore the marginal effect of individual cell-type

proportions on survival, we performed Kaplan–Meier analysis by separating patients into two groups based on K-means clustering (Fig. 4c). We observe that among the cell types, endocrine cell type exhibits a significant hazard ratio, predicting a good survival outcome (Fig. 4d;  $p$ -value = 0.0091; log-rank test), which is supported by the literature as PNETs are mostly benign [43]. However, ductal cell type 2 is associated with poor survival outcome, which is expected as pancreatic adenocarcinoma are aggressive ( $p$ -value = 0.03; log-rank test). On the other hand, Topic 1 and Topic 5 indicate poorer survival, with hazard ratios of 4 and 40 respectively (Additional file 1: Fig. S21) although their roles are unclear since they do not cluster with any of the CTS topics (Fig. 4b). This reiterates the usefulness of using CTS topics in conjunction with de novo bulkRS topics to enable their interpretation wherever possible.

A similar analysis conducted using a separate scRNA-seq reference from healthy pancreatic subjects only reveals similar results (Additional file 1: Fig. S22). However, in this analysis, only two of the de novo bulkRS topics clearly correspond to CTS topics (acinar and ductal), suggesting that the usage of an appropriate scRNA-seq reference with matched tumor environment is preferred.

#### Deconvolving breast cancer transcriptomes revealed subtype-specific markers

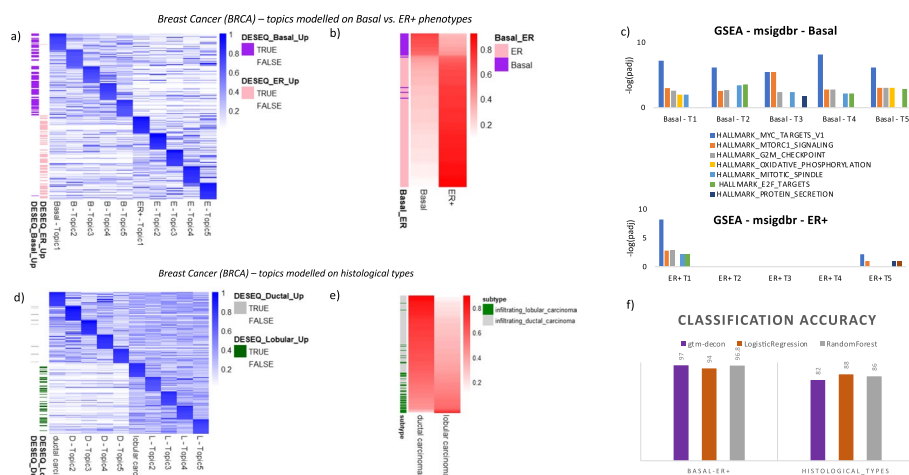
To capture specific subtypes of breast tumor samples from the TCGA data, we trained GTM-decon on an scRNA-seq reference data from 26 primary tumors of breast cancer (BRCA) patients with three major clinical subtypes of BRCA, including 11 ER+, 5 HER2+ and 10 TNBC [44]. The data consists of about 1 million cells and covers 7 major cell types and 29 minor cell types. This served as a high-resolution reference for annotating the TCGA-BRCA tumor samples ( $n = 1212$ ) based on the major subtypes from the scRNA-seq dataset, namely Cancer-Basal, Cancer-Her2, Cancer-LumA, and Cancer-LumB (Additional file 1: Fig. S23a). As expected, significantly higher proportions of endothelial and myoepithelial cell types are found in the normal-like samples, in comparison to the cancer subtypes (Fig. S23a). Furthermore, basal subtype is enriched for cancer-basal cells (Fig. S23b, shown in dotted brown rectangle), and the cancer-associated fibroblasts (CAFs) are enriched in almost all the samples (Fig. S23b, shown in blue rectangles). Similar to the above analyses, some of the de novo bulkRS topics from sparsified samples overlap with the most represented cell types. For example, Topic 6 resembles myofibroblast-like cancer-associated fibroblasts (myCAF-like), and Topic 8 resembles LumA subtype cancer cells.

Deconvolution using the scRNA-seq reference from the healthy individuals also captures the cell type of origin for the different breast cancer subtypes implicated in the bulk samples (Additional file 1: Fig. S24). For this analysis, using highly variable genes is more discriminatory than all genes (Additional file 1: Fig. S25). The inferred CTS topic distributions from normal breast tissue recapitulate several marker genes from CellmarkerDB and PanglaoDB (Fig. S24a). Furthermore, GTM-decon-inferred cell-type proportions clearly distinguish TCGA breast tumor samples from GTEx normal breast tissues (Fig. S24b). Moreover, the basal cell type is predicted to have the highest proportion in the basal subtype defined by the PAM50 classification in comparison to other subtypes (Fig. S24c) [45]. Also, we observed higher predicted proportion for Luminal\_2 cell type in both LumA and LumB subtypes as expected [46] (Fig. S24c). Among the de novo bulkRS

topics from sparsified samples, Topic 5 is highly enriched in LumB, while Topic 7 and Topic 8 are enriched in basal and LumA subtypes (Fig. S24c). The guided topic score for the basal subtype also correlates with higher proliferation score as expected [45]. Specifically, the basal cell type is enriched in this subtype (Fig. S24d; enclosed in green rectangle), whereas Luminal\_2 is depleted (Fig. S24d; enclosed in green dashed rectangle). In contrast, ER + samples appear to be enriched for Luminal\_2 cells (Wilcoxon test  $p$ -value =  $4.5e - 10$ ; Fig. S26), as expected [46, 47]. Furthermore, Topic 7 clearly captures the basal subtype (enclosed in brown rectangle in Fig. S24d), whereas there is no clear topic capturing the ER + phenotype.

**GTM-decon learns phenotype-guided gene topics specific to BRCA subtypes**

Our guided topic mechanism is not limited to inferring CTS topics, but can be extended to inferring phenotype-specific topics (i.e., topics capturing phenotype-specific gene expression) (Fig. 1b), thereby discovering gene signatures of subtypes or different cancer stages. We applied this approach to study the differences between basal and ER + BRCA subtypes (Fig. 5a–c) and the difference between ductal carcinoma and lobular carcinoma



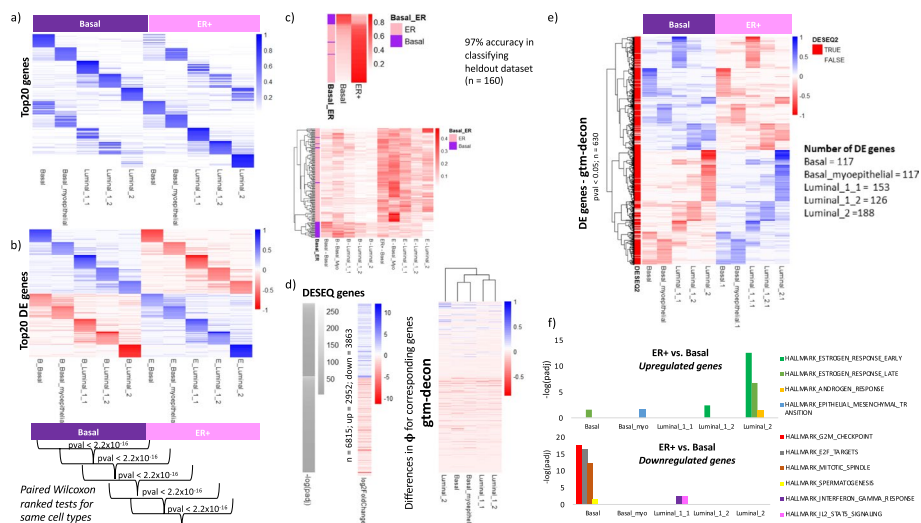
**Fig. 5** Phenotype-guided topic modeling of bulk RNA-seq data of breast cancer. **a** Predicted top genes for the phenotype-guided topics for basal and ER + breast cancer subtypes. GTM-decon was trained on the sparsified TCGA-BRCA bulk RNA-seq data with the basal and ER + cancer subtypes as the guide. Five topics were used per subtype and therefore 10 topics in total. The heatmap illustrates the topic probabilities of the top 20 genes from each topic. As a comparison, the genes were also labeled as up- or downregulated if they were deemed differentially expressed by the DESeq2 analysis. **b** Classification of basal and ER + subtypes based on phenotype-guided topic scores. The 5 topics for the same subtype were summed to obtain the overall score for basal and ER + subtype. Subjects in the rows were sorted by their basal topic scores. **c** GSEA analysis of the basal and ER + subtype topics. Significantly enriched MSigDb HALLMARK pathways were identified for each topic and displayed as barplots. The heights of the bar indicate the  $-\log_{10}$  adjusted  $p$ -values and the colors indicate enriched pathways. **d** Predicted top genes for the phenotype-guided topics for histological subtypes. Same as in panel **a** but for ductal and lobular subtypes. **e** Classification of histological subtypes. Same as in panel **b** but for ductal and lobular subtypes. **f** Evaluation of the subtype classification accuracy on the test breast tumor samples. We trained the phenotype-guided GTM-decon separately on 80% of the sparsified TCGA-BRCA tumor samples using basal/ER + and histological types as the guides and evaluated its phenotype prediction accuracy on the 20% held sparsified samples. As a comparison, we also trained and evaluated logistic regression and random forest on the same training and test split, respectively. The classification accuracy on the test set by each method were displayed in the barplots

(Fig. 5d–f). We modeled each phenotype using 5 topics each, based on a sparsified matrix of bulk RNA-seq data (“Methods”), resulting in a genes-by-phenotypes matrix  $\phi$  with 5 topics per phenotype (Fig. 5a). We then ranked genes by the topic scores under each topic. Almost all the genes identified by our approach were also deemed as differentially expressed (DE) genes by DESeq2 differential analysis [48] (Fig. 5a). GSEA of the topics shows differences between the basal and ER+ phenotypes, although there is not much difference among the 5 topics for basal (Fig. 5c). Moreover, the trained GTM-decon confers accurate phenotype classification with 97% accuracy for discriminating basal and ER+ (Fig. 5b) and 82% for discriminating ductal and lobular subtypes (Fig. 5d), which is comparable to the traditional supervised learning methods namely logistic regression and random forest (Fig. 5f). We also evaluated phenotype classification accuracy as a function of sparsification rate (described in Additional file 1 Section S7 and illustrated in Additional file 1: Fig. S27) suggesting the importance of sparsification for inferring topics from bulk RNA-seq data using GTM-decon. We also trained GTM-decon only on the highly variable genes in the unsparified bulk RNA-seq samples. This results in lower classification accuracy (76% for ductal-lobular samples as compared to 82%, and 95% for basal-ER+ samples as compared to 97%), which might be due to information loss (HVG = 1391 for basal-ER+ samples, HVG = 445 for ductal-lobular samples). Training using DE genes between the two phenotypes identified by DESeq2 also results in lower classification accuracy of 83% and 75%, respectively. These experiments suggest that GTM-decon can utilize more informative genes to discriminate the breast cancer subtypes than the traditional differential analysis approach.

#### **Nested-guided topic modeling identifies CTS DE genes in breast cancer subtypes**

scRNA-seq can facilitate molecular understanding of genes and pathways in specific cell types with respect to phenotypic states. This level of detail is absent in bulk RNA-seq data, which profiles only the averaged gene expression from all cell types in the tissue. However, due to the cost, scRNA-seq profiles at the patient cohort size are rare. To take advantage of both types of data, we sought a way to identify cell-type-specific gene expression differences corresponding to the phenotypes observed for the bulk samples. Briefly, we took a pretrained GTM-decon on a scRNA-seq reference data and then updated its CTS topics based on the corresponding phenotypes from the sparsified bulk data (Fig. 1c). This is equivalent to treating the phenotype as level 1 and cell types as level 2 in a two-stage nested factor design in statistics.

For the TCGA-BRCA data, in particular, we first initialized the genes-by-topic matrix  $\phi$  with the *pretrained* guided topics learned from the scRNA-seq reference for normal breast tissue. We then *fine-tuned* 5 CTS topics for each phenotype (i.e., ER+ or basal) from sparsified bulk data, resulting in a 10-topic model. During the fine-tuning, all CTS topics corresponding to the patient’s phenotype are assigned a prior value of 0.9. Applying this approach to 798 sparsified samples from TCGA-BRCA, corresponding to basal ( $N_1 = 140$ ) and ER+ ( $N_2 = 658$ ) led to a new genes-by-topic matrix  $\phi^*$ . First of all, the top genes for the topics corresponding to the same cell type are similar in both phenotypes, as they are expected to capture CTS signatures (Fig. 6a). However, differential analysis of genes between two phenotypes revealed upregulated genes in one phenotype being downregulated in the other (Fig. 6b). These differences between basal and ER+ are



**Fig. 6** Identification of cell-type-specific differentially expressed genes from bulk RNA-seq data for basal vs. ER+ subtypes. **a** Top cell-type-specific gene signatures for basal and ER+. GTM-decon was pretrained on a scRNA-seq reference dataset from normal breast tissue to infer the expression distribution of 5 cell types, namely basal, basal myoepithelial, luminal 1–1, luminal 1–2, and luminal 2. The resulting genes-by-cell-type estimates were then used as the initial topic distributions for another GTM-decon, which is guided by the basal and ER+ cancer subtypes in modeling the sparsified TCGA-BRCA bulk data. This led to a 10-topic distribution, each of which was specifically tailored for a combination of cell type and cancer subtype. The heatmap displays the probabilities of the top 20 genes for each topic. The left half displays the cell-type-specific topic distribution for basal and the right half for ER+. **b** Predicted differentially expressed (DE) genes for each cell type between basal and ER+. The top DE genes for basal in contrast to ER+ were identified by subtracting the gene topic scores for ER+ from the gene topic score for basal under the same cell type. The resulting DE scores were shown in the top half of the heatmap. The bottom half displays the DE scores of the top genes for ER+ in contrast to basal. The pairwise Wilcoxon signed-rank tests were performed to compare the gene topic scores across all genes between the two subtypes for the same cell type. All tests yielded  $p$ -values lower than  $2.2e^{-16}$ . **c** Classification of basal and ER+ based on the phenotype probabilities. As a validation for our nested phenotype-cell-type guided approach, we evaluated the classification accuracy on the 160 held-out sparsified breast tumor samples. For each subtype, we summed the cell-type-specific topic probabilities from bottom heatmap for each sample to obtain the phenotype scores, which are shown in the top heatmap. **d** Comparison of DE genes detected by our approach and by DESeq2. DESeq2 was applied to the bulk RNA-seq gene expression data to compare gene expression between ER+ and basal samples. In total, 6815 DE genes were deemed significant by adjusted  $p$ -value  $< 0.05$  (Wald test) with 2952 upregulated and 3863 downregulated genes in ER+ relative to basal. The grey bar and the heatmap on the left display the  $-\log$  adjusted  $p$ -value for all of the upregulated genes (top half) and the downregulated genes (bottom half). Genes were ordered in decreasing order of the absolute test statistic for each half. The corresponding  $\log_2$  fold-change of ER+ over basal was also shown as heatmap. The heatmap on the right displays the change of gene topic score from basal subtype to ER+ subtype. **e** Cell-type-specific DE genes identified by nested-guided topic approach. The top and bottom part of the heatmap displays the topic scores for the upregulated and downregulated genes in basal relative to ER+, respectively ( $p$ -value  $< 0.05$ ; permutation test). Genes that were also detected by DESeq2 were labeled in the color bar. **f** ORA was applied to the differential topic scores of upregulated and downregulated genes in ER+ relative to basal. MSigDb HALLMARK pathway gene sets were used in ORA. The  $-\log p$ -values for the significant pathways were shown in the bar plot

statistically significant across all genes for all the cell types based on paired Wilcoxon signed-rank tests. Next, we evaluated the ability of the fine-tuned GTM-decon to classify the held-out test set and observed a 97% classification accuracy (Fig. 6c, top panel). These higher-resolution deconvolved CTS breast cancer profiles reveal that ER+ samples are enriched for Luminal-2 cell type, whereas the basal subtype is depleted for that cell type (Fig. 6c).



Next, we identified the DE genes per cell type by subtracting the genes-by-topic entries in the  $\phi$  matrix for phenotype  $d$  (e.g., basal) from phenotype  $d'$  (e.g., ER+) under the same cell type (e.g., Luminal 2). To evaluate the consistency of our DE genes, we compared them against the DE genes identified by DESeq2 (Fig. 6d). We observe that all upregulated and downregulated DE genes nominated by DESeq2 in ER+ versus basal comparison were also deemed upregulated and downregulated by our approach, respectively. We visualize the expression of the 630 statistically significant DE genes (adjusted  $p$ -value < 0.05; permutation test) in the BRCA tumor samples (Fig. 6e). We found that most of our DE genes do not only agree with those by DESeq2 but also exhibit CTS patterns. Notably, the most DE genes correspond to the Luminal\_2 cell type (Fig. 6e), which exhibit the larger difference between ER+ and basal (Additional file 1: Fig. S26). Over representation analysis (ORA) of these CTS DE genes against Hallmark pathways from MSigDb revealed several meaningful pathways. As expected, estrogen response (early) and (late) pathways are highly upregulated in the ER+ phenotype, mainly in the Luminal\_2 cell type (Fig. 6f). Similarly, most DE genes in the basal phenotype are upregulated for typical pathways involved in cancer, such as G2M checkpoint, E2F targets, and mitotic spindle. This reflects the aggressive nature of basal cell type as the origin for the cancer subtype (Fig. 6f). The contribution of sparse genes (i.e., genes with zero counts due to sparsification) to CTS topics and DE analysis is described in Additional file 1 Section S7 and illustrated in Additional file 1: Fig. S28. We obtained similar results comparing the 2 histological subtypes—ductal carcinoma and lobular carcinoma (Additional file 1: Fig. S29).

To further demonstrate the phenotype-guided and phenotype+cell-type-guided functionality, we applied GTM-decon to the same scRNA-seq data from breast cancer tumors using phenotypes (i.e., ER+ and TNBC) and cell-type labels as the guides and then used the inferred topics to deconvolve the bulk TCGA breast tumor transcriptome data. Detailed analyses are presented in Additional file 1 Section S8 and Fig. S30. This was feasible because we have scRNA-seq references that were collected from similar tissue sites from patients of the same disease phenotypes as the target bulk transcriptomes.

## Discussion

In this study, we developed a Bayesian approach called *GTM-decon* to infer CTS gene topic distributions from scRNA-seq reference data. During the topic inference of each cell, we introduce the guidance by setting the topic hyperparameter for the cell type of that cell to be relatively larger than the hyperparameters for other topics. This enables us to anchor each topic to a specific cell type and subsequently guide the inference of the global topic distributions over genes to automatically prioritize cell-type marker genes. The resulting topic distributions can then be used to infer the relative cell-type proportions from bulk RNA-seq datasets (i.e., cell-type deconvolution).

Through our analysis of the pancreatic and breast tissue datasets, we observe that for those cell types, where marker gene information is available, most of the top genes under the CTS topics correspond to known marker genes (Fig. 3a; Additional file 1: Fig. S24a). Because *GTM-decon* infers a distribution over all the genes under each cell-type-guided topic, it can be used to not only quantify the contribution of the known marker genes but also score novel marker genes. In terms of cell-type

deconvolution, GTM-decon confers comparable performance to the existing state-of-the-art methods (Fig. 2; Additional file 1: Fig. S10-15). The deconvolved cell-type proportions can be used to distinguish healthy and diseased samples as shown in the case of diabetic patients (Fig. 3c,e) as well as cancer subtypes from pancreatic and breast tumors (Fig. 4 and Additional file 1: Fig. S23). This enables investigation of molecular contribution to the phenotypic differences. Phenotypic differences between healthy and diabetic patients were captured even when the scRNA-seq reference datasets were generated from different platforms (e.g., Smart-seq and Drop-seq) (Fig. 3d, f).

Using GTM-decon, we revisited two cancer datasets from TCGA, namely pancreatic adenocarcinoma (PAAD) and breast cancer (BRCA). To dissect the heterogeneous tumor microenvironment, we deconvolved these datasets based on scRNA-seq reference sets from pancreatic cancer and breast cancer, respectively. As a result, we identified the ductal and acinar origin of PAAD, the endocrine origin of pancreatic neuroendocrine tumors (PNET) (Fig. 4b), and enrichment of subtype-specific cells for the BRCA subtypes (Additional file 1: Fig. S23a). Interestingly, using scRNA-seq references from pancreatic and breast tissues of healthy individuals is also sufficient to identify the cell type of origin for all subtypes in pancreatic cancer (Additional file 1: Fig. S22), as well as the basal origin of basal subtype and the luminal origin of ER+ breast cancer (Additional file 1: Fig. S24). However, using a cancer-specific scRNA-seq dataset improves the resolution of deconvolution in identifying more subtypes (Additional file 1: Fig. S23). By combining the de novo topics inferred directly from the bulk RNA-seq data with the scRNA-guided topics, we identified putative prognostic biomarkers that correlate with survival time (Fig. 4c; Additional file 1: Fig. S21, S22c).

We further extended GTM-decon by modeling the sparsified bulk RNA-seq data using the patient phenotype labels as the guide rather than cell types. In contrast to the traditional differential analysis approach, the phenotype-guided GTM-decon provides a different way to investigate gene signatures and molecular pathways underpinning the phenotypes of interest (Fig. 5). To leverage the scRNA-seq reference data, we further extended this framework to a nested-guided topic model by fine-tuning a dedicated set of CTS topic distributions for each phenotypic state (e.g., BRCA subtypes) of the patients from the bulk RNA-seq data. This enables learning not only intra-phenotype changes of cell-type distributions but inter-phenotype changes of gene expression. The latter allowed us to identify CTS DE genes directly from bulk RNA-seq data (Fig. 6). We extended this approach to infer phenotype and cell-type guided topics from single-cell breast cancer RNA-seq data with cancer subtypes as the phenotype guide. This led to accurate deconvolution of cancer subtypes from bulk TCGA-BRCA data, as well as identification of phenotype-specific and CTS DE genes from the single-cell data (Additional file 1: Fig. S30). Only a few methods can perform both deconvolution and CTS DE analysis. For example, CIBERSORTx [18] uses a non-negative matrix factorization approach based on partial observations to identify CTS DE genes across phenotypes. Other methods such as TOAST [24] and bMIND [23] that can estimate CTS DE require precomputed cell-type proportions by an external deconvolution method.

As future works, we can adapt GTM-decon to leverage other single-cell omic data such as scATAC-seq for reference and deconvolve the equivalent omic in bulk samples. It can also be extended to work with different cell states, like in BayesPrism [22], by modeling the topics based on cell states instead of cell types. Like all other deconvolution approaches, GTM-decon alone is unable to identify novel cell types from bulk RNA-seq datasets (i.e., cell types that are not present in the reference scRNA-seq data GTM-decon is trained on), while it is able to capture perturbed cell types based on the phenotypes (e.g., Fig. 6c). This can be addressed by training GTM-decon on an atlas-level scRNA-seq data such as the Human Cell Landscape [49] and Tabula Sapiens [50], which comprehensively covers most of the cell types in the primary tissues. The resulting model can then be used to deconvolve bulk samples of any given target tissue. Lastly, cell types are not independent entities but rather form a lineage. One future direction is to exploit the relations among the cell types, which may better capture the underlying phenotypic states of the subjects. A more recently developed method called CeDAR [51] uses known cell-type hierarchy as prior to infer CTS expression in bulk data as opposed to our de novo sub-cell-type inference and will leave a more detailed comparison as future work. Moreover, we can also extend GTM-decon to modeling multi-omic single-cell data to identify multi-omic CTS topic distributions and then use them to deconvolve multi-omic bulk data. To this end, while several multi-omic modeling methods have been developed [52–55], their benefits in deconvolution are not fully realized. Lastly, deep-learning-based methods such as Scaden [22] can train on simulated or real bulk RNA-seq datasets to predict cell-type proportions. While Scaden can confer accurate deconvolution results, it compromises interpretability because of its non-linear distributed representation of the gene expression features. Some recently developed variational autoencoder and embedded topic modeling frameworks [56, 57] may be extended to strike a balance between deconvolution accuracy and model interpretability.

## Conclusions

Computational cell-type deconvolution of heterogeneous bulk transcriptome is highly cost-effective in revealing the underlying phenotypic states and identifying CTS differentially expressed genes. GTM-decon represents a significant advance in the deconvolution methods with 3 prominent contributions: (1) automatic inference of interpretable CTS topic distributions by directly modeling large single-cell RNA-seq reference data without using known marker genes, therefore providing a principled and amenable reference map for the subsequent deconvolution; (2) identifying sub-CTS gene expression distributions by inferring multiple topics anchored at the same cell type, leading to a finer resolution of the deconvolution results; (3) detecting CTS DE genes directly from bulk or single-cell samples via an extended nested-guided topic design leveraging both the phenotype states and cell-type label information. Our comprehensive experiments on pancreatic and breast datasets demonstrated the utilities of all 3 contributions. Together, GTM-decon is an efficient marker-free deconvolution method that takes the full advantage of the single-cell RNA-seq reference data for cell-type deconvolution.

## Methods

### Modeling scRNA-seq reference data via topic inference guided by cell-type labels

We adapted GTM-decon from the MixEHR-Guided (Mixture of Electronic Health Records – Guided) [58] (which was in turn inspired by MixEHR (Mixture of Electronic Health Records) [59] and sureLDA (Surrogate-guided ensemble Latent Dirichlet Allocation) [60]) to model gene expression from scRNA-seq data. We assume that the scRNA-seq data are generated from the following data generative process. Each cell indexed by  $m \in \{1, \dots, M\}$  is a mixture of the  $K$  cell types. The cell-type mixture  $\theta_m$  is sampled from a  $K$ -dimensional Dirichlet distribution  $\theta_m \sim \text{Dir}(\alpha_m + 0.1)$  with the hyperparameter  $\alpha_m$  being specific to the cell. The key assumption here that separates GTM-decon from the standard LDA [25] is the use of noisy cell-type label  $y_m \in \{1, \dots, K\}$  for the cell. The hyperparameter corresponding to the cell-type label  $y_m = k$  has higher value ( $\alpha_{m,k} = 0.9$  by default) in contrast to the rest of the hyperparameters  $\alpha_{m,k'}$  for  $k' \neq k$ , which are randomly set to small values between 0.01 and 0.1. Note that having non-zero  $\alpha_{m,k'}$  allows other cell types to be assigned to the cell (i.e.,  $\theta_{m,k'} \geq 0$ ) and therefore  $\theta_m$  reflects the statistical uncertainty of the cell-type label  $y_m$ , which can be error prone due to various technical and data preprocessing aspects of the scRNA-seq data. To not clutter the notation, we omit the baseline value 0.1 in the following model description and use the more general form of  $\theta_m \sim \text{Dir}(\alpha_m)$  instead. Following the above default setting, this is equivalent to setting  $\alpha_{m,k} = 1$  for the observed cell type and  $\alpha_{m,k'} \in [0.11, 0.2]$  for other cell types. For each cell  $m$ , each scRNA-sequenced read  $i \in \{1, \dots, N_m\}$  originates from one of the  $K$  cell types with the probabilities dictated by its cell-type mixture ( $\theta_m$ ):  $z_{i,m} \sim \text{Cat}(\theta_m)$ . Given the cell type  $z_{i,m} = k$ , the  $i$ th read maps to a specific gene indexed by  $x_{i,m}$  with the probabilities dictated by the CTS topic distribution over all genes ( $\phi_k$ ):  $x_{i,m} \sim \text{Cat}(\phi_k)$ , where  $\phi_k$  follows a  $G$ -dimensional Dirichlet distribution with a fixed hyperparameter  $\beta$  across all  $K$  dimensions:  $\phi_k \sim \text{Dir}(\beta)$ . Since the topic mixture  $\theta_m$  are softly clamped to specific cell types via the  $K$ -dimensional hyperparameter  $\alpha_m$ , by following the above data generative process, it is straightforward to see that the  $K$  sets of topic distributions  $\phi_k$ 's are also CTS.

The posterior distribution for the latent variables  $\theta_m$ 's,  $z_{i,m}$ 's and  $\phi_k$ 's conditioned on the scRNA-seq reference data can be either approximated by Gibbs sampling [61] or by collapsed mean-field variational inference [62]. Specifically, for algorithmic convenience, we can leverage the conjugacy of the Dirichlet to categorical distribution by integrating out  $\theta_m$  and  $\phi_k$  resulting in two Dirichlet-Multinomial distributions [63]:

$$p(z|\alpha) = \int p(\theta|\alpha)p(z|\theta)d\theta = \prod_m \frac{\Gamma(\sum_k \alpha_{m,k})}{\prod_k \Gamma(\alpha_{m,k})} \frac{\prod_k \Gamma(\alpha_{m,k} + n_{.,m,k})}{\Gamma(\sum_k \alpha_{m,k} + n_{.,m,k})}$$

$$p(x|z) = \int p(\phi|\beta)p(x|z, \phi)d\phi = \prod_k \frac{\Gamma(G\beta)}{\prod_g \Gamma(\beta)} \frac{\prod_g \Gamma(\beta + n_{g.,k})}{\Gamma(G\beta + \sum_g n_{g.,k})}$$

Note that we can recover the expected values for  $\theta_m$ 's and  $\phi_k$ 's given the posterior estimates of  $z_{i,m}$ 's as they are proportional to the unnormalized counts  $n_{.,m,k} = \sum_i [z_{i,m} = k]$  and  $n_{g.,k} = \sum_m [z_{i,m} = k, x_{i,m} = g]$ , respectively.

The conditional distribution of the topic assignment for read  $i$  and cell  $m$  has a closed form expression:

$$\begin{aligned}
 p(z_{i,m} = k | x_{i,m} = g, z^{-(i,m)}, x^{-(i,m)}) &\propto p(z_{i,m} = k, z^{-(i,m)} | \alpha_{m,k}) p(x_{i,m} = g, x^{-(i,m)} | z_{i,m} = k, z^{-(i,m)}) \\
 &\propto \frac{\prod_k \Gamma(\alpha_{m,k} + n_{.,m,k})}{\Gamma(\sum_k \alpha_{m,k} + n_{.,m,k})} \prod_k \frac{\prod_g \Gamma(\beta + n_{g.,k})}{\Gamma(G\beta + \sum_g n_{g.,k})} \\
 &\propto \prod_{k' \neq k} \Gamma(\alpha_{m,k'} + n_{.,m,k'}^{(-i)}) \Gamma(\alpha_{m,k} + n_{.,m,k}^{(-i)} + 1) \prod_{k' \neq k} \frac{\prod_g \Gamma(\beta + n_{g.,k}^{-(i,m)})}{\Gamma(G\beta + \sum_g n_{g.,k}^{-(i,m)})} \frac{\Gamma(\beta + n_{g.,k}^{-(i,d)} + 1)}{\Gamma(G\beta + \sum_g n_{g.,k}^{-(i,d)} + 1)} \\
 &\propto \prod_k \Gamma(\alpha_{m,k'} + n_{.,m,k'}^{(-i)}) (\alpha_{m,k} + n_{.,m,k}^{(-i)}) \prod_k \frac{\prod_g \Gamma(\beta + n_{g.,k}^{-(i,m)})}{\Gamma(G\beta + \sum_g n_{g.,k}^{-(i,m)})} \frac{\beta + n_{g.,k}^{-(i,d)}}{G\beta + \sum_g n_{g.,k}^{-(i,d)}} \\
 &\propto (\alpha_{m,k} + n_{.,m,k}^{(-i)}) \left( \frac{\beta + n_{g.,k}^{-(i,m)}}{G\beta + \sum_g n_{g.,k}^{-(i,m)}} \right)
 \end{aligned}$$

where the second last equality exploits the property of the Gamma function, i.e.,  $\Gamma(x + 1) = \Gamma(x)x$ . Here  $n_{.,m,k}^{(-i)}$  is the total number of scRNA-seq reads allocated for topic  $k$  for cell  $m$  without counting the current  $i$ th read, and  $n_{g.,k}^{-(i,m)}$  is the total read counts for gene  $g$  under topic  $k$  across all of the  $M$  cells, without counting the current  $i$ th read in the  $m$ th cell:

$$\begin{aligned}
 n_{.,m,k}^{(-i)} &= \sum_{i' \neq i}^{N_m} [z_{i',m} = k] \\
 n_{g.,k}^{-(i,m)} &= \sum_{m' \neq m, i' \neq i} [z_{i',m'} = k, x_{i',m'} = g]
 \end{aligned}$$

From here, the topic inference can be done by collapsed Gibbs sampling from  $p(z_{i,m} = k | z^{-(i,m)}, x)$ , while fixing the topic assignments for all other reads [61]. For a large number of cells in the scRNA-seq data, the collapsed Gibbs sampling approach tends to be slow in reaching an equilibrium state. Therefore, we took a deterministic mean-field variational inference approach, known as the collapsed variational Bayes (CVB) [62]. Specifically, we approximate the posterior distribution of the cell-type assignment  $p(z_{i,m} | \mathbf{x}, \alpha_m)$  via the variational categorical distribution  $q(z_{i,m} | x_{i,m} = g) = \prod_k \gamma_{i,m,k}^{[z_{i,m}=k]}$ , where

$$\gamma_{i,m,k} \propto (\alpha_{m,k} + n_{.,m,k}^{(-i)}) \left( \frac{\beta + n_{g.,k}^{-(i,m)}}{G\beta + \sum_g n_{g.,k}^{-(i,m)}} \right)$$

Using the variational parameters for the topic assignments, the above sufficient statistics are replaced by the soft counts:

$$n_{.,m,k}^{(-i)} = \sum_{i' \neq i}^{N_m} \gamma_{i',m,k}$$

$$n_{g,,k}^{-(i,m)} = \sum_{m' \neq m \text{ or } i' \neq i} [x_{i',m'} = g] \alpha_{m,k} \gamma_{i',m',k}$$

Here in updating the global gene distribution in the second equation, we further make use of the CTS prior  $\alpha_m$  to obtain more interpretable results.

The above topic inference formulation operates at the level of read. For computational efficiency, our actual implementation of the topic inference was simplified to operate at the level of gene instead of the level of read:

$$\gamma_{g,m,k} \propto \left( \alpha_{m,k} + n_{.,m,k}^{(-g)} \right) \left( \frac{\beta + n_{g,,k}^{(-m)}}{G\beta + \sum_g n_{g,,k}^{(-m)}} \right)$$

Similar to the above sufficient statistics,  $n_{.,m,k}^{(-g)} = \sum_{g' \neq g} \gamma_{g',m,k}$  is the total number of scRNA-seq reads allocated for topic  $k$  for cell  $m$ , and  $n_{g,,k}^{(-m)} = \sum_{m' \neq m} \alpha_{m,k} \gamma_{g,m',k}$  is the total read counts for gene  $g$  under topic  $k$  across all of the  $M$  cells, without counting the  $g$ th gene for cell  $m$ . This is equivalent to a reasonable assumption that all the reads from the same gene for the same cell are originated from the same cell type, i.e.,  $\forall_{i,i'} z_{i,m} = z_{i',m}$  if  $x_{i,m} = x_{i',m}$ .

Together, the inference algorithm alternates between two simple steps: (1) for each cell and each gene, perform coordinate ascent by computing  $\gamma_{g,m,k}$  while fixing the variational parameters  $\gamma_{g',m,k}$  for other genes  $g'$ ; (2) update the sufficient statistics. This algorithm maximizes the evidence lower bound  $ELBO = E_q[\log p(x|z)] + E_q[\log p(z|\alpha)] - E_q[\log q(z|\gamma)]$  [62]. The model is deemed converged when ELBO stops improving by a small threshold ( $1e^{-6}$  by default).

Upon convergence, the expected values for  $\theta_{m,k}$  and  $\phi_{g,k}$  are:

$$\mathbb{E}_q[\theta_{m,k}] = \frac{\alpha_{m,k} + n_{.,m,k}}{\sum_k \alpha_{m,k} + n_{.,m,k}} \equiv \hat{\theta}_{m,k}; \mathbb{E}_q[\phi_{g,k}] = \frac{\beta + n_{g,,k}}{G\beta + \sum_g n_{g,,k}} \equiv \hat{\phi}_{g,k}$$

where  $n_{.,m,k} = \sum_{i=1}^{N_m} \gamma_{i,m,k}$  and  $n_{g,,k} = \sum_{m=1}^M [x_{i,m} = g] \alpha_{m,k} \gamma_{i,m,k}$ . Here  $\{\hat{\theta}_{m,k}\}_{M \times K} = \hat{\Theta}$  can be used to assess the ‘‘purity’’ of the single cells as a quality control step and  $\{\hat{\phi}_{g,k}\}_{G \times K} = \hat{\Phi}$  probabilities are used as the CTS topics in the subsequent deconvolution step.

### Inferring multiple topics per cell type

Suppose we use  $L$  topics per cell type for  $K$  cell types, the hyperparameter  $\alpha_m$  for cell  $m$  can be formulated into a  $K \times L$  matrix. Given that the cell-type label for cell  $m$  is  $y_m = k$ , we set the  $k$ th row of  $\alpha_m$  to relatively high values and the rest of the  $K-1$  rows to relatively low values. For example, suppose we have 5 cell types and 3 sub-topics per cell type. If cell  $m$  is labeled with cell type 2, then the topic prior hyperparameter matrix  $\alpha_m$  can be set to the following values: [0.01, 0.01, 0.01; 0.9, 0.9, 0.9; 0.01, 0.01, 0.01; 0.01, 0.01, 0.01; 0.01, 0.01, 0.01], where comma separates the columns and semicolons separate the rows. Here, the 3 topic prior values in the second row corresponding to cell type 2 are set to 0.9 and the values in the remaining rows are set to 0.01.

The data generative process is identical to the basic GTM-decon except having  $K \times L$  topics instead of  $K$  topics. Specifically, to sample the cell-type mixture  $\theta_m$ , we flatten the

$K \times L$  matrix for  $\alpha_m$  to have a row vector of  $1 \times (K \times L)$  so that the  $\theta_m$  will have relatively high expected value for the  $y_m^{th}$  consecutive  $L$  values that correspond to the labeled cell type and relatively low expected value for the rest of the entries. We experimented modeling each cell type using  $L \in \{2, 3, 4, 5\}$  topics per cell type, with the hyperparameters  $\alpha_{m,k}$  for each cell of cell type  $k$  set to 0.45, 0.3, 0.22, and 0.18, respectively. The prior values for the remaining topics were assigned with random values between 0.001 and 0.01. These priors were heuristically chosen based on the hyperparameter of value 0.9 for one topic divided by the number of topics per cell type.

**Inferring mixing cell-type proportions in bulk transcriptome**

We assume a similar data generative process of the bulk transcriptome as the single-cell transcriptome described above. In particular, each bulk sample  $j \in \{1, \dots, D\}$  is a mixture of  $K$  cell types. Its cell-type mixture  $\theta_j$  is sampled from a  $K$ -dimensional symmetric Dirichlet distribution  $\theta_j \sim \text{Dir}(\alpha)$  with the hyperparameter fixed at a constant value across all  $K$  cell types (default:  $\alpha_k = 0.1 \forall k$ ). The flat hyperparameter value is used here since we typically do not have any prior information about the cell-type mixtures in the bulk RNA-seq data. For  $N_j$  total RNA-seq reads of bulk sample  $j$ , each read  $i \in \{1, \dots, N_j\}$  originates from one of the  $K$  cell types with the categorical rates set to be  $\theta_j$ ;  $z_{i,j} \sim \text{Cat}(\theta_j)$ , where  $z_{i,j} \in \{1, \dots, K\}$ , and maps to a specific gene indexed by  $x_{i,j} \in \{1, \dots, G\}$  with a known CTS categorical rates  $\hat{\phi}_{z_{i,j}}$ ;  $x_{i,j} \sim \text{Cat}(\hat{\phi}_{z_{i,j}})$ .

Performing deconvolution on a bulk transcriptome profile is equivalent to inferring the posterior distribution of the CTS topic mixture given its gene expression and the CTS topic distributions:  $p(\theta_j | \mathbf{x}_j, \hat{\Phi})$ . To this end, we used the  $G \times K$  genes-by-CTS-topic estimates  $\hat{\Phi}$  inferred from the scRNA-seq reference data (described in section “[Modeling scRNA-seq reference data via topic inference guided by cell-type labels](#)”) and perform variational inference to infer the cell-type mixture  $\theta_j$  for the  $j$ th bulk RNA-seq profile. Implementation-wise, similar to the scRNA-seq topic modeling, we also use the simplified topic inference algorithm at the gene level  $g \in \{1, \dots, G\}$  as opposed to at read level. This involves alternating between the topic assignment inference for  $\gamma_{g,j,k}$  and computing the sufficient statistics  $n_{.j,k}^{(-g)}$  while fixing  $\hat{\Phi}$ . Algorithmically, for each gene, we infer the topic assignments:

$$\gamma_{g,j,k} \propto \left( \alpha + n_{.j,k}^{(-g)} \right) \hat{\phi}_{g,k}$$

where  $n_{.j,k}^{(-g)} = \sum_{g' \neq g} \gamma_{g',j,k}$ . Upon convergence, we compute the expected CTS-topic mixture:

$$\mathbb{E}_q[\theta_{j,k}] = \frac{\alpha + n_{.j,k}}{K\alpha + \sum_k n_{.j,k}}$$

The bulk sample is transformed in the same way as the scRNA-seq reference data as described in the “[Preprocessing the reference scRNA-seq data](#)” section.

We also show that although GTM-decon infers RNA fractions instead of cell-type fractions, the correlation between the RNA fractions and cell-fractions is strong across cell types in multiple datasets, suggesting that RNA fractions per cell type can serve as a good surrogate to the cell fraction per cell type despite the potential

differences in cell sizes among cell types (described in Additional file 1: Section S9 and illustrated in Fig. S31).

#### Sparsification of the bulk data to directly infer topic distributions from them

In order to identify de novo topics that are not present in the reference scRNA-seq data, we applied standard LDA using the CVB implementation from MixEHR [59] to the sparsified bulk RNA-seq data. However, bulk RNA-seq data is a dense matrix with most genes having non-zero entries while topic models excel at modeling sparse matrices (e.g., scRNA-seq data). To make the bulk RNA-seq data amenable to our approach, we sparsified the dense matrix by setting all values below the 75th percentile to 0. This cut-off was derived from scRNA-seq datasets, where on average 25% of genes in a cell have non-zero count values. While sparsifying works well in our applications, as a caveat, we acknowledge that it will lead to information loss. We also experimented with two other approaches to sparsify the matrix: (1) training GTM-decon only using HVG genes in bulk RNA-seq data; (2) training GTM-decon using differentially expressed genes identified between the phenotypes using DESeq2 [64].

Note that the sparsification procedure was done on the bulk data only when we directly inferred gene topic distributions from them, which pertains to the de novo topic inference from bulk samples, phenotype-guided topic inference from the bulk, and nest-guided CTS-phenotype topic inference. All of the deconvolution experiments, where we first inferred topics from a single-cell reference dataset and then applied the inferred topics to deconvolve bulk data, does not involve the sparsification procedure (i.e., deconvolving the original bulk transcriptomes as they are).

#### Identifying statistically significant marker genes per cell type

After GTM-decon topic inference, marker genes for each cell type were identified using permutation test. Specifically, for gene  $g$  under topic  $k$ , we computed the difference of its topic score  $\phi_{g,k}$  from the average topic score over the rest of the  $K-1$  topics  $\phi_{g,k' \neq k}$ . For example, for 3 cell types, the test statistic for cell type 1 is calculated as  $\phi_{g,1} - (\phi_{g,2} + \phi_{g,3})/2$ . More generally, for  $K$  cell types and  $T$  topics per cell topic, the test statistic for cell type  $k$  and gene  $g$  is

$$\frac{1}{T} \left( \sum_{t=1} \phi_{k,g}^{\{t\}} - \frac{1}{K-1} \sum_{k' \neq k} \sum_{t=1} \phi_{k',g}^{\{t\}} \right)$$

The significance of the observed statistic is compared against the same statistic calculated from 100,000 permutations. The empirical  $p$ -value is computed as fraction of permutations, where the test statistic is greater than the observed value.

#### Evaluation of deconvolution accuracy using simulated and real bulk data

We evaluated the deconvolution accuracy by comparing the inferred cell-type proportions against the ground truth values for the datasets. First, we simulated bulk RNA-seq data using the scRNA-seq dataset. We summed the gene expression counts of each sample from the scRNA-seq data to represent the bulk data of that sample (artificial bulk transcriptome). The ground-truth cell-type proportions are the fraction of cells for each



cell type. To avoid information bias in the evaluation, we performed leave-one-out cross-validation (LOOCV) by inferring topics from scRNA-seq data for  $N - 1$  subjects as the training examples and deconvolving the held-out subject as the validation example. The results are shown in Additional file 1: Fig. S10.

In addition, we also used real bulk data with known cell-type proportion to benchmark our methods. To conduct different benchmarking and qualitative experiments, we used scRNA-seq datasets from several tissues for training GTM-decon (Additional file 1: Table S1): **a** human pancreas (E-MTAB-5061, GSE81433) (with 4 type 2 diabetes and 3 healthy subjects) and mouse (GSE81433); **b** Peripheral blood mononucleocytes (PBMC) of healthy human (GSE132044); **c** human blood cells (HBC) of healthy humans (GSE149938); **d** post-mortem brain tissue of frontal cortex from adult human (GSE97930); **e** Breast tissue from healthy individuals (GSE113197). Apart from healthy individuals, scRNA-seq references from patients with **a** pancreatic cancer (CRA001160) and **b** breast cancer (GSE176078) were used to deconvolve cancer bulk RNA-seq datasets.

Bulk RNA-seq datasets with known ground truth values were chosen for benchmarking (Additional file 1: Table S2). Datasets with known ground truth proportions from flow cytometry include **a** whole blood from 12 individuals (obtained from CIBERSORTx webpage), **b** PBMC from a cohort of 13 individuals (GSE107011), **c** PBMC from a cohort of 346 individuals (SDY67). For brain tissue, ground truth proportions from immunohistochemistry were available for a cohort of 41 individuals from the Religious Orders Study / Memory and Aging Project (ROSMAP) study (CortexCellDeconv). A unique dataset with both bulk RNA-seq data and scRNA-seq data from pancreas was accessed from E-MTAB-5060 and E-MTAB-5061, respectively, with the cell type proportions from scRNA-seq considered as ground truth values.

For the 3 immune (SDY67, whole blood, PBMC S13 cohort) and brain prefrontal cortex (ROSMAP) bulk datasets, we used two independent single-cell references with the most closely matched tissues of origin (Additional file 1: Table S1). For the pancreatic dataset (Segerstolpe), paired single-cell and bulk data from the same individuals were used in the LOOCV manner, where we used the cell-type proportions from the single-cell dataset to compare the estimated proportions from the paired bulk data in the same held-out subject.

To evaluate the concordance between the ground truth  $y_m$  and our inferred cell-type mixture  $\theta_m$  for each test sample  $m$ , we computed four common metric scores: (1) Pearson correlation coefficient (PCC), (2) Spearman correlation (SCC), (3) cross entropy (CE), and (4) residual mean squared error (RMSE) (Fig. 2a). Moreover, we also evaluated the deconvolution performance at each cell type using PCC and RMSE across samples (Fig. 2b; Additional file 1: Figs. S11-15).

For the qualitative analyses on the cancer datasets, the bulk RNA-seq datasets for pancreatic cancer (PAAD) and breast cancer (BRCA) datasets from TCGA were downloaded from the GDC data portal (<https://portal.gdc.cancer.gov/>). Similarly, bulk RNA-seq data from the pancreas of a cohort of 89 normal and diabetic individuals was obtained from GSE50244.

### Implementation of the existing deconvolution methods

We compared GTM-decon with other state-of-the-art methods, including BSEQ-sc, BISQUE, MuSiC, CIBERSORTx, and BayesPrism. For BSEQ-sc, marker genes were selected from CellmarkerDB for the brain and immune datasets, while the built-in pancreatic marker genes were used for the Segerstolpe dataset. BISQUE requires at least two paired reference bulk and reference single-cell samples, which is the case for the Segerstolpe dataset. When using human blood cell (HBC) scRNA-seq data as a reference, in order to use BISQUE, artificial bulk samples were constructed using single-cell data. For the other 3 bulk datasets, since we used only one scRNA-seq reference, we left out BISQUE from the evaluation. For CIBERSORTx, all genes were provided to the web portal (<https://cibersortx.stanford.edu/>) for Signature Matrix generation and Cell Fraction imputation. For BayesPrism, all genes were provided to the web portal (<https://www.bayesprism.org/>) for cell type composition estimation, with the metadata column for tumor status set to 0 for all cells.

### Preprocessing the reference scRNA-seq data

The gene expression profiles of scRNA-seq data were used as training data. Since the performance of GTM-decon may vary depending on how the scRNA-seq count data are processed, we explored different gene selection and transformations of the scRNA-seq data. Specifically, the following gene selection were considered:

- (i) *all*: all genes without removal of any gene;
- (ii) *pp*: preprocessed genes to remove uninformative genes (frequently expressed genes—found in  $\geq 80\%$  of cells, infrequently expressed genes—found in  $\leq 5$  cells);
- (iii) *hvg*: highly variable genes identified using the *highly\_variable\_genes* (HVG) function of scanpy [65].

Furthermore, for each of these gene sets, the following transformation were considered:

- (i) *raw count* (all / pp / hvg)—scRNA-seq read counts per gene;
- (ii) *normr* (all\_normr / pp\_normr / hvg\_normr)—normalize counts (while excluding highly expressed genes for calculating the normalization factor) as counts divided by the sum of counts per cell multiplied by a scaling factor of 10,000 (a commonly used factor for scRNA-seq data [66]), and round the values to their nearest integer to make the input suitable for topic modeling;
- (iii) *normr\_log1p* (all\_normr\_log1p / pp\_normr\_log1p / hvg\_normr\_log1p): log-transform normalized counts in (ii), and round the values to their nearest integer to make the input suitable for topic modeling.

### Gene set enrichment analysis (GSEA)

GSEA was performed using the *fgsea* package from R, on two different gene sets: (a) gene sets corresponding to specific cell types from CellMarkerDB [30]; (b) gene set corresponding to the HALLMARK pathways from MSigDb using *msigdb* [67].

### Survival analysis

For TCGA datasets, after cell-type proportions are inferred using GTM-decon, survival analysis is performed using the survival and survminer R packages in order to assess if there is any correlation between survival probability of the cancer subjects and cell-type proportions for the cell types. Kaplan–Meier curves were generated for those cell types, where the association is deemed statistically significant at  $p$ -value  $< 0.05$  based on log-rank test.

### Learning phenotype-specific gene signatures from bulk RNA-seq data

To identify phenotype-specific genes, we applied GTM-decon to directly infer phenotype-topics from the sparsified bulk RNA-seq data (Fig. 1b). To evaluate phenotype prediction accuracy, we randomly split the bulk RNA-seq samples into training and test datasets in 80:20 ratio. The training dataset is guided in a manner similar to GTM-decon for the scRNA-seq datasets. Briefly, for each subject, the hyperparameter(s) for the topic(s) corresponding to the observed phenotype was (were) set to the prior value of 0.9, and the hyperparameters for the rest of the topics were set to a small value between 0.01 and 0.1. Since a model of 5 topics per cell type worked well for several scRNA-seq datasets, we opted to model each phenotype by 5 topics as well. For the subjects in the 20% test data, we summed up the inferred phenotype mixture  $\theta_{d,k}$  values for all topics corresponding to the same phenotype. The predicted phenotype was the one with the highest summed up topic score. As a comparison, we also evaluated two other common machine learning methods namely logistic regression and random forest. We used their scikit-learn implementations with the default settings [68].

### Identifying CTS DE genes from bulk transcriptomes from nested-guided topics

In identifying phenotype-specific differentially expressed genes, we developed a way to leverage the CTS topics inferred from the scRNA-seq reference data by fine-tuning the CTS topics based on gene expression from the sparsified bulk RNA-seq data for different phenotypes (Fig. 1c). Specifically, we initialized the genes-by-topic matrix  $\Phi_d$  for each phenotype  $d$ . For example, for  $D = 3$  phenotypes in a tissue with  $K = 5$  cell types, we will have 15 topics, comprising of 3 sets of CTS  $\Phi_d$  matrices of 5 columns each. These matrices are guided by the phenotype labels during the variational inference.

For subject  $j$ , suppose his/her phenotype label is  $d \in \{1, \dots, D\}$ . The topic hyperparameters  $\alpha_{j,d,k}$ 's for all the  $K$  topics corresponding to the phenotype label  $d$  were set to 0.9, and the  $\alpha_{j,d',k}$  for the other phenotypes were set to a small value between 0.01 and 0.1. These values were propagated throughout the variational inference, guiding the phenotypic inference appropriately. The inference algorithm is identical to the one described in “Modeling scRNA-seq reference data via topic inference guided by cell-type labels” section for modeling the scRNA-seq data. This approach also works with single-cell transcriptomes from a disease study such as the single-cell breast cancer study, where we not only have the cell type information but also phenotype states of the subjects (i.e., ER + vs TNBC) (Additional file 1: Section S8 and Fig. S30).

From the learned matrix, the test statistic of a DE candidate gene  $g$  between phenotypes  $d$  and  $d' \neq d$  per cell type  $k$  were identified by subtracting the  $\phi_{g,d,k}$  value for

phenotype  $d$  from the  $\phi_{g,d',k}$  value for phenotype  $d'$ . We identify statistically significant genes with a  $p$ -value  $< 0.05$  by permutation test calculated from the 100,000 randomly shuffled  $\Phi_d$  matrices. The  $p$ -value is the fraction of times the null statistic from the permutations is greater than the observed test statistic. As a reference, we compared our approach against the DE genes identified from the same samples using DESeq2 [64]. We identified hallmark pathways from MSigDb which were enriched for the DE genes by Over Representation Analysis (ORA) using the ClusterProfiler package from R [69].

## Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-023-03034-4>.

**Additional file 1.** Supplementary Section S1 (Evaluation on gene selection strategies); S2 (Evaluation on raw count and transformation strategies); S3 (Experimenting hyperparameters  $\alpha_{m,k}$  for cell-type mixture prior); S4 (Experimenting hyperparameter  $\beta$  for CTS topics); S5 (Experimenting number of topics per cell type); S6 (Benchmark time and memory usage); S7 (Effect of sparsification on phenotype classification); S8 (Phenotype-CTS topic modeling of single-cell breast cancer transcriptomes for TCGA-BRCA bulk deconvolution); S9 (Effect of cell size on inference of cell-type deconvolution); Table S1-S3; Figure S1-S31.

**Additional file 2.** Review history.

## Acknowledgements

We thank Yixuan Li for helping troubleshoot the GTM-decon software at its later stage.

## Review history

The review history is available as Additional file 2.

## Peer review information

Andrew Cosgrove was the primary editor of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

## Authors' contributions

Y.L. conceived of the study. S.S. and Y.L. analyzed and interpreted the data, wrote the manuscript, and wrote the code for GTM-decon. M.H. ran the experiments on benchmarking the deconvolutional accuracy. Y.L. supervised the project. All authors approved the final manuscript.

## Funding

Y.L. is supported by Canada Research Chair (Tier 2) in Machine Learning for Genomics and Healthcare, New Frontier Research Fund – Exploration (NFRFE-2019–00980), and Canada First Research Excellence Fund Healthy Brains for Healthy Life (HBHL) initiative New Investigator start-up award (G249591).

## Availability of data and materials

The datasets analyzed during the current study were obtained from publicly available repositories or data portals. All the scRNA-seq datasets [29, 70–81] used for training in this study are listed in Table S1, and all the bulk RNA-seq datasets with known ground truth proportions used in this study are listed in Table S2. The Human pancreatic islet single-cell datasets from Segerstolpe et al. used in this study is available in the ArrayExpress database under the accession code E-MTAB-5061 [82], with the corresponding bulk datasets accessible using E-MTAB-5060 [83]. The Human and Mouse pancreatic islet datasets from Baron et al. used in this study are available in the GEO database under the accession code GSE84133 [84]. The cancerous pancreatic tissue dataset is available in the Genome Sequence Archive database under the accession code CRA001160 [85]. The normal breast tissue dataset is available in the GEO database under the accession code GSE113197 [86]. The cancerous breast tissue dataset is available in the GEO database under the accession code GSE176078 [87]. scRNA-seq dataset for PBMC is available from the Single Cell Portal as PBMC2 or at GSE132044 [88], and the dataset for human blood cells at GSE149938 [89]. Post-mortem brain frontal cortex from adult human is available under the accession code GSE97930 [90]. The RA synovium dataset is available in the ImmPort database under the accession code SDY998 [91]. The bulk RNA-seq dataset for normal and type 2 diabetes patients is available in the GEO database under the accession code GSE50244 [92]. The bulk RNA-seq datasets for pancreatic cancer (PAAD) and breast cancer (BRCA) datasets from TCGA can be downloaded from the GDC data portal (<https://portal.gdc.cancer.gov/>). Bulk RNA-seq datasets for PBMC can be accessed using the accession code GSE107011 [93] from GEO database, and SDY67 [94] from Immport database. The bulk dataset for purified immune cells can be accessed using the GEO identifier GSE64655 [95].

GTM-decon source codes have been deposited at the GitHub repository (<https://github.com/li-lab-mcgill/gtm-decon>) [96]. The repository is licensed under the open-source GPL-3.0. The source code was also deposited at Zenodo: <https://doi.org/10.5281/zenodo.8200316> [97].

## Declarations

### Ethics approval and consent to participate

Not applicable.

### Consent for publication

Not applicable.

### Competing interests

The authors declare that they have no competing interests.

Received: 22 December 2022 Accepted: 9 August 2023

Published online: 18 August 2023

## References

1. Cho J-H, Kim J-W, Shin J-A, Shin J, Yoon K-H.  $\beta$ -cell mass in people with type 2 diabetes. *J Diab Investig*. 2011;2:6–17.
2. Sasaki H, Saisho Y, Inaishi J, Watanabe Y, Tsuchiya T, Makio M, Sato M, Nishikawa M, Kitago M, Yamada T, Itoh H. Reduced beta cell number rather than size is a major contributor to beta cell loss in type 2 diabetes. *Diabetologia*. 2021;64:1816–21.
3. van Galen P, Hovestadt V, Wadsworth li MH, Hughes TK, Griffin GK, Battaglia S, Verga JA, Stephansky J, Pastika TJ, Lombardi Story J, et al. Single-cell RNA-Seq reveals AML hierarchies relevant to disease progression and immunity. *Cell*. 2019;176:1265–1281.e1224.
4. Chen G, Ning B, Shi T. Single-cell RNA-Seq technologies and related computational data analysis. *Front Genet*. 2019;10:317–317.
5. Haque A, Engel J, Teichmann SA, Lönnberg T. A practical guide to single-cell RNA-sequencing for biomedical research and clinical applications. *Genome medicine*. 2017;9:75–75.
6. Hwang B, Lee JH, Bang D. Single-cell RNA sequencing technologies and bioinformatics pipelines. *Exp Mol Med*. 2018;50:1–14.
7. Van den Berge K, Hembach KM, Sonesson C, Tiberi S, Clement L, Love MI, Patro R, Robinson MD. RNA sequencing data: Hitchhiker's guide to expression analysis. *Ann Rev Biomed Data Sci*. 2019;2:139–73.
8. Barkley D, Rao A, Pour M, França GS, Yanai I. Cancer cell states and emergent properties of the dynamic tumor system. *Genome Res*. 2021;31:1719–27.
9. Davis-Marcisak EF, Deshpande A, Stein-O'Brien GL, Ho WJ, Laheru D, Jaffee EM, Fertig EJ, Kagohara LT. From bench to bedside: single-cell analysis for cancer immunotherapy. *Cancer Cell*. 2021;39:1062–80.
10. Conesa A, Madrigal P, Tarazona S, Gomez-Cabrero D, Cervera A, McPherson A, Szczesniak MW, Gaffney DJ, Elo LL, Zhang X, Mortazavi A. A survey of best practices for RNA-seq data analysis. *Genome Biol*. 2016;17:13–13.
11. Clough E, Barrett T. The Gene Expression Omnibus Database. *Methods Mol Biol (Clifton, NJ)*. 2016;1418:93–110.
12. Edgar R, Domrachev M, Lash AE. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Res*. 2002;30:207–10.
13. Consortium GT. The Genotype-Tissue Expression (GTEx) project. *Nat Genet*. 2013;45:580–5.
14. Hoadley KA, Yau C, Hinoue T, Wolf DM, Lazar AJ, Drill E, Shen R, Taylor AM, Cherniack AD, Thorsson V, et al. Cell-of-origin patterns dominate the molecular classification of 10,000 tumors from 33 types of cancer. *Cell*. 2018;173:291–304.e296.
15. Shen-Orr SS, Gaujoux R. Computational deconvolution: extracting cell type-specific information from heterogeneous samples. *Curr Opin Immunol*. 2013;25:571–8.
16. Newman AM, Liu CL, Green MR, Gentles AJ, Feng W, Xu Y, Hoang CD, Diehn M, Alizadeh AA. Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods*. 2015;12:453–7.
17. Baron M, Veres A, Wolock SL, Faust AL, Gaujoux R, Vetere A, Ryu JH, Wagner BK, Shen-Orr SS, Klein AM, et al. A single-cell transcriptomic map of the human and mouse pancreas reveals inter- and intra-cell population structure. *Cell Syst*. 2016;3:346–360.e344.
18. Newman AM, Steen CB, Liu CL, Gentles AJ, Chaudhuri AA, Scherer F, Khodadoust MS, Esfahani MS, Luca BA, Steiner D, et al. Determining cell type abundance and expression from bulk tissues with digital cytometry. *Nat Biotechnol*. 2019;37:773–82.
19. Wang X, Park J, Susztak K, Zhang NR, Li M. Bulk tissue cell type deconvolution with multi-subject single-cell expression reference. *Nat Commun*. 2019;10:380–380.
20. Racle J, Gfeller D. EPIC: a tool to estimate the proportions of different cell types from bulk gene expression data. *Methods Mol Biol (Clifton, NJ)*. 2020;2120:233–48.
21. Jew B, Alvarez M, Rahmani E, Miao Z, Ko A, Garske KM, Sul JH, Pietiläinen KH, Pajukanta P, Halperin E. Accurate estimation of cell composition in bulk expression through robust integration of single-cell information. *Nat Commun*. 2020;11:1971–1971.
22. Chu T, Wang Z, Pe'er D, Danko CG. Cell type and gene expression deconvolution with BayesPrism enables Bayesian integrative analysis across bulk and single-cell RNA sequencing in oncology. *Nature Cancer*. 2022;3:505–17.
23. Wang J, Roeder K, Devlin B. Bayesian estimation of cell type-specific gene expression with prior derived from single-cell data. *Genome Res*. 2021;31:268722.268120.
24. Wu H. TOAST: improving reference-free cell composition estimation by cross-cell type differential analysis. *Genome Biol*. 2019;20:1–17.
25. Blei DM, Ng AY, Jordan MI. Latent dirichlet allocation. *J Mach Learn Res*. 2003;3:993–1022.
26. Teh YW, Jordan MI, Beal MJ, Blei DM. Hierarchical dirichlet processes. 2006;101:1566–81.

27. Geering B, Stoeckle C, Conus S, Simon H-U. Living and dying for inflammation: neutrophils, eosinophils, basophils. *Trends Immunol.* 2013;34:398–409.
28. Hoffman MD, Blei DM, Wang C, Paisley J. Stochastic variational inference. *J Mach Learn Res.* 2013;14:1303–47.
29. Segerstolpe Å, Palasantza A, Eliasson P, Andersson E-M, Andréasson A-C, Sun X, Picelli S, Sabirsh A, Clausen M, Bjursell MK, et al. Single-cell transcriptome profiling of human pancreatic islets in health and type 2 diabetes. *Cell Metab.* 2016;24:593–607.
30. Zhang X, Lan Y, Xu J, Quan F, Zhao E, Deng C, Luo T, Xu L, Liao G, Yan M, et al. Cell Marker: a manually curated resource of cell markers in human and mouse. *Nucleic Acids Res.* 2019;47:D721–8.
31. Franzén O, Gan LM, Björkregren JL. PanglaoDB: a web server for exploration of mouse and human single-cell RNA sequencing data. *Database.* 2019;2019:baz046.
32. Fadista J, Vikman P, Laakso EO, Mollet IG, Esguerra JL, Taneera J, Storm P, Osmark P, Ladenvall C, Prasad RB, et al. Global genomic and transcriptomic analysis of human pancreatic islets reveals novel genes influencing glucose metabolism. *Proc Natl Acad Sci.* 2014;111:13924–9.
33. Cunningham F, Allen JE, Allen J, Alvarez-Jarreta J, Amode MR, Armean Irina M, Austine-Orimoloye O, Azov Andrey G, Barnes I, Bennett R, et al. Ensembl 2022. *Nucleic Acids Res.* 2022;50:D988–95.
34. Anderson NM, Simon MC. The tumor microenvironment. *Curr Biol CB.* 2020;30:R921–5.
35. Lei X, Lei Y, Li J-K, Du W-X, Li R-G, Yang J, Li J, Li F, Tan H-B. Immune cells within the tumor microenvironment: biological functions and roles in cancer immunotherapy. *Cancer Lett.* 2020;470:126–33.
36. Peng J, Sun B-F, Chen C-Y, Zhou J-Y, Chen Y-S, Chen H, Liu L, Huang D, Jiang J, Cui G-S, et al. Single-cell RNA-seq highlights intra-tumoral heterogeneity and malignant progression in pancreatic ductal adenocarcinoma. *Cell Res.* 2019;29:725–38.
37. Xu Y, Liu J, Nipper M, Wang P. Ductal vs. acinar? Recent insights into identifying cell lineage of pancreatic ductal adenocarcinoma. *Ann Pancreat Cancer.* 2019;2:11. <https://doi.org/10.21037/apc.2019.06.03>.
38. Peran I, Madhavan S, Byers SW, McCoy MD. Curation of the pancreatic ductal adenocarcinoma subset of the cancer genome atlas is essential for accurate conclusions about survival-related molecular mechanisms. *Clin Cancer Res.* 2018;24:3813–9.
39. Di Domenico A, Pipinikas CP, Maire RS, Bräutigam K, Simillion C, Dettmer MS, Vassella E, Thirlwell C, Perren A, Marinoni I. Epigenetic landscape of pancreatic neuroendocrine tumours reveals distinct cells of origin and means of tumour progression. *Commun Biol.* 2020;3:740–740.
40. Whittle MC, Hingorani SR. Fibroblasts in pancreatic ductal adenocarcinoma: biological mechanisms and therapeutic targets. *Gastroenterology.* 2019;156:2085–96.
41. Garcia PE, Scales MK, Allen BL, Pasca di Magliano M. Pancreatic Fibroblast Heterogeneity: From Development to Cancer. *Cells.* 2020;9(11):2464. <https://doi.org/10.3390/cells9112464>.
42. Alkasalias T, Moyano-Galceran L, Arsenian-Henriksson M, Lehti K. Fibroblasts in the Tumor Microenvironment: Shield or Spear? *Int J Mol Sci.* 2018;19(5):1532. <https://doi.org/10.3390/ijms19051532>.
43. Starzyńska T, Karczmarski J, Paziewska A, Kulecka M, Kuśnierz K, Żeber-Lubecka N, Ambrożkiewicz F, Mikula M, Kos-Kudła B, Ostrowski J. Differences between Well-Differentiated Neuroendocrine Tumors and Ductal Adenocarcinomas of the Pancreas Assessed by Multi-Omics Profiling. *Int J Mol Sci.* 2020;21(12):4470. <https://doi.org/10.3390/ijms21124470>.
44. Wu SZ, Al-Eryani G, Roden DL, Junankar S, Harvey K, Andersson A, Thennavan A, Wang C, Torpy JR, Bartonicek N, et al. A single-cell and spatially resolved atlas of human breast cancers. *Nat Genet.* 2021;53:1334–47.
45. Bertucci F, Finetti P, Birnbaum D. Basal breast cancer: a complex and deadly molecular subtype. *Curr Mol Med.* 2012;12:96–110.
46. Nguyen QH, Pervolarakis N, Blake K, Ma D, Davis RT, James N, Phung AT, Willey E, Kumar R, Jabart E, et al. Profiling human breast epithelial cells using single cell RNA sequencing identifies cell diversity. *Nat Commun.* 2018;9:2028–2028.
47. Dai X, Cheng H, Bai Z, Li J. Breast cancer cell line classification and its relevance with breast tumor subtyping. *J Cancer.* 2017;8:3131–41.
48. Love M, Anders S, Huber W. Differential analysis of count data—the DESeq2 package. *Genome Biol.* 2014;15:10–1186.
49. Han X, Zhou Z, Fei L, Sun H, Wang R, Chen Y, Chen H, Wang J, Tang H, Ge W, et al. Construction of a human cell landscape at single-cell level. *Nature.* 2020;581:303–9.
50. Consortium\* TS, Jones RC, Karkanas J, Krasnow MA, Pisco AO, Quake SR, Salzman J, Yosef N, Bulthaupt B, Brown P, et al. The Tabula Sapiens: a multiple-organ, single-cell transcriptomic atlas of humans. *Science.* 2022;376:eabl4896.
51. Chen L, Li Z, Wu H. CeDAR: incorporating cell type hierarchy improves cell type-specific differential analyses in bulk omics data. *Genome Biol.* 2023;24:37.
52. Hao Y, Hao S, Andersen-Nissen E, Mauck WM, Zheng S, Butler A, Lee MJ, Wilk AJ, Darby C, Zager M. Integrated analysis of multimodal single-cell data. *Cell.* 2021;184:3573–87.
53. Efremova M, Teichmann SA. Computational methods for single-cell omics across modalities. *Nat Methods.* 2020;17:14–7.
54. Zhou M, Zhang H, Bai Z, Mann-Krzisnik D, Wang F, Li Y. Single-cell multi-omic topic embedding reveals cell-type-specific and COVID-19 severity-related immune signatures. *bioRxiv.* 2023;2023.2001.2031:526312.
55. Argelaguet R, Cuomo ASE, Stegle O, Marioni JC. Computational principles and challenges in single-cell data integration. *Nat Biotechnol.* 2021;39:1202–15.
56. Zhao Y, Cai H, Zhang Z, Tang J, Li Y. Learning interpretable cellular and gene signature embeddings from single-cell transcriptomic data. *Nat Commun.* 2021;12:5261–5261.
57. Choi Y, Li R, Quon G. siVAE: interpretable deep generative models for single-cell transcriptomes. *Genome Biol.* 2023;24:29.
58. Ahuja Y, Zou Y, Verma A, Buckeridge D, Li Y. MixEHR-Guided: a guided multi-modal topic modeling approach for large-scale automatic phenotyping using the electronic health record. *J Biomed Inform.* 2022;134:104190–104190.
59. Li Y, Nair P, Lu XH, Wen Z, Wang Y, Dehaghi AA, et al. Inferring multimodal latent topics from electronic health records. *Nat Commun.* 2020;11(1):2536. <https://doi.org/10.1038/s41467-020-16378-3>.

60. Ahuja Y, Zhou D, He Z, Sun J, Castro VM, Gainer V, Murphy SN, Hong C, Cai T. sureLDA: A multidisease automated phenotyping method for the electronic health record. *J Am Med Inform Assoc.* 2020;27:1235–43.
61. Griffiths TL, Steyvers M. Finding scientific topics. *Proc Natl Acad Sci USA.* 2004;101(Suppl 1):5228–35.
62. Teh YW, Newman D, Welling M. A collapsed variational Bayesian inference algorithm for latent Dirichlet allocation. *Advances in neural...* 2006.
63. Minka T. Estimating a Dirichlet distribution. Technical report, MIT. 2000.
64. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* 2014;15:550–550.
65. Wolf FA, Angerer P, Theis FJ. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol.* 2018;19:15–15.
66. Satija R, Farrell JA, Gennert D, Schier AF, Regev A. Spatial reconstruction of single-cell gene expression data. *Nat Publ Group.* 2015;33:495–502.
67. Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, Tamayo P. The Molecular Signatures Database (MSigDB) hallmark gene set collection. *Cell Syst.* 2015;1:417–25.
68. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J. Scikit-learn: Machine learning in Python. *J Mach Learn Res.* 2011;12:2825–30.
69. Yu G, Wang LG, Han Y, He QY. clusterProfiler: an R package for comparing biological themes among gene clusters. *OMICS.* 2012;16:284–7.
70. Ding J, Adiconis X, Simmons SK, Kowalczyk MS, Hession CC, Marjanovic ND, Hughes TK, Wadsworth MH, Burks T, Nguyen LT, et al. Systematic comparison of single-cell and single-nucleus RNA-sequencing methods. *Nat Biotechnol.* 2020;38:737–46.
71. Xie X, Liu M, Zhang Y, Wang B, Zhu C, Wang C, Li Q, Huo Y, Guo J, Xu C, et al. Single-cell transcriptomic landscape of human blood cells. *Natl Sci Rev.* 2021;8:nwaa180.
72. Lake BB, Ai R, Kaeser GE, Salathia NS, Yung YC, Liu R, Wildberg A, Gao D, Fung HL, Chen S, et al. Neuronal subtypes and diversity revealed by single-nucleus RNA sequencing of the human brain. *Science.* 2016;352:1586–90.
73. Nguyen QH, Pervolarakis N, Blake K, Ma D, Davis RT, James N, Phung AT, Willey E, Kumar R, Jabart E, et al. Profiling human breast epithelial cells using single cell RNA sequencing identifies cell diversity. *Nat Commun.* 2028;2018:9.
74. Steen CB, Liu CL, Alizadeh AA, Newman AM. Profiling cell type abundance and expression in bulk tissues with CIBER-SORTx. *Methods Mol Biol.* 2020;2117:135–57.
75. Monaco G, Lee B, Xu W, Mustafah S, Hwang YY, Carre C, Burdin N, Visan L, Ceccarelli M, Poidinger M, et al. RNA-Seq signatures normalized by mRNA abundance allow absolute deconvolution of human immune cell types. *Cell Rep.* 2019;26(1627–1640):e1627.
76. Zimmermann MT, Kennedy RB, Grill DE, Oberg AL, Goergen KM, Ovsyannikova IG, Haralambieva IH, Poland GA. Integration of immune cell populations, mRNA-Seq, and CpG methylation to better predict humoral immunity to influenza vaccination: dependence of mRNA-Seq/CpG methylation on immune cell populations. *Front Immunol.* 2017;8:445.
77. Hoek KL, Samir P, Howard LM, Niu X, Prasad N, Galassie A, Liu Q, Allos TM, Floyd KA, Guo Y, et al. A cell-based systems biology assessment of human blood to monitor immune responses after influenza vaccination. *PLoS ONE.* 2015;10:e0118528.
78. Patrick E, Taga M, Ergun A, Ng B, Casazza W, Cimpean M, Yung C, Schneider JA, Bennett DA, Gaiteri C, et al. Deconvolving the contributions of cell-type heterogeneity on cortical gene expression. *PLoS Comput Biol.* 2020;16:e1008120.
79. Schulze A, Oshi M, Endo I, Takabe K. MYC Targets Scores Are Associated with Cancer Aggressiveness and Poor Survival in ER-Positive Primary and Metastatic Breast Cancer. *Int J Mol Sci.* 2020;21(21):8127. <https://doi.org/10.3390/ijms21218127>.
80. Oshi M, Takahashi H, Tokumaru Y, Yan L, Rashid OM, Nagahashi M, Matsuyama R, Endo I, Takabe K. The E2F Pathway Score as a Predictive Biomarker of Response to Neoadjuvant Therapy in ER+/HER2- Breast Cancer. *Cells.* 2020;9(7):1643. <https://doi.org/10.3390/cells9071643>.
81. Oshi M, Takahashi H, Tokumaru Y, Yan L, Rashid OM, Matsuyama R, Endo I, Takabe K. G2M Cell Cycle Pathway Score as a Prognostic Biomarker of Metastasis in Estrogen Receptor (ER)-Positive Breast Cancer. *Int J Mol Sci.* 2020;21(8):2921. <https://doi.org/10.3390/ijms21082921>.
82. Sandberg R, Palasantza A, Segerstolpe A. Single-cell RNA-seq analysis of human pancreas from healthy individuals and type 2 diabetes patients. *ArrayExpress.* 2016. <https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-5061>.
83. Palasantza A, Sandberg R, Clausen M. Whole-islet RNA-sequencing analysis of human pancreas from healthy individuals and type 2 diabetes patients. *ArrayExpress.* 2016. <https://www.ebi.ac.uk/biostudies/arrayexpress/studies/E-MTAB-5060>.
84. Veres A, Baron M. A single-cell transcriptomic map of the human and mouse pancreas reveals inter- and intra-cell population structure. *Gene Expression Omnibus.* 2016. <https://identifiers.org/geo:GSE84133>.
85. Ying Y: GSA-PDAC. *Genome Sequence Archive.* 2018. <https://ngdc.cnc.ac.cn/gsa/browse/CRA001160>.
86. Kessenbrock K. Single cell RNA sequencing of adult human breast epithelial cells. . *Gene Expression Omnibus.* 2018. <https://identifiers.org/geo:GSE113197>.
87. Swarbrick A, Wu S, Al-Eryani G, Roden D. A single-cell and spatially resolved atlas of human breast cancers. . *Gene Expression Omnibus.* 2021. <https://identifiers.org/geo:GSE176078>.
88. Ding J, Adiconis X, Simmons S, Kowalczyk M, Hession C, Marjanovic N, Hughes T, Wadsworth M, Burks T, Nguyen L, et al. Systematic comparative analysis of single cell RNA-sequencing methods. *Gene Expression Omnibus.* 2019. <https://identifiers.org/geo:GSE132044>.
89. Zhu P, Cheng T. Single-cell transcriptomic landscape of human blood cells. *Gene Expression Omnibus.* 2020. <https://identifiers.org/geo:GSE149938>.
90. Lake B, Chen S, Sos B, Fan JB, Yung Y, Chun J, Kharchenko P, Zhang K. Integrative single-cell analysis of transcriptional and epigenetic states in the human adult brain. *Gene Expression Omnibus.* 2017. <https://identifiers.org/geo:GSE97930>.

91. Anolik J, Bykerk V, Moreland L, Holers M, McGeachy M, Seifert J, Filer A, Pitzalis C, Gregersen P, Firestein G, et al. AMP Rheumatoid Arthritis Phase 1. *Immport*. 2018. <https://doi.org/10.21430/M3KXJHSP4T>.
92. Fadista J, Groop L. Global transcriptomic analysis of human pancreatic islets reveals novel genes influencing glucose metabolism. *Gene Expression Omnibus*. 2014. <https://identifiers.org/geo:GSE50244>.
93. Monaco G, Lee B, Xu W, Hwang Y, Poidinger M, Poidinger M, de Magalhães J, Larbi A. RNA-Seq profiling of 29 immune cell types and peripheral blood mononuclear cells. *Gene Expression Omnibus*. 2019. <https://identifiers.org/geo:GSE107011>.
94. Poland G. Bioinformatics Approach to 2010–2011 TIV Influenza A/H1N1 Vaccine Immune Profiling. *Immport*. 2015. <https://doi.org/10.21430/M3OYWCJHO1>.
95. Hoek K, Link A. A cell-based systems biology assessment of human blood to monitor immune responses after influenza vaccination. *Gene Expression Omnibus*. 2015. <https://identifiers.org/geo:GSE64655>.
96. Swapna LS, Huang M, Li Y: GTM-decon: Guided Topic Modeling for Deconvolution of cell types from bulk RNA-seq data. *Github*. 2023. <https://github.com/li-lab-mcgill/gtm-decon>.
97. Swapna LS, Huang M, Li Y. Source package and associated scripts for GTM-decon: guided-topic modelling of single-cell transcriptomes enables sub-cell-type and disease-subtype deconvolution of bulk transcriptomes. 2023. *Zenodo*. <https://doi.org/10.5281/zenodo.8200316>.

### Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Ready to submit your research? Choose BMC and benefit from:**

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

**At BMC, research is always in progress.**

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

