

**Performance of tumour microenvironment deconvolution methods in breast cancer
using single-cell simulated bulk mixtures**

Supplementary Information

Khoa A. Tran^{1,2}, Venkateswar Addala¹, Rebecca L. Johnston¹, David Lovell^{3,4}, Andrew Bradley⁵, Lambros T. Koufariotis¹, Scott Wood¹, Sunny Z. Wu^{6,7}, Dan Roden^{6,7}, Ghamdan Al-Eryani^{6,7}, Alexander Swarbrick^{6,7}, Elizabeth D. Williams^{2,8}, John V. Pearson¹, Olga Kondrashova^{1*}, Nicola Waddell^{1,2*##}

*Jointly supervised and contributed to the work

¹ Medical Genomics and Genome Informatics, QIMR Berghofer Medical Research Institute, Brisbane, 4006, Australia

² School of Biomedical Sciences at the Translational Research Institute (TRI), Queensland University of Technology (QUT), Brisbane, 4000, Australia

³ School of Computer Science, Queensland University of Technology, Brisbane, 4000, Australia

⁴ QUT Centre for Data Science, Brisbane, 4000, Australia

⁵ Faculty of Engineering, Queensland University of Technology, Brisbane, 4000, Australia

⁶ Cancer Ecosystems Program, Garvan Institute of Medical Research, Darlinghurst, NSW 2010, Australia

⁷ School of Clinical Medicine, Faculty of Medicine and Health, UNSW Sydney, Kensington, NSW 2052, Australia

⁸ Australian Prostate Cancer Research Centre – Queensland (APCRC-Q) and Queensland Bladder Cancer Initiative (QBCI), Brisbane, 4000, Australia

Corresponding author: Dr Nicola Waddell, Cancer Program, QIMR Berghofer Medical Research Institute, 300 Herston Road, Brisbane, 4006, Australia. Tel: +61 7 3845 3538. email: Nic.waddell@qimrberghofer.edu.au. ORCID: 0000-0002-3950-2476

Supplementary Note 1

Technical overviews of deconvolution methods

BayesPrism

BayesPrism¹ is a Bayesian method which uses single-cell reference profiles as the prior to estimate a joint posterior of cell-type fractions and cell-type-specific gene expressions. With each bulk mixture, the algorithm updates the posterior distribution using Gibbs sampling². This process is executed for each cell state, i.e. cell subtype, present in the single-cell reference, then cell-state fractions and gene expressions are summed over all cell states of each cell type. Furthermore, BayesPrism also factors in the presence of cancer cells and splits the single-cell reference into a malignant reference (used for deconvolving cancer population), and a non-malignant reference (used for deconvolving non-cancer populations).

bisqueRNA (Bisque)

Bisque³ is designed specifically to address issues that may originate from the proportionality assumptions between single-cell reference and bulk mixtures expression values. That is, many regression-based deconvolution methods rely on the assumption that single-cell reference and bulk mixture expression values are from the same distributions. However, this is not always true, as the methods that generate bulk and single-cell sequencing data, such as library preparation and mRNA capture methods, can differ greatly. Bisque overcomes these potential biases by transforming gene-specific bulk expression values into the same distribution space with gene-specific expression values of the single-cell reference data. Bisque achieves this by first averaging gene counts within each cell type across the single-cell reference matrix and then multiplying the results with cell type proportions of each patient (normalised counts of cells within each cell type), producing a “pseudo-bulk” reference sample for each individual. If single-cell and bulk expression values come from the same patient, expression of each gene in bulk mixtures will be transformed to the same distribution space as the “pseudo-bulk” reference via linear regression.

CIBERSORT_x (CBX)

CBX⁴ is an improvement of an earlier method called CIBERSORT⁵ (CB). Both CB and CBX use the same linear Support Vector Regression (SVR) algorithm to impute cell fractions using the same formula as other regression methods:

$$(1) B * F_i = M_i$$

where B is a signature matrix containing n discriminatory marker genes across c cell types, F_i is fractions of c cell types in mixture j , and M_i is bulk expression of n genes in bulk mixture. While CB relies on the signature matrix LM22⁵, CBX has built-in functionality to build a custom signature matrix using any scRNA-seq dataset. Essentially, this method filters out lowly expressed genes, i.e. those with transcripts per million (TPM) < 0.75 in log₂ scale, randomly selects 50% of available expression profiles for each cell type with count > 3, then sums their gene expression values and normalises gene expression values in non-log space into TPM. This process is repeated five times (default setting of CBX), with gene expression values averaged across duplicates to produce a final signature matrix.

Cell Population Mapping (CPM)

CPM⁶ performs deconvolution following a sequence of steps that traverse through a cell state space that captures the essence of gene-regulation variation among the reference single cells. This cell state space can be generated using dimensionality reduction techniques such as tSNE or UMAP.

By integrating the cell state space and reference single-cell expression values, CPM performs deconvolution by repeating many iterations of:

- 1) randomly sampling N number of reference subsets (each containing N_s cells) from the provided reference single cell (N is decided so that each single cell is randomly sampled at least N_r times),
- 2) predicting abundance of each cell in the N_s subset using the same linear SVR model in CBX⁴ and its earlier version CIBERSORT⁵,

- 3) averaging the abundance of each individual single cell over N runs (as one single cell can be sampled for a N_s subset at least N_r times)
- 4) computing smoothed abundance of each individual single cell as the averaged abundance of its N_d nearest-neighbour cell. This is to generate a smooth distribution of cell abundance across the entire cell space,
- 5) estimating cell-type-level abundance by averaging single-cell abundances across cell types.

Dampened Weighted Least Squares (DWLS)

DWLS⁷ relies on a single cell reference matrix to build its internal signature matrix in a 3-step process. First, DWLS clusters cells in the single-cell reference into different cell types using a rare-cell-type-sensitive clustering method and the provided unique cell types. After that, different expression analyses are conducted between cell clusters (using the MAST library) to detect distinct clusters containing biologically relevant cell types and identify marker genes for each cell type. DWLS defines marker genes as those with False Discovery Rate (FDR) adjusted p -value < 0.01 and \log_2 mean fold change > 0.5 . Next, 151 candidate signature matrices are created by selecting between 50-200 (hence 151) marker genes from each cell type. Within each candidate matrix, gene expression values are averaged across cell types. Similar to CIBERSORT, the candidate matrix with the lowest condition number is chosen.

DWLS then performs an improved version of the Ordinary Least Squares (OLS) method to deconvolute cell proportions. Traditional OLS methods model bulk expressions (t) as a function of signature gene expression values (S) and proportions of cells in the bulk mixture (x) in a linear relationship:

$$(2) t = S * x$$

By iterating over the provided data, DWLS (and any least-squares-based models) attempts to produce a set of x that minimizes the error function:

$$(3) \text{ Error} = t - B * x$$

This function can be refactored as a multiplication of relative percent error, S and x . This leads to two limitations of traditional OLS: 1) rare cell types with very small x have little to no impact on the error function, and 2) genes pertaining to prevalent cell types usually dominate the error term, regardless of whether they are differentially expressed. The second limitation also implies that differentially expressed genes with low expression play very minimal roles in the error terms. To address these two issues, DWLS introduces a weight to "weight down" the impact of abundant cell types such that

$$(4) \text{ Error} = W * (t - B * x)$$

Estimating the Proportion of Immune and Cancer cells (EPIC)

EPIC⁸ relies on gene markers for the cell types it can deconvolve to estimate populations of these cell types. The rest of the cell population is characterised as uncharacterized cell types, where the expression of signature genes is either very low or non-existent. This uncharacterised cell population can be referred to as cancer population.

hybrid-scale proportions estimation (hspe)

hspe⁹ is the updated version of dtangle¹⁰ benchmarked in Cobos et al¹¹. The inspiration for dtangle and hspe is that most existing deconvolution methods are regression-based models which require gene expression values either in linear or logarithmic scales. Both dtangle and hspe assert that while log-scale gene expression values are biologically questionable, linear-scale gene expression values have unrealistic error assumptions. Both methods rely on the use of a hybrid linear-logarithm approach to overcome these limitations. In our benchmarking study, hspe and dtangle produced identical results. We, therefore, report only results of hspe representative of both methods.

MUlti-Subject SIngle Cell deconvolution (MuSiC)

Similar to other single-cell-based methods benchmarked in this study, MuSiC¹² uses single-cell reference profiles to build its own set of marker genes. The method, however, implemented an innovation in its marker gene selection algorithm and another innovation in its quantification of cell-type abundance. The first innovation comes from its multi-subject gene-weighting approach. MuSiC dynamically emphasizes the importance of genes with low cross-subject variances by increasing their weights compared to the weights of genes with low cross-subject variances. This technique essentially incorporates cell-type-specific gene expression into the method's algorithm. MuSiC's second innovation was to deal with high collinearity of gene expression values among closely related cell types by incorporating a hierarchical deconvolution structure. The algorithm first groups cell types into clusters and estimate cluster-specific proportions. Within each cluster, MuSiC recursively repeats the same process to arrive at cell-type-specific proportions.

Single cell–assisted deconvolutional DNN (Scaden)

Scaden¹³ is an artificial intelligence method that uses an ensemble of three deep neural networks to infer the cellular composition of bulk tissues. By using gene expression values of simulated bulk mixtures to predict cell compositions, the three neural networks in Scaden “learn” how to adjust the weight of each neuron to get their predictions as close to the ground truth as possible. Each of the three neural networks in Scaden produce their own set of predictions, which are averaged to produce the method's final prediction.

Supplementary Note 2

Library size of simulated bulk mixtures

The recommended total RNA input for bulk RNA-seq library preparation can be as low as 10ng (SMART-seq) or 25ng (Illumina Stranded mRNA prep), which translates to 500 or 1,250 cells given the assumption that a typical mammalian cell contains 10-30 pg of total RNA. We chose to use 500 cells per pseudobulk to ensure we had enough cells for the simulated mixtures, particularly for the cancer purity experiments (where cancer cells made up to 95% of the mixture). As some samples contained between 500-1,000 cells for the most prevalent cell type that was used as the cell number for SMOTE oversampling, we were unable to SMOTE more cells than this number. We wanted to avoid sampling some cells twice, which would have been required if we used a higher number of cells per pseudo-bulk.

In terms of the library size, the average number of reads per cell is around 7,000 in Wu et al¹⁴. This adds up to around 3.5 million reads per pseudobulk mixture, which is less than 25-50 million reads for a typical real bulk mixture. We did, however, normalise scRNA-seq read counts to counts-per-10,000, which was comparable to using transcripts-per million normalised read counts in real bulk. We note that the smaller total library size means genes with lower expression would be less represented compared to bulk RNA-seq, which is a potential limitation of our study design. However, due to cell number constraints outlined above, we are unable to adjust the design.

Supplementary Figures

Supplementary Fig. 1: Number of cells per cell type before and after data oversampling using SMOTE in each breast cancer patient.

Supplementary Fig. 2: UMAP visualisations of synthesised and original cells from 26 breast cancer patients.

Supplementary Fig. 3: Cell counts for nine major cell types across training simulated variable tumour purity mixtures.

Supplementary Fig. 4: Impact of variable tumour purity levels on deconvolution for nine methods.

Supplementary Fig. 5: Performance of CPM across tumour purity levels.

Supplementary Fig. 6: Direction of mis-predictions of each cell type for nine computational methods.

Supplementary Fig. 7: Overall generalisation of impact of tumour purity on deconvolution to simulated bulk mixtures generated using scRNA-Seq from Bassez et al and Pal et al.

Supplementary Fig. 8: Cell-type-specific generalisation of impact of tumour purity on deconvolution to simulated bulk mixtures generated using scRNA-Seq from Bassez et al and Pal et al.

Supplementary Fig. 9: Histogram of predicted cancer proportions.

Supplementary Fig. 10: Studying the impact of tumour purity on deconvolution to simulated bulk mixtures generated using scRNA-Seq from Bassez et al and Pal et al, and single-cell reference from Wu et al.

Supplementary Fig. 11: Cell-type specific deconvolution performance for nine methods with tumour and normal epithelial cells grouped as epithelial.

Supplementary Fig. 12: The performance of the nine deconvolution methods assessed by false positive and false negative rates.

Supplementary Fig. 13: Cross-lineage performance of BayesPrism and DWLS by patient.

Supplementary Fig. 14: Raw prediction errors of BayesPrism and DWLS across immune lineages.

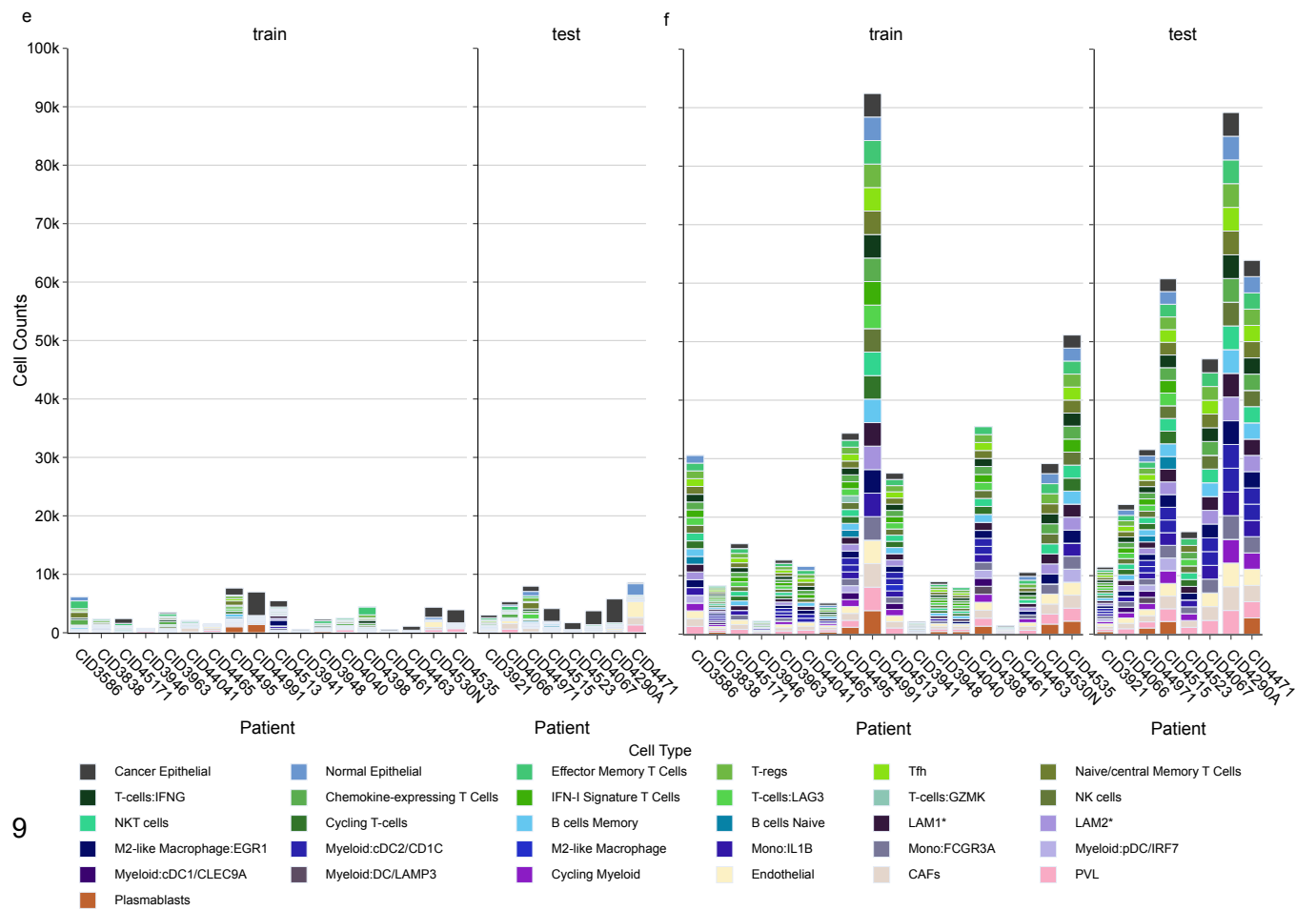
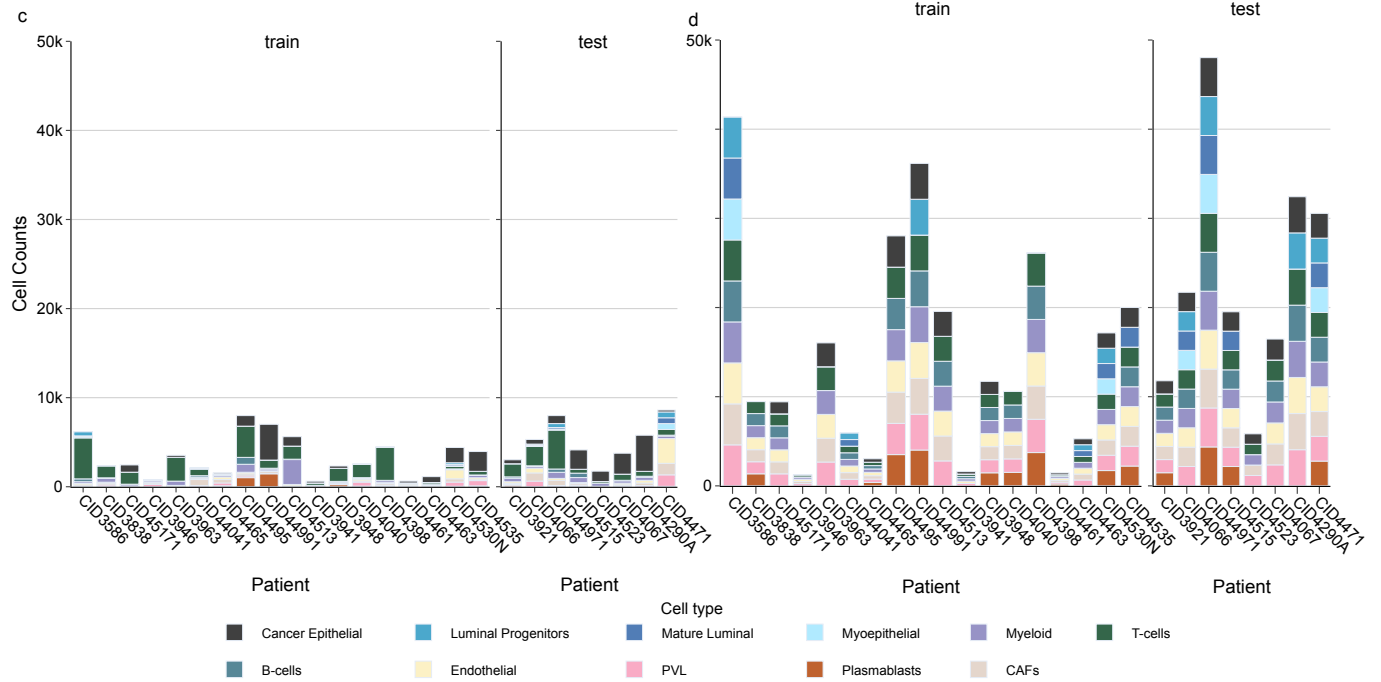
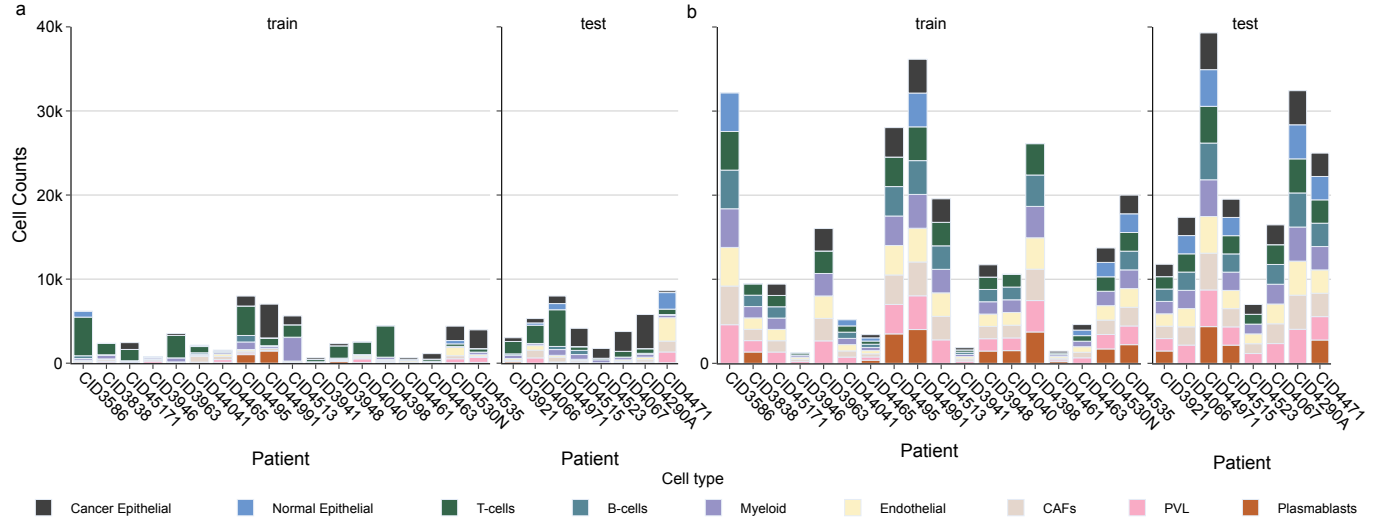
Supplementary Fig. 15: The performance of BayesPrism and DWLS methods for minor and subset immune cell types assessed by false positive and false negative rates.

Supplementary Fig. 16: Comparison of nine deconvolution approaches to four alternative methods to predict tumour purity in TCGA breast data.

Supplementary Fig. 17: Performance of BayesPrism before and after Gibbs sampling on artificial bulk mixtures generated using scRNA-seq data from Wu et al, Bassez et al and Pal et al.

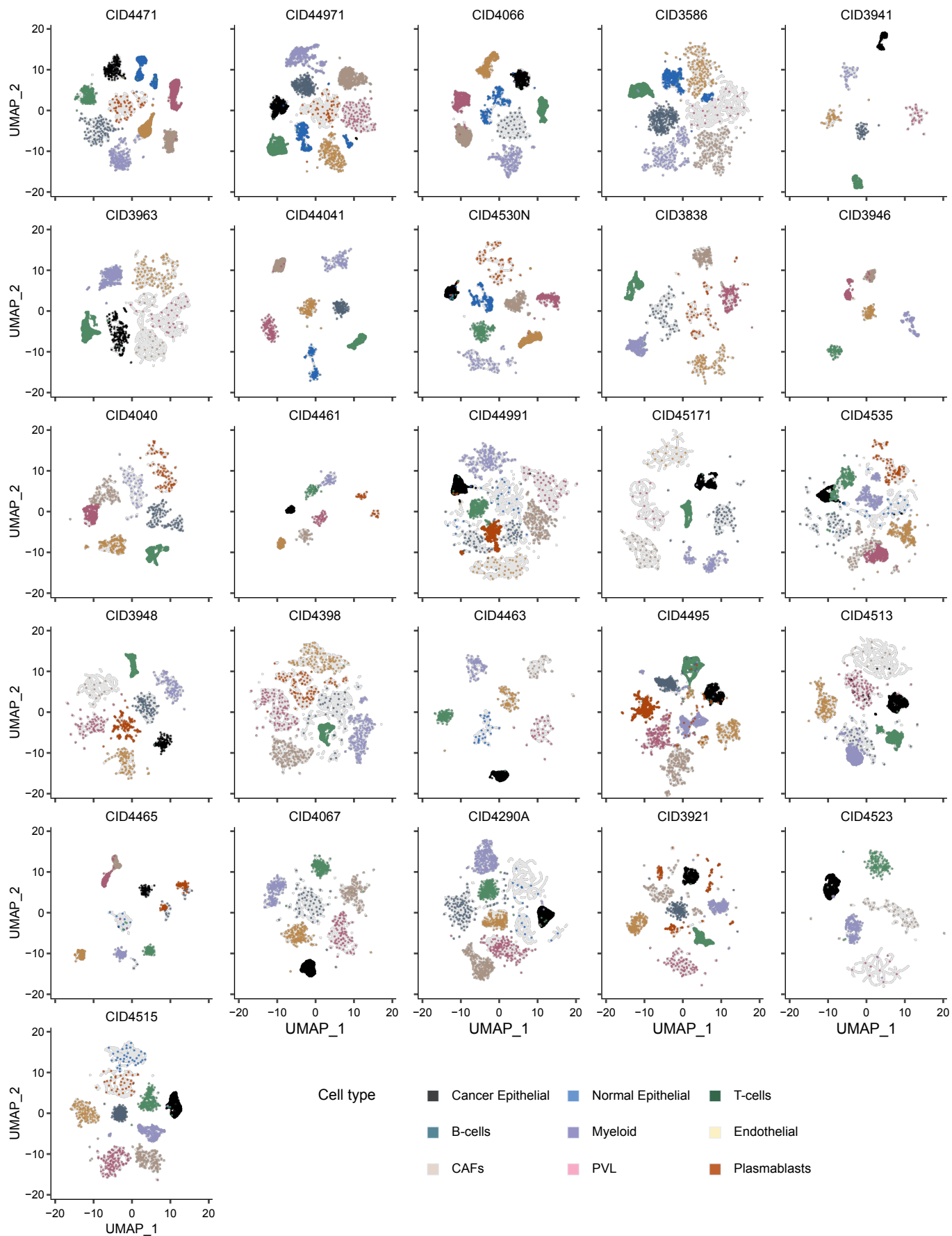
Supplementary Fig. 18: Performance comparison of different versions of BayesPrism, Bisque, MuSiC, hspe and DWLS.

9



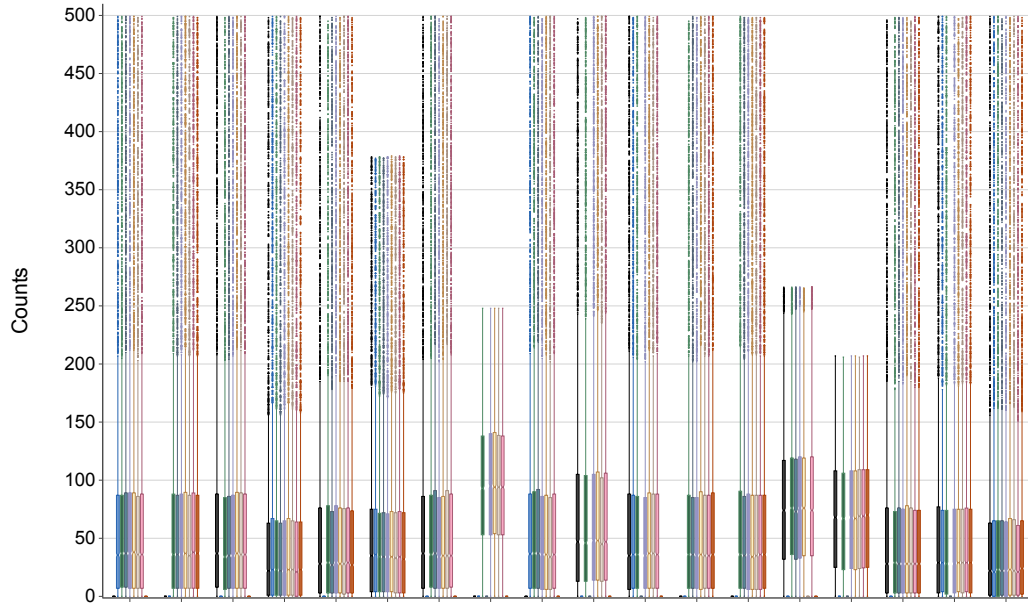
9

Supplementary Fig. 1: Number of cells per cell type before and after data oversampling using SMOTE in each breast cancer patient. Number of cells for each of the nine major cell types per patient before (a) and after data oversampling used to assess impact of tumour purity levels on deconvolution (b). Number of cells for each of the six non-normal-epithelial cell types and three minor normal epithelial cell types per patient before (c) and after data oversampling used to assess impact of normal cell lineages on deconvolution (d). Number of cells for each of the five non-immune cell types, as well as 17 subset cell types and six minor cell types of T-cells, B-cells and Myeloid for each patient before (e) and after data oversampling used to assess impact of immune cell lineages on deconvolution (f). All original cell counts before data oversampling were obtained from Wu *et al*, Nature Genetics 2021. Source data are provided as a Source Data file.

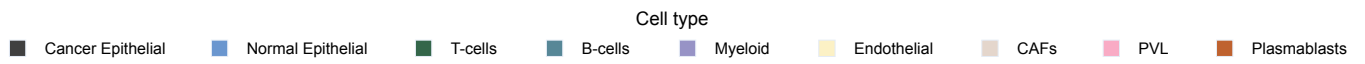
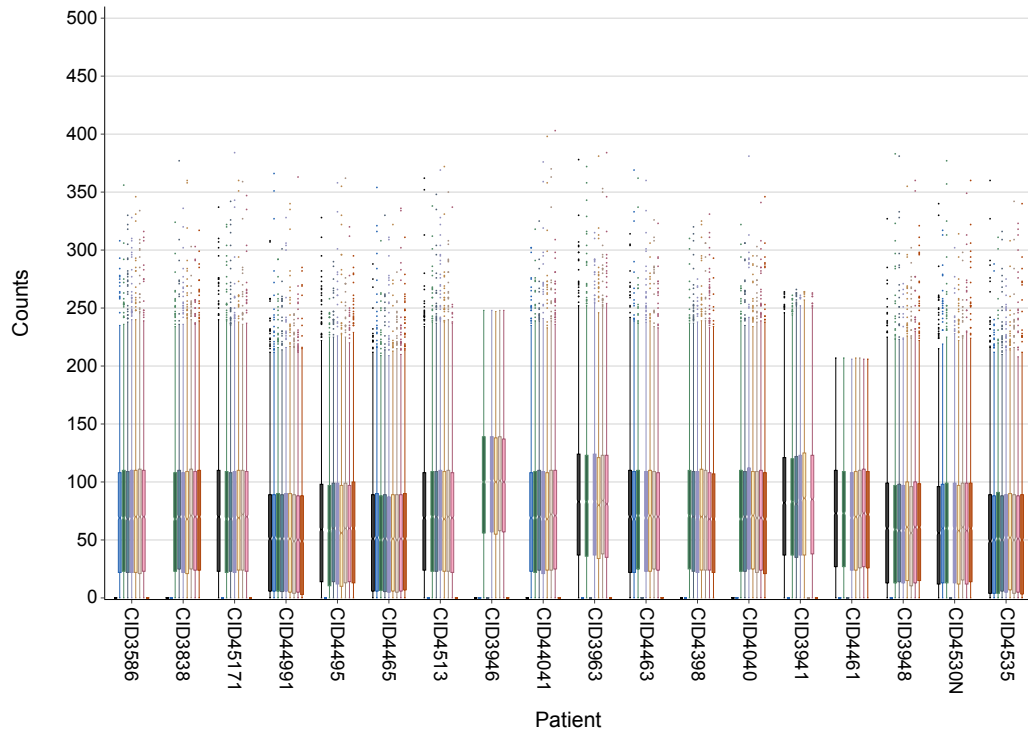


Supplementary Fig. 2: UMAP visualisations of synthesised and original cells from 26 breast cancer patients. UMAP was run on gene expression values of both synthesised and original cells for each individual patient. Within each plot, grey points represent SMOTE synthesised cells across all cell types, and points of other colours represent original cells from each of the nine major cell types. Source data are provided as a Source Data file.

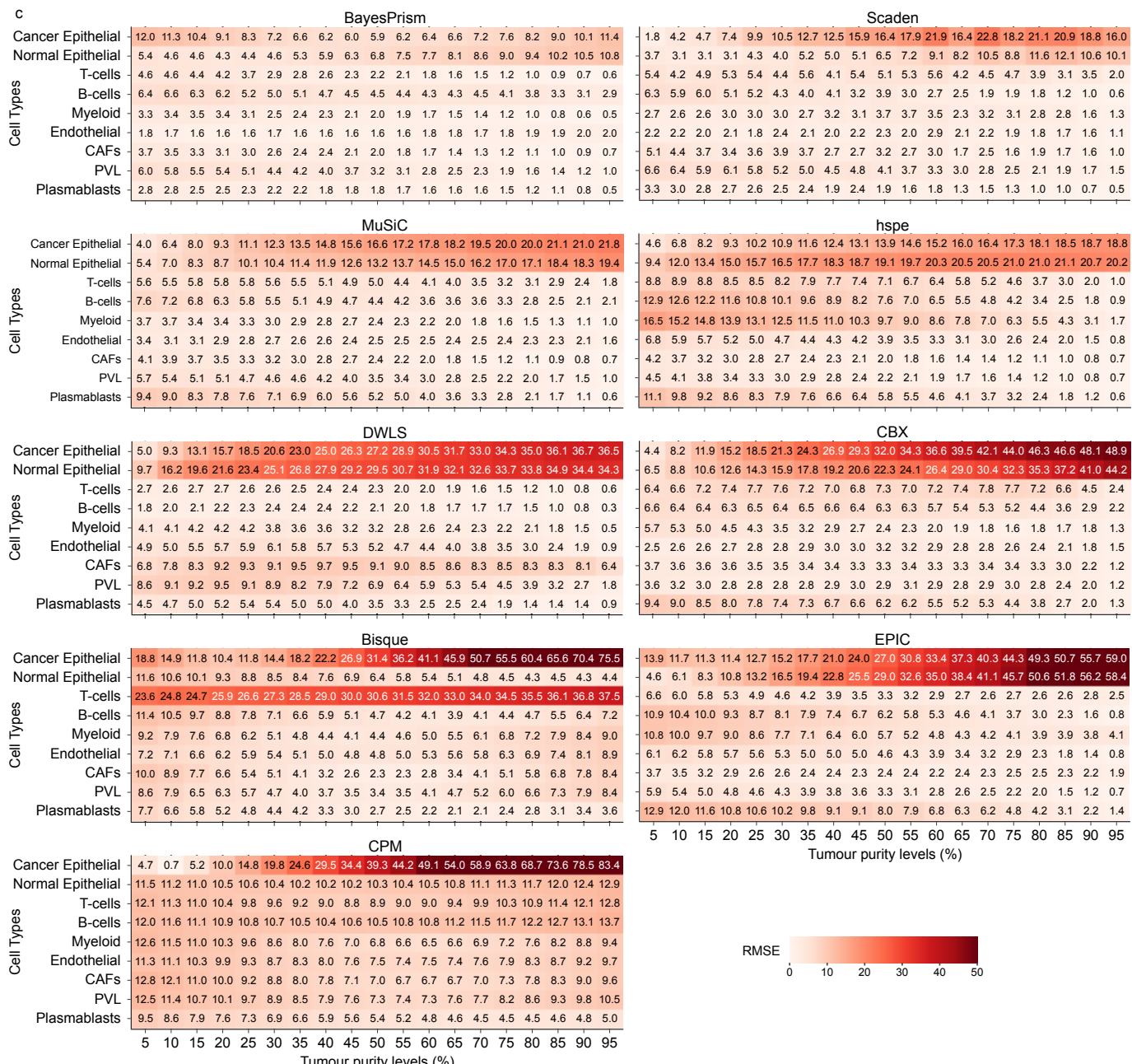
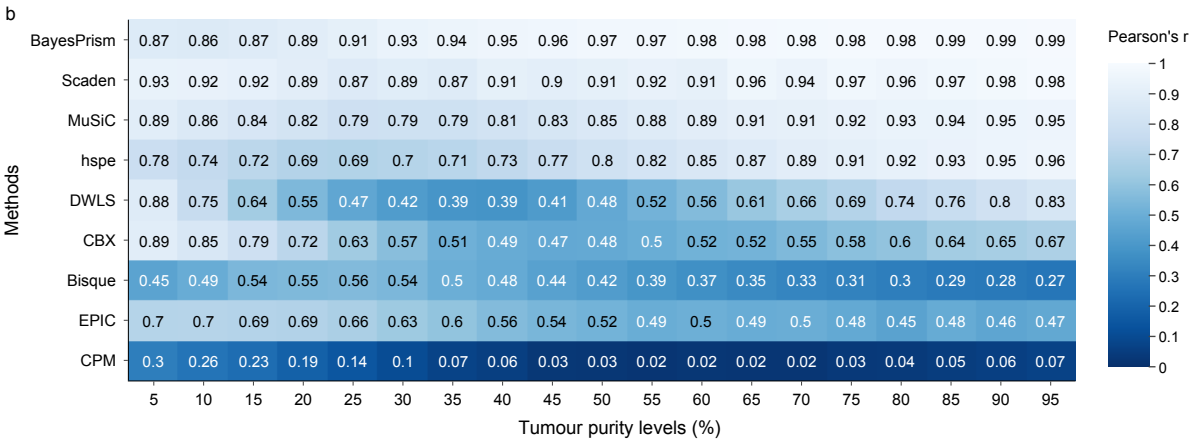
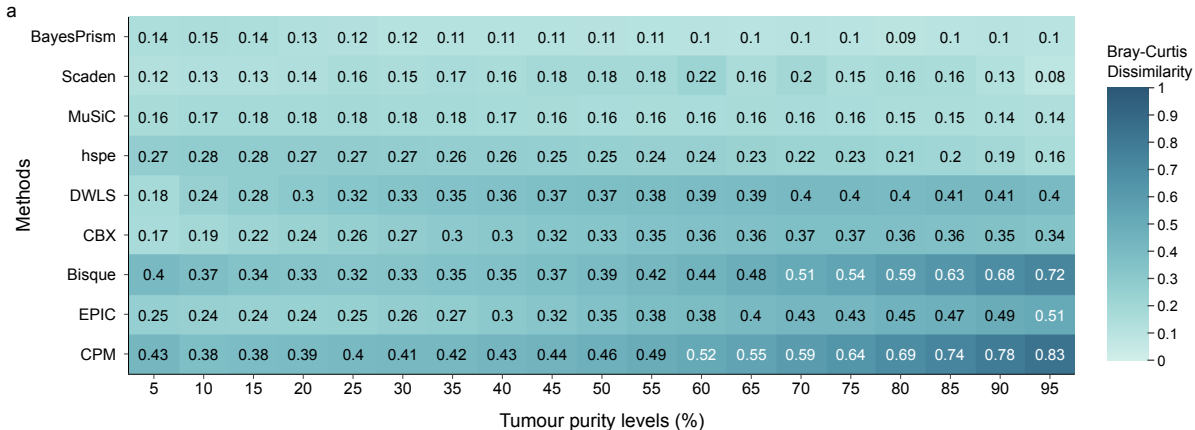
a



b

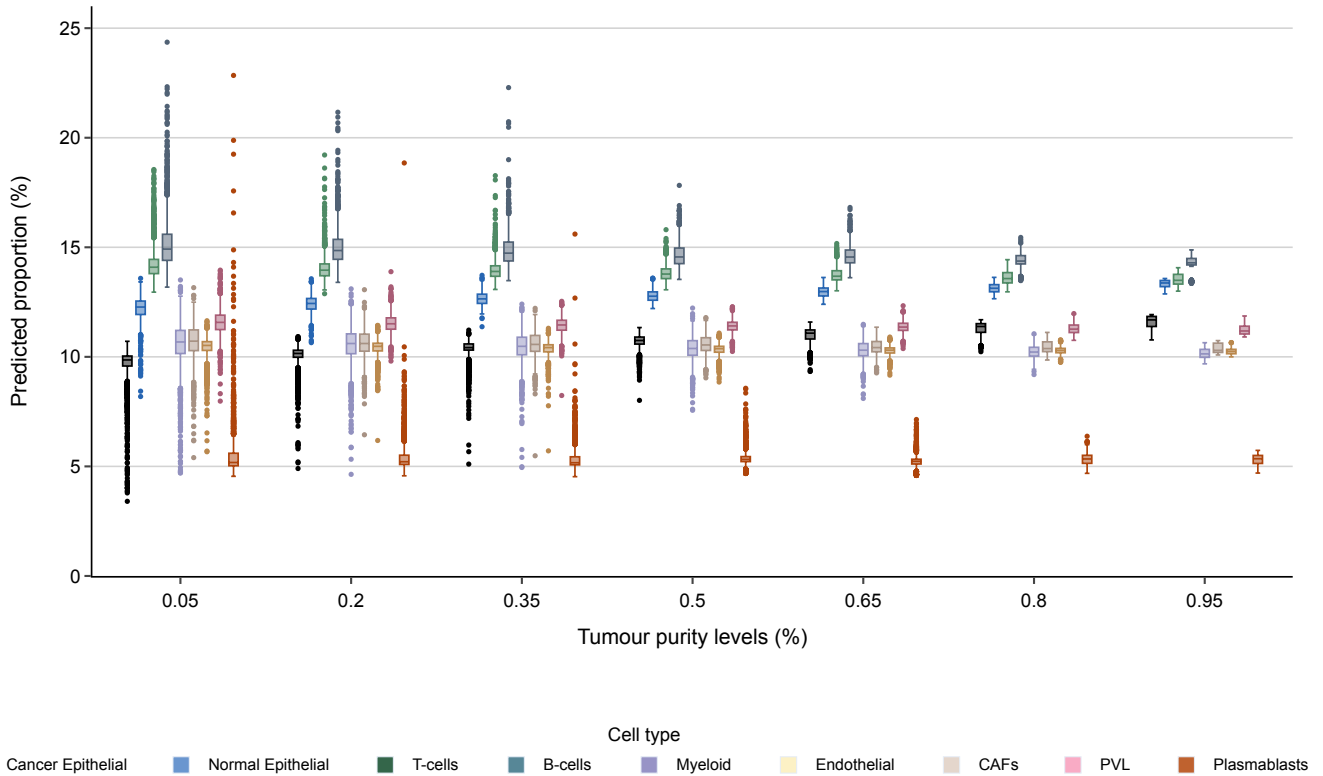


Supplementary Fig. 3: Cell counts for nine major cell types across training simulated variable tumour purity mixtures. The total of 38,000 mixtures across 19 tumour purity levels for each patient in the training dataset is shown (n=38,000 per boxplot). Boxplots represent cell counts in mixtures generated using the sparse method which was utilised in this study **(a)**, and cell counts in mixtures simulated using method from Menden *et al*, Science Advances 2020 **(b)**. All simulated mixtures had 500 cells, therefore a point closer to 500 on y-axis indicates a cell type accounting for the majority of a simulated mixture, while a point closer to 0 on the y-axis indicates a cell type with low representation in a simulated mixture. The sparse method enabled a more diverse proportion range across all cell types. Upper and lower whiskers depict cell counts outside of the centre 50%. Each box represents the middle 50% of cell-type counts per mixture, which includes the first quartile (Q1), the median, and the third quartile (Q3). Upper and lower whiskers depict maxima and minima of cell-type counts, excluding outliers. Outliers are cell-type counts that are more than 1.5x the interquartile range from either Q1 or Q3. Source data are provided as a Source Data file.

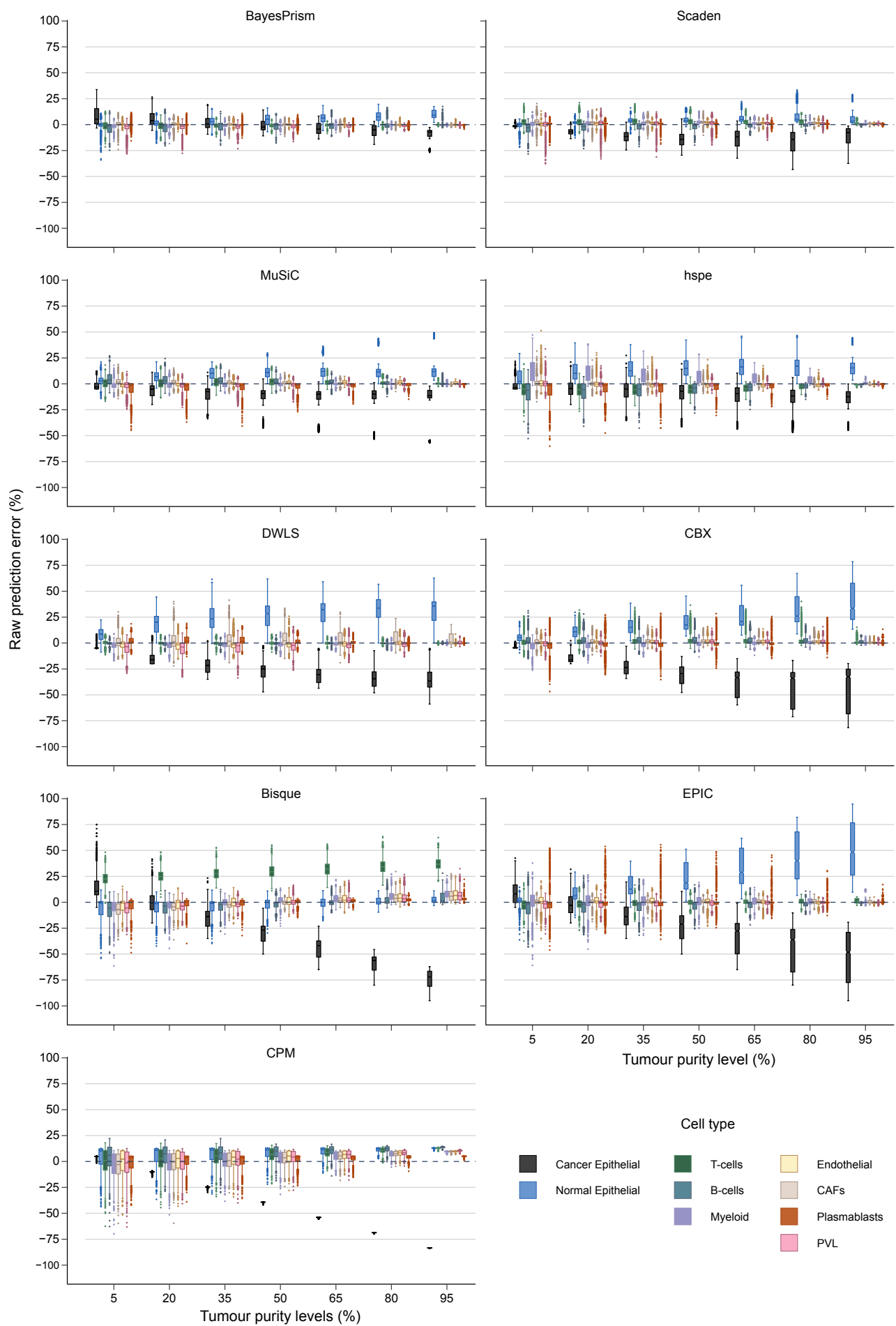


Supplementary Fig. 4: Impact of variable tumour purity levels on deconvolution for nine methods.

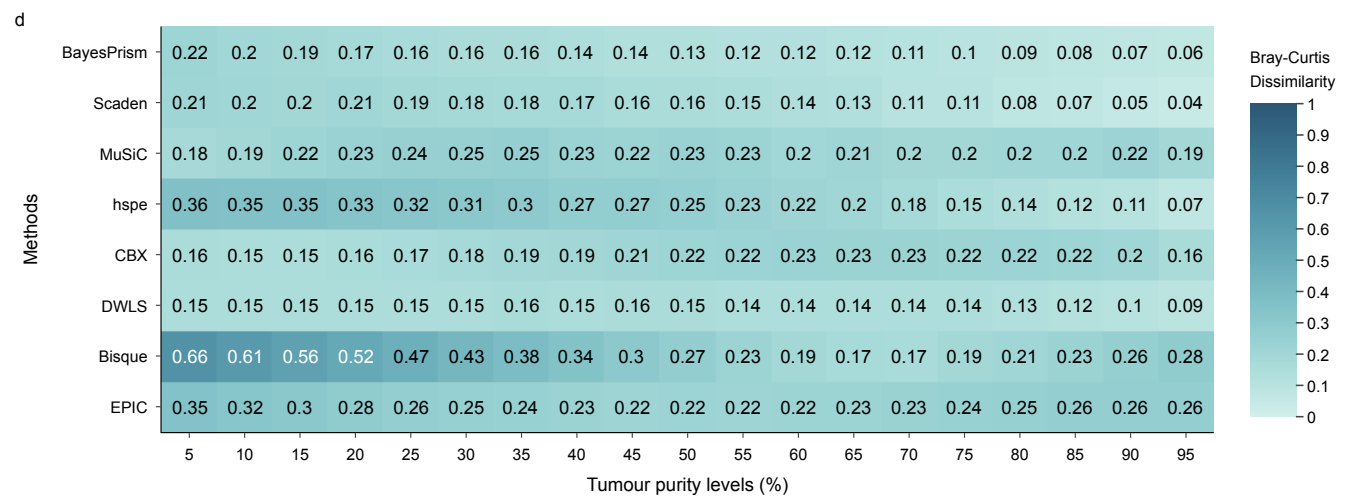
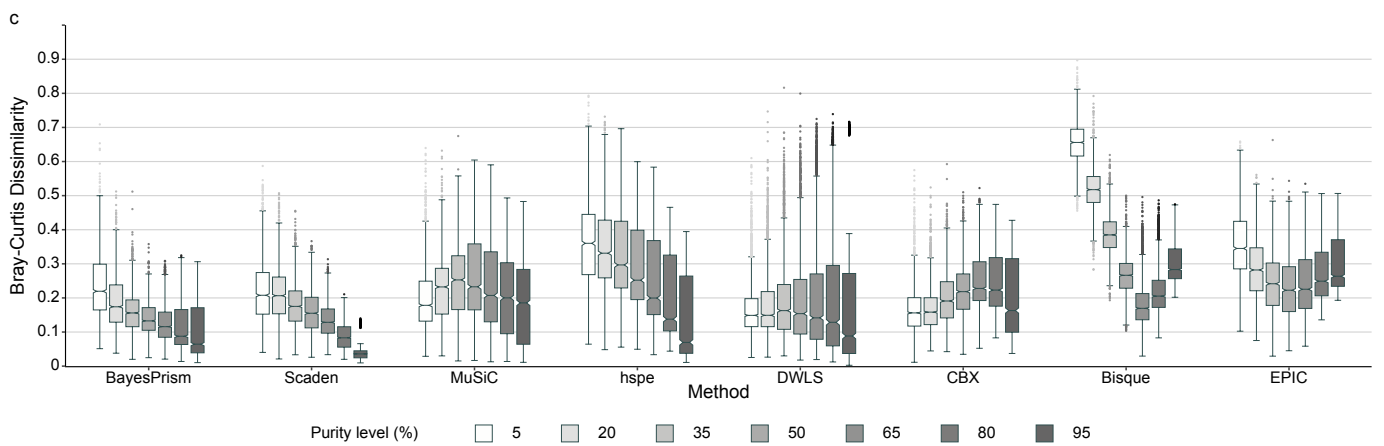
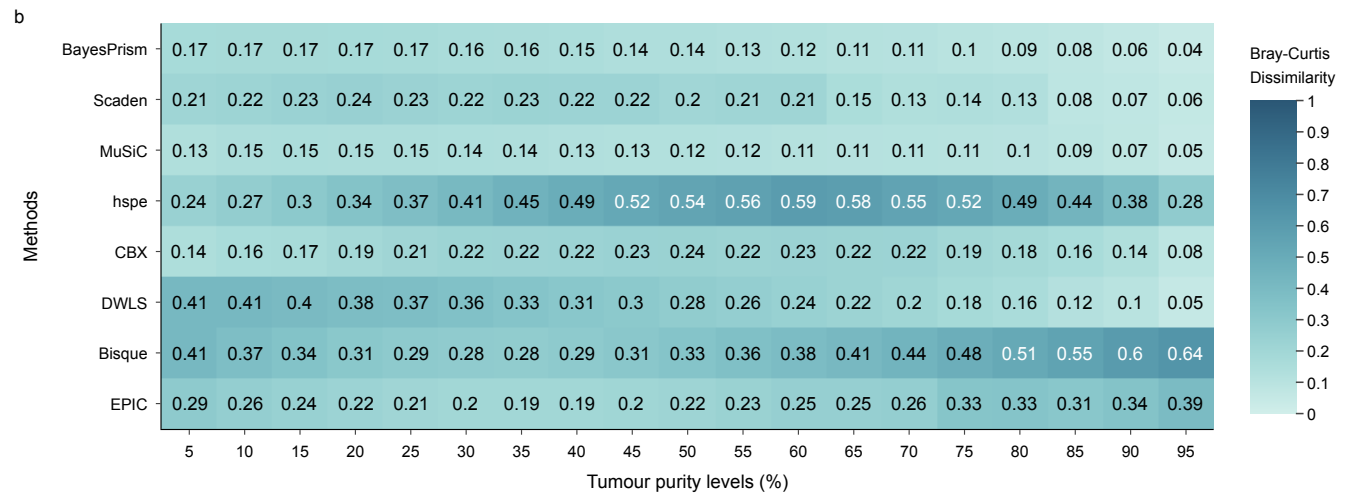
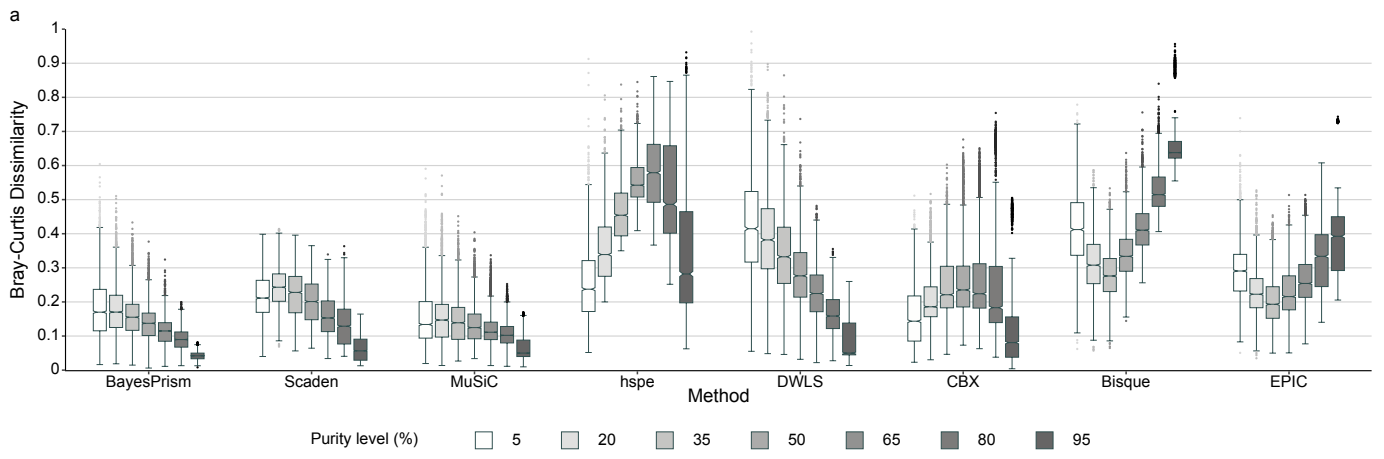
a) Median Bray-Curtis dissimilarity and **b)** Median Pearson's r correlation coefficients between predicted and actual cell compositions across 19 tumours purity levels (samples with tumour cells ranging from 5% to 95%, at 5% intervals) which are shown as proportions (x-axis) for nine computational methods (y-axis). There were 2,000 artificial bulk mixtures per purity level. Darker shade of teal represents higher Bray-Curtis dissimilarity (worse performance), and darker shade of blue represents lower Pearson's r values (worse performance). The numeric values are shown. The computational methods are organised in order of decreasing performance according to their median Bray-Curtis dissimilarity values. Upper and lower whiskers depict Bray-Curtis values outside of the centre 50%. **c)** Median RMSE of each of the nine major cell types across all 2,000 artificial mixtures and 19 tumour purity levels (from 5% to 95%, 5% interval). Darker shades of red indicate higher RMSE values and worse deconvolution performance. The numeric RMSE values are shown. Methods are ordered in the same order as **(a)**. RMSE: Root Mean Square Error. Source data are provided as a Source Data file.



Supplementary Fig. 5: Performance of CPM across tumour purity levels. Predicted proportions of the CPM method across 7 tumour purity levels (from 5% to 95%, 15% interval). At each tumour purity level, proportions are aggregated into nine major cell types. Upper and lower whiskers depict predicted cell-type proportions outside of the centre 50%. Each box represents the middle 50% of predicted cell-type proportions, which includes the first quartile (Q1), the median, and the third quartile (Q3). Upper and lower whiskers depict maxima and minima of cell-type proportions, excluding outliers. Outliers are cell-type proportions that are more than 1.5x the interquartile range from either Q1 or Q3. Source data are provided as a Source Data file.

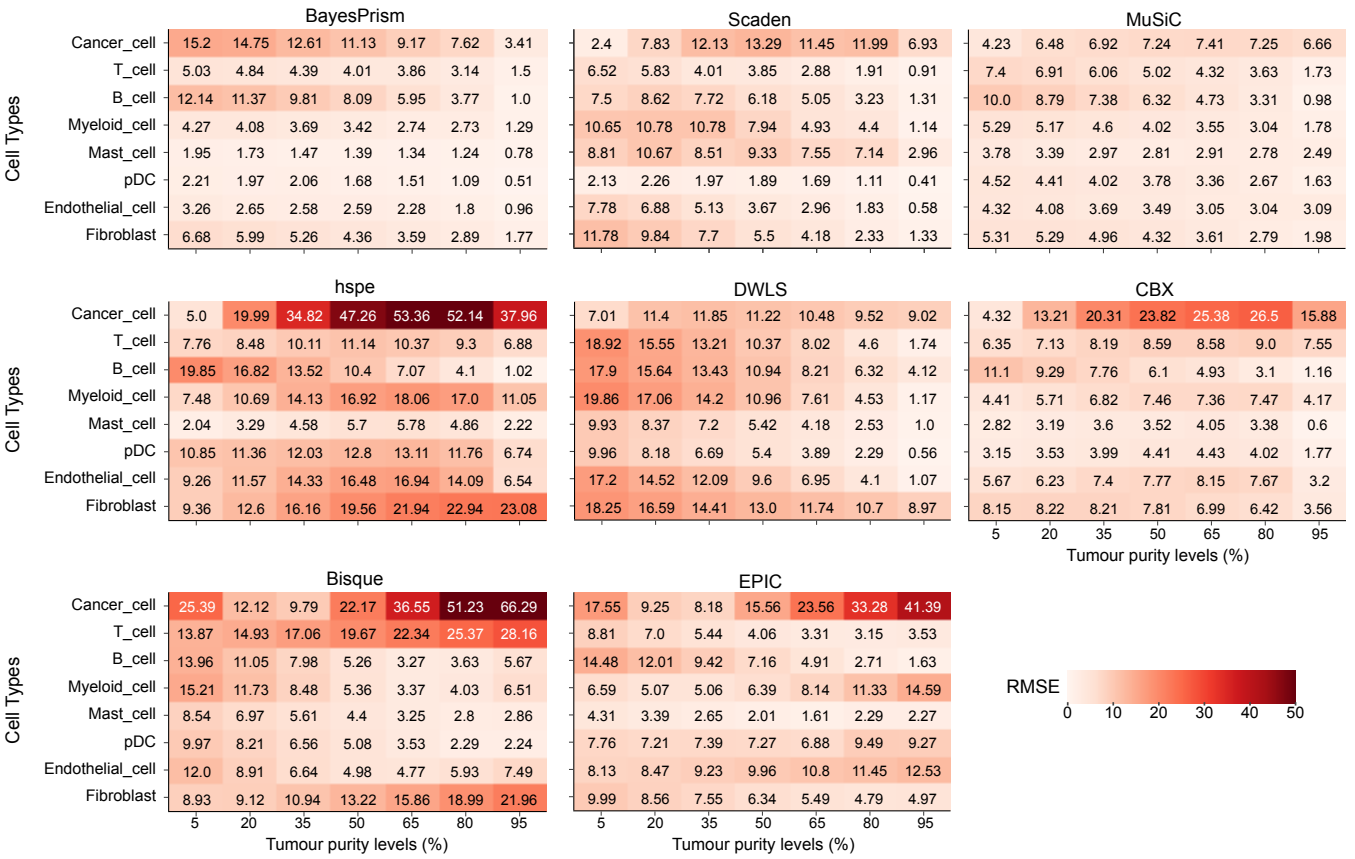


Supplementary Fig. 6: Direction of mis-predictions of each cell type for nine computational methods. Raw prediction errors (Predicted – Actual Cell Fractions) between predicted and ground truth cell compositions of nine major cell types across 7 tumour purity levels (from 5% to 95%, 15% interval). Higher positive and lower negative raw prediction errors represent worse performance. Each of the nine plots represents raw prediction errors for one deconvolution method. Negative raw prediction errors represent under-estimation (predicted fraction lower than ground truth), positive raw prediction errors represent over-estimation (prediction fraction higher than ground truth). Upper and lower whiskers depict raw prediction errors outside of the centre 50%. Each box represents the middle 50% of raw prediction errors, which includes the first quartile (Q1), the median, and the third quartile (Q3). Upper and lower whiskers depict maxima and minima of raw prediction errors, excluding outliers. Outliers are raw prediction errors that are more than 1.5x the interquartile range from either Q1 or Q3. Source data are provided as a Source Data file.

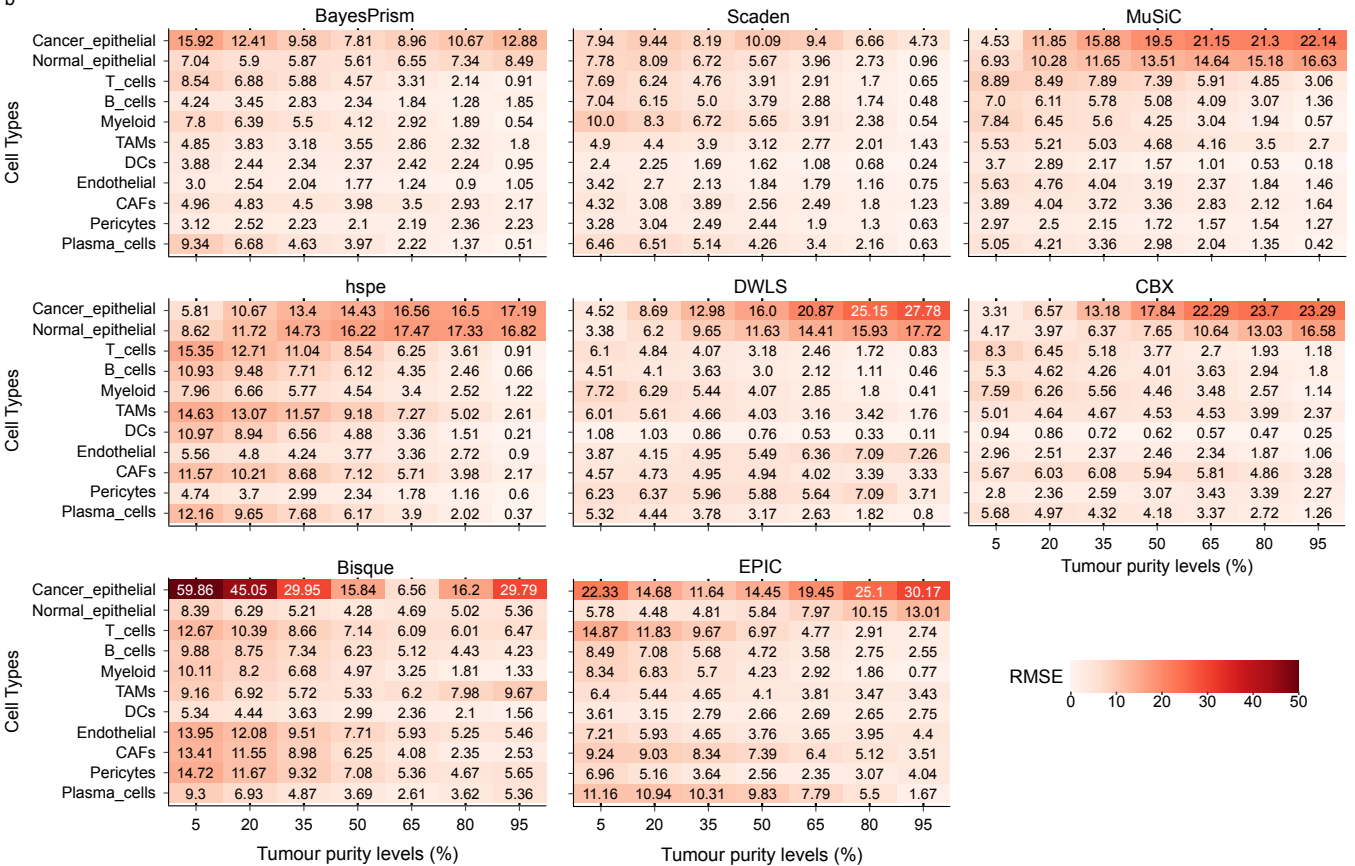


Supplementary Fig. 7: Overall generalisation of impact of tumour purity on deconvolution to simulated bulk mixtures generated using scRNA-Seq from Bassez et al and Pal et al. **a)** Bray-Curtis dissimilarity values predicted and ground truth cell compositions across 7 tumour purity levels (from 5% to 95%, 15% interval) using scRNA-Seq from Bassez et al. n=2,000 artificial bulk at each purity level. Each box represents the middle 50% of Bray-Curtis values, which includes the first quartile (Q1), the median, and the third quartile (Q3). Upper and lower whiskers depict maxima and minima of Bray-Curtis values, excluding outliers. Outliers are Bray-Curtis values that are more than 1.5x the interquartile range from either Q1 or Q3. Higher Bray-Curtis dissimilarity indicates poorer performance. **b)** Median Bray-Curtis dissimilarity between predicted and actual cell compositions across 19 tumour purity levels (from 5% to 95%, at 5% intervals) which are shown as proportions (x-axis) for nine computational methods (y-axis). There were 2,000 artificial bulk mixtures from Bassez et al per purity level. Darker shade of teal represents higher Bray-Curtis dissimilarity (worse performance). The numeric values are shown. Upper and lower whiskers depict Bray-Curtis values outside of the centre 50%. **c)** Bray-Curtis dissimilarity predicted and ground truth cell compositions across 7 tumour purity levels (from 5% to 95%, 15% interval) using scRNA-Seq from Pal et al. n=2,000 artificial bulk at each purity level. Each box represents the middle 50% of Bray-Curtis values, which includes the first quartile (Q1), the median, and the third quartile (Q3). Upper and lower whiskers depict maxima and minima of Bray-Curtis values, excluding outliers. Outliers are Bray-Curtis values that are more than 1.5x the interquartile range from either Q1 or Q3. Higher Bray-Curtis dissimilarity indicates poorer performance. **d)** Median Bray-Curtis dissimilarity between predicted and actual cell compositions across 19 tumour purity levels (samples with tumour cells ranging from 5% to 95%, at 5% intervals) which are shown as proportions (x-axis) for nine computational methods (y-axis). There were 2,000 artificial bulk mixtures from Pal et al per purity level. Darker shade of teal represents higher Bray-Curtis dissimilarity (worse performance). The numeric values are shown. Upper and lower whiskers depict Bray-Curtis values outside of the centre 50%. For all sub-figures, deconvolution methods, excluding CPM, are organised in the same order as Fig. 2a and Supplementary Fig. 4 for comparison. Source data are provided as a Source Data file.

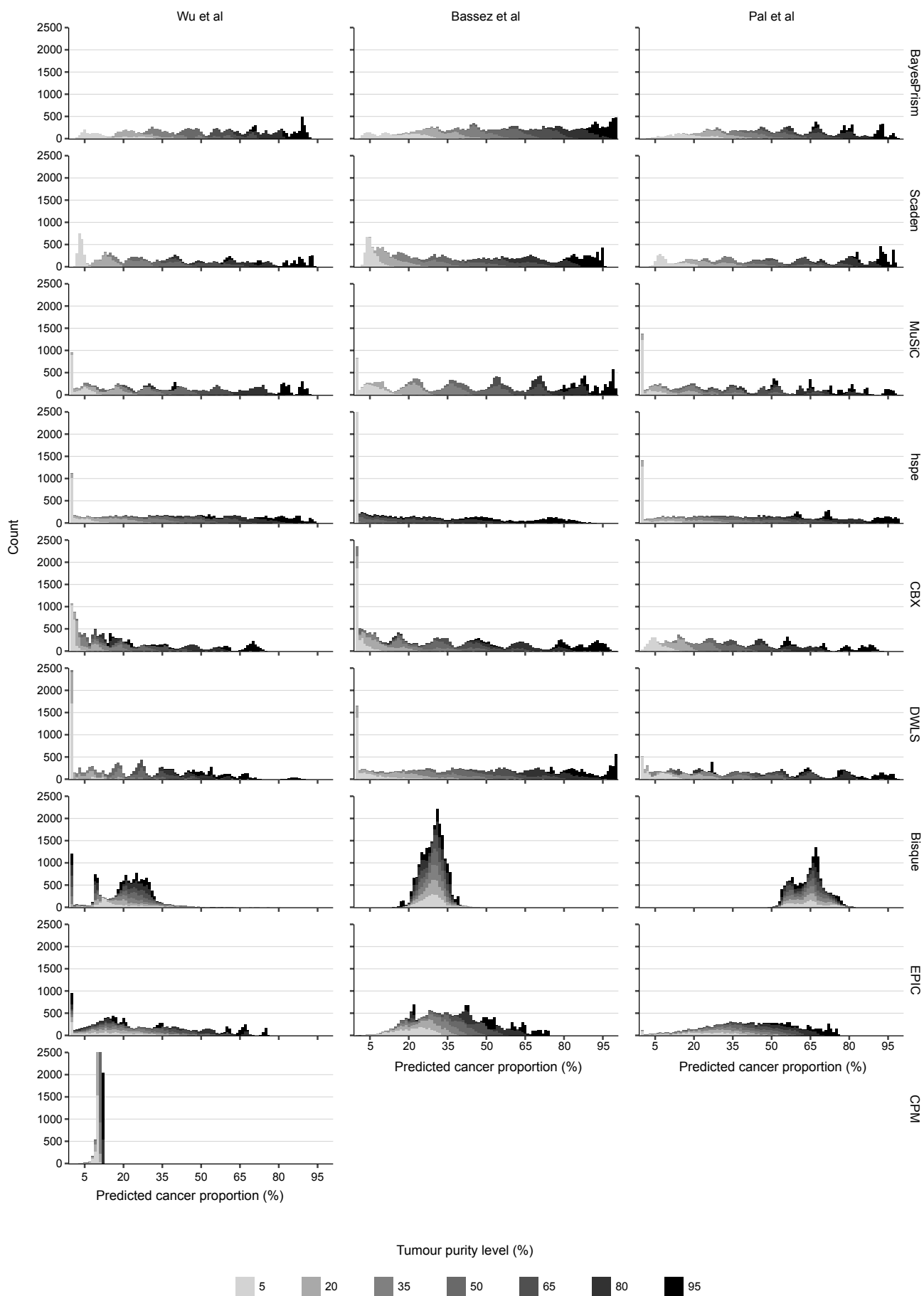
a



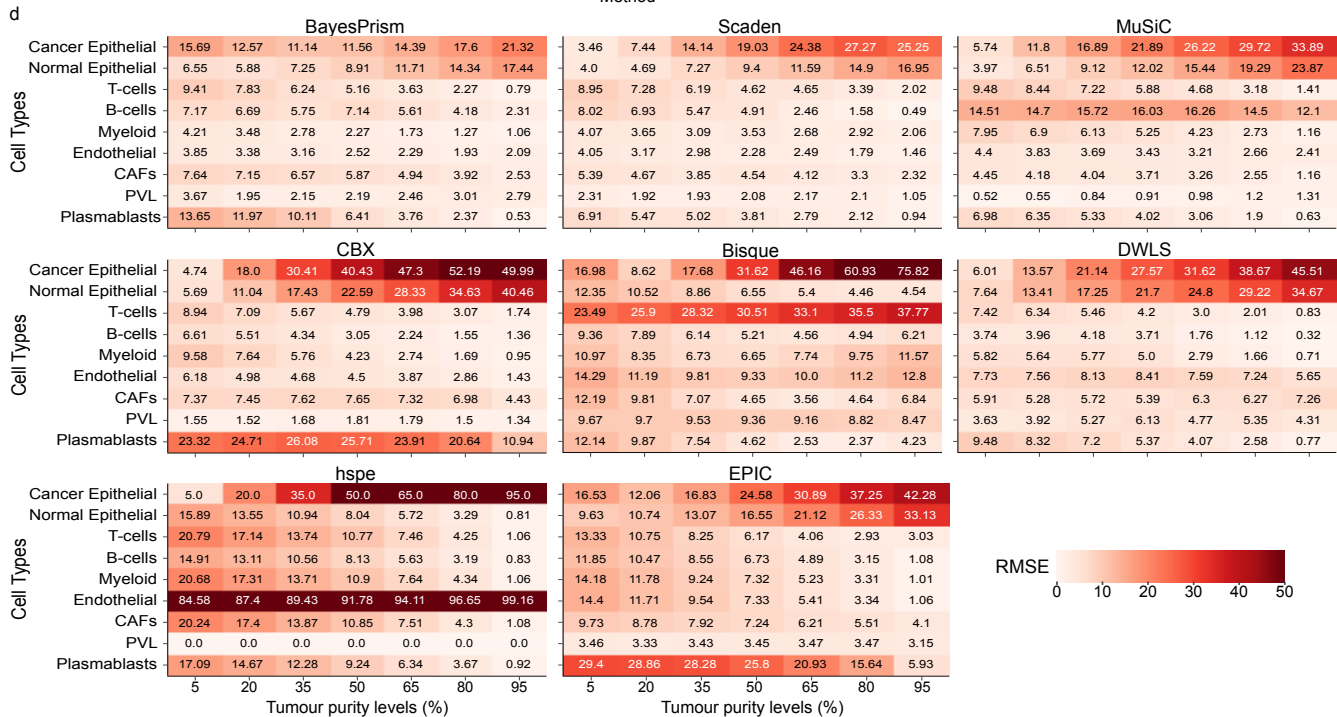
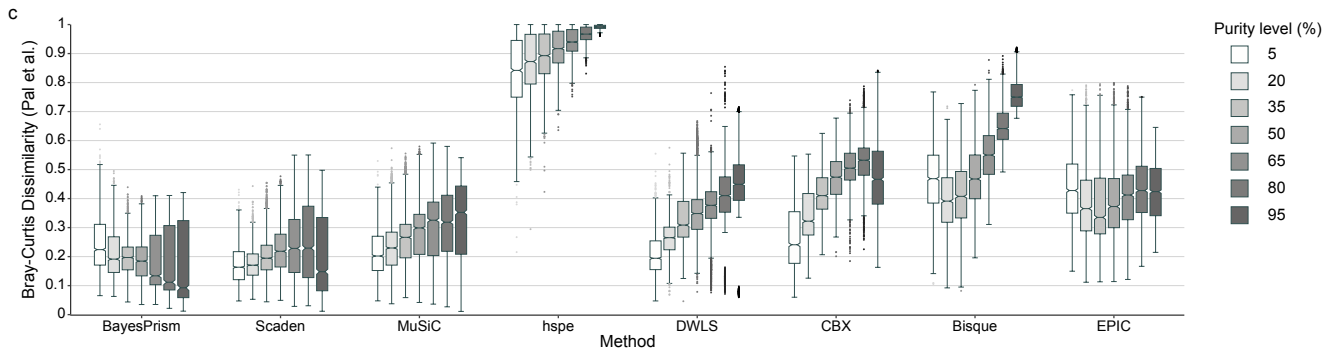
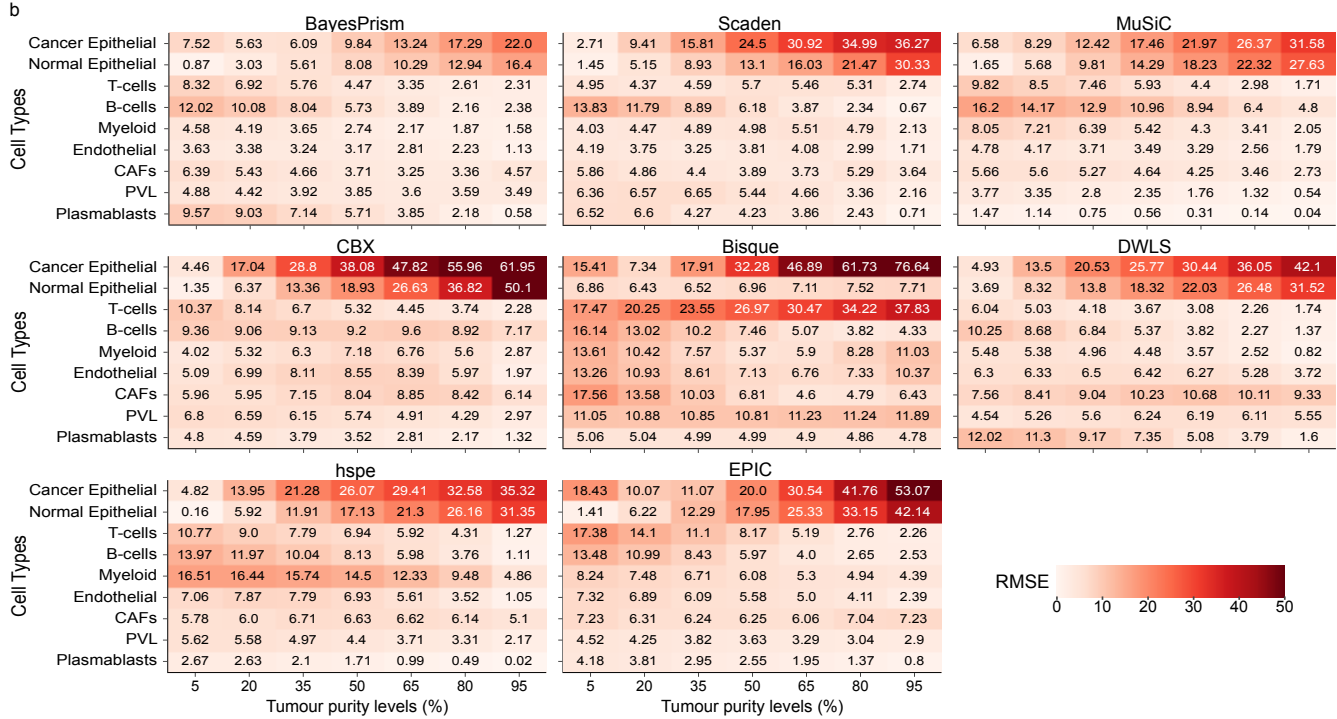
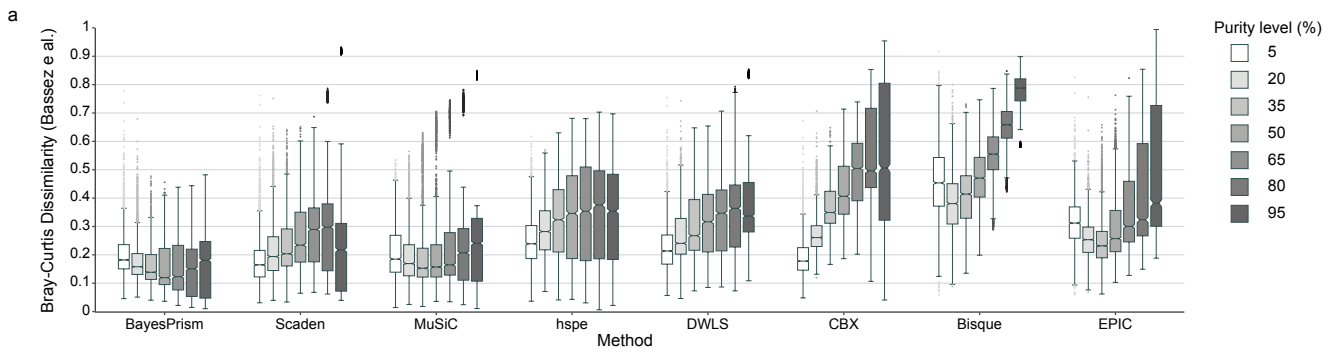
b



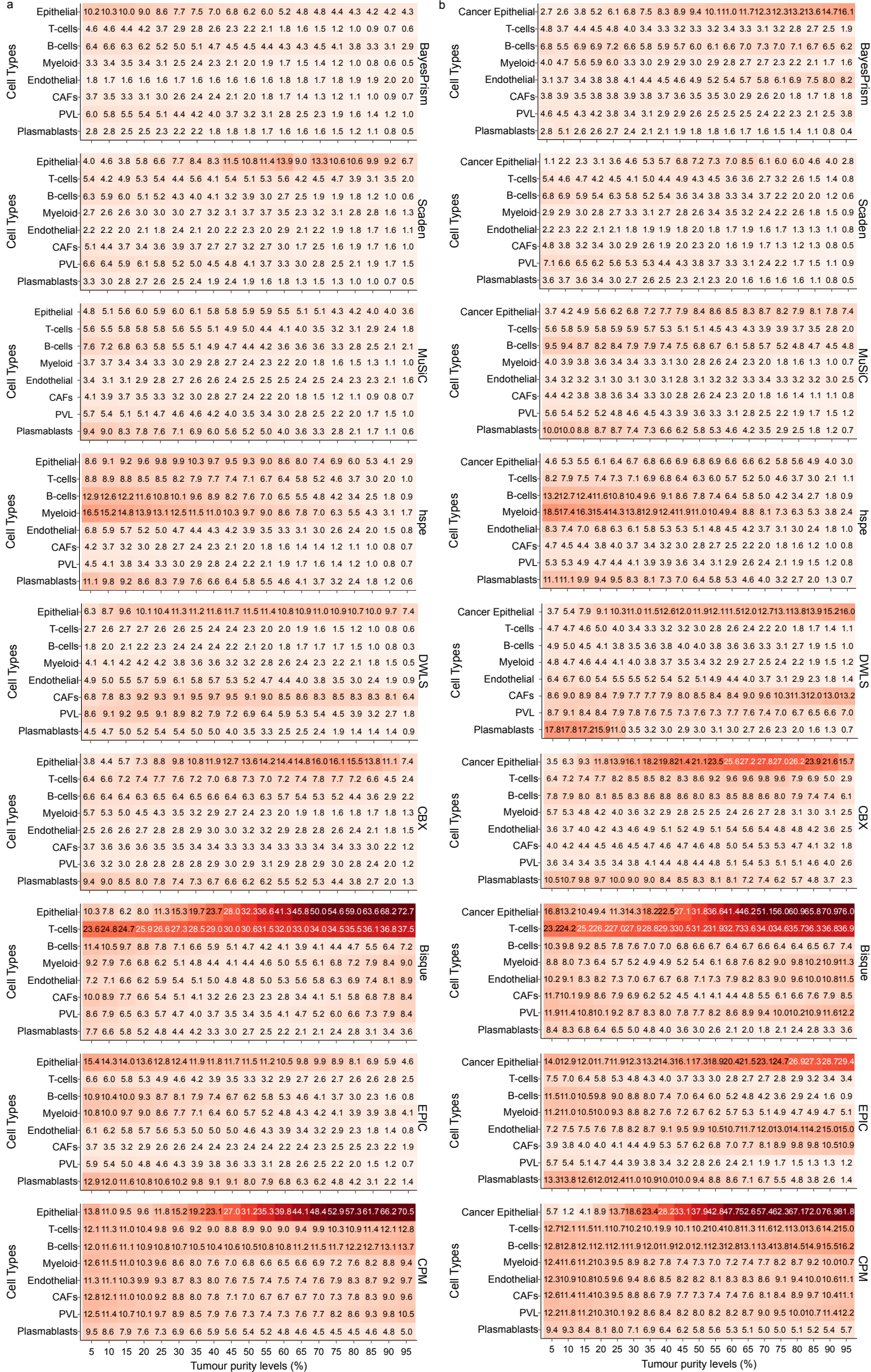
Supplementary Fig. 8: Cell-type-specific generalisation of impact of tumour purity on deconvolution to simulated bulk mixtures generated using scRNA-Seq from Bassez et al and Pal et al. **a)** Median RMSE between predicted and actual cell compositions of simulated bulk mixtures generated using scRNA-Seq from Bassez et al, aggregated by cell type. Seven tumour purity levels are shown (from 5% to 95%, 15% interval). Darker shade of red represents higher RMSE values (worse performance), with numeric RMSE values shown. Major cell types (y-axis) are organised into three categories: cancer (Cancer_cell), immune (T_cell, B_cell, Myeloid_cell, Mast_cell and pDC), and stromal cells (Endothelial_cell and Fibroblast). **b)** RMSE between predicted and actual cell compositions of simulated bulk mixtures generated using scRNA-Seq from Bassez et al, aggregated by cell type. Seven tumour purity levels are shown (from 5% to 95%, 15% interval). Darker shade of red represents higher RMSE values (worse performance), with numeric RMSE values shown. Major cell types (y-axis) are organised into three categories: epithelial (Cancer_epithelial and Normal_epithelial), immune (T_cells, B_cells, Myeloid, TAMs and DCs), and stromal cells (Endothelial, CAFs, Pericytes and Plasma_cells). TAMs: tumour-associated macrophage, DCs: dendritic cells, CAFs: cancer-associated fibroblast, pDC: plasmacytoid dendritic cell, RMSE: Root Mean Square Error, scRNA-Seq: Single-cell RNA Sequencing. Source data are provided as a Source Data file.



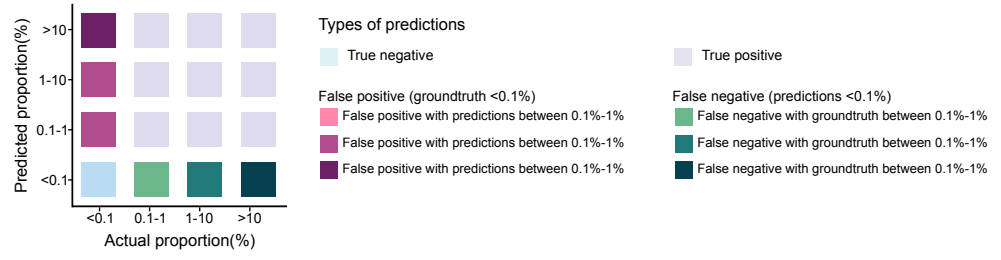
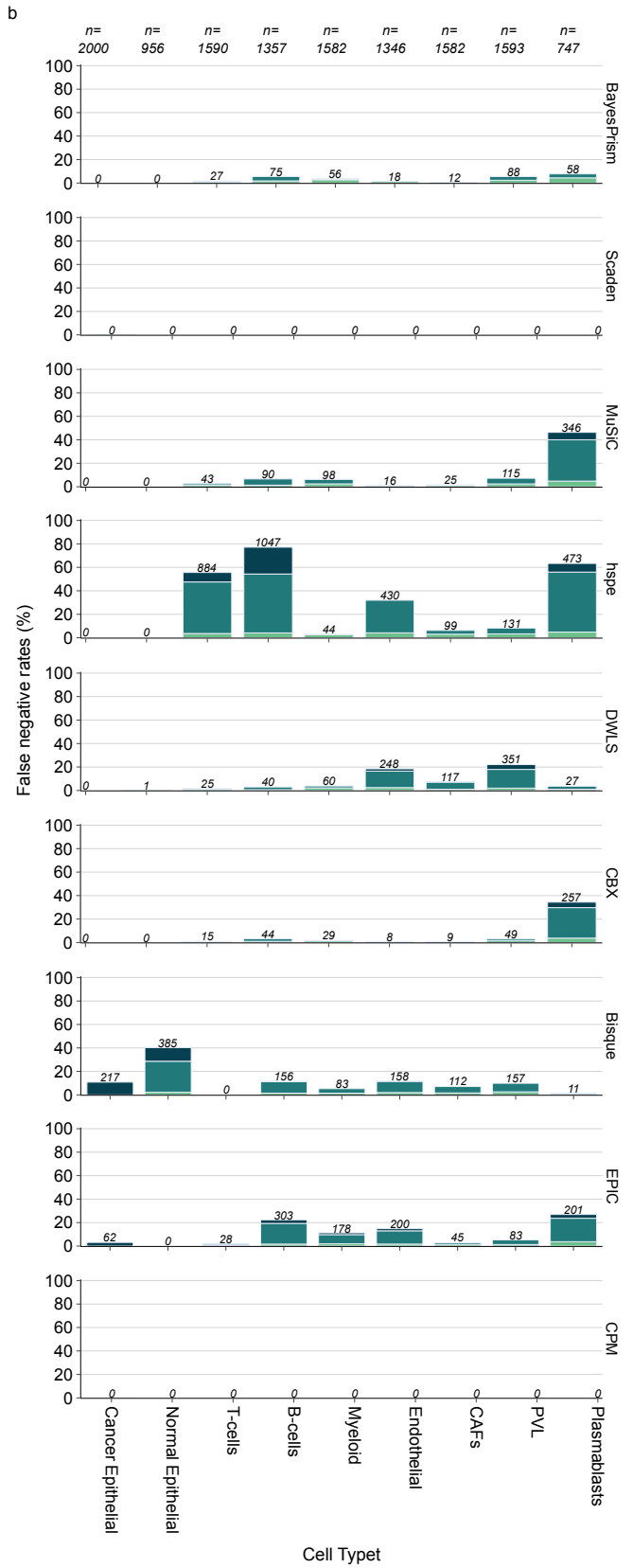
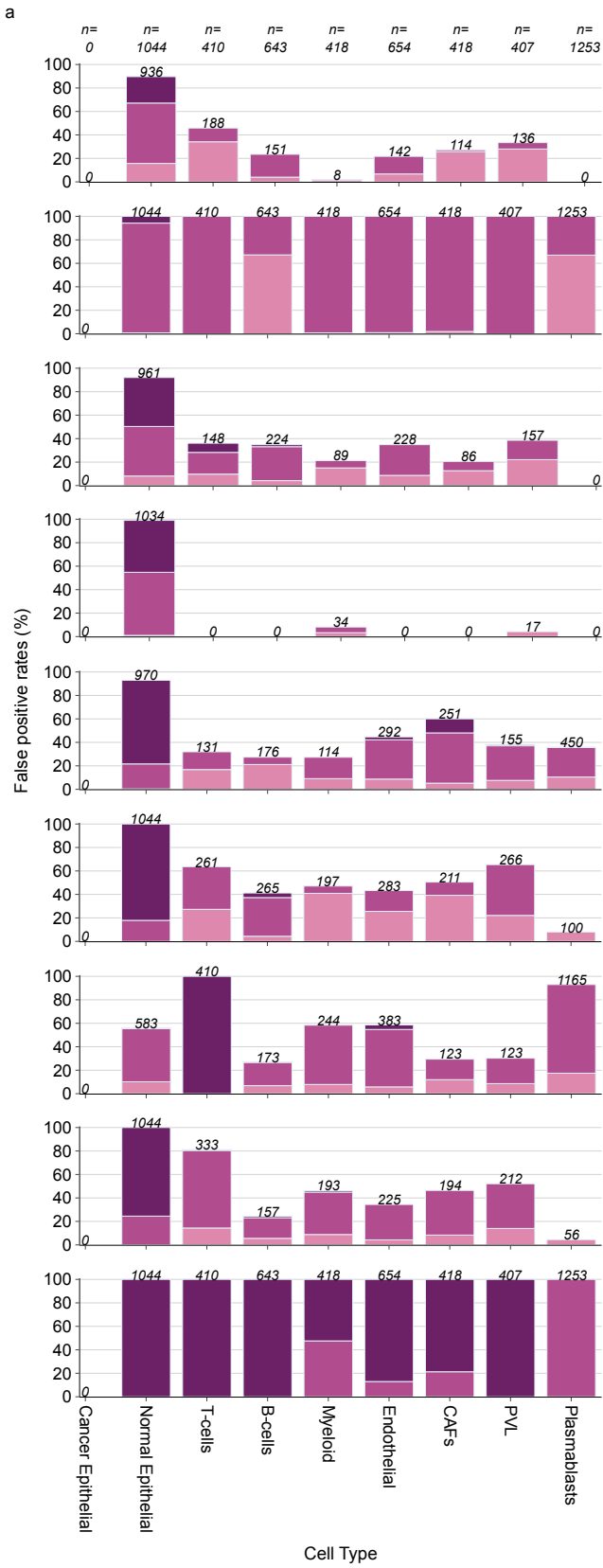
Supplementary Fig. 9: Histogram of predicted cancer proportions. Histogram distributions (*bin_size*=100) of predicted cancer proportions from simulated mixtures generated using single cells from Wu et al (left column), Bassez et al (middle column) and Pal et al (right column). Each histogram is aggregated and coloured by 7 tumour purity levels (from 5% to 95%, 15% interval, depicted as tick values on x-axes), and represents performance of each method in estimating cancer populations. x-axes represent predicted cancer proportions by each method, aggregated into 100 bins between 0% and 100% (i.e. bin width is 1%). y-axes represent count of the predicted bin on x-axis. The more accurate the prediction, the more each coloured histogram will be centred at the corresponding tumour purity level. A perfectly predicting model would produce 7 straight lines on top of the 7 annotated tumour purity levels. Source data are provided as a Source Data file.



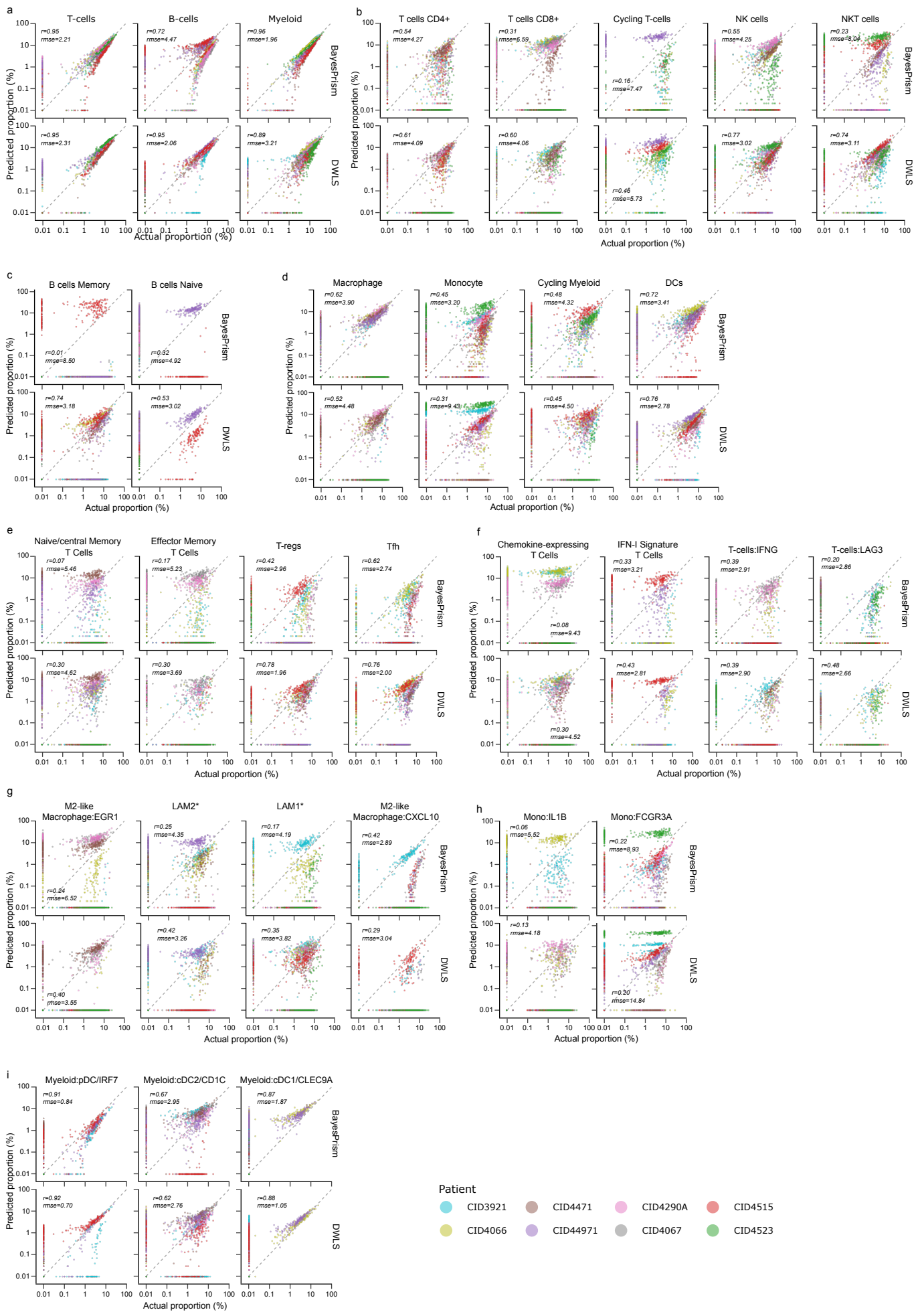
Supplementary Fig. 10: Studying the impact of tumour purity on deconvolution to simulated bulk mixtures generated using scRNA-Seq from Bassez et al and Pal et al, and single-cell reference from Wu et al. **a)** Bray Curtis dissimilarity predicted and ground truth cell compositions across 7 tumours purity levels (from 5% to 95%, 15% interval) using scRNA-Seq from Bassez et al. n=2,000 artificial bulk at each purity level. Each box represents the middle 50% of Bray-Curtis values, which includes the first quartile (Q1), the median, and the third quartile (Q3). Upper and lower whiskers depict maxima and minima of Bray-Curtis values, excluding outliers. Outliers are Bray-Curtis values that are more than 1.5x the interquartile range from either Q1 or Q3. Higher Bray-Curtis dissimilarity indicates poorer performance. **b)** RMSE between predicted and actual cell compositions of artificial bulk mixtures from Bassez et al, aggregated by cell type. Seven tumour purity levels are shown (from 5% to 95%, 15% interval). Darker shade of red represents higher RMSE values (worse performance), with numeric RSME values shown. Major cell types (y-axis) are organised into three categories: epithelial (Cancer_epithelial and Normal_epithelial), immune (T_cells, B_cells, Myeloid, TAMs and DCs), and stromal cells (Endothelial, CAFs, Pericytes and Plasma_cells). **c)** Bray Curtis dissimilarity predicted and ground truth cell compositions across 7 tumours purity levels (from 5% to 95%, 15% interval) using scRNA-Seq from Pal et al. Deconvolution methods, excluding CPM, are organised in the same order as Fig. 2a for comparison. n=2,000 artificial bulk at each purity level. Each box represents the middle 50% of Bray-Curtis values, which includes the first quartile (Q1), the median, and the third quartile (Q3). Upper and lower whiskers depict maxima and minima of Bray-Curtis values, excluding outliers. Outliers are Bray-Curtis values that are more than 1.5x the interquartile range from either Q1 or Q3. Higher Bray-Curtis dissimilarity indicates poorer performance. **d)** RMSE between predicted and actual cell compositions of artificial bulk mixtures from Pal et al, aggregated by cell type. Seven tumour purity levels are shown (from 5% to 95%, 15% interval). Darker shade of red represents higher RMSE values (worse performance), with numeric RSME values shown. Major cell types (y-axis) are organised into three categories: epithelial (Cancer_epithelial and Normal_epithelial), immune (T_cells, B_cells, Myeloid, TAMs and DCs), and stromal cells (Endothelial, CAFs, Pericytes and Plasma_cells). TAMs: tumour-associated macrophage, DCs: dendritic cells, CAFs: cancer-associated fibroblast, RMSE: Root Mean Square Error, scRNA-Seq: Single-cell RNA Sequencing. Source data are provided as a Source Data file.



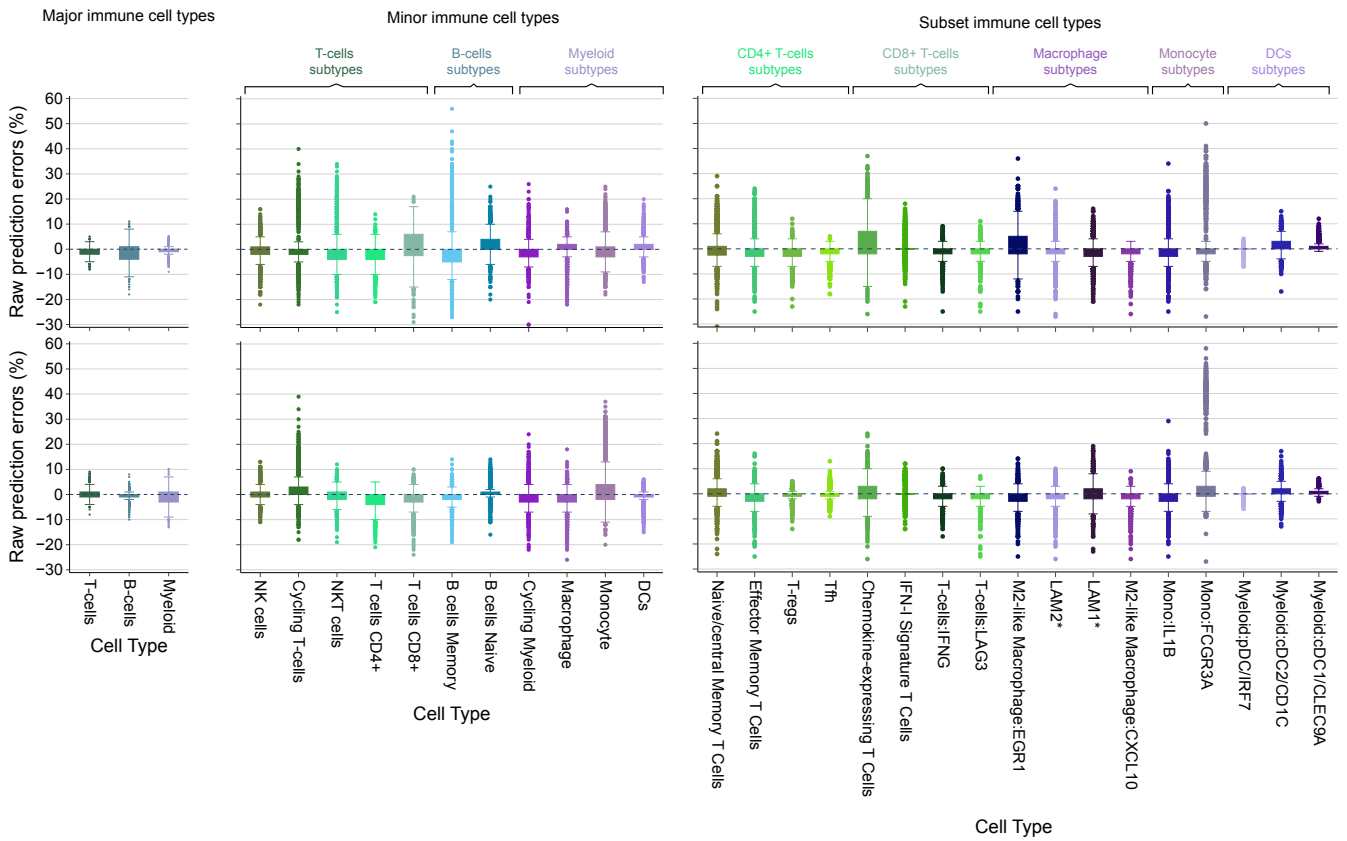
Supplementary Fig. 11: Cell-type specific deconvolution performance for nine methods with tumour and normal epithelial cells grouped as epithelial. a) Heatmap of RMSE of eight cell types across 2,000 mixtures and 19 tumour purity levels (from 5% to 95%, 5% interval) with normal epithelial and cancer epithelial predictions collapsed into 'epithelial' cell type. **b)** Heatmap of RMSE of eight cell types across 2,000 artificial mixtures and 19 tumour purity levels (from 5% to 95%, 5% interval) with no normal epithelial cells. For both **a)** and **b)**, darker shade of red represents higher RMSE values and worse deconvolution performance. The numeric RMSE values are shown. RMSE: Root Mean Square Error. Source data are provided as a Source Data file.



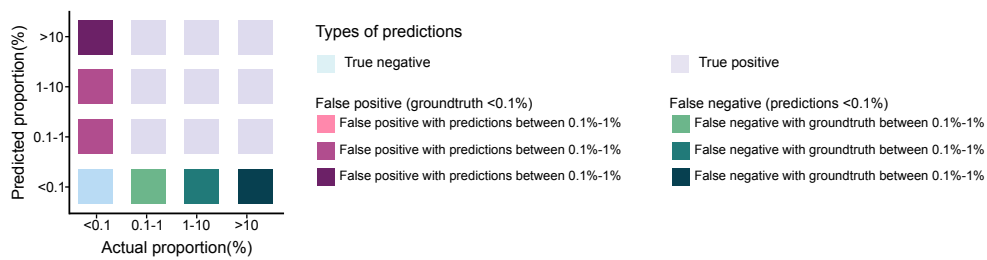
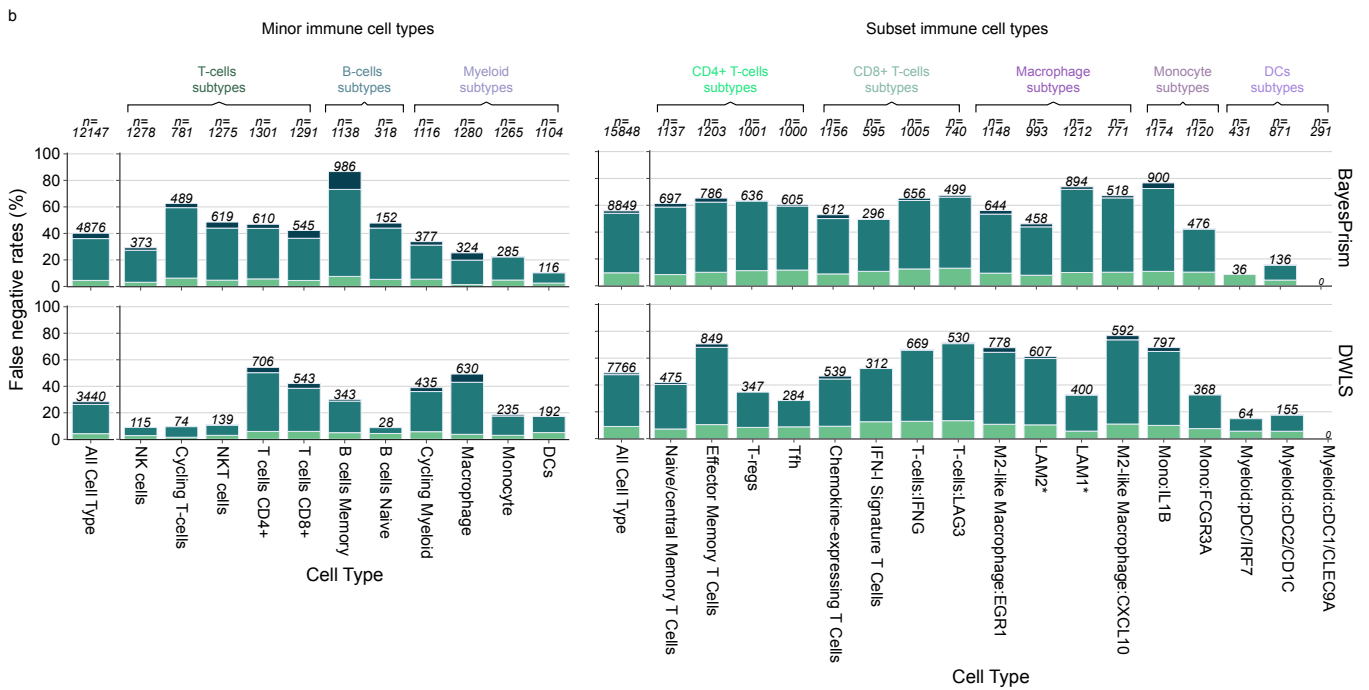
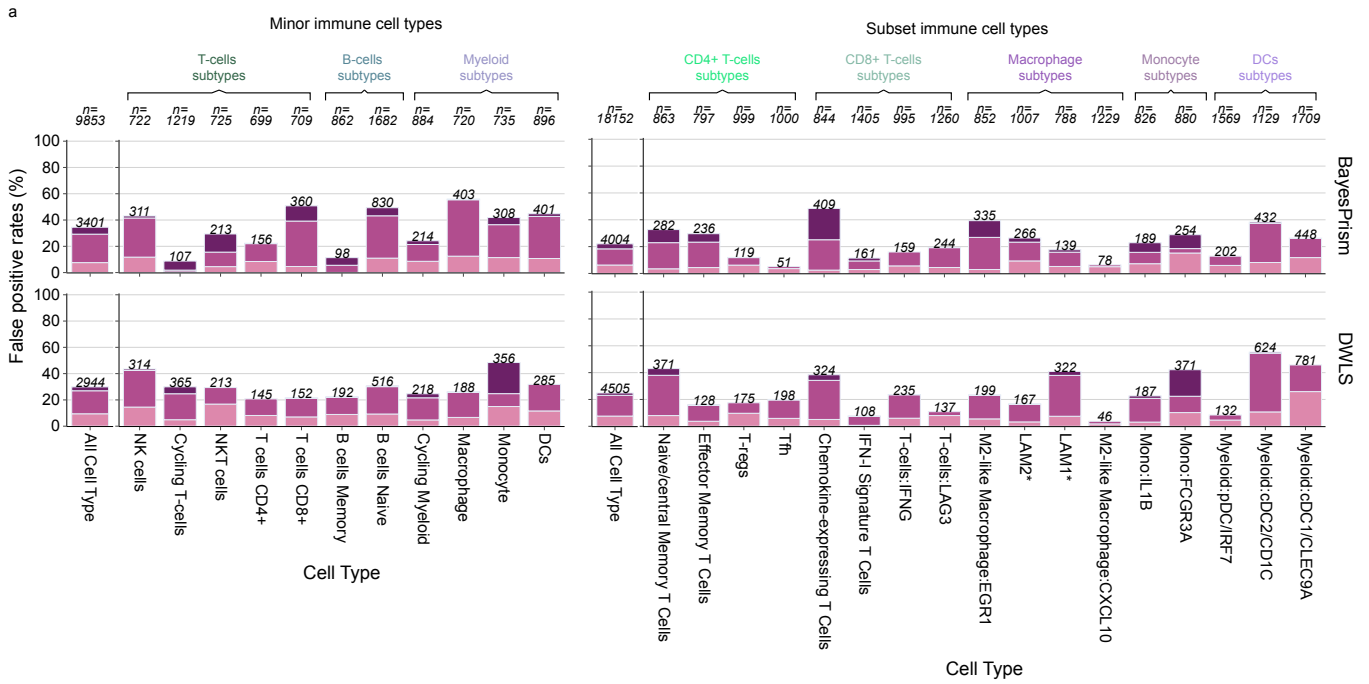
Supplementary Fig. 12: The performance of the nine deconvolution methods assessed by false positive and false negative rates. **a)** Percentages of the three levels of false positives out of the total number of false positives and true negatives (actual proportion less than 0.1%). Percentages and counts (shown above each bar) of false positives are shown for each individual cell type. The total number of missing cell type components used to determine false positive rates for all methods is shown at the top with the “*n*=” prefix. **b)** Percentages of the three levels of false negatives out of total number of all false negatives and true positives (predicted proportions less than 0.1%). Percentages and counts (shown above each bar) of false negatives are shown for each individual cell type. The total number of cell type components used to determine false negative rates for all methods is shown at the top with the “*n*=” prefix. Figure legend illustrates definitions of true negative, false positive, true positive, and false negative predictions. Source data are provided as a Source Data file.



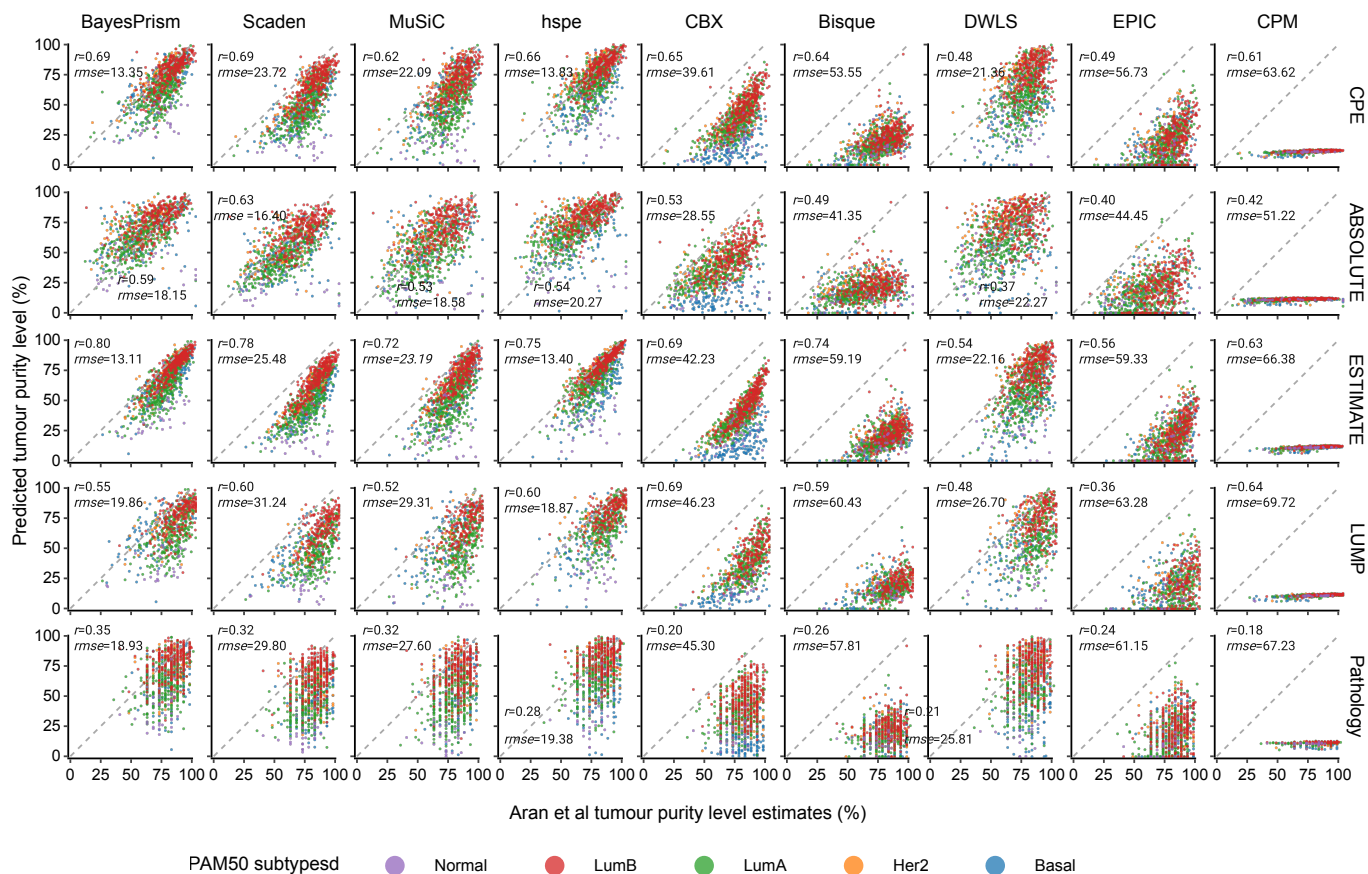
Supplementary Fig. 13: Cross-lineage performance of BayesPrism and DWLS by patient. Scatter plots of predicted (y-axis) versus actual (x-axis) cell compositions of BayesPrism and DWLS for major **(a)**, minor **(b-d)** and subset immune cell types **(e-i)**. Each point represents a test mixture component, with its colour representing one of eight test patients (from which single cells were used to generate the mixture). Dotted 45-degree diagonal line represents perfect prediction where predicted proportions match actual proportions. Cell types represented in each plot are three major immune cell types **(a)**, five minor cell types of T-cells **(b)**, two minor cell types of B-cells **(c)**, four minor myeloid cell types **(d)**, four subset cell types of CD8+ T-cells **(e)**, four subset cell types of CD4+ T-cells **(f)**, four subset macrophage cell types **(g)**, two subset monocyte cell types **(h)**, and three subset cell types of dendritic cells **(i)**. Each scatterplot is annotated with the associated Pearson's r correlation coefficient and RMSE values. RMSE: Root Mean Square Error. Source data are provided as a Source Data file.



Supplementary Fig. 14: Raw prediction errors of BayesPrism and DWLS across immune lineages. Raw prediction errors (Predicted – Actual Cell Fractions) between predicted and ground truth cell compositions for BayesPrism (top) and DWLS (bottom). $n=2,000$ artificial bulk mixtures with 50% tumour purity. Higher positive and lower negative raw prediction errors represent poorer performance. Raw prediction error values are aggregated into major (left sub-plot), minor (middle sub-plot), and subset (right sub-plot) immune cell types. Minor cell types include five subtypes of T-cells, two subtypes of B-cells, and four subtypes of Myeloid. Subset cell types include eight subtypes of T-cells (four subtypes of CD4+ T cells and four subtypes of CD8+ T-cells), and nine subtypes of Myeloid (four subtypes of macrophages, two subtypes of monocyte, and three subtypes of DCs). Zero line indicates a perfect match between prediction and ground truth. Each box represents the middle 50% of raw prediction errors, which includes the first quartile (Q1), the median, and the third quartile (Q3). Upper and lower whiskers depict maxima and minima of raw prediction errors, excluding outliers. Outliers are raw prediction errors that are more than 1.5x the interquartile range from either Q1 or Q3. Source data are provided as a Source Data file.



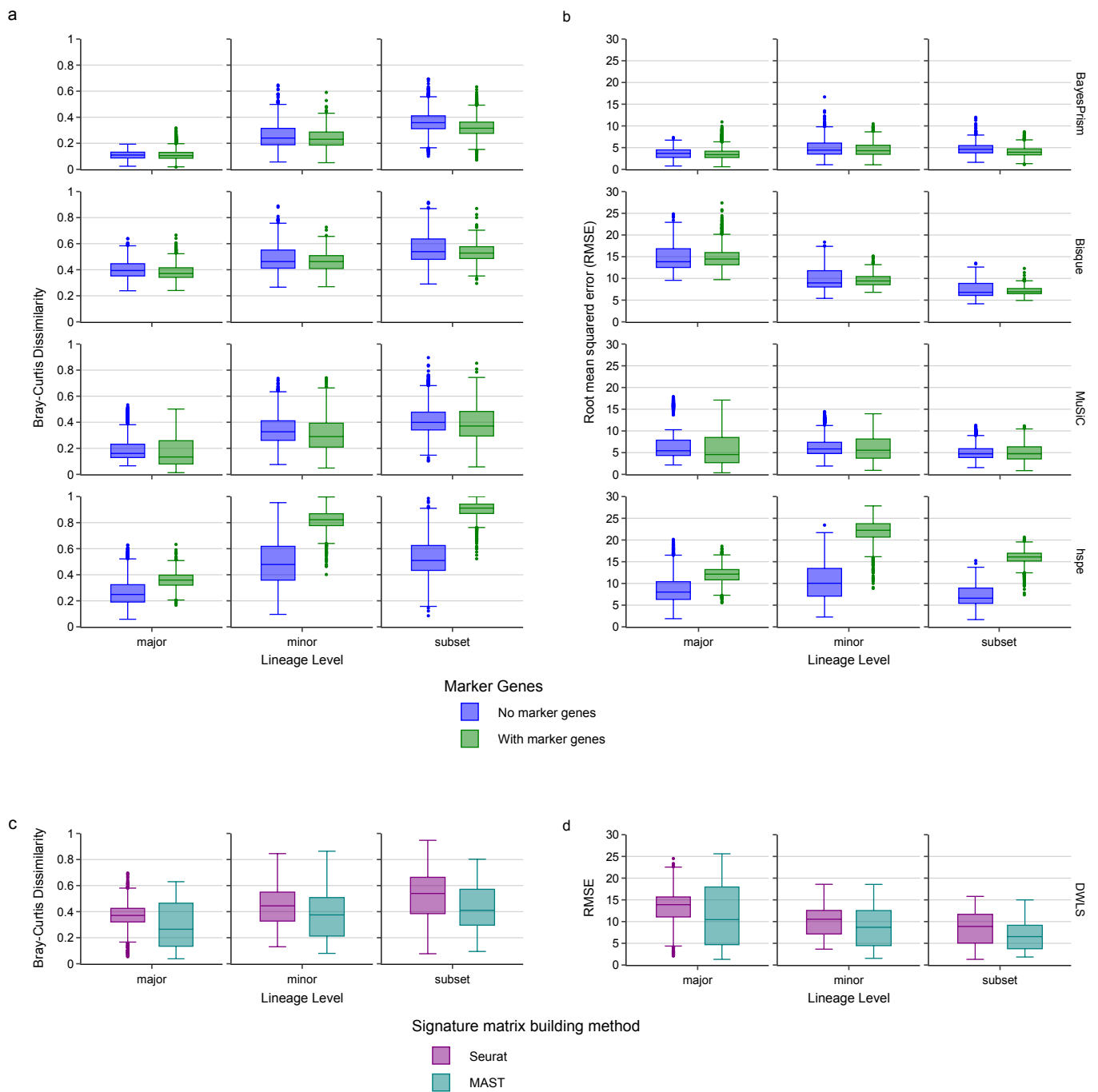
Supplementary Fig. 15: The performance of BayesPrism and DWLS methods for minor and subset immune cell types assessed by false positive and false negative rates. **a)** Percentages of the three levels of false positives out of the total number of false positives and true negatives across minor (left) and subset (right) immune cell types. Percentages and counts (shown above each bar) of false positives are either aggregated across all cell types (first stacked bar) or shown for each individual cell type. The total number of missing cell type components used to determine false positive rates for both methods is shown at the top with the “*n*=” prefix. **b)** Percentages of the three levels of false negatives out of total number of all false negatives and true positives across minor (left) and (subset) immune cell types. Percentages and counts (shown above each bar) of false negatives are either aggregated across all cell types (first stacked bar) or shown for each individual cell type. The total number of cell type components used to determine false negative rates for all methods is shown at the top with the “*n*=” prefix. Figure legend illustrates definitions of true negative, false positive, true positive, and false negative predictions. Source data are provided as a Source Data file.



Supplementary Fig. 16: Comparison of nine deconvolution approaches to four alternative methods to predict tumour purity in TCGA breast data. Scatter plot comparing purity levels determined by the nine deconvolution methods (columns) with the four alternative methods (ABSOLUTE, ESTIMATE, LUMP and Pathology) and their Consensus Purity Estimates (CPE) provided by Aran et al. The data comprised $n=968$ breast cancers from TCGA, with each sample coloured by the PAM50 subtype. The y-axis for each plot is the percentage of tumour cells predicted from the relevant deconvolution tool and the x-axis is the percentage of tumour cells predicted from the orthogonal method. Each scatterplot is annotated with the associated Pearson's r correlation coefficient and RMSE values. RMSE: Root Mean Square Error. Source data are provided as a Source Data file.



Supplementary Fig. 17: Performance of BayesPrism before and after Gibbs sampling on artificial bulk mixtures generated using scRNA-seq data from Wu et al, Bassez et al and Pal et al. RMSE between predicted and actual cell compositions of artificial bulk mixtures using scRNA-Seq from Wu et al (a, b), Bassez et al (c, d) and Pal et al (e, f), aggregated by cell type. Single-cell reference profiles and train simulated mixtures (only for Scaden) were from Wu et al, with intersecting genes with Bassez et al and Pal et al when appropriate. Left column (a, c, e) depict BayesPrism performance before Gibbs sampling, while right column (b, d, f) depict its performance after Gibbs sampling. Seven tumour purity levels are shown (from 5% to 95%, 15% interval). Darker shade of red represents higher RMSE values (worse performance), with numeric RMSE values shown. CAFs: cancer-associated fibroblast, RMSE: Root Mean Square Error, scRNA-Seq: Single-cell RNA Sequencing. Source data are provided as a Source Data file.



Supplementary Fig. 18: Performance comparison of different versions of BayesPrism, Bisque, MuSiC, hspe and DWLS. Bray-Curtis dissimilarity (**a, c**) and Root-mean-squared errors (RMSE) (**b, d**) between predicted and ground truth cell compositions of BayesPrism, Bisque, MuSiC and hspe when marker genes are used (green) or not used (blue), and of DWLS when either Seurat (purple) or MAST (teal) is used for building its internal signature matrix. $n=2,000$ artificial bulk mixtures generated at 50% tumour purity and either major (top row), minor (middle row), or subset (bottom row) immune cell types. Each box represents the middle 50% of Bray-Curtis/RMSE values, which includes the first quartile (Q1), the median, and the third quartile (Q3). Upper and lower whiskers depict maxima and minima of Bray-Curtis/ RMSE values, excluding outliers. Outliers are Bray-Curtis/ RMSE values that are more than 1.5x the interquartile range from either Q1 or Q3. Higher Bray-Curtis dissimilarity and RMSE indicates poorer performance. Source data are provided as a Source Data file.

Supplementary references

1. Chu, T., Wang, Z., Pe'er, D. & Danko, C. G. Cell type and gene expression deconvolution with BayesPrism enables Bayesian integrative analysis across bulk and single-cell RNA sequencing in oncology. *Nat. Cancer* **3**, 505–517 (2022).
2. Casella, G. & George, E. I. Explaining the Gibbs Sampler. 9 (2022).
3. Jew, B. *et al.* Accurate estimation of cell composition in bulk expression through robust integration of single-cell information. *Nat. Commun.* **11**, 1971 (2020).
4. Newman, A. M. *et al.* Determining cell type abundance and expression from bulk tissues with digital cytometry. *Nat. Biotechnol.* **37**, 773–782 (2019).
5. Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue expression profiles. *Nat. Methods* **12**, 453–457 (2015).
6. Frishberg, A. *et al.* Cell composition analysis of bulk genomics using single-cell data. *Nat. Methods* **16**, 327–332 (2019).
7. Tsoucas, D. *et al.* Accurate estimation of cell-type composition from gene expression data. *Nat. Commun.* **10**, 2975 (2019).
8. Racle, J., de Jonge, K., Baumgaertner, P., Speiser, D. E. & Gfeller, D. Simultaneous enumeration of cancer and immune cell types from bulk tumor gene expression data. *eLife* **6**, e26476 (2017).
9. Gregory J. Hunt & Johann A. Gagnon-Bartsch. The role of scale in the estimation of cell-type proportions. *Ann. Appl. Stat.* **15**, 270–286 (2021).
10. Hunt, G. J., Freytag, S., Bahlo, M. & Gagnon-Bartsch, J. A. dtangle: accurate and robust cell type deconvolution. *Bioinformatics* **35**, 2093–2099 (2019).
11. Cobos, F. A., Alquicira-Hernandez, J., Powell, J. E., Mestdagh, P. & Preter, K. D. Benchmarking of cell type deconvolution pipelines for transcriptomics data. *Nat. Commun.* **11**, 14 (2020).
12. Wang, X., Park, J., Susztak, K., Zhang, N. R. & Li, M. Bulk tissue cell type deconvolution with multi-subject single-cell expression reference. *Nat. Commun.* **10**, 380 (2019).
13. Menden, K. *et al.* Deep learning–based cell composition analysis from tissue expression profiles. *Sci. Adv.* **6**, (2020).
14. Wu, S. Z. *et al.* A single-cell and spatially resolved atlas of human breast cancers. *Nat. Genet.* **53**, 1334–1347 (2021).

15. Lambrechts, D. *et al.* Phenotype molding of stromal cells in the lung tumor microenvironment. *Nat. Med.* **24**, 1277–1289 (2018).
16. Aran, D., Sirota, M. & Butte, A. J. Systematic pan-cancer analysis of tumour purity. *Nat. Commun.* **6**, 8971 (2015).