

1 **Contents**

2 Supplementary2

3 A key distinction from traditional deconvolution tasks2

4 Further validation on ROSMAP cohort2

5 Methods3

6 Mapping METABRIC samples to TCGA-defined niche subtypes3

7 Simulation of cfDNA methylation profiles and prediction of niche abundances4

8 Compatibility constraints of methods with proteomics data4

9 Details of Comparison Methods5

10

11

12 Supplementary

13 A key distinction from traditional deconvolution 14 tasks

15 Specifically, in the NanoString CosMx spatial single-cell dataset of non-small cell lung cancer
16 (NSCLC), the proportion of each cell's k -nearest neighbors sharing the same niche label (niche kNN
17 purity) and the proportion sharing the same cell-type label (cell-type kNN purity) were computed.
18 The results (Fig. 1a) showed that local niche consistency was substantially lower than that of cell
19 types (niche kNN purity [1] = 50.96%; cell-type kNN purity = 83.66%). This discrepancy suggests
20 that the dominant axes of variation captured by gene expression are more strongly aligned with cell-
21 type-associated transcriptional programs than with niche-level signals, which represent
22 multicellular cooperative units. Because a given niche typically comprises multiple cell types
23 coexisting within defined spatial neighborhoods, niche labels inherently correspond to mixed
24 expression distributions at the single-cell level across cell types, rather than a pseudo-cell type that
25 forms a tight cluster in expression space and can be separated by simple linear unmixing.

26 Therefore, recovering niche abundances from conventional tissue-level data is not a trivial
27 substitution of label sets, but a substantially more challenging identification problem: it requires
28 disentangling cross-cell-type coordinated functional patterns—driven by spatial neighborhood
29 constraints and intercellular interactions—from mixed signals dominated by cell-type variation.
30 Moreover, most existing deconvolution methods do not explicitly leverage structural priors carried
31 by spatial single-cell data, such as local continuity, spatial autocorrelation, and neighborhood
32 relationships. As a consequence, the strategy of first annotating niches and then applying
33 conventional deconvolution is face intrinsic challenges in statistical identifiability and robustness.
34 Motivated by these considerations, NicheDECODE is proposed to explicitly integrate spatial
35 structural evidence and functional module priors during both training and inference, improving the
36 separability of niche-level signals and enabling reliable estimation of niche abundances from
37 conventional tissue-level datasets.

38 [1] To quantify local label consistency, we defined k -nearest neighbor (kNN) purity based on the
39 neighborhood graph constructed from gene expression profiles. For each cell, we identified its k
40 nearest neighbors from the connectivity graph and computed the proportion of neighbors sharing
41 the same label (e.g., niche label or cell-type label) as the query cell. The kNN purity was then defined
42 as the average of these proportions across all cells. A higher kNN purity indicates stronger local
43 clustering of cells with the same label in expression space, reflecting greater separability and
44 homogeneity of the corresponding biological category.

45

46 Further validation on ROSMAP cohort

47 To further evaluate robustness, NicheDECODE models were independently trained by two
48 distinct spatial reference datasets (Chen et al. and Maynard et al.) and subsequently applied to
49 ROSMAP cohort. Notably, both training settings yielded highly consistent trends: white matter
50 (WM) abundance was positively correlated with Alzheimer's disease (AD) progression, whereas

51 gray matter abundance showed a negative correlation with Braak stage (Fig S3). These results
52 indicate that NicheDECODE reliably recovers spatial structural alterations from tissue-level data
53 and remains stable across reference sources.

54 Using the DLPFC single-cell atlas from Song *et al.* as reference, we performed cell-type
55 deconvolution of the ROSMAP cohort using CIBERSORTx (nine cell types; Fig. S5a). To assess
56 predictive relevance, we trained multi-layer perceptron (MLP) models to predict Braak stage using
57 three feature sets separately: (i) bulk gene expression, (ii) inferred cell-type proportions, and (iii)
58 inferred niche proportions (Fig. S5b). Across multiple hyperparameter settings, models based on
59 niche proportions consistently achieved the lowest mean squared error (MSE), indicating that niche-
60 level composition provides stronger predictive signal for AD progression than either gene
61 expression or cell-type proportions alone.

62 Furthermore, linear regression analysis (linregress function of SciPy) revealed that none of the
63 individual cell-type proportions were significantly associated with Braak stage (Fig S5c),
64 reinforcing the notion that niche-level organization captures disease-related variation beyond cell-
65 type abundance alone.

66 Methods

67 Mapping METABRIC samples to TCGA-defined niche subtypes

68 First, a sample-level niche feature matrix was constructed in the TCGA-BRCA cohort using
69 the niche abundance profiles of individual samples. For patients with multiple samples, the
70 abundance of each niche component was averaged across samples to obtain a single patient-level
71 profile. The resulting TCGA niche features were then standardized using z-score transformation. K-
72 means clustering was performed on the standardized profiles to classify TCGA patients into two
73 niche-defined subtypes. The number of clusters was set to ($k=2$), the random seed was fixed at 2026,
74 and multiple random initializations were used to improve clustering stability.

75 To map the METABRIC cohort onto the two niche subtypes defined in TCGA, the METABRIC
76 niche abundance profiles were first standardized using the means and standard deviations estimated
77 from the TCGA cohort. The standardized METABRIC profiles were then supplied to the K-means
78 model fitted in TCGA. When the prediction yielded both subtypes, each METABRIC sample was
79 assigned a TCGA-like niche subtype according to its nearest TCGA niche centroid.

80 In practice, direct nearest-centroid assignment classified all METABRIC samples into a single
81 TCGA niche subtype. We therefore used a centroid-axis projection strategy as an alternative
82 mapping approach. In the standardized TCGA niche feature space, the vector extending from the
83 TCGA C1 centroid to the TCGA C2 centroid was defined as the C1-to-C2 projection axis. Each
84 standardized METABRIC sample was projected onto this axis to obtain a continuous projection
85 score. Because the axis was oriented from C1 toward C2, a higher projection score indicated a more
86 TCGA C2-like niche state, whereas a lower score indicated a more TCGA C1-like state. Finally,
87 METABRIC samples were divided into C1-like and C2-like groups using the median projection
88 score within the METABRIC cohort as the cutoff. Samples with projection scores below the median
89 were classified as C1-like, whereas those with scores greater than or equal to the median were
90 classified as C2-like.

91 Simulation of cfDNA methylation profiles and prediction of niche 92 abundances

93 We first identified niche-associated methylation features using paired samples from the TCGA-
94 BRCA cohort. Samples with available RNA-sequencing data were used to match tissue HM450
95 methylation beta-value profiles with the seven niche abundances inferred by NicheDECODE. For
96 each niche state, samples were ranked according to their inferred abundance, and the top and bottom
97 quartiles were defined as the high- and low-abundance groups, respectively. Differentially
98 methylated CpG sites between the two groups were identified using a two-sided Wilcoxon rank-
99 sum test. CpG sites with ($P < 0.05$) and an absolute difference in mean beta value greater than 0.1
100 were retained. For each niche state, up to 5,000 positively associated and 5,000 negatively
101 associated CpG sites were selected.

102 To generate simulated cfDNA methylation profiles, the selected tumour tissue methylation
103 profiles were linearly mixed with healthy cfDNA methylation backgrounds. The healthy-
104 background fraction was randomly sampled between 0.2 and 0.6, and the tumour fraction was
105 defined as its complement. Each healthy-background profile was generated by randomly combining
106 three healthy cfDNA samples and introducing mild multiplicative noise. The resulting mixed CpG
107 beta values were used as model inputs, whereas the niche abundance profiles inferred from the
108 corresponding TCGA transcriptomic samples were retained as reference labels.

109 A multilayer perceptron model was implemented in PyTorch to predict the relative abundances
110 of the seven niche states from the simulated cfDNA methylation profiles. Samples were divided
111 equally into training and test sets using stratified sampling based on the previously defined TCGA
112 ($k=2$) niche subtypes. Model performance was evaluated in the independent test set by calculating
113 Pearson correlation coefficients between predicted and reference abundances for each niche state.
114 The predicted niche abundance profiles in the test set were subsequently standardized and grouped
115 into two clusters using spectral clustering. The prognostic relevance of the resulting groups was
116 evaluated using Kaplan–Meier analysis, the log-rank test and Cox proportional-hazards regression.

117 Compatibility constraints of methods with proteomics data

118 In both proteomics datasets, CIBERSORTx failed to detect significant differential proteins
119 across niches. This is likely attributable to the comparatively weaker inter-niche molecular contrast
120 in proteomic data relative to transcriptomic data, together with the substantially smaller number of
121 available molecular features.

122 RCTD could not be successfully applied to the CITE-seq dataset for similar reasons, as
123 insufficient discriminative molecular features were identified. In the CODEX dataset, RCTD was
124 further incompatible because the measured protein intensities are not provided in integer count
125 format, which is required by the method.

126 UCDSselect, as a large pretrained model developed on transcriptomic data, relies on modality-
127 specific learned representations that do not readily generalize to proteomics data.

128 For these reasons, the three methods were excluded from comparative analyses in the
129 proteomics benchmarking experiments.

130 Details of Comparison Methods

131 **ScpDeconv**

132 ScpDeconv is a deep learning-based method that employs autoencoders to effectively explore the
133 relationships between tissue proteomics data and single-cell proteomics data. It can interpolate
134 features with high information content that are not present in the single-cell reference, thereby
135 enhancing the reliability of the deconvolution results. In this study, when using ScpDeconv, the
136 training data were generated from pseudo-tissue samples randomly drawn from a uniform
137 distribution. The tissue data is scaled to the range of [0, 1] using the MinMaxScaler() class from the
138 Sklearn preprocessing module. When employing this method, refer to the example workflow
139 provided in its GitHub repository and adhere to its default parameters.

140 **CIBERSORTx**

141 CIBERSORTx is a deconvolution method based on support vector regression that can derive gene
142 expression features from single-cell reference data and estimate cell type proportions.
143 CIBERSORTx is a user-friendly web tool that requires feature alignment of the single-cell reference
144 matrix and the data to be deconvoluted. Users should structure the data according to the tutorial
145 examples and run the method using the default parameters.

146 **Tangram**

147 Tangram is a deep learning-based spatial mapping deconvolution algorithm that constructs a
148 probability matrix for mapping cells to spatial locations, aligning single-cell and spatial
149 transcriptomic expression using cosine similarity loss to achieve deconvolution of cell type
150 composition in spatial data. For configuration, we use the deconvolution examples provided by
151 Tangram on GitHub and apply cluster mode (mode='clusters') to map and output normalized cell
152 type density matrices. The single-cell data is normalized using Scanpy's pp.normalize_total, and
153 cell types with fewer than 2 samples are filtered out. We obtain the top 200 marker genes for each
154 cell type using Scanpy's tl.rank_genes_groups, then intersect these with the spatial data to get
155 common genes. Genes with no information are removed through tg.pp_adata. During mapping, we
156 use the default parameters from the Tangram authors' tutorial without Squidpy to ultimately output
157 the deconvolution results of spatial points and niche types.

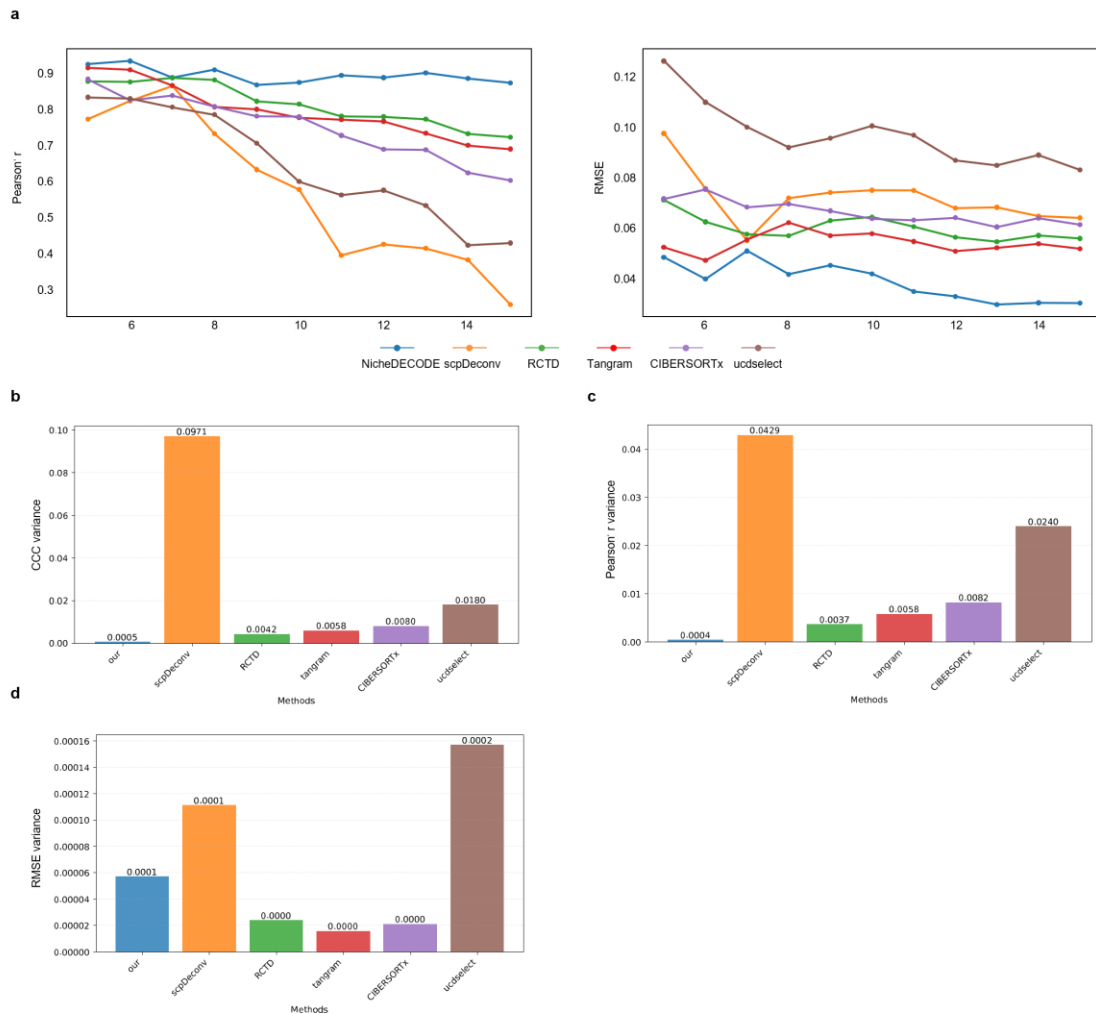
158 **Ucdselect**

159 UCD is a deconvolution algorithm based on a pre-trained model. The model is pre-trained using 10
160 million pseudotissue samples from the author's integrated scRNA-seq training database. When using
161 UCD, the first step is to match the gene expression features of the input reference single-cell data
162 and spatial mixed data with the feature dictionary of the pre-trained data, ensuring that they share
163 common features with the pre-trained model. Since the deconvolution results of the UCD_base
164 mode refer to cell types from the pre-trained dataset, this paper employs the UCD_select mode. In
165 this mode, niche type pairs required for UCD_select. Subsequently, the API is called to invoke the
166 cloud-based pre-trained model, which performs deconvolution on the input data and returns the
167 abundance matrix of niche types.

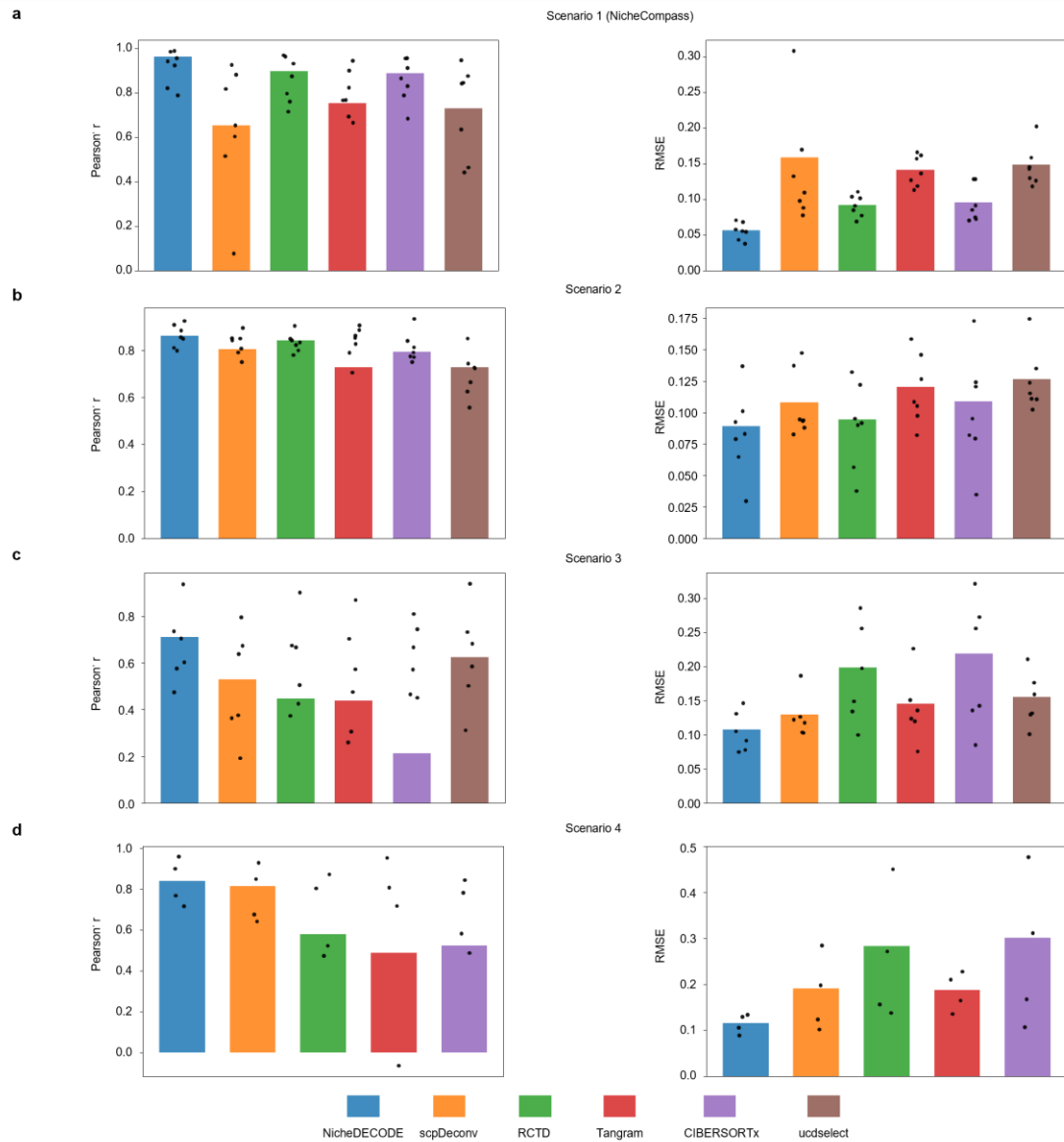
168 **RCTD**

169 RCTD is a spatial transcriptomics deconvolution method implemented in R, based on a mixed linear
170 model framework that corrects for batch effects. It optimizes model parameters through maximum
171 likelihood iterative optimization to compute the posterior probability distribution of niche type
172 composition at each spatial point, decomposing the gene counts of that spatial point into a weighted

173 sum of the contributions from each niche type. This paper follows the Bioconductor tutorial,
 174 constructing a puck that combines spatial expression data and spatial location information.
 175 The createRCTD() function is then used to preprocess the puck data and single-cell reference. This
 176 function filters spots and genes based on thresholds such as UMI counts and identifies differentially
 177 expressed genes. It creates niche type profiles from the reference,
 178 setting CELL_MIN_INSTANCE=20 in the createRCTD() function to filter out low-abundance
 179 niche types. Finally, it specifies doublet_mode="full" in the run.RCTD() function to enable a
 180 complete mode that allows for an unrestricted number of niche types, resulting in the deconvolution
 181 outcomes.



182
 183 **Fig. S1| Additional results on CosMx NSCLC data.** (a) Line plots illustrating Pearson's r (left)
 184 and RMSE (right) of NicheDECODE, scpDeconv, RCTD, Tangram, CIBERSORTx, and ucdselect
 185 across niche resolutions ranging from 5 to 15 clusters. (b-d) Variance of CCC (b), Pearson's r (c),
 186 and RMSE (d) across niche resolutions (5 - 15 clusters) for each method.



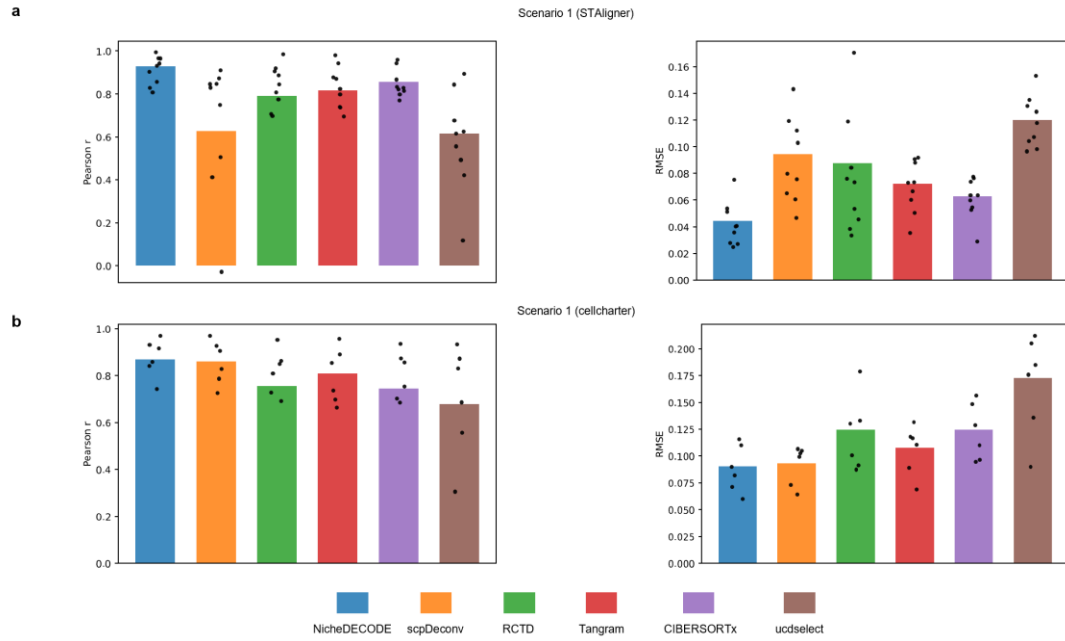
187

188 **Fig. S2| Pearson's r and RMSE of different methods across four generalization scenarios.**

189 Pearson's r and RMSE results for NicheDECODE, scpDeconv, RCTD, Tangram, CIBERSORTx, and

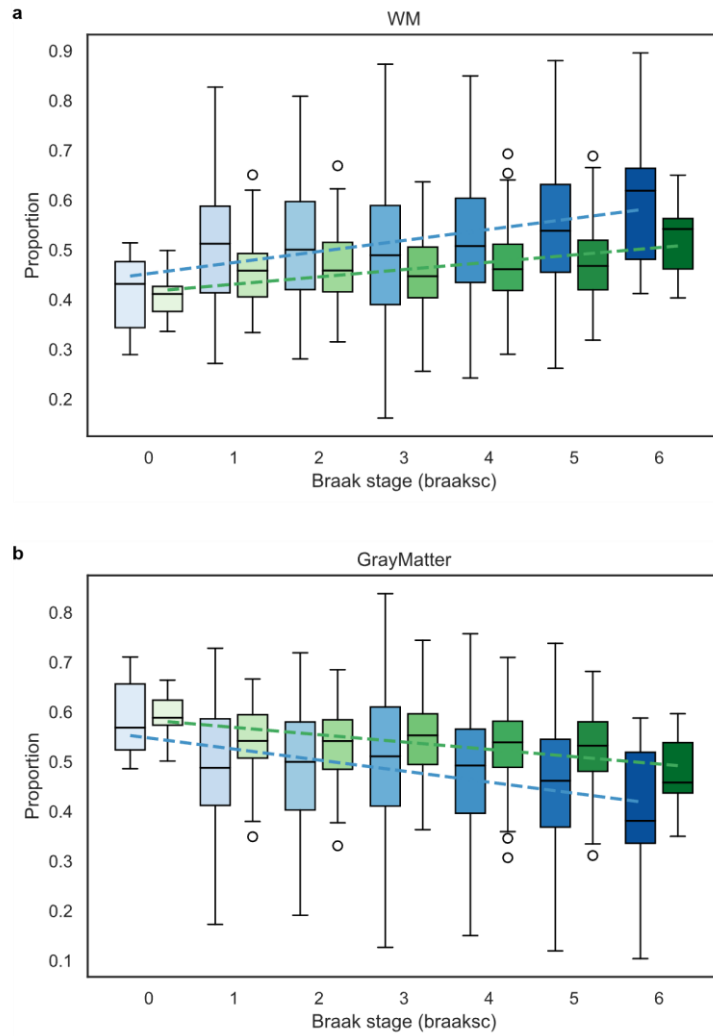
190 ucselect across Scenario 1 (a), Scenario 2 (b), Scenario 3 (c), and Scenario 4 (d).

191



192
 193
 194
 195
 196

Fig. S3| Pearson's r and RMSE of different methods across scenario 1. Pearson's r and RMSE results for NicheDECODE, scpDeconv, RCTD, Tangram, CIBERSORTx, and ucselect across Scenario 1 STAligner (a), Scenario 1 cellcharter (b).

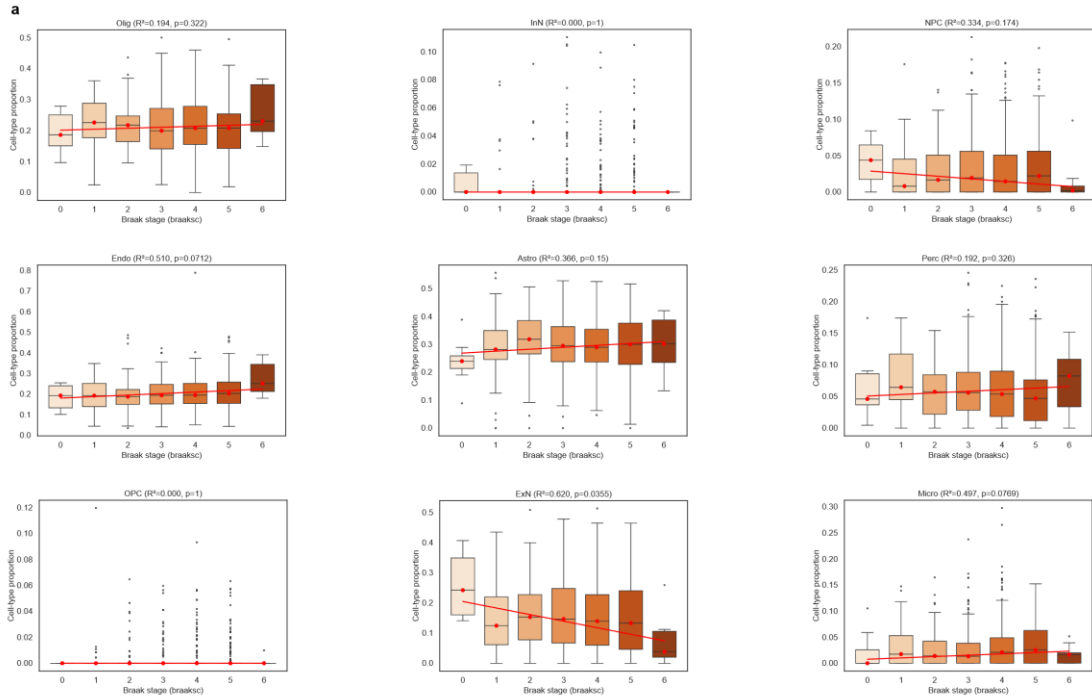


197

198 **Fig. S4| Box plots showing changes in white matter (a) and gray matter (b) abundances across**
 199 **Alzheimer's disease Braak stages.** AD cortex (blue) and healthy cortex (green) references.

200 Notations: Box plots show the median (centre line), the 25th–75th percentiles (bounds of the box),
 201 and the minimum and maximum values (whiskers); dots indicate outliers.

202

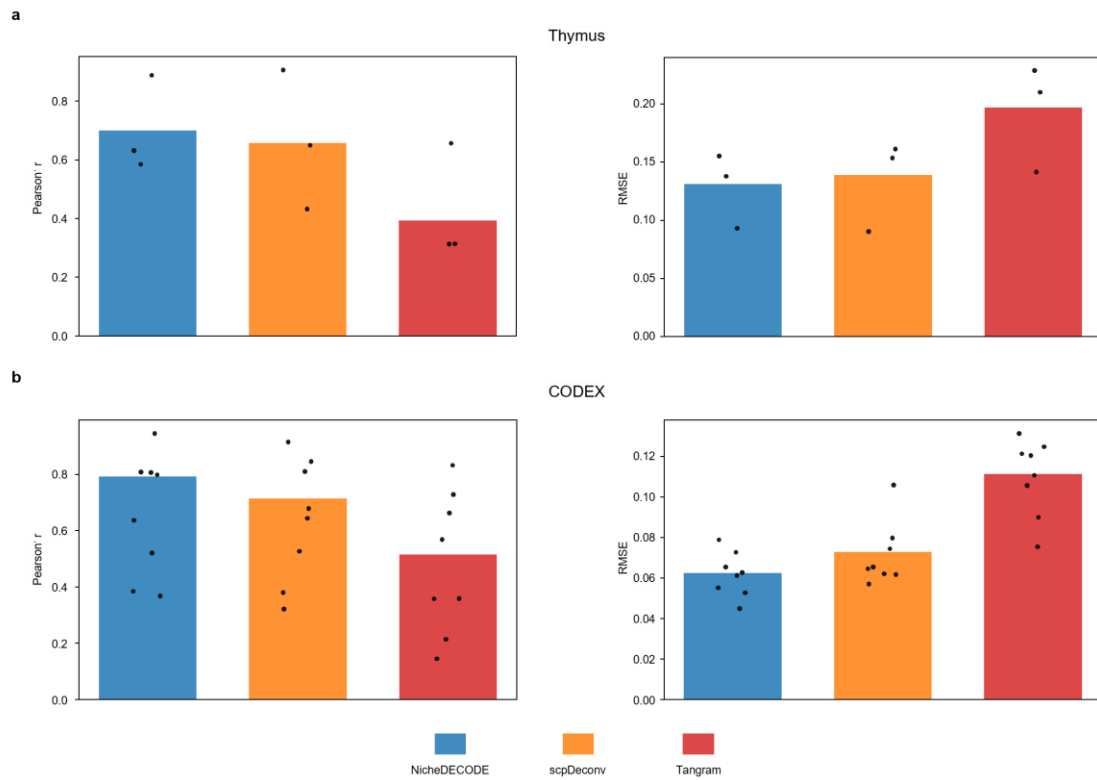


203

204

Fig. S5| Cell-type deconvolution results on the ROSMAP cohort. (a) Box plots of cell-type abundances in the ROSMAP cohort inferred by CIBERSORTx. The x-axis represents samples from different AD Braak stages, and the y-axis represents the cell-type proportions. (stage-wise $n = 10, 53, 76, 201, 208, 275,$ and 9). Regression lines were fitted by ordinary least squares using the median abundance within each Braak stage. Accordingly, both R^2 and the p value for the slope (from the stage-median regression) were computed based on these stage medians using the linregress-function from the SciPy library. Notations: Box plots show the median (centre line), the 25th–75th percentiles (bounds of the box), and the minimum and maximum values (whiskers); dots indicate outliers.

212



213
 214
 215
 216
 217

Fig. S6| Pearson's r and RMSE of different methods across proteomics data. Pearson's r and RMSE results for NicheDECODE, scpDeconv, and Tangram on the Thymus (a) and CODEX (b) datasets.

