

Identification of spatial co-expression patterns and intra-tissue heterogeneity in spatially resolved transcriptomics by region-specific denoising

Linhua Wang

Baylor College of Medicine <https://orcid.org/0000-0002-6717-860X>

Zhandong Liu (✉ zhandong.liu@bcm.edu)

Baylor College of Medicine <https://orcid.org/0000-0002-7608-0831>

Article

Keywords: Spatial transcriptomics, cell spatial organization, gene-expression

Posted Date: November 8th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-647777/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: There is **NO** Competing Interest.

Version of Record: A version of this preprint was published at Nature Communications on November 14th, 2022. See the published version at <https://doi.org/10.1038/s41467-022-34567-0>.

1 **Identification of spatial co-expression patterns and intra-tissue heterogeneity in**
2 **spatially resolved transcriptomics by region-specific denoising**

3 Linhua Wang¹, Zhandong Liu^{2, 3*}

4 1. Graduate Program in Quantitative and Computational Biosciences, Baylor College of Medicine,
5 Houston, USA

6 2. Jan and Dan Duncan Neurological Research Institute at Texas Children's Hospital, Houston, USA

7 3. Department of Pediatrics, Baylor College of Medicine, Houston, USA

8 * Corresponding author: zhandonl@bcm.edu

9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

34 **Abstract**

35 We are pleased to introduce a first-of-its-kind tool that combines *in-silico* region detection and
36 missing value estimation for spatially resolved transcriptomics. Spatial transcriptomics by 10X
37 Visium (ST) is a new technology used to dissect gene and cell spatial organization. Analyzing
38 this new type of data has two main challenges: automatically annotating the major tissue
39 regions and excessive zero values of gene expression due to high dropout rates. We developed
40 a computational tool—MIST—that addresses both challenges by automatically identifying tissue
41 regions and estimating missing gene-expression values for each detected region. We validated
42 MIST detected regions across multiple datasets using manual annotation on the histological
43 staining images as references. We also demonstrated that MIST can accurately recover ST's
44 missing values through hold-out experiments. Furthermore, we showed that MIST could identify
45 intra-tissue heterogeneity and recover spatial gene-gene co-expression signals. We therefore
46 strongly encourage using MIST before downstream ST analysis because it provides unbiased
47 region annotations and enables accurately denoised spatial gene-expression profiles.

48

49 **Introduction**

50 To understand the biological mechanisms underlying diseases, it is essential to delineate cell
51 and gene spatial organizations; the scientific community has therefore invested significant time
52 and effort in detecting positional gene-expression^{1,2}. Of all the spatial gene-expression profiling
53 techniques that have been developed, 10X Visium Spatial Transcriptomics (ST) gained its most
54 popularity due to its whole-genome scalability and cost-efficiency^{2,3}. So far, ST has been used
55 to study many tissues' spatial gene-expression, including breast cancer, prostate tumors,
56 melanoma tumors, human and mouse brains, and more³⁻⁸.

57

58 While the location of every sequenced tissue domain (spot) in ST is assigned, regions of the
59 tissue are not directly provided. Routinely, researchers circled out anatomical regions by looking
60 at the histological staining image that is aligned with ST. However, in many cases, pathological
61 regions are not detectable through human eyes. Moreover, assigning each spot to a region is
62 labor-intensive and might be erroneous. Therefore, it is of great significance to automatically
63 and unbiasedly extract regions within ST tissues.

64

65 Another challenge researchers faced when analyzing ST data is the sparsity of the gene
66 expression profile caused by technical dropouts, which leads to excessive zero values in the
67 gene expression data³. As a consequence, it drastically reduces ST's signal-to-noise ratio and

68 prevents accurate co-expression calculation, cluster detection, and other downstream analyses.
69 Therefore, accurately estimating dropped-out expression values in ST data is vital to rescuing
70 such signals and facilitating more accurate downstream analyses.

71

72 While we lack computational methods specifically designed to estimate missing values in ST,
73 plenty of such methods have been developed for single-cell RNA sequencing (scRNA-seq)
74 data^{9–14}. Many of these methods—like MAGIC⁹, PRIME¹⁰ and knnSmoothing¹¹—smooth the
75 missing gene expression values by taking the weighted average of the gene expression values
76 from cells with similar molecular profiles. None of these methods, however, interpret ST's spatial
77 connectivity information to determine the functionally and physically similar neighbors to learn
78 from. Other kinds of methods—like McImpute¹² and ALRA¹⁴, estimate the missing values by
79 approximating a low rank of the gene expression matrix, assuming a relative number of cell
80 types within the tissue. However, they are performed at the entire tissue and lack regional
81 specificity that could allow us to rescue regional heterogeneity.

82

83 We realized that these two challenges: (1) difficulty automatically annotating tissue regions and
84 (2) excessive zero gene expression values due to technical dropouts are not independent of
85 each other. On the contrary, solving the first challenge could facilitate region-specific denoising.

86

87 An intuitive way of denoising ST using spatial information is to use spatially adjacent neighbors
88 to estimate missing values. However, using adjacent neighbors introduces errors regarding
89 spots on functionally different tissue regions' boundaries. In such cases, detecting the
90 boundaries is vital to avoiding spatial-information misuse when denoising ST data and allows
91 ST to be grouped into functionally similar regions. Denoising within such functional and spatial
92 regions, therefore, avoids errors at tissue boundaries.

93

94 We hypothesize that the number of cell types within a functional region in the tissue is limited,
95 leading us to believe that a more accurate denoising result will be achieved by region-specific
96 denoising through low-rank matrix completion. Therefore, we introduce our new tool – Missing-
97 value Imputation for Spatially resolved Transcriptomics (MIST) that addresses the *in-silico*
98 detection and region-specific denoising tasks in a two-step pipeline.

99

100 **Results**

101 **The MIST Algorithm**

102 MIST provides a two-step tool to address the challenges in ST's analyses including
103 automatically detecting ST's functional regions and denoising ST data at a region-specific
104 schema (Fig. 1).

105
106 In the first step, MIST automatically detects functional regions in the sample of interest. MIST
107 embeds ST as a 2D graph where every spot will be presented as a node. Every pair of adjacent
108 nodes will have an edge connecting each other with a weight defined by the molecular similarity
109 between the pair of nodes.

110
111 To simulate the region boundaries, MIST filtered out low-weight edges with a threshold and
112 extract the connected components within the remaining graph. To avoid bias in selecting the
113 filtering threshold, MIST searches for the optimal value by assessing different thresholds'
114 accuracy in recovering the hold-out values of a set of highly expressed genes.

115
116 In the second step, MIST estimates the missing values in each detected region by averaging the
117 outcomes from multiple runs of a low-rank approximation algorithm¹². Under the assumption that
118 each detected region should have a small number of cell types, we expect the denoised region-
119 specific expression matrix to have a low rank. To achieve this goal, we used a low-rank matrix
120 completion approach that estimates the missing values by minimizing the singular values of the
121 denoised matrix.

122
123 To increase the reliability of the estimated values and avoid an extremely singular gene-
124 expression matrix, before denoising, MIST adds contrast to each region by augmenting the
125 region with random spots from other regions. To remove the bias from randomly selected spots
126 and stabilize the estimated values, MIST runs multiple repeats of low-rank approximation on the
127 augmented regions with different sets of augmented spots and takes the average of the
128 outcomes as the estimated values for region-specific spots.

129 130 **MIST detected functional tissue regions that agree with manual annotations**

131 To prove that MIST could faithfully detect functional regions within tissues, we tested MIST on
132 two examples including a melanoma sample and a mouse brain sample (Fig. 2a, b). We
133 showed that under different filtering thresholds, different shapes of regions were detected.
134 However, regions using the threshold automatically selected by MIST mostly resembles the
135 manual annotation provided by human experts.

136

137 While higher filtering thresholds on the edges will find strongly correlated spots within regions,
138 they also lead to a significant number of isolated spots with no connected edges (Fig. 2c). But,
139 while small thresholds will include most spots, they will lose the specificity that allows us to
140 accurately detect tissue regions (Fig. 2c). Balancing the proportion of non-isolated spots and
141 region-detection accuracy is, therefore, crucial to successful region detection and region-
142 specific denoising.

143

144 To demonstrate that MIST strikes this balance automatically and accurately, we used Adjusted
145 Rand Index (ARI) to evaluate the consistency between MIST-detected regions and the human
146 experts' annotations. An ARI of zero shows the concordance of two random clustering results
147 while a score of one indicates a perfect match of two clustering results. Using the mouse wild-
148 type brain sample, we validated that the threshold that MIST automatically selected found a
149 balance point that maximizes the ARI while minimizing the proportion of isolated spots (Fig. 2c).
150 Using the selected threshold, the mouse brain regions detected by MIST agreed with the
151 anatomical regions (Fig. 2b) with an ARI value of 0.64.

152

153 To confirm the biological relevance of MIST-detected regions in the melanoma sample, we
154 looked at the gene profiles in these regions. Specifically, we identified significantly activated
155 genes in the tumor region relative to the lymphoid region and vice versa using the Wilcoxon
156 rank-sum test. By selecting genes with a log₂ fold-change greater than 0.58 and an adjusted P-
157 value less than 10^{-5} , we found 143 and 49 marker genes for the tumor and lymphoid regions,
158 respectively (Fig. 2d). We observed that some well-known melanoma marker genes, including
159 *MLANA*¹⁵, *MCAM*¹⁶, *SPP1*¹⁷, and *HSP90AA1*¹⁸, topped the list of the melanoma-activated genes
160 that we identified, which demonstrates our detected region's accuracy. We then performed gene
161 set enrichment analysis on the significantly activated genes and observed that the immune
162 response gene ontology term was significantly enriched (adjusted P-value = 9×10^{-8} , Fig. 2e)
163 for the detected lymphoid region. Taken together, the region automatically detected from the
164 Melanoma ST reflected the hidden functional regions within the tissue.

165

166 **MIST accurately recovers hold-out values across multiple data sets**

167 To assess MIST's accuracy when estimating missing values, we performed 5-fold random hold-
168 out experiments in which we withheld a random set of the observed non-zero values and used
169 these as ground truth to evaluate the models' performances. By withholding some of the

170 observed values, we simulated cases in which non-zero expression values have dropped out.
171 To consider the heterogeneity of datasets that might lead to biased performance, we tested four
172 datasets from samples including a wild-type mouse brain sample, an AD mouse brain sample, a
173 melanoma tumor sample, and a prostate tumor sample that vary in the number of genes,
174 number of spots, and sparsity level (Table 1). For each dataset, we selected genes that are
175 expressed across at least half of all spots in the sample to generate hold-out test data sets. For
176 each gene, we partitioned the non-zero expression values into five non-overlapping sets. Then,
177 we iteratively held out one-fold of the values and assessed MIST and other methods' accuracy
178 in recovering the held-out gene expression values.

179
180 To demonstrate MIST's supremacy in missing value recovery, we compared MIST with state-of-
181 the-art scRNA-seq methods¹⁹, including MAGIC⁹, knn-smoothing¹¹, McImpute¹² and SAVER¹³,
182 and a baseline k-nearest neighbor method we constructed (spKNN) that estimates missing
183 values by averaging spatially adjacent neighbors.

184
185 To evaluate the accuracy of missing value estimation in hold-out experiments, we used two
186 metrics including Rooted Mean Square Error (RMSE) and Pearson Correlation Coefficient
187 (PCC) where RMSE represents the error and PCC shows the agreement between the ground
188 truth and estimated values. Better imputation methods are expected to have lower RMSE and
189 higher PCC scores.

190
191 We found that MIST consistently outperformed other methods across all datasets with higher
192 PCC and lower RMSE scores during hold-out value evaluation (Fig. 3a, b). For the spots that
193 are assigned to functional regions, MIST had an average RMSE improvement (decreasing
194 values) of 12% (P-value= 3×10^{-4} compared with McImpute, and 14% (P-value= 10^{-14}) and
195 35% (P-value= 7×10^{-10}) improvement compared with the baseline spKNN algorithm.
196 Moreover, MIST's mean PCC is also 8% (P-value= 9×10^{-6}), 4% (P-value= 2×10^{-4}), and 16%
197 (P-value= 6×10^{-11}) higher than McImpute, MAGIC and spKNN respectively. Knn-smoothing
198 and SAVER consistently performed substantially worse than the other four methods (Ext. Fig.
199 3.1, 3.2).

200
201 To investigate the impact of genes' sparsity on denoising accuracy, we further stratified the
202 performance assessment at a per-gene level grouped by the sparsity level (zero-value
203 proportion). While MAGIC and spKNN's estimation error monotonically increased with sparsity

204 level, MIST's performance was not vulnerable to gene sparsity (Fig. 3c). While McImpute's
205 performance was also not influenced by gene sparsity, MIST has superior performance than
206 McImpute at every gene sparsity level (Fig. 3c).

207

208 To demonstrate that MIST recovers the gene-expression patterns, we then visualized and
209 showed that MIST can faithfully recover the gene-expression spatial pattern for gene *GAPDH*
210 after denoising using the Melanoma tissue sample (Fig. 3d). With the hold-out input, MIST
211 accurately estimated the original expression values by increasing Spearman's correlation
212 coefficient from 0.65 to 0.96 (Fig. 3e). When evaluating all genes across the four tested data
213 sets, after denoising, the median correlation had significantly improved from 0.64 to 0.88 (P-
214 value ≈ 0 , Fig 3f, Ext. Fig. 3.3, 3.4).

215

216 **MIST discovered intra-cortex heterogeneity within an Alzheimer's Disease (AD) mouse** 217 **brain**

218 To demonstrate that MIST could improve the clustering results of ST data, we used Uniform
219 Manifold Approximation and Projection (UMAP)²⁰ to reduce the dimension of Mouse brain ST
220 samples' gene expression data and visualize the clustering structures. First, we performed
221 UMAP on the original and denoised wild-type mouse brain sample. After denoising, MIST
222 enhanced the heterogeneity with most of the spots within the cortex, hippocampus, and
223 thalamus, forming individual clusters (Fig. 4a, Ext. Fig. 4.1).

224

225 While similar enhanced clustering patterns were also identified in the Mouse AD brain sample
226 after denoising, in addition, we observed that the cortex region was separated into two individual
227 clusters, something we did not detect using the WT mouse brain nor in the original Mouse AD
228 data (Fig. 4b, Ext. Fig. 4.2). Further cortex analysis revealed a clear separation of the cortex
229 region into two spatially separable parts (Fig. 4c, Ext. Fig. 4.3). The first cluster consisted of the
230 cortical subplate, olfactory, entorhinal, ectrorhinal, temporal association, and perirhinal areas.
231 The second cluster contained the auditory, primary somatosensory, posterior parietal
232 association, and retrosplenial areas. When these two clusters were mapped to the anatomical
233 reference, cluster 1 occupied the upper quadrant while cluster 2 occupied the lower quadrant
234 (Fig. 4d-f). Interestingly, such heterogeneity was only detected in the MIST-denoised AD cortex
235 but not the WT cortex.

236

237 To understand the biological difference of these two clusters in AD progression, we further
238 performed differential gene analysis to extract AD activated genes for these two clusters,
239 respectively. By selecting upregulated genes in the AD sample with a fold change greater than
240 50% and adjusted P-value < 0.01, we identified 55 markers for cluster 1 and 41 markers for
241 cluster 2 (Ext. Fig. 4.4, 4.5). We found only 21 AD-activated genes, such as *Clu*, are shared for
242 these two clusters (Fig. 4d). 34 genes, such as *Hap1*, are AD-upregulated only in cluster 1 (Fig.
243 4e) and 20 genes, such as *Sez6*, that are only AD-upregulated in cluster 2 (Fig. 4f). We showed
244 that the spatial clusters identified using MIST-denoised cortex data play different roles in AD
245 progression by activating different sets of genes, demonstrating that MIST allows the discovery
246 of biological heterogeneity from ST data.

247

248 **MIST recovers spatial gene-gene co-expression patterns**

249 Spatial gene-gene co-expression plays an important role in understanding gene interactions
250 across 2D space. The dropouts within the ST data undermine the correlation analysis's power
251 and cause inaccurate estimation of gene-gene spatial correlation, which is the fundamental
252 element in many analyses such as weighted correlation network analysis (WGCNA)²¹.

253

254 To demonstrate that MIST recovers the spatial co-expression patterns, we examined two pairs
255 of genes: *Cldn11-Arhgef10* and *Gfap-Aqp4*.

256

257 The first pair of genes *Cldn11-Arhgef10* showed a high spatial correlation score of 0.97 based
258 on the reference Allen Brain Atlas²² (Fig. 5a). However, in the original ST data, the correlation
259 score between *Cldn11* and *Arhgef10* is only 0.15 (P-value = 2×10^{-3} , Fig. 5b). After denoising
260 by MIST, the correlation score was improved to 0.5 (P-value = 3.5×10^{-29} , Fig. 5c). To visualize
261 the gene expression patterns, we plotted the heatmap of log-scale expression values and
262 showed that *Cldn11* and *Arhgef10* have similar gene expression patterns only after MIST
263 denoising (Fig. 5d-g).

264

265 To show that MIST can recover co-expressed gene pairs that are not significantly correlated
266 with the original data, we examined the second pair of genes: *Gfap-Aqp4*, which have a
267 moderately good spatial correlation with a score of 0.74 in the reference Allen Brain Atlas
268 database (Fig. 5h). Before denoising, we observed an insignificant correlation with a score of
269 0.08 (P-value = 0.56, Fig. 5i). After denoising, we could recover a significant spatial correlation
270 with a score of 0.35 (P-value = 9×10^{-12} , Fig. 5j). Similar to the first pair of genes, the visual

271 pattern improvement of *Gfap-Aqp4* could also be observed only after denoising by MIST (Fig.
272 5k-n).

273

274 By showing these two cases, we demonstrated that MIST could rescue the spatial correlation of
275 gene pairs whose co-expression patterns are either lessened or lost in the original ST data.
276 Given co-expression estimation is vital in many downstream analyses such as identifying gene
277 modules²¹, we see the significance and importance of using MIST to denoise ST data before
278 carrying out such analyses.

279

280 **Conclusions**

281 In this study, we tackled two major problems encountered in ST data analyses – (1) in-silico
282 region detection and (2) missing value estimation. We developed a computational tool, MIST,
283 that solved both problems.

284

285 We solved the first challenge by combining the molecular similarity and spatial connectivity
286 between spots and enabled detecting tissue regions automatically. Before our work,
287 researchers align the histological images with ST coordinates to manually assign regions to
288 every spot. This procedure is laborious and might be inaccurate in cases where pathological
289 changes are not visually detectable through human eyes. With MIST, we could bypass manual
290 annotation and allows unbiased region assignment to every spot. We proved that MIST could
291 accurately detect regions by comparing MIST-detected regions with manually assigned regions
292 at the per-spot level. Addressing this challenge allows us to analyze ST data at a region-specific
293 level and improve the specificity for denoising.

294

295 We solved the second challenge by approximating a low rank of region-specific gene
296 expression matrix through singular value decomposition. This is based on a simple yet
297 interpretable assumption that the number of cell types for every region is small, which has been
298 adopted by many other single-cell RNA-sequencing denoising methods such as Mclmpute¹² and
299 ALRA¹⁴. Compared with Mclmpute and ALRA, MIST utilized spatial information to define regions
300 and improved the denoising specificity.

301

302 To demonstrate that MIST's accuracy in estimating the missing values, we developed hold-out
303 experiments to simulate dropout events and assessed MIST's performance in the hold-out trials.

304 As a result, MIST could accurately recover the hold-out values and outperformed the state-of-
305 the-art single-cell denoising methods and a baseline spatial-information-based approach.

306

307 While addressing these two major challenges, MIST allows researchers to drastically recover
308 biological signals that are missing in the original ST data. Such biological signals include spatial
309 co-expression patterns, intra-region heterogeneity, and more.

310

311 We demonstrated that MIST could improve or recover spatial patterns of co-expressed genes
312 that are highly correlated in reference atlas but either poorly or insignificantly correlated in the
313 original ST. Since many downstream analyses such as WGCNA²¹ are based on co-expression
314 estimation, denoising ST by MIST is rather essential before such analyses because it avoids
315 false conclusions due to inaccurate co-expression estimation.

316

317 It is also of great significance to understand the heterogeneous responses of different tissue
318 regions to pathology, which will shed light on the treatment of diseases. Therefore, it is critical to
319 recovering such heterogeneity that was lost in the original ST data.

320

321 We showed that MIST enables recovering the spatial heterogeneity within the mouse cortex
322 during AD progression. By extracting the differentially expressed genes between AD and wild-
323 type mouse brain samples, we identified genes that are AD-upregulated only in specific regions
324 of the mouse cortex.

325

326 We therefore unequivocally recommend using MIST for ST analyses. Our *sui generis*
327 *in-silico* region detection enables analyzing the ST at a brand-new level that combines
328 anatomical connectivity and molecular similarity. It will be useful in many areas, including
329 identifying local subregions within tumors whose heterogeneity is hard to see through
330 pathological staining images. Second, MIST allows researchers to recover important biological
331 signals in downstream analyses, such as when identifying spatial gene-gene co-expression
332 patterns. While the original ST data provided by 10X Visium might hinder ST analyses'
333 accuracy, MIST will accurately rescue the missing values and drastically increase the signal-to-
334 noise ratio.

335

336 **Methods**

337 **Data collection and preprocessing**

338 The ST data sets we used in this study include a 12-month wild-type mouse brain sample, a 12-
339 month Alzheimer's Disease (AD) mouse brain sample, a Melanoma tumor sample, and a
340 prostate tumor sample^{5,6,8} (Table 1). Every dataset has a raw mRNA count matrix form where
341 rows indicate spots and columns indicate genes.

342

343 To filter out low-quality genes that might otherwise introduce noise to the pipeline, we kept
344 genes with raw counts ≥ 3 in more than 2 spots within each sample.

345

346 To account for the different library sizes of every spot due to variance in sequencing depth and
347 the number of cells, we normalized the raw mRNA count matrix using count per million ($CPM =$

348
$$\frac{\text{Raw count} * 10^6}{\text{Library size}}).$$

349

350 ***In-silico* region detection**

351 Graph embedding

352 Suppose the ST expression matrix has M spots and N genes, the spatial gene-expression
353 profile can be defined as $Y \in R_*^{M \times N}$, where Y is the observed gene expression matrix, and R_*
354 denotes non-negative real matrices with M rows and N columns. The M spots in an ST slide can
355 form a lattice graph, $G = \langle V, E \rangle$, where V is the node-set and E is the edge-set. Every pair
356 of adjacent (u, v) spots are connected with edge $E(u, v)$.

357

358 Weight calculation

359 To infer the weights for every connected edge $E(u, v)$, we calculated the Pearson correlation
360 coefficient between the gene expression profile of spot u and v . To remove the noise in high-
361 dimensional gene expression data while keeping the major signals, we used Principal
362 Component Analysis (PCA)²³ to reduce the dimension of the gene expression matrix and kept
363 the first p principal components that could explain 80% of the variance. The weight for edge E
364 (u, v) is then inferred using the correlation score between the first p principal components of
365 spot u and v .

366

367 Edge removal and parameter selection

368 To simulate the boundaries between functionally dissimilar regions, we removed edges whose
369 weights are less than a threshold ε .

370

371 To automatically select the threshold and avoid bias, we performed a hold-out experiment on a
372 set of genes that are expressed in at least 80% of all spots. We randomly held out 25% of the
373 gene expression values for every such gene and used them as ground truth to test MIST's
374 denoising error defined as $RMSE = \sqrt{\frac{\sum_{i=1}^n (X_i - X'_i)^2}{n}}$, where X is the hold-out non-zero values, X'
375 represents the MIST denoised gene expression matrix using certain ε and n denotes the number
376 of non-zero hold-out values. To determine the optimal threshold, we did a grid search with
377 values ranged from 0.1 to 0.9 and selected the one with the minimal RMSE.

378

379 Region detection by extracting connected components

380 To simulate the functional regions within tissues, we used a depth-first search algorithm²⁴ to
381 identify all the connected components in graph G after edge removal. Each connected
382 component with greater than 5 spots is predicted as an independent tissue region.

383

384 **Region-specific denoising through multiple runs of low-rank matrix completions**

385 Low-rank matrix completion

386 Suppose the observed gene expression matrix for region C is Y_C where Y_C is a sparse matrix
387 with m rows (spots) and n columns (genes). The task is to estimate X_C , which represents the
388 denoised gene expression matrix for region C . We adapted the low-rank-matrix completion
389 algorithm through singular value decomposition used by McImpute¹².

390

391 Given the assumption that the number of cell types within a functional region is small, we expect
392 X_C to have a low rank. To achieve this goal, the task is turned into a low-rank matrix completion
393 problem by solving the following objective function:

$$394 \min_{X_C} \|Y_C - A(X_C)\|^2 + \lambda * rank(X_C) \quad (1)$$

395

396 The first term in equation (1) minimizes the error between the non-missing gene expression
397 values of X_C and Y_C with a projection function A that returns values in X_C at the indices of non-
398 missing values in Y_C . The second term in equation (1) minimizes the rank of the denoised gene
399 expression matrix. The objective function is a linear combination of these two terms regularized
400 a non-zero tuning parameter λ . Theoretically, a larger λ will give us a lower-ranked denoised
401 gene expression matrix whose values on the non-missing indices might deviate from the ground

402 truth. On the other hand, a small λ will result in a relatively high-rank denoised matrix with a
403 lower error on the non-missing indices.

404

405 However, minimizing the rank of a matrix is non-convex. To transform it to a convex problem
406 with a globally optimal solution, we relaxed equation 1 as:

$$407 \min_{X_C} \|Y_C - A(X_C)\|^2 + \lambda * \|X_C\|_{nuc} \quad (2),$$

408 , where we transformed the term 2 in equation 1 as a nuclear norm of X_C , which can be
409 calculated by summing up the singular values obtained through singular value decomposition.

410 Specifically, equation (2) can be further transformed as

$$411 \min_{X_C} \|B - X_C\|^2 + \lambda * \|X_C\|_{nuc} \quad (3),$$

412 where $B_{k+1} = X_{k,C} + \frac{1}{\alpha} A^T(Y_C - A(X_{k,C}))$. Using the inequality $\|B - X\|_2 \geq \|s_B - s_X\|$, where s_X

413 denotes the singular value vector for matrix X, equation (3) can be rewritten as

$$414 \min_{X_C} \|s_B - s_{X_C}\|^2 + \lambda * \|s_{X_C}\|_1 \quad (4).$$

415

416 Therefore, by taking the derivative of (4), the solution of s_{X_C} is achieved by soft thresholding the
417 singular values of s_B with a threshold equals to $\lambda/2$.

418

419 To tune the parameter λ that strikes the balance between low matrix singularity and low error on
420 non-missing indices, we find the maximal λ that achieves a fixed low error (10^{-12}) calculated by
421 the sum of the absolute difference between denoised and observed values on non-missing
422 indices.

423

424 Compared with the original solution provided by McImpute, MIST also forces the observed
425 values to be unchanged. By doing so, MIST guaranteed that it will only touch the missing values
426 and keep the observed non-zero expression values unharmed. Moreover, MIST's matrix
427 completion was implemented in a free and open-source language—Python, as compared to
428 MATLAB, which was used by McImpute.

429

430 Region augmentation

431 To avoid an extremely singular denoised gene expression matrix that might otherwise introduce
432 instability, we added contrast to every detected region before low-rank matrix completion.

433 Specifically, we augmented the region C with r randomly selected spots from other regions

434 where $r = \min(\frac{\text{size}(\text{Other regions})}{10}, \text{size}(C))$. The low-rank matrix completion algorithm will be
435 performed on the augmented region C^a with input X^a .

436

437 To include more diversity and reduce the randomness in region augmentation, we repeated the
438 previous step T times. For every spot $s \in C$, we take the average of denoised gene expression
439 values $Y^a(s)$ from outcomes of T repeats. T is set to be 10 as a default to avoid high
440 computational load while including as much signals as possible.

441

442 To demonstrate the essentiality of region augmentation, we tested MIST's performance in the
443 hold-out experiments using detected regions and augmented regions, respectively. We found
444 that region augmentation significantly improved MIST's denoising accuracy (RMSE P-value =
445 1.6×10^{-7} , PCC P-value = 4.4×10^{-7} , Ext. Fig. 6).

446

447 **Quantifying region detection accuracy**

448 Adjusted rand index is used to quantify the agreement between the manual spot-level
449 annotation and MIST-predicted regions. A score of one means a perfect agreement between
450 two grouping results while zero score means the predicted regions are similar to the results from
451 random guessing. Adjusted rand index scores are calculated using the Python *scikit-learn*
452 package²⁶.

453

454 **Differential gene expression analysis**

455 Wilcoxon rank-sum test provided by the Python package *Scipy*²⁵ was used to infer the
456 significant level of differentially expressed genes. Fold change of genes between condition and
457 control was calculated based on the difference of average gene expression within groups. To
458 get the significantly upregulated genes in the MIST-predicted melanoma region relatively to
459 lymphoid region, we selected genes with a fold change greater than 50% and adjusted P-value
460 less than 10^{-5} . To get the AD upregulated genes in the mouse cortex regions, we changed the
461 adjusted P-value threshold to 10^{-2} because of smaller sample sizes.

462

463 **Gene set enrichment analysis**

464 The R package *clusterProfiler*²⁷ was used to perform gene ontology enrichment analysis using
465 the list of differentially expressed genes in the melanoma sample.

466

467 **Data generation for hold-out experiments**

468 To improve the diversity in the hold-out experiments, we tested samples including a mouse wild-
469 type brain sample, a mouse AD brain sample, a melanoma tumor sample, and a prostate tumor
470 sample. To extract good-quality genes to simulate the dropout events, we removed genes that
471 were expressed in less than 50% of the spots.

472

473 To generate the hold-out data, we used a 5-fold cross-validation schema for the non-zero
474 values. Specially, we first randomly partitioned every gene's nonzero expression values into 5
475 groups. In each hold-out, we created missing values by setting one group to zero and performed
476 imputation based on the remaining values. The held-out values served as ground truth for
477 evaluating the imputation algorithms' accuracy.

478

479 **Evaluating hold-out experiments' performance**

480 To quantifying the accuracy in recovering the held-out values, we reported two metrics including
481 RMSE and PCC, where RMSE measures the estimation error and PCC measures the linear
482 correlation between the true expression values and the estimated values. RMSE is defined as

483 $RMSE = \sqrt{\frac{\sum_{i=1}^n (X_i - Y_i)^2}{n}}$ and PCC is defined as $PCC = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$, where Y is the hold-

484 out non-zero values, X represents the MIST estimated values and n denotes the number of
485 hold-out values.

486

487 To quantify the recovery of gene expression patterns after denoising the hold-out data, we used
488 Spearman's rank correlation test implemented by *Scipy*²⁵ to assess the correlation and the
489 corresponding significance level between the original ST and denoised gene expression values.

490

491 **Co-expression analysis using Allen Brain Atlas as a reference**

492 We obtained mouse brain (coronal section) regional expression values for gene *Cldn11*
493 (experiment: *RP_070116_01_G04*), *Arhgef10* (experiment: *RP_070116_01_B05*), *Gfap*
494 (experiment: *RP_Baylor_253913*) and *Aqp4* (experiment: *RP_040324_01_F07*) from Allen
495 Brain Atlas²² as references. Gene expression values provided by Allen Brain Atlas are at the
496 log₂ scale.

497

498 Correlation scores between gene pairs in both reference and ST data are represented by
499 Spearman's correlation coefficient calculated using the Python package *Scipy*²⁵.

500

501

Table 1. Data summary for hold-out experiments

	<i>MouseWT</i>	<i>MouseAD</i>	<i>Melanoma</i>	<i>Prostate</i>
<i>#Genes (preprocessed)</i>	10178	11241	4498	8073
<i>#Genes (<50% sparsity)</i>	3024	4081	952	1718
<i>#Spots</i>	447	488	293	406
<i>Non-zero values (%)</i>	22	26	15	21

502

503 Source Code Availability

504 The MIST algorithm is implemented in Python and is available at

505 <https://github.com/linhuawang/MIST.git>.

506 All the source code and raw data for this manuscript can be downloaded from

507 https://github.com/LiuzLab/MIST_manuscript.git.

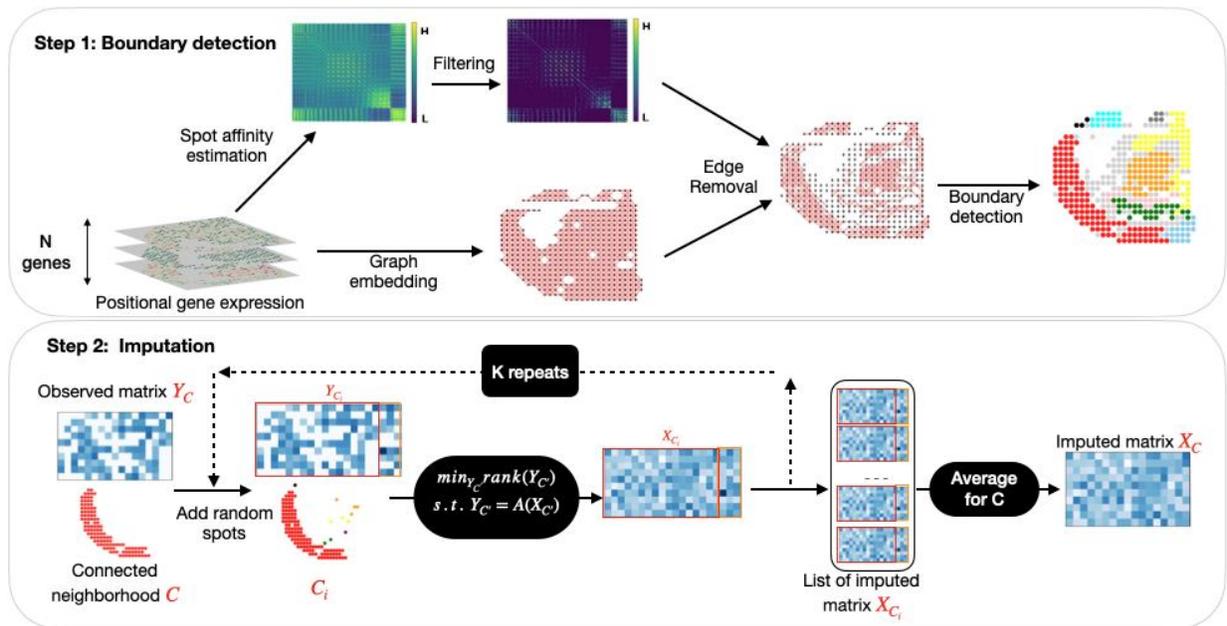
508

509 References

- 510 1. Marx, V. Method of the Year: spatially resolved transcriptomics. *Nat. Methods* **18**, 9–14
511 (2021).
- 512 2. Asp, M., Bergenstr hle, J. & Lundeberg, J. Spatially resolved transcriptomes—next
513 generation tools for tissue exploration. *BioEssays* **42**, 1900221 (2020).
- 514 3. St hl, P. L. *et al.* Visualization and analysis of gene expression in tissue sections by
515 spatial transcriptomics. *Science (80-.)*. **353**, 78–82 (2016).
- 516 4. He, B. *et al.* Integrating spatial gene expression and breast tumour morphology via deep
517 learning. *Nat. Biomed. Eng.* 1–8 (2020).
- 518 5. Berglund, E. *et al.* Spatial maps of prostate cancer transcriptomes reveal an unexplored
519 landscape of heterogeneity. *Nat. Commun.* **9**, 1–13 (2018).
- 520 6. Thrane, K., Eriksson, H., Maaskola, J., Hansson, J. & Lundeberg, J. Spatially resolved
521 transcriptomics enables dissection of genetic heterogeneity in stage III cutaneous
522 malignant melanoma. *Cancer Res.* **78**, 5970–5979 (2018).
- 523 7. Maynard, K. R. *et al.* Transcriptome-scale spatial gene expression in the human
524 dorsolateral prefrontal cortex. *Nat. Neurosci.* 1–12 (2021).
- 525 8. Chen, W.-T. *et al.* Spatial transcriptomics and in situ sequencing to study Alzheimer’s
526 disease. *Cell* **182**, 976–991 (2020).
- 527 9. Van Dijk, D. *et al.* Recovering gene interactions from single-cell data using data diffusion.

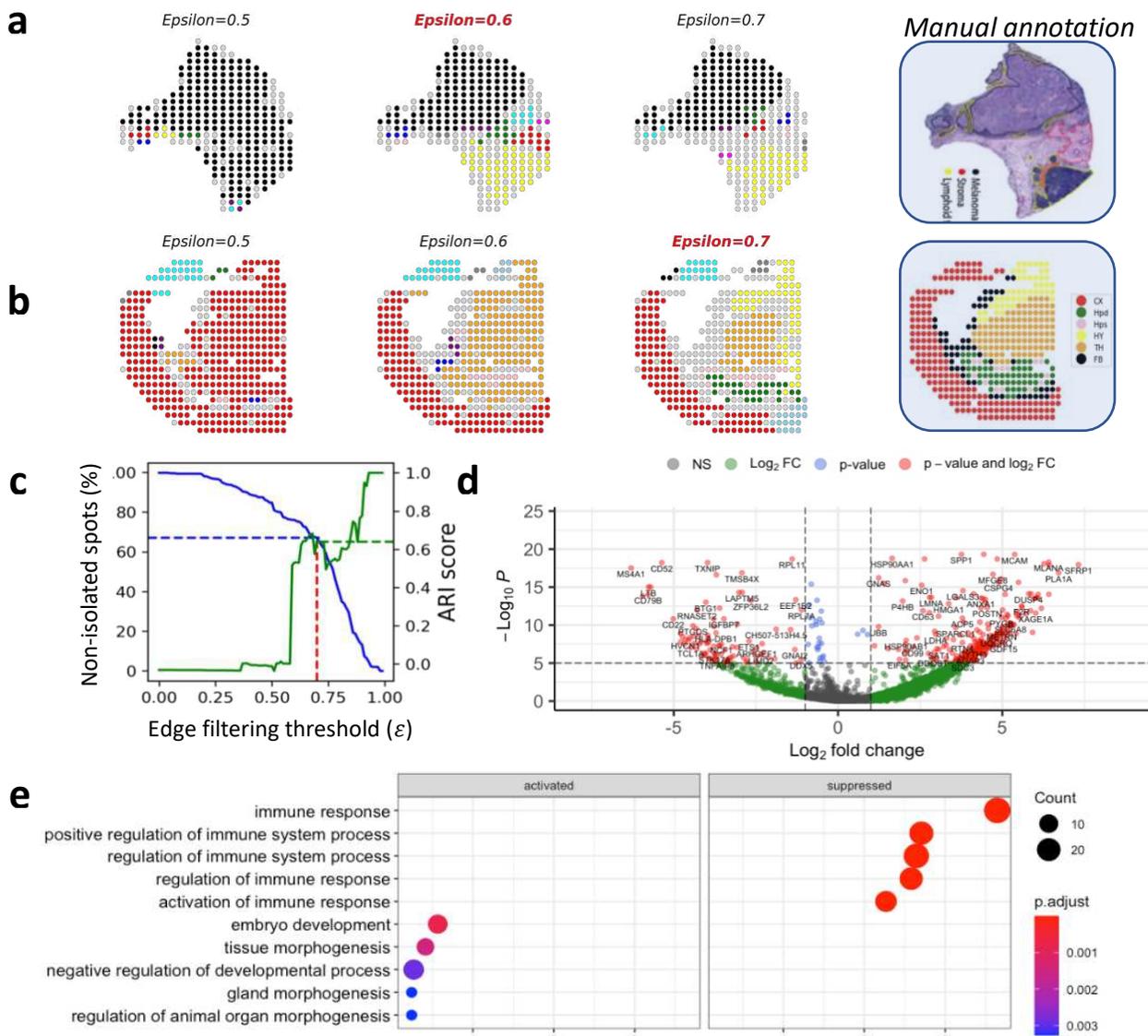
- 528 *Cell* **174**, 716–729 (2018).
- 529 10. Jeong, H. & Liu, Z. PRIME: a probabilistic imputation method to reduce dropout effects in
530 single-cell RNA sequencing. *Bioinformatics* **36**, 4021–4029 (2020).
- 531 11. Wagner, F., Yan, Y. & Yanai, I. K-nearest neighbor smoothing for high-throughput single-
532 cell RNA-Seq data. *BioRxiv* 217737 (2017).
- 533 12. Mongia, A., Sengupta, D. & Majumdar, A. McImpute: matrix completion based imputation
534 for single cell RNA-seq data. *Front. Genet.* **10**, 9 (2019).
- 535 13. Huang, M. *et al.* SAVER: gene expression recovery for single-cell RNA sequencing. *Nat.*
536 *Methods* **15**, 539–542 (2018).
- 537 14. Linderman, G. C., Zhao, J. & Kluger, Y. Zero-preserving imputation of scRNA-seq data
538 using low-rank approximation. *bioRxiv* 397588 (2018).
- 539 15. Chen, Y.-T. *et al.* Serological analysis of Melan-A (MART-1), a melanocyte-specific
540 protein homogeneously expressed in human melanomas. *Proc. Natl. Acad. Sci.* **93**,
541 5915–5919 (1996).
- 542 16. Xie, S. *et al.* Expression of MCAM/MUC18 by human melanoma cells leads to increased
543 tumor growth and metastasis. *Cancer Res.* **57**, 2295–2303 (1997).
- 544 17. Zhou, Y. *et al.* Osteopontin expression correlates with melanoma invasion. *J. Invest.*
545 *Dermatol.* **124**, 1044–1052 (2005).
- 546 18. Shomali, N. *et al.* Heat shock proteins regulating toll-like receptors and the immune
547 system could be a novel therapeutic target for melanoma. *Curr. Mol. Med.* **21**, 15–24
548 (2021).
- 549 19. Hou, W., Ji, Z., Ji, H. & Hicks, S. C. A systematic evaluation of single-cell RNA-
550 sequencing imputation methods. *Genome Biol.* **21**, 1–30 (2020).
- 551 20. McInnes, L., Healy, J. & Melville, J. Umap: Uniform manifold approximation and
552 projection for dimension reduction. *arXiv Prepr. arXiv1802.03426* (2018).
- 553 21. Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network
554 analysis. *BMC Bioinformatics* **9**, 1–13 (2008).
- 555 22. Jones, A. R., Overly, C. C. & Sunkin, S. M. The Allen brain atlas: 5 years and beyond.
556 *Nat. Rev. Neurosci.* **10**, 821–828 (2009).
- 557 23. Wold, S., Esbensen, K. & Geladi, P. Principal component analysis. *Chemom. Intell. Lab.*
558 *Syst.* **2**, 37–52 (1987).
- 559 24. Tarjan, R. Depth-first search and linear graph algorithms. *SIAM J. Comput.* **1**, 146–160
560 (1972).
- 561 25. Virtanen, P. *et al.* {SciPy} 1.0: Fundamental Algorithms for Scientific Computing in

- 562 Python. *Nat. Methods* **17**, 261–272 (2020).
- 563 26. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in {P}ython. *J. Mach. Learn. Res.* **12**,
564 2825–2830 (2011).
- 565 27. Yu, G., Wang, L.-G., Han, Y. & He, Q.-Y. clusterProfiler: an R package for comparing
566 biological themes among gene clusters. *Omi. a J. Integr. Biol.* **16**, 284–287 (2012).
- 567

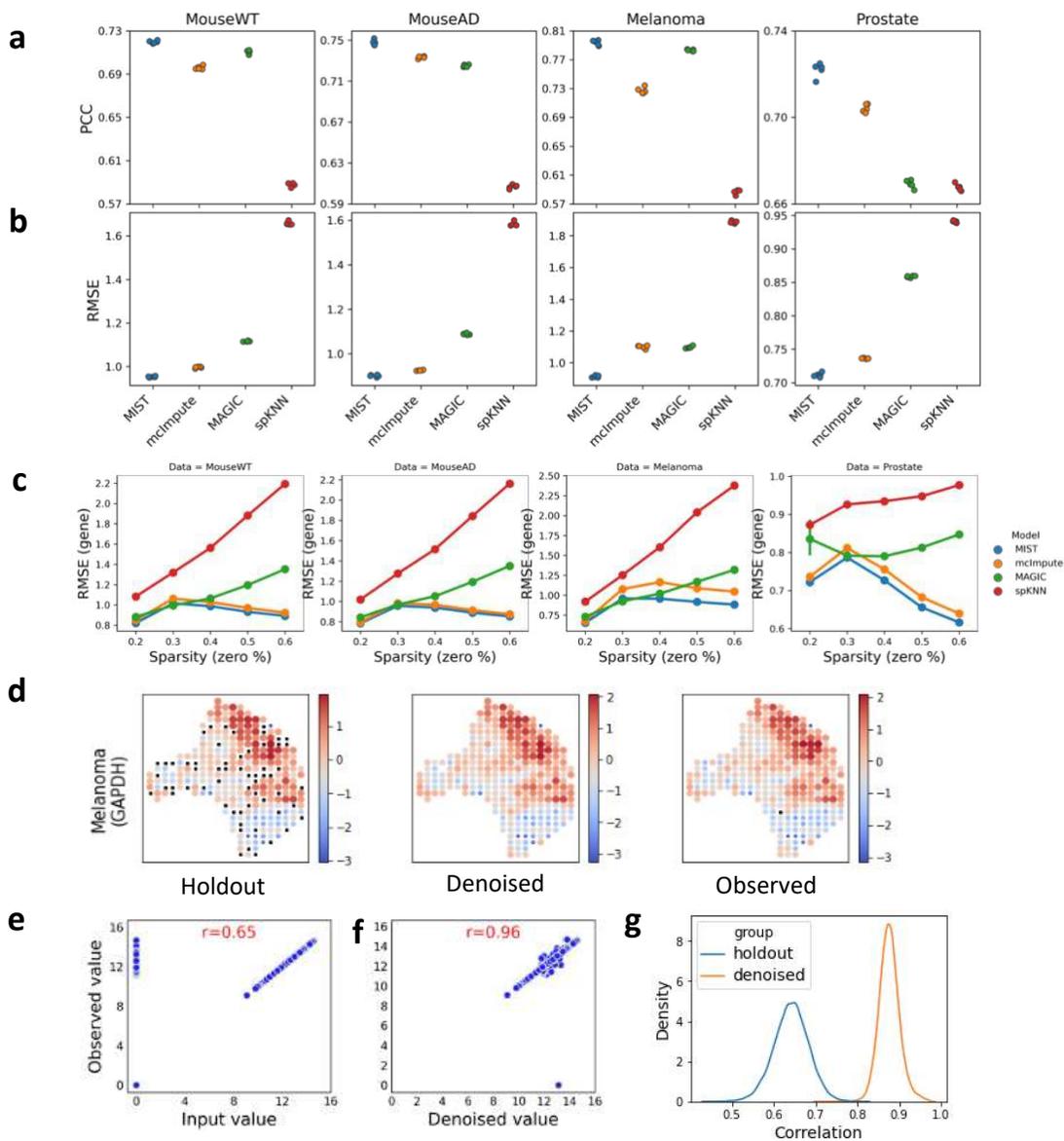


568 **Figure 1. | The MIST pipeline.** In step 1, region boundaries are detected by extracting locally
 569 connected components through graph embedding and edge filtering. In step 2, missing values
 570 for every detected region are estimated using the average of multiple runs of low-rank matrix
 571 completion algorithm.

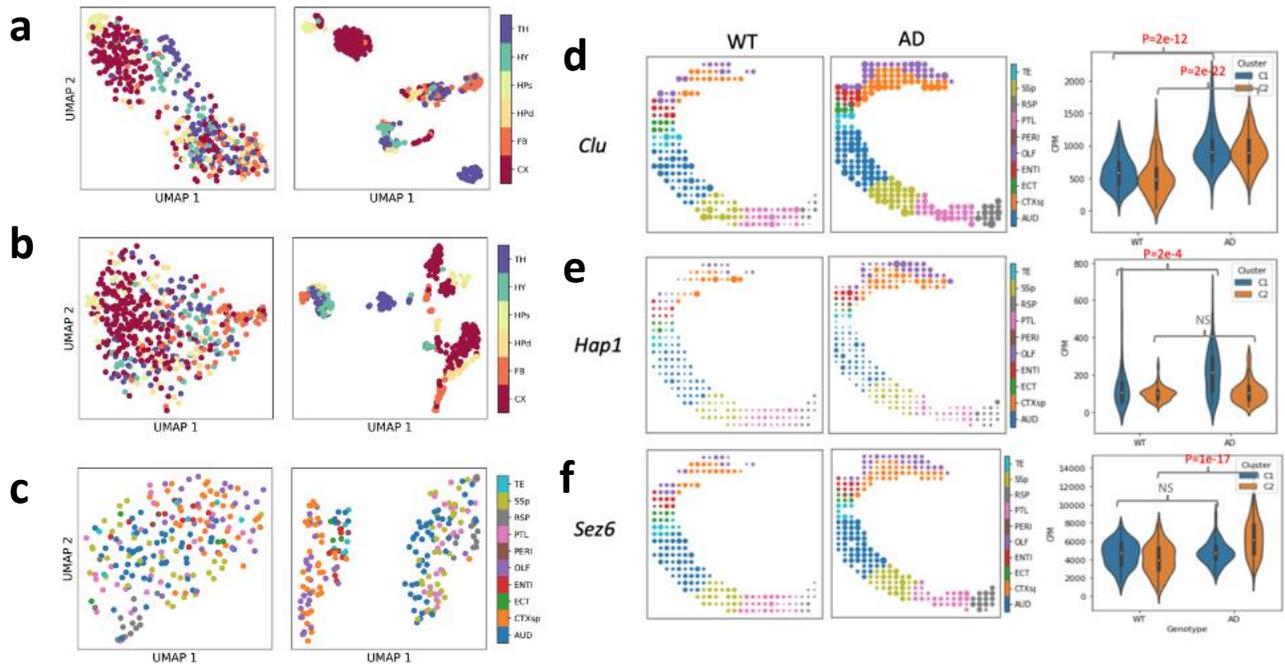
572



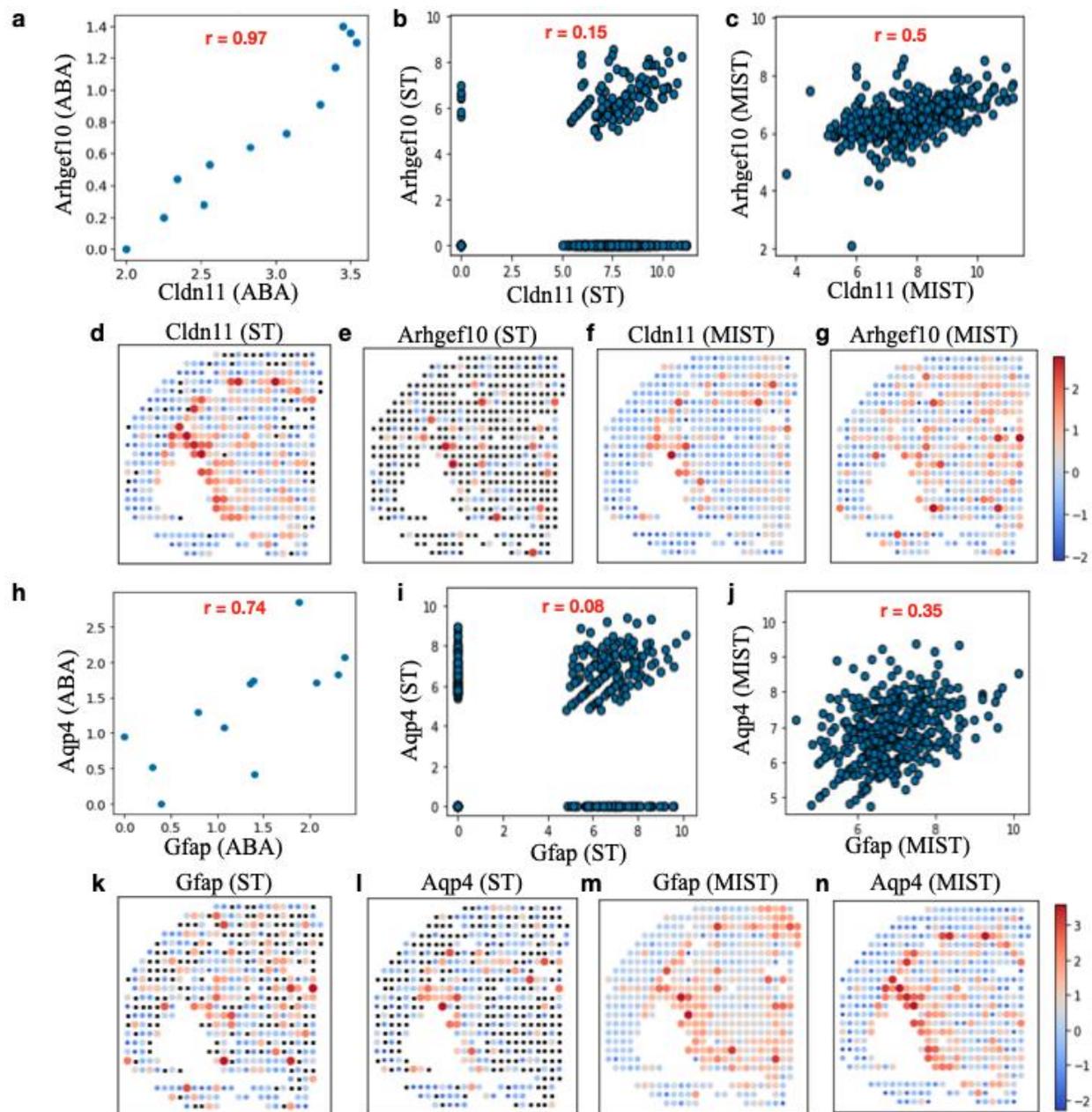
573 **Figure 2. | MIST detected regions agree with manual annotations.** **a**, Melanoma tissue regions
574 detected using different epsilon values. Column 2 with epsilon of 0.6 was selected by MIST.
575 Column 4 is the pathological annotation on the aligned histological staining image. **b**, Mouse
576 brain tissue regions detected using different epsilon values. Column 3 with epsilon of 0.7 was
577 selected by MIST. Column 4 is the spot-level manual annotation from the original study. **c**,
578 Percentage of non-isolated spots (left y-axis, blue curve) and adjusted rand index (ARI, right y-
579 axis, green curve) as functions of edge-filtering parameter epsilon. **d**, Volcano plot of
580 differential expressed genes in tumor region contrasted to lymphoid region of the melanoma
581 sample. **e**, Gene ontologies activated and suppressed for MIST-detected melanoma tumor
582 region compared relatively to lymphoid region.
583



584 **Figure 3. | MIST's outperforms other imputation methods in holdout experiments. a-b,**
585 **Holdout experiment performance across multiple data sets using metrics a, PCC and b, RMSE.**
586 **Each column is an individual data set. Points in the same color indicate the 5 non-overlap fold of**
587 **tests for each model. c, Gene-level performance of each model represented by the RMSE (y-**
588 **axis) as a function of zero-value percentage of the gene expression values (sparsity, x-axis). d,**
589 **Expression patterns recovered for gene GAPDH in the Melanoma sample. From left to right**
590 **shows the spatial pattern of the input to MIST with holdouts, the denoised gene expression**
591 **pattern and the original observed gene expression pattern. Color and size indicate relative gene**
592 **expression abundance. e, the correlation between holdout input and observed GAPDH**
593 **expression values. f, the correlation between MIST denoised and observed GAPDH expression**
594 **values. g, distribution of the gene-level correlation between original non-zero gene expression**
595 **values and the holdout-input (blue), and denoised values (orange), respectively.**
596



597 **Figure 4. | MIST identified intra-cortex heterogeneity within Alzheimer's Disease (AD) mouse**
 598 **brain. a,** UMAP of Mouse WT brain using original observed ST data (left) and MIST-
 599 denoised ST data (right). **b,** UMAP of Mouse AD brain using original observed ST data (left) and MIST-
 600 denoised ST data (right). **c,** UMAP of Mouse AD cortex region using original observed ST data
 601 (left) and MIST-denoised ST data (right). **d-f,** Examples of AD up-regulated genes that are
 602 activated in both clusters (**d**), only significantly upregulated in Cluster 1 (**e**, CTXsp: cortical
 603 subplate, OLF: olfactory, ENTI: entorhinal, TE: temporal association, ECT: ectorhinal, and PERI:
 604 perirhinal areas) and Cluster 2 (**f**, AUD: auditory, PTL: posterior parietal association area, RSP:
 605 retrosplenial area, SSp: primary somatosensory). Left, spatial expression pattern in WT mouse
 606 brain cortex; middle, spatial expression pattern in AD mouse brain cortex; right, boxplot of CPM
 607 grouped by genotype and spatial clusters.
 608
 609



610 **Figure 5. | MIST recovers spatially co-expressed gene pairs.** a-c, spatial correlation of *Cldn11*-
611 *Arhgef10* in the (a) Allen Brain Atlas reference, (b) original ST data and (c) MIST-denoised ST
612 data. d, expression patterns of gene *Cldn11* using the original ST data. e, expression patterns of
613 gene *Arhgef10* using the original ST data. f, expression patterns of gene *Cldn11* using the MIST-
614 denoised data. g, expression patterns of gene *Arhgef10* using the MIST-denoised data. h-j,
615 spatial correlation of *Gfap*-*Aqp4* in the (h) Allen Brain Atlas reference, (i) original ST data and (j)
616 MIST-denoised ST data. k, expression patterns of gene *Gfap* using the original ST data. l,
617 expression patterns of gene *Aqp4* using the original ST data. m, expression patterns of gene
618 *Gfap* using the MIST-denoised data. n, expression patterns of gene *Aqp4* using the MIST-
619 denoised data.

620

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [extendedfigures.docx](#)