

1 Spectral sparsification helps restore the spatial 2 structure at single-cell resolution

3 Jingwan Wang^{1,†}, Shiyong Li^{1,†}, Lingxi Chen¹, and Shuai Cheng Li^{1,*}

4 ¹Department of Computer Science, City University of Hong Kong, 83 Tat Chee Ave, Kowloon Tong, Hong Kong,
5 China

6 [†]These authors contributed equally to this work.

7 ABSTRACT

Single-cell RNA sequencing thoroughly quantifies the individual cell transcriptomes but renounces the spatial structure. Conversely, recently emerged spatial transcriptomics technologies capture the cellular spatial structure but skimp cell or gene resolutions. Cell-cell affinity estimated by ligand-receptor interactions can partially reconstruct the quasi-structure of cells but falsely include the pseudo affinities between distant or indirectly interacting cells. Here, we develop a software package, STORM, to reconstruct the single-cell resolution quasi-structure from the spatial transcriptome with diminished pseudo affinities. STORM first curates the representative single-cell profiles for each spatial spot from a candidate library, then reduces the pseudo affinities in the intercellular
8 affinity matrix by partial correlation, spectral graph sparsification, and spatial coordinates refinement. STORM embeds the estimated interactions into a low-dimensional space with the cross-entropy objective to restore the intercellular quasi-structures, which facilitates the discovery of dominant ligand-receptor pairs between neighboring cells at single-cell resolution. STORM reconstructed structures achieved shape Pearson correlations ranging from 0.91 to 0.97 on the mouse hippocampus and human organ tumor microenvironment datasets. Furthermore, STORM can solely *de novo* reconstruct the quasi-structures at single-cell resolution, *i.e.*, reaching the cell-type proximity correlations 0.68 and 0.89 between reconstructed and immunohistochemistry-informed spatial structures on a human developing heart dataset and a tumor microenvironment dataset, respectively.

9 Introduction

10 Revealing the spatial context and molecular abundance of cells and tissue is critical for understanding the
11 composition and functions of complex tissues. Single-cell RNA sequencing (scRNA-seq) technologies quantify the
12 single-cell transcriptome by a high sequencing depth with whole-transcriptome coverage¹. The thorough scope of
13 single-cell transcriptome enables investigations on cell heterogeneities, subpopulations, and interactions^{2,3}. However,
14 the isolation procedure renounces the spatial context of these cells.

15 Spatial transcriptomics (ST) technologies have been developed to acquire spatial context and expression profiles
16 simultaneously. High-plex RNA imaging technologies⁴⁻⁶ only localize dozens to hundreds of genes, and spatial
17 barcoding technologies such as 10X Visium, Slide-Seq, and HDST⁷⁻⁹ yield a greater magnitude. However, they

18 have achieved unsatisfied abundances or inadequate cell resolution, which restricts the potential of ST data for
19 downstream analyses.

20 Except for wet-lab approaches, researchers also proposed computational methods to restore the spatial structure
21 from the scRNA-seq data. NovoSpaRc¹⁰ assigns cells to tissue locations by probability. Its premise only considers
22 the similarity in gene expression as the neighboring factor, neglecting the heterogeneity of, for instance, the transition
23 areas¹¹ or immune cell infiltration regions¹². CSOmap reconstructs the intercellular proximity based on the contact-
24 required ligand-receptor (LR) interactions^{13,14}. Specifically, CSOmap estimates the affinity of two cells by the
25 mRNA expression summation of the interacting LR pairs, forming a k -nearest neighbor affinity graph simulating
26 cell-cell interactions. However, the pseudo affinities in the affinity graph remain untended, leading to a defective
27 reconstruction of spatial structure.

28 Researchers also started to integrate the ST data with the scRNA-seq data. Early attempts for integration focus
29 on reconstructing cellular spatial structure based on spatial references such as immunohistochemistry (IHC) or
30 fluorescence *in situ* hybridization (FISH)^{15,16}. Spatial barcoding presents a new aspect for integrating scRNA-seq
31 and spatial data, leading to two primary integration approaches: deconvolution and mapping¹⁷. One objective of
32 deconvolution methods is to infer the proportion of cell types from each ST capture location or spot in the ST
33 data. Provided with a labeled scRNA-seq dataset, non-negative least squares and dampened weighted least squares
34 linear regression can deconvolute the captured spot mixtures^{18,19}. Alternatively, deconvolution can be accomplished
35 by fitting a model of negative binomial distribution or Poisson distribution to the scRNA-seq expression with
36 the empirical data of ST spot as a prior. Subsequently, maximized posterior yields an estimation of the cell-type
37 distribution^{20–22}. Moreover, several studies on the tumor microenvironment (TME) map subgroups of single-cell
38 to specific subregions in ST data by the enrichment score^{23,24}. These mappings improve the resolution on the
39 subpopulation level but require prior clustering and annotation on both data types, which is inaccurate when
40 mapping tissue regions comprised of mixed cell types. SpaOTsc²⁵ maps cells by minimizing the gene expression
41 dissimilarity between single-cell data and ST with the optimal transport distance, neglecting the heterogeneity in
42 spot.

43 Here, we present a software package, STORM, that recapitulates the single-cell resolution cell quasi-structure
44 of the spatial transcriptome from a sparsified affinity graph where the pseudo affinities are reduced by partial
45 correlation²⁶, spectral sparsification²⁷, and spatial coordinates refinement. Instead of solely delivering cell-type
46 acknowledgment, STORM locates single-cell expression profiles in spots from a candidate library, hence enabling the
47 exploration of the spatial intercellular communication mechanisms at single-cell resolution.

48 Results

49 Overview of STORM algorithm: reconstructing spatial organization at single-cell resolution from the spatial 50 transcriptome

51 STORM provides a preprocessing module for ST datasets which select and aggregate single-cell profiles representing
52 the expression profile of each spot. For a spot of the spatial data, the module derives the quantities of each cell type
53 by deconvoluted cell type proportions produced by the stereoscope²⁰ and a prespecified parameter ℓ_s representing
54 the average number of cells in a spot (Figure 1a). The module then aggregates a set of single cells agreeing with
55 the derived quantities and maximizing the correlation between the aggregated cell expression profile and the ST
56 spot. Note that if the paired single-cell data are unavailable, we can use a labeled single-cell candidate library of the
57 similar tissue to create aggregations (Figure 1b).

58 Cells interact with proximal cells, and in this work, we use the term *affinity* as the measurement for the interaction
59 strengths between interacting cells. We can build a cellular spatial configuration, termed *quasi-structure*, from the
60 affinity values. We first assume that the cell-cell affinity can be estimated by the concentration of LR complexes
61 which can be approximated by their mRNA abundance. Furthermore, we assume that cells compete for space because
62 of the limitation of biological constraints. STORM has no prior knowledge of cell proximity when forming the initial
63 affinity matrix. It calculates the affinity value between any two cells. Therefore, the approximated affinities based
64 on the first assumption contain pseudo affinities between distant or indirect interacting cells. Following the above
65 assumptions, STORM reconstructs the quasi-structure from scRNA-seq data with four steps: (a) establishing the
66 initial affinity matrix by the LR expression profiles, which falsely includes the pseudo affinities between distant or
67 indirect interacting cells; (b) constructing an affinity graph regards cells as vertices and the initial affinity matrix as
68 the adjacency matrix; (c) reducing the underlying pseudo affinities in the initial affinity graph by partial correlation,
69 spectral graph sparsification, and spatial coordinates refinement; and (d) embedding the sparsified affinity graph
70 into a low-dimensional space as the quasi-structure in the single-cell resolution.

71 STORM approximates the cell-cell affinity by the mRNA abundance of interacting LR pairs (Figure 1c). For
72 initial affinities of high variances, STORM replaces the initial cell-cell affinity matrix with the precision matrix
73 to reduce the indirect correlations for subsequent procedures (Figure 1d). STORM reduces the pseudo affinities
74 from the initial affinity matrix by imposing spectral graph sparsification and spatial coordinates refinement on the
75 affinity matrix (Figure 1d). STORM adopts a local fuzzy set (LFS) embedding method to embed the processed
76 affinity matrix to a low-dimensional space. The LFS step first builds a fuzzy topological representation from the
77 processed affinity matrix, limiting the number of neighbors required by the second assumption (Figure 1e, top
78 panel). Subsequently, the LFS step optimizes the representation in the low-dimensional space by minimizing the
79 fuzzy set cross-entropy between the two representations (Figure 1e, bottom panel). STORM can take the curated
80 single-cell aggregates, yielding the reconstructed quasi-structure for downstream analyses (Figure 1f). The embedding

81 result, that is, the reconstructed quasi-structure by STORM, facilitates further evaluation of discovering dominant
82 ligand-receptor pairs between neighboring cells at single-cell resolution (Figure 1g). Furthermore, with proper
83 sparsification, STORM is capable of *de novo* reconstruction from the single-cell transcriptome. In the head and neck
84 cancer (HNC) scRNA-seq dataset, STORM recapitulates the quasi-structure features which are commonly observed
85 in the partial epithelial to mesenchymal transition (p-EMT) process: (a) p-EMT cells locating at the interface
86 between malignant cells and cancer-associated fibroblasts (CAF) cells; (b) CAF-1 cells presenting at closer proximity
87 to the p-EMT cells compared to CAF-2 cells; (c) malignant cells showing minimum interactions with immune cells
88 due to immune evasion (Figure 1h).

89 **Assessing the performance of STORM in processing ST datasets**

90 We demonstrate the performance of STORM on simulated and real-world ST datasets.

91 ***In silico evaluation of STORM on processing the ST datasets***

92 We assess the validity of STORM in processing ST data coupled datasets by simulated datasets generated from
93 the scRNA-seq data of the mouse hippocampus²⁸. Since neurons, oligodendrocytes, and astrocytes are the main
94 constituents of the hippocampus, we prepare two distinct scRNA-seq candidate libraries and coupled ST data
95 for the simulated datasets: library A consisting of astrocytes and neuron cluster 1, and library B consisting of
96 oligodendrocytes and neuron cluster 2. Moreover, the average number of cells per spot varies according to the
97 tissue density and the spot diameter^{23,29,30}. Therefore, we simulate ST data with the parameter number of cells
98 per spot set as 10, 20, 30, and 40 to test the adaptability of the preprocessing module. Meanwhile, we perform
99 five simulations for each parameter and candidate library to assess the robustness of STORM. Every simulated ST
100 dataset consists of 30 spots. For each spot in the dataset, we arbitrarily sample the designated number of cells from
101 each candidate library and regard the aggregated expression profile of these selected cells as the spot expression
102 simulating the ST profile.

103 The preprocessing module of STORM selects 300 ($\ell_s = 10$), 600 ($\ell_s = 20$), 900 ($\ell_s = 30$) and 1200 ($\ell_s = 40$) cells
104 respectively from each candidate library, constituting 30 single-cell aggregates to represent the expression profile of
105 ST spots. The aggregated expression profiles of each single-cell aggregate regarding various parameters ℓ_s achieve
106 an average Pearson correlation coefficient $r = 0.97$ with their corresponding ST profiles (Figure 2a). Moreover,
107 we perform a paired t-test on the expression correlation of simulations across different cell-number parameters in
108 each candidate library. In the best simulation of each candidate library, that is, the simulation with the highest
109 average expression correlation, we observe that in candidate library A, the expression correlation differences between
110 parameter ten and other parameters are significant. Yet, in candidate library B, the differences between parameters
111 are not statistically significant (Figure 2b).

112 Subsequently, STORM reconstructs the quasi-structure from the selected single-cell aggregates. The quasi-
113 structure of each simulation reached a high shape correlation of $r = 0.94$, $\ell_s = 10$ with the simulated spot organization

114 (Figure 2c). The quasi-structure has a coincidental interspot organization as the cells originating from the same
115 spot remain in the same compartment as illustrated by the Voronoi partition (Figure 2c, right). Furthermore, we
116 compare the shape correlation between various parameters, and the average correlations decrease with the increase
117 in cell number per spot (Figure 2d).

118 We also calculate the expression correlation of each gene between ST and single-cell aggregates. *Nckap5l* and
119 *Smdt1*, achieve high expression correlations, $r = 0.82$, $r = 0.67$, between ST spots and single-cell aggregates (Figure 2e).
120 The high-quality single-cell aggregates and the quasi-structure demonstrate the accuracy and robustness of the
121 preprocessing module of STORM.

122 **STORM reconstructed a high-quality quasi-structure for the mouse hippocampus dataset**

123 We apply STORM to reconstruct the single-cell resolution quasi-structure for the mouse hippocampus dataset.
124 The spatial data provided by stereoscope²⁰ contains 609 spots, and the single-cell library from the *mousebrain.org*
125 contains 8,449 cells, re-clustered and annotated by stereoscope, covering 56 subtypes across seven major groups.

126 The preprocessing module of STORM selects 6,071 cells that constitute 609 single-cell aggregates ($\ell_s = 10$) from
127 the single-cell candidate library to represent the expression profile of ST spots. Then STORM reconstructs the
128 quasi-structure from the selected single-cell aggregates (Figure 3a). The reconstructed quasi-structure of STORM
129 achieves a 0.97 Pearson correlation with its coupling ST spots in the pairwise distance (Supplementary information,
130 Table S1).

131 Furthermore, we calculate the expression correlation of each gene between ST and single-cell aggregates. *Cnp*,
132 *Plp1* and *Ppp3ca*, achieve high expression correlations, $r = 0.71$, $r = 0.70$, $r = 0.65$, between ST spots and single-cell
133 aggregates (Figure 3b, Supplementary information, Table S3). Meanwhile, the aggregated expression profile of each
134 single-cell aggregate achieves a median Pearson correlation coefficient $r = 0.66$ with their corresponding ST profiles
135 (Supplementary information, Fig. S1).

136 The cell-type proximity summarized by cell locations is vital for downstream analyses. Thus, the recapitulation
137 of such information should also be a metric for evaluating the reconstructed quasi-structure. Specifically, we
138 use Kullback-Liebler (KL) divergence to assess the difference of the cell-type proximity between the original and
139 quasi-structure. The quasi-structure achieves a low KL divergence, 0.067, in the cell-type proximity (Figure 3c,
140 Supplementary information, Table S2). Moreover, we assess the effectiveness of each step in STORM by comparing
141 the KL divergence with different combinations of embedding and sparsification methods (Figure 3c). Comparing the
142 LFS embedding that STORM utilizes with constrained t-SNE used by CSOmap, the lower median KL divergence
143 in the combination of LFS embedding with a sparsification method is demonstrated. For sparsification methods,
144 spectral graph sparsification partially reduces the pseudo affinities in the cell-cell affinity matrix, hence achieving a
145 smaller median KL divergence compared to the hard-filtering method of keeping the top fifty high-affinity edges for
146 each node. The additional distance metric provided by spatial information effectively reduces more pseudo affinities

147 in the cell-cell affinity graph, leading to a smaller median KL divergence. The smallest KL divergence, 0.067, is
148 acquired in the combination of LFS embedding and dual sparsification, which suggests the validity of each step in
149 STORM.

150 The well-captured neighboring information in the reconstructed quasi-structure enables identifying the driver LR
151 pairs mediating interactions between cell types. In the reconstructed quasi-structure, we observe that the interactions
152 between lipoprotein receptor-related protein 1 (*Lrp1*) and apolipoprotein E (*apoE*) is the leading interactions
153 among neurons, vascular cells, and astrocytes (Figure 3d, Supplementary information, Fig. S2). LRP1 mediates
154 the metabolism of Amyloid-beta ($A\beta$), whose accumulation is a vital pathogenic element of Alzheimer's disease.
155 Yet apoE can block the LRP1-mediated pathway in astrocytes, hindering the clearance of $A\beta$ ³¹. Hence, certain
156 immunotherapy targeting apoE has been applied on APP/PS1 mice to meliorate the accumulation of $A\beta$ ³². The
157 reveal of the fundamental interaction between *Lrp1* and apoE in our quasi-structure consolidates the validity of
158 STORM and, therefore, its capability of providing valuable biological insights.

159 ***STORM uncovers the metastasis-promoting effect of HMGB1-SDC1 interaction in the human squamous cell carcinoma*** 160 ***dataset***

161 High-quality reconstructions of STORM help reveal the underlying molecular mechanisms of human diseases. We
162 apply STORM on the human squamous cell carcinoma (SCC) dataset of patient 02 in Andrew *et al.*'s work²⁴. The
163 ST and scRNA-seq data are collected from the same malignant skin tissue. The spatial data contains 666 spots, and
164 the matching scRNA-seq data contains 2,689 cells across 14 cell types.

165 The preprocessing module of STORM curates 6,625 cells with replacement regarding the SCC scRNA-seq data
166 as the candidate library ($\ell_s = 10$), forming single-cell aggregates to represent the expression profile of 666 spots in
167 the spatial data. Subsequently, STORM rebuilds the quasi-structure from the curated single-cell aggregates. The
168 reconstructed quasi-structure has high consistency, $r = 0.91$, with its coupling ST spot structure, regarding the
169 pairwise distance (Figure 4a, Supplementary information, Table S1).

170 Furthermore, we calculated the expression correlation of each gene between ST and single-cell aggregates
171 (Supplementary information, Table S3). Several cell-type marker genes annotated in Andrew *et al.*'s work, *e.g.*,
172 *CALML5*, *SPRR1B*, *KRT2*, achieve high expression correlations, $r = 0.79$, $r = 0.65$, $r = 0.61$, between ST spots and
173 single-cell aggregates (Figure 4b). Moreover, the aggregated expression profile of each single-cell aggregate has a
174 median Pearson correlation coefficient $r = 0.72$ with their corresponding ST profiles (Supplementary information,
175 Fig. S1).

176 The quasi-structure achieves a low KL divergence of 0.42 in the cell-type proximity between the original and the
177 quasi-structure (Figure 4c, Supplementary information, Table S2). When comparing across different combinations of
178 embedding and sparsification methods, STORM also achieves the smallest median KL divergence while combining
179 dual sparsification and LFS embedding, which emphasizes the stability of STORM on cancer datasets.

180 Leveraging the high-quality quasi-structure STORM reconstructed, we identify the LR pair HLA-B-CANX as
181 a driving force behind the interaction of T cells, constituting about 29% of the T cell affinities. Our finding is
182 supported by a report regarding an impaired CD8+ T cell-mediated immune response due to the disturbance in
183 HLA-B-CANX interaction in colorectal cancer³³. We investigate the dominating LR pairs facilitating the crosstalk
184 between T and epithelial cells. We identify that the interaction between HMGB1 and SDC1 contributes around
185 30% to the affinity between T and epithelial cells. HMGB1 and SDC1 have been reported to associate with the
186 drug resistance in glioma³⁴. Furthermore, the increase in HMGB1 promotes tissue invasion and metastasis of
187 cancer³⁵, and SDC1 influences the migration of mouse keratinocytes³⁶. Our finding connects HMGB1 with SDC1,
188 indicating that the reported promotion of metastasis may result from the interaction between HMGB1 and SDC1.
189 The discovery demonstrates that the high-quality quasi-structure reconstructed by STORM facilitates disclosing the
190 decisive LR interaction underneath the cell-cell communications.

191 **STORM reveals different dominating LR pairs in two types of cancer cells from the high-quality quasi-structure**

192 Tumor heterogeneity has been an obstacle to cancer therapy since mutant clones escape and thrive from the targeted
193 therapy. Our spatially informed single-cell transcriptome can characterize the driver interactions between distinct
194 subpopulations. We apply STORM on the patient PDAC-A of the pancreatic ductal adenocarcinoma (PDAC)
195 dataset in Moncada *et al.*'s work²³. Three tissue sections of PDAC-A were sequenced. We use the spatial data of
196 replica 1. The ST and scRNA-seq data are processed from the same malignant tissue. The spatial data contains 428
197 spots, and the scRNA-seq data contains 1,926 cells annotated by 17 cell types.

198 Regarding the PDAC scRNA-seq data as the candidate library, the preprocessing module of STORM curates
199 4,289 cells with replacement ($\ell_s = 10$), constructing single-cell aggregates to represent the expression profile of 428
200 spots in the spatial data. Given the high variance in the affinity values of the PDAC dataset, STORM reconstructs
201 the quasi-structure of the curated single-cell aggregates with the precision matrix form of affinity matrix. The
202 reconstructed quasi-structure achieves high similarity, $r = 0.93$, of the pairwise distance with its coupling spatial
203 data (Figure 5a, Supplementary information, Table S1).

204 Moreover, we calculated the expression correlation of each gene between ST and single-cell aggregates (Supple-
205 mentary information, Table S3). The feature gene of the main regions identified in Moncada *et al.*'s work, *CRISP3*,
206 *PRSS1*, *TM4SF1*, also express in the corresponding regions in the quasi-structure (Figure 5b).

207 Furthermore, the quasi-structure achieves a low KL divergence of 0.13 in the cell-type proximity between the
208 original and the quasi-structure (Figure 5c, Supplementary information, Table S2). The KL-divergence between
209 ST and reconstructed quasi-structure decreases after progressively reducing the pseudo affinities by spectral graph
210 sparsification and spatial coordinates refinement. The smallest median KL divergence is also achieved with the
211 combination of LFS embedding and dual sparsification.

212 Subsequently, by evaluating the LR contribution to the cell-cell affinity, we observe that the interaction between

213 HLA-A and APLP2 contributes around 16% to the overall interaction potential in both *TM4SF1*- and *S100A4*-
214 expressing cancer cells (Figure 5d). APLP2 can cause a reduction in the expression of the total cell surface major
215 histocompatibility complex (MHC) class I³⁷, which is a crucial molecule for cancer cell recognition and elimination.
216 The high interaction between HLA-A and APLP2 observed in the quasi-structure indicates a potential immune
217 escape mechanism adopted by both *TM4SF1*- and *S100A4*-expressing cancer cells. Expect for the mutual LR
218 interactions, we also found distinct dominating LR pairs in these two cancer types (Figure 5d). The LR pair
219 ITGB1-SPP1 is a major contributing factor to the interaction between *TM4SF1*-expressing cancer cells between
220 myeloid dendritic cell (mDC) and macrophage (Figure 5e). SPP1 has been proved to abet immune escape in
221 lung adenocarcinoma through its mediation on macrophage polarization³⁸. Experiments have also revealed how
222 ITGB1-SPP1 interaction incites the cancer progression in ovarian cancer³⁹. Our finding suggests that the interaction
223 between ITGB1 and SPP1 potentially triggers the immune escape of PDAC. However, in *S100A4*-expressing cancer
224 cells, the interaction between ITGA3 and CALR is more prevalent (Figure 5d, right). ITGA3 has been identified
225 as a biomarker for diagnosing and prognostic predicting pancreatic cancer⁴⁰. The LR pair ITGA3-CALR has also
226 been predicted as a poor-prognostic LR pair by other datasets from the same tissue in the recent work of Suzuki *et*
227 *al.*⁴¹. These discoveries demonstrate that researchers can characterize the tumor heterogeneity with the high-quality
228 quasi-structure by revealing the driver interactions between distinct subpopulations.

229 **Evaluating the effectiveness of STORM on *de novo* reconstruction of single-cell datasets**

230 We have demonstrated that the quasi-structure can be reconstructed from cell-cell affinity with proper sparsification.
231 Therefore, we further evaluate the validity of STORM in reconstructing the spatial organization of scRNA-seq data
232 without prior spatial structure.

233 ***STORM outperforms CSOmap on the hepatocellular carcinoma (HCC) dataset***

234 We apply STORM on the HCC dataset consisting of 1,329 cells from Ren *et al.*'s work, for which the reconstruction
235 of CSOmap obtains a Spearman correlation of $r = 0.69$ in the cell-type proximity with the IHC image of the same
236 tumor sample. Given the large variance in the initial affinity values of the HCC dataset, STORM rebuilds the
237 quasi-structure with the precision matrix form of affinity matrix. Compared with CSOmap, the reconstructed
238 quasi-structure of STORM is visually less compact (Figure 6a) and achieves higher cell-type proximity, that is, a
239 Spearman correlation of $r = 0.89$, with its IHC image reference (Figure 6b).

240 Subsequently, we evaluate the performance of combinations in embedding and sparsification methods regarding the
241 cell-type proximity similarity (Figure 6b, Supplementary information, Table S4). A higher correlation is observed in
242 the combination of LFS embedding and any sparsification method when comparing LFS embedding with constrained
243 t-SNE. Moreover, spectral graph sparsification reduces the pseudo affinities, achieving a higher correlation than the
244 hard-filtering method. The comparison between different combinations reveals the collaborative contribution of LFS
245 embedding and spectral graph sparsification for reconstruction.

246 The high-quality reconstructed structure enables investigations on intercellular regulatory mechanisms. The
247 interaction between regulatory T cells (Tregs) and CD8+ T cells suggests an ongoing suppression of the immune
248 response⁴², during which Treg cells induce the p38 and ERK1/2 signaling pathways in effective T cells, which
249 initiate DNA damage, resulting in cell senescence⁴³. Consistent with the previous study, we observe an increase in
250 the mRNA expression of ERK1 in the Treg-CD8+ T cell interacting area, indicating the potential of STORM in
251 discovering the immune response signals hidden in the scRNA-seq data.

252 Furthermore, the well-captured cell-type proximity in the quasi-structure enables the analysis of the dominating
253 LR pairs contributing to the cell-cell affinity. We analyze the main LR pairs between any two cell types. Specifically,
254 we identified the difference in the dominating LR pair between Tregs and CD8+ T cells as well as between Treg and
255 exhausted T cells, which indicates a distinct regulation mechanism of Treg in these two types of cells. *CCL5* is one of
256 the signature genes identified in exhausted T cells⁴⁴. The contribution of CXCR3-CCL5 increases in the interaction
257 between Treg and exhausted T cells compared with CD8+ T cells. Indicating that the Tregs originated expression
258 of CXCR3 may trigger the exhaustion.

259 The discovery demonstrates that the high-quality quasi-structure reconstructed *de novo* by STORM promotes
260 the reveal of the LR interaction underneath the cell-cell regulatory mechanism.

261 ***STORM recapitulates the signal transmission process in the developing human heart.***

262 We apply STORM on a human developing heart dataset consisting of 3,717 cells from the 6.5 post-conception weeks
263 (PCW) heart⁴⁵. We apply STORM to reconstruct the quasi-structure of the heart dataset. The 3D quasi-structure
264 of the developing human heart demonstrates a compact structure (Figure 7a, left). The atrial cardiomyocytes are
265 spatially segregated from ventricular cardiomyocytes (Figure 7a, middle), which is consistent with the separation of
266 the atrium and the ventricle in anatomy (Figure 7a, right). Moreover, we evaluate the cell-type proximity similarity
267 between the quasi-structure and the *in situ* sequencing data. The quasi-structure achieves a high normalized
268 Spearman correlation of $r = 0.68$ in the cell-type proximity.

269 When comparing the different combinations of embedding and sparsification methods, Figure 7b demonstrates that
270 the reconstructed quasi-structure rebuilt by the combination of spectral sparsification and LFS embedding achieves
271 the highest resemblance in cell-type proximity (Supplementary information, Table S4). The cell-type proximity
272 STORM recapitulated includes fibroblasts and cardiac cells (Figure 7c), enabling fibroblasts to modify gene and
273 protein expression, and ultimately cardiac function⁴⁶. Ang II activates the paracrine secretion of TGF- β 1 (*TGFB1*,
274 transforming growth factor- β) and endothelin-1 (*EDN1*) in fibroblasts, leading to the cardiac myocyte hypertrophy
275 (Figure 7d)⁴⁷. Angiotensinogen (*AGT*) is a precursor for angiotensin I, which will be eventually converted to
276 Ang II for further activities⁴⁸. Therefore, we inspect the proximity of *AGT* high-express cell and *TGFB1*, *EDN1*
277 high-express cell through the neighboring cell pair numbers between these cells in the quasi-structure (Figure 7d).
278 We consider a pair of cells are neighboring if the distance is less than the median distance between any cell to

279 its third-nearest neighbor. The proximity between cells that express critical signaling genes provides conditions
280 for signaling through paracrine, consistent with the experimentally validated signaling pathway. This consistency
281 indicates the effectiveness of the quasi-structure rebuild by STORM to reveal the local signal transmission process in
282 the tissue.

283 Discussion

284 The combination of the spatial context and expression profile of each cell enables our understanding of the
285 intercellular regulation mechanism of tissue homeostasis and pathogenesis. The scRNA-seq discards the spatial
286 context, and ST technologies skimp the cell resolution. Therefore, current technologies are inadequate to produce
287 the spatial structure of tissues with single-cell resolution. In this work, we presented STORM to reconstruct the
288 single-cell resolution spatial structure from the spatial and/or single-cell transcriptome. STORM rebuilds the
289 quasi-structure of cells by embedding the sparsified affinity graph to a low-dimensional space. The reconstruction
290 accuracy of STORM has been demonstrated in the mouse hippocampus, human heart, and tumor microenvironment
291 of different organs in expression similarity, shape similarity, and cell-type proximity.

292 Although STORM relies on a comprehensive and valid LR pair database, extensive tests across different organisms
293 and diseases demonstrate a consistent performance of STORM. The recapitulation of literature-supported major LR
294 interactions in TMEs and immune responses also shows the effectiveness of the default LR datasets in providing valid
295 biological observations. However, STORM can delineate a broader range of interactions with a higher accuracy if a
296 more extensive LR pair network is expected with future developments. In addition, the preprocessing module benefits
297 from a comprehensive single-cell candidate library. It is therefore subjected to the influence of sequencing depth of
298 ST data, the imbalanced sizes, inconsistent cell-type constitution, and batch effects between ST and scRNA-seq
299 data, and the accuracy of the estimated cell numbers per spot. Nevertheless, our evaluations consistently show that
300 STORM produces high-correlation quasi-structures across various paired and unpaired datasets with different library
301 sizes. In particular, we recommend using paired datasets for disease studies to ensure an accurate reconstruction
302 against high heterogeneity among samples. In contrast, unpaired datasets have little influence on normal tissues
303 with smaller divergence in mRNA expression across different samples.

304 Previous deconvolution methods^{18–22} failed to achieve a single-cell resolution, integrative methods either fall
305 short in dealing with heterogeneous tissue^{10,23,24} or omit single-cell datasets without spatial reference²⁵, and
306 LR-based reconstruction¹³ neglected the pseudo affinities of distant or indirect interacting cells. Unlike previous
307 methods, STORM utilizes the single-cell transcriptome, spatial transcriptome, and LR interactions to reconstruct
308 a quasi-structure of cells in single-cell resolution by a curated affinity graph. A limitation of the preprocessing
309 module is that the actual number of cells in each spot varies according to spots and tissues. For instance, tissue like
310 the lung, which contains many alveoli, leaves plenty of cavities in the tissue section⁴⁹. Therefore hard to estimate
311 the cell number in each spot accurately. For future development, we intend to include an algorithm for accurate

312 quantification of cell numbers per spot by the high-resolution histological image of the tissue section.

313 STORM reconstructs the spatial structure in single-cell resolution, utilizing the spatial context of each cell. The
314 quasi-structure facilitates the acquisition of the dominating LR pairs in each cell pair, leading to the discovery of
315 subpopulations based on dominating LR since cell talk subdivides cell functions. With a precise reconstruction,
316 STORM reveals the co-occurrence of different types of cells and divergent colonization of subpopulations, which
317 cannot be detected solely by scRNA-seq or ST technologies. Besides, STORM can acquire the dominating LR
318 pairs in each cell pair, leading to the discovery of novel subpopulations based on dominating LR since cell talk
319 subdivides cell functions. These abilities shed light on the studies on tumor heterogeneity and immune therapy.
320 For instance, identifying the disparity of immune microenvironment around different cancer subpopulations could
321 guide medication and metastatic evaluation. Furthermore, the quantification of intercellular interactions between
322 the cancer cell and immune cell can predicate the prognosis of patients with clinical information.

323 **Materials and Methods**

324 **Constructing the single-cell aggregates to reproduce ST expression profiles**

325 We propose a preprocessing module to integrate ST data with scRNA-seq data. The module takes two parameters,
326 the cell number ℓ_s and cell proportion $p_{s,t}$, $t \in T$ for a ST spot s , T denotes the set of cell types. The parameter ℓ_s
327 denotes the average number of cells in a spot. The number of cells captured in a spot varies according to sequencing
328 methods and tissue density; our module allows users to specify it. The software package *stereoscope*²⁰ can infer
329 $p_{s,t}$. Stereoscope assumes a negative binomial distribution model on single-cell and ST data, building a reference
330 expression profile of each cell type from the scRNA-seq data, then maximizing the posterior estimation to obtain the
331 approximate cell proportion at every spot of the ST dataset.

332 Let $k_{s,t}$ denote the cell number of type t at ST spot s , then $k_{s,t} \approx \ell_s \times p_{s,t} = f_{s,t}$. Note that $f_{s,t}$ can be fractional.
333 Here, we round on $f_{s,t}$ randomly⁵⁰ to acquire the integer number of $k_{s,t}$ while stabilizing the expectation of ℓ_s .
334 Denoting the decimal part of $f_{s,t}$ as $\{f_{s,t}\} \in [0,1)$, $f_{s,t}$ randomly rounds up or down to $k_{s,t}$ according to the
335 probability $P(k_{s,t} = \lceil f_{s,t} \rceil) = \{f_{s,t}\}$.

336 The preprocessing module chooses cells from a predefined library to reproduce the single-cell resolution for the
337 ST data. The summed expression profile of all chosen cells in \mathbf{M}_s termed the *aggregated expressions* $E(\mathbf{M}_s)$. It
338 curates the single-cell aggregates set \mathbf{M}_s by maximizing the Pearson correlation between $E(\mathbf{M}_s)$ and the expression
339 $E(s)$ of spot s ; that is, by the following objective function.

$$\text{maximize } \sum_{s \in S, \mathbf{M}_s \subset \mathcal{L}} \rho(E(\mathbf{M}_s), E(s)) \quad \text{s.t. } k_{s,t} = |\{c \in \mathbf{M}_s | t(c) = t\}| \quad \forall t \in T \quad (1)$$

340 The number of chosen cells from each type in \mathbf{M}_s should be the same as the value of $k_{s,t}$. where $\mathcal{L} \in \mathbb{R}^{m \times n}$ is the
341 expression matrix of the single-cell library composed of m cells and n genes.

342 The module adopts a heuristic method of two steps, initialization and *swapping* to optimize the objective function.
343 The initialization selects top $k_{s,t}$ cells of type t for spot s according to the Pearson correlation coefficients between

344 the spot and the cell from the single-cell candidate library.

345 If a better objective value is obtained, the swapping step swaps a cell in aggregates with a cell from the library.
346 The process is repeated until convergence, or a predefined maximum number of iterations is achieved. The swapping
347 process can be time-consuming, and we adopted a local sensitive hash (LSH) strategy to accelerate the swapping
348 step⁵¹. During the swapping procedure, the module removes one cell from the aggregate \mathbf{M}_s at spot s randomly,
349 denoting the aggregate after the removal as \mathbf{M}_s' . The module chooses a new cell \mathbf{m} in each iteration to further
350 increases the $\rho(E(\mathbf{M}_s' \cup \{\mathbf{m}\}), E(s))$. It can be chosen by querying a cell in LSH that has the highest correlation
351 with $E(s) - E(\mathbf{M}_s')$.

352 The module performs feature selection⁵² on the single-cell candidates to reduce the noise introduced by sequencing
353 and low variable genes by choosing the top 3,000 highly variable genes and 80% highly variable LR genes to maintain
354 the capability to infer the intercellular affinity.

355 **Measuring the intercellular affinity by ligand-receptor interactions in single-cell profiles**

356 We denote the single-cell expression matrix as $\mathbf{T} \in \mathbb{R}^{r \times n}$ consisting of r cells and n genes. With n_{lr} ligand-receptor
357 (LR) pairs, we define the ligand and receptor expression matrices as \mathbf{T}_L and $\mathbf{T}_R \in \mathbb{R}^{r \times n_{lr}}$, whose columns are the
358 corresponding LR pairs' ligand and receptor expressions, respectively. The multiplication of the two expression
359 matrices yields the affinity between each pair of cells suggested by the co-expression of each LR pair. As a
360 cell can simultaneously express both ligand and receptor genes, we have two symmetric terms $\mathbf{A}_1 = \mathbf{T}_L \mathbf{T}_R^T$ and
361 $\mathbf{A}_2 = \mathbf{T}_R \mathbf{T}_L^T$ representing two possible LR orders in each cell pair. We formulate the *initial affinity matrix* \mathbf{W} as
362 $\mathbf{A}_1 + \mathbf{A}_2 = \mathbf{T}_L \mathbf{T}_R^T + \mathbf{T}_R \mathbf{T}_L^T$ of size $r \times r$.

363 **Reducing the pseudo affinities to refine the affinity matrix by sparsification**

364 The initial affinity matrix includes pseudo affinities between distant or indirectly interacting cells. Here we present
365 three different approaches for diminishing the pseudo affinities, that is, partial correlation, spectral graph sparsification,
366 and spatial coordinates refinement for ST coupled datasets.

367 We first adopt partial correlation to reduce the pseudo affinities for initial affinities of high variance²⁶. Partial
368 correlation identifies the latent variables representing direct causation and removes indirect relationships among
369 entities^{53, 54}. While a covariance matrix represents the relations between any two entities, the inverse of a covariance
370 matrix, also known as the precision matrix, approximates the partial correlations among entities⁵⁵. For the block
371 expression matrix $\mathbf{T}_{LR} = \begin{pmatrix} \mathbf{T}_L \\ \mathbf{T}_R \end{pmatrix}$, we denote its covariance matrix as \mathbf{K} , that is, $\mathbf{K} = \mathbf{T}_{LR} \mathbf{T}_{LR}^T$. In particular, we have
372 the block form of $\mathbf{K} = \begin{pmatrix} \mathbf{S}_L & \mathbf{A}_1 \\ \mathbf{A}_2 & \mathbf{S}_R \end{pmatrix}$, where $\mathbf{A}_1 = \mathbf{T}_L \mathbf{T}_R^T$, $\mathbf{A}_2 = \mathbf{T}_R \mathbf{T}_L^T$, and $\mathbf{S}_L = \mathbf{T}_L \mathbf{T}_L^T$ and $\mathbf{S}_R = \mathbf{T}_R \mathbf{T}_R^T$ represent
373 the ligand and receptor gene expression similarity between any two cells. We could distinguish direct and indirect
374 LR interactions among cells and keep the direct ones by using the precision matrix of \mathbf{K} , *i.e.*, $\mathbf{K}^{-1} = \begin{pmatrix} \mathbf{K}_{11}^{-1} & \mathbf{K}_{12}^{-1} \\ \mathbf{K}_{21}^{-1} & \mathbf{K}_{22}^{-1} \end{pmatrix}$.
375 Therefore, we have $\mathbf{W} = \mathbf{K}_{11}^{-1} + \mathbf{K}_{12}^{-1} + \mathbf{K}_{21}^{-1} + \mathbf{K}_{22}^{-1}$ representing direct LR interactions.

376 We build the affinity graph \mathbf{G} by regarding the cells as vertices and the cell-cell affinity as the edge weight. When
377 the context is clear, we also refer \mathbf{W} as the adjacency matrix for \mathbf{G} for notation simplicity. We further denote the
378 Laplacian matrix of \mathbf{G} as \mathbf{L} . Therefore, we apply the Spielman-Srivastava spectral graph sparsification algorithm⁵⁶
379 to remove pseudo affinities. Spectral graph sparsification aims to find a sparse approximation of the original graph
380 while maintaining high spectral similarity between two graphs²⁷. In the Spielman-Srivastava algorithm, the effective
381 resistance, *i.e.*, the distance between two vertices connecting by an edge is proportional to the reciprocal of its edge
382 weight. In the sparsification step, edges are sampled by the probabilities proportional to their effective resistances.
383 The algorithm preserves the spectrum of the graph Laplacian, *i.e.* the eigenspaces spanned by eigenvalues, and their
384 relations by requiring high similarity between the two Laplacian matrices, while some previous works only maintain
385 the span of the dominant eigenvectors^{57,58}. We define the effective resistance between two cells u and v as

$$\text{Reff}(u, v) = (\delta_u - \delta_v)^T \mathbf{L}^{-1} (\delta_u - \delta_v) \quad (2)$$

386 where $\delta_u \in \{0, 1\}^r$ is the indicator vector of vertex u . Following the definition, the sparse graph preserves the crucial
387 edges of the original graph. We sample the edge (u, v) by the probability $p_{u,v} = \min\{1, C \cdot (\log r) W_{u,v} \cdot \text{Reff}(u, v) / \epsilon^2\}$,
388 where C is some constant and ϵ is the approximation parameter. We further adjust the weight of the sampled edge
389 (u, v) as $W_{u,v} / p_{u,v}$. We determine the value of the term C / ϵ^2 by the user-defined proportion of preserved edges
390 $q = 2 \sum_{u,v} p_{u,v} / r(r-1)$. Since the expected number of chosen edges can be bounded by

$$\sum_{u,v} p_{u,v} = \sum_{u,v} \min\{1, C \cdot (\log r) W_{u,v} \cdot \text{Reff}(u, v) / \epsilon^2\} \leq \frac{Cr \log r}{\epsilon^2} \quad (3)$$

391 where $\frac{C}{\epsilon^2} \geq \frac{\sum_{u,v} p_{u,v}}{r \log r} = \frac{q(r-1)}{2 \log r}$, thus by adjusting the parameter q we can control the percentage of preserved edges.

392 Moreover, we utilize the spot coordinates in the coupled spatial data as one sparsification approach. If two cells
393 belong to nonadjacent spots, the affinity between them is considered to be pseudo affinities.

394 **Reconstructing the quasi-structure with fuzzy set cross-entropy embedding**

395 The embedding of a cell-cell affinity graph to a low-dimensional space consists of two stages: (a) forming a topological
396 representation \mathbb{W} of sparsified the cell-cell affinity \mathbf{W} ; and (b) finding an embedding \mathbb{E} in the low-dimensional
397 space of the topological representation to minimize the discrepancy between the embedding and the representation.
398 A reliable topological representation of \mathbf{W} should maintain the affinity relations while restricting the number of
399 neighbors for each cell. Here, we maintain the top k_n affinities in \mathbf{W} for each cell while setting other values to
400 be zeros. Subsequently, we perform min-max normalization on the remaining affinities to obtain the membership
401 strength in the range of $[0, 1]$, denoting the matrix as \mathbb{W} . The fuzzy simplicial set expands the classical binary
402 definition of membership by allowing continuous membership strength in the range of $[0, 1]$ ⁵⁹, and the union of the
403 fuzzy simplicial sets⁶⁰ yields the fuzzy topological representation. Hence, \mathbb{W} is the fuzzy topological representation
404 of \mathbf{W} .

405 Subsequently, we apply strategies from UMAP⁶¹ to minimize the fuzzy set cross-entropy between the embedding

406 \mathbb{E} and the topological representation \mathbb{W} , that is,

$$CE(\mathbb{E}, \mathbb{W}) = P(\mathbb{E}) \log \frac{P(\mathbb{E})}{Q(\mathbb{W})} + (1 - P(\mathbb{E})) \log \frac{1 - P(\mathbb{E})}{1 - Q(\mathbb{W})} \quad (4)$$

407 where $P(\mathbb{E})$ and $Q(\mathbb{W})$ represent the normalized adjacency matrices of \mathbb{E} and \mathbb{W} , respectively. We use a spectral
408 layout, that is, the Laplacian matrix of \mathbb{W} to as the initial Cartesian coordinates of \mathbb{E} ⁶². By regarding edges as
409 attractive forces and vertices as repulsive forces, we alternatively apply the attractive and repulsive forces until
410 $CE(\mathbb{E}, \mathbb{W})$ converges to a local minimum.

411 **Evaluating the reconstruction performance of STORM**

412 A major metric for assessing the quality of the reconstructed spatial structure is its reproduction of the spatial
413 characteristics of the tissue. Given a spatial structure of cells, we construct a fixed-volume neighbor graph, where
414 the radius is the median distance between any cell to its third-nearest neighbor. According to the fixed-volume
415 neighbor graph, we quantify the spatial characteristics as the number of neighboring pairs between any two cell
416 types, indicating whether the two are enriched or depleted near each other. Therefore, we evaluate the cell type
417 enrichment or depletion discrepancy by the Kullback-Leibler (KL) divergence⁶³ of the neighboring pair numbers for
418 any two cell types between a given spatial structure and the embedding structure. To further evaluate the statistical
419 significance of observed possible enrichment or depletion, we compare the number of neighboring pairs with 1000
420 random permutations of the cell type labels. We test the enrichment hypothesis, that is, the observed number of
421 neighboring pairs is larger than the random expectation by p -values from both the right-tailed and left-tailed tests.
422 We further adjust the p -values following the Benjamini-Hochberg procedure⁶⁴ and obtain the q -values with a cutoff
423 of 0.05 for significance.

424 **Revealing the dominating LR pairs contributing to intercellular affinity**

425 Given a pair of cell expression profiles \mathbf{E}_i and \mathbf{E}_j , the contribution from the k -th LR pair to the total cell-cell
426 interacting affinity can be formulated as:

$$b_k^{ij} = \frac{\mathbf{E}_{L_k}^i \mathbf{E}_{R_k}^j T + \mathbf{E}_{R_k}^i \mathbf{E}_{L_k}^j T}{\mathbf{E}_L^i \mathbf{E}_R^j T + \mathbf{E}_R^i \mathbf{E}_L^j T} \quad (5)$$

427 The contribution of each LR pair between two cell types t_1 and t_2 is calculated by:

$$b_k^{t_1, t_2} = \frac{1}{N} \sum_{i \in t_1, j \in t_2} b_k^{ij} \quad (6)$$

428 where N is the number of neighboring cell pairs between t_1 and t_2 .

429 **Data availability**

430 The ST and scRNA-seq data we use has been previously published^{13, 20, 23, 24, 28, 45} and are available online
431 mousebrain.org, <https://github.com/almaan/stereoscope>. The PDAC, HNC, and SCC datasets are deposited at the
432 Gene Expression Omnibus under GSE111672, GSE103322, and GSE144240. The count matrix of the developing
433 human heart is available at <https://www.spatialresearch.org> with the erythrocytes and immune cells removed, and
434 the labels we use in this work remain consistent with the original publication. The HCC dataset CSOmap used is

435 deposited at EGA with accession number EGAS00001003449.

436 Code availability

437 The software implementation and analysis notebooks of STORM are available at <https://github.com/deepomicslab/STORM>.

438 References

- 439 1. Wu, A. R. *et al.* Quantitative assessment of single-cell RNA-sequencing methods. *Nat. methods* **11**, 41–46
440 (2014).
- 441 2. Kolodziejczyk, A. A., Kim, J. K., Svensson, V., Marioni, J. C. & Teichmann, S. A. The technology and biology
442 of single-cell RNA sequencing. *Mol. cell* **58**, 610–620 (2015).
- 443 3. Hwang, B., Lee, J. H. & Bang, D. Single-cell RNA sequencing technologies and bioinformatics pipelines. *Exp.*
444 *É molecular medicine* **50**, 1–14 (2018).
- 445 4. Ke, R. *et al.* In situ sequencing for RNA analysis in preserved tissue and cells. *Nat. methods* **10**, 857–860 (2013).
- 446 5. Lubeck, E., Coskun, A. F., Zhiyentayev, T., Ahmad, M. & Cai, L. Single-cell in situ RNA profiling by sequential
447 hybridization. *Nat. methods* **11**, 360–361 (2014).
- 448 6. Eng, C.-H. L. *et al.* Transcriptome-scale super-resolved imaging in tissues by RNA seqFISH+. *Nature* **568**,
449 235–239 (2019).
- 450 7. Ståhl, P. L. *et al.* Visualization and analysis of gene expression in tissue sections by spatial transcriptomics.
451 *Science* **353**, 78–82 (2016).
- 452 8. Rodriques, S. G. *et al.* Slide-seq: A scalable technology for measuring genome-wide expression at high spatial
453 resolution. *Science* **363**, 1463–1467 (2019).
- 454 9. Stickels, R. R. *et al.* Highly sensitive spatial transcriptomics at near-cellular resolution with Slide-seqV2. *Nat.*
455 *biotechnology* **39**, 313–319 (2021).
- 456 10. Moriel, N. *et al.* NovoSpaRc: flexible spatial reconstruction of single-cell gene expression with optimal transport.
457 *Nat. Protoc.* 1–24 (2021).
- 458 11. Thrane, K., Eriksson, H., Maaskola, J., Hansson, J. & Lundeberg, J. Spatially resolved transcriptomics enables
459 dissection of genetic heterogeneity in stage III cutaneous malignant melanoma. *Cancer research* **78**, 5970–5979
460 (2018).
- 461 12. Vestweber, D. How leukocytes cross the vascular endothelium. *Nat. Rev. Immunol.* **15**, 692–704 (2015).
- 462 13. Ren, X. *et al.* Reconstruction of cell spatial organization from single-cell RNA sequencing data based on
463 ligand-receptor mediated self-assembly. *Cell research* **30**, 763–778 (2020).
- 464 14. Wells, A. & Wiley, H. S. A systems perspective of heterocellular signaling. *Essays biochemistry* **62**, 607–617
465 (2018).

- 466 **15.** Achim, K. *et al.* High-throughput spatial mapping of single-cell RNA-seq data to tissue of origin. *Nat. Biotechnol.*
467 **33**, 503–509, [10.1038/nbt.3209](https://doi.org/10.1038/nbt.3209) (2015).
- 468 **16.** Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell gene
469 expression data. *Nat. Biotechnol.* **33**, 495–502, [10.1038/nbt.3192](https://doi.org/10.1038/nbt.3192) (2015).
- 470 **17.** Longo, S. K., Guo, M. G., Ji, A. L. & Khavari, P. A. Integrating single-cell and spatial transcriptomics to
471 elucidate intercellular tissue dynamics. *Nat. Rev. Genet.* 1–18 (2021).
- 472 **18.** Elosua-Bayes, M., Nieto, P., Mereu, E., Gut, I. & Heyn, H. SPOTlight: seeded NMF regression to deconvolute
473 spatial transcriptomics spots with single-cell transcriptomes. *Nucleic acids research* **49**, e50–e50 (2021).
- 474 **19.** Dong, R. & Yuan, G.-C. SpatialDWLS: accurate deconvolution of spatial transcriptomic data. *Genome biology*
475 **22**, 1–10 (2021).
- 476 **20.** Andersson, A. *et al.* Single-cell and spatial transcriptomics enables probabilistic inference of cell type topography.
477 *Commun. biology* **3**, 1–8 (2020).
- 478 **21.** Kleshchevnikov, V. *et al.* Comprehensive mapping of tissue cell architecture via integrated single cell and spatial
479 transcriptomics. *bioRxiv* (2020).
- 480 **22.** Cable, D. M. *et al.* Robust decomposition of cell type mixtures in spatial transcriptomics. *Nat. Biotechnol.*
481 1–10 (2021).
- 482 **23.** Moncada, R. *et al.* Integrating microarray-based spatial transcriptomics and single-cell RNA-seq reveals tissue
483 architecture in pancreatic ductal adenocarcinomas. *Nat. biotechnology* **38**, 333–342 (2020).
- 484 **24.** Ji, A. L. *et al.* Multimodal analysis of composition and spatial architecture in human squamous cell carcinoma.
485 *Cell* **182**, 497–514 (2020).
- 486 **25.** Cang, Z. & Nie, Q. Inferring spatial and signaling relationships between cells from single cell transcriptomic
487 data. *Nat. communications* **11**, 1–13 (2020).
- 488 **26.** de la Fuente, A., Bing, N., Hoeschele, I. & Mendes, P. Discovery of meaningful associations in genomic data
489 using partial correlation coefficients. *Bioinformatics* **20**, 3565–3574, [10.1093/bioinformatics/bth445](https://doi.org/10.1093/bioinformatics/bth445) (2004).
490 <https://academic.oup.com/bioinformatics/article-pdf/20/18/3565/521954/bth445.pdf>.
- 491 **27.** Spielman, D. A. & Teng, S.-H. Spectral sparsification of graphs. *SIAM J. on Comput.* **40**, 981–45 (2011).
492 Copyright - Copyright] © 2011 Society for Industrial and Applied Mathematics; Last updated - 2021-09-11.
- 493 **28.** Zeisel, A. *et al.* Molecular architecture of the mouse nervous system. *Cell* **174**, 999–1014 (2018).
- 494 **29.** Ståhl, P. L. *et al.* Visualization and analysis of gene expression in tissue sections by spatial transcriptomics.
495 *Science* **353**, 78–82 (2016).

- 496 **30.** Vickovic, S. *et al.* High-definition spatial transcriptomics for in situ tissue profiling. *Nat. methods* **16**, 987–990
497 (2019).
- 498 **31.** Liu, C.-C. *et al.* Astrocytic LRP1 mediates brain A β clearance and impacts amyloid deposition. *J. Neurosci.*
499 **37**, 4023–4031 (2017).
- 500 **32.** Kim, J. *et al.* Anti-apoe immunotherapy inhibits amyloid accumulation in a transgenic mouse model of A β
501 amyloidosis. *J. Exp. Medicine* **209**, 2149–2156 (2012).
- 502 **33.** Zheng, J. *et al.* mir-148a-3p silences the CANX/MHC-I pathway and impairs CD8⁺ T cell-mediated immune
503 attack in colorectal cancer. *The FASEB J.* **35**, e21776 (2021).
- 504 **34.** Yuan, D., Tao, Y., Chen, G. & Shi, T. Systematic expression analysis of ligand-receptor pairs reveals important
505 cell-to-cell interactions inside glioma. *Cell Commun. Signal.* **17**, 1–10 (2019).
- 506 **35.** Seidu, R. A., Wu, M., Su, Z. & Xu, H. Paradoxical role of high mobility group box 1 in glioma: a suppressor or
507 a promoter? *Oncol. Rev.* **11** (2017).
- 508 **36.** Stepp, M. A. *et al.* Reduced migration, altered matrix and enhanced TGF β 1 signaling are signatures of mouse
509 keratinocytes lacking *sdcl1*. *J. cell science* **120**, 2851–2863 (2007).
- 510 **37.** Tuli, A. *et al.* Amyloid precursor-like protein 2 association with HLA class I molecules. *Cancer Immunol.*
511 *Immunother.* **58**, 1419 (2009).
- 512 **38.** Zhang, Y., Du, W., Chen, Z. & Xiang, C. Upregulation of PD-L1 by SPP1 mediates macrophage polarization
513 and facilitates immune escape in lung adenocarcinoma. *Exp. cell research* **359**, 449–457 (2017).
- 514 **39.** Zeng, B., Zhou, M., Wu, H. & Xiong, Z. SPP1 promotes ovarian cancer progression via Integrin β 1/FAK/AKT
515 signaling pathway. *OncoTargets therapy* **11**, 1333 (2018).
- 516 **40.** Jiao, Y., Li, Y., Liu, S., Chen, Q. & Liu, Y. ITGA3 serves as a diagnostic and prognostic biomarker for
517 pancreatic cancer. *OncoTargets therapy* **12**, 4141 (2019).
- 518 **41.** Suzuki, S. R., Kuno, A. & Ozaki, H. Cell-to-cell interaction analysis of prognostic ligand-receptor pairs in
519 human pancreatic ductal adenocarcinoma. *Biochem. biophysics reports* **28**, 101126 (2021).
- 520 **42.** Chen, M.-L. *et al.* Regulatory T cells suppress tumor-specific CD8 T cell cytotoxicity through TGF- β signals in
521 vivo. *Proc. Natl. Acad. Sci.* **102**, 419–424 (2005).
- 522 **43.** Liu, X. *et al.* Regulatory T cells trigger effector T cell DNA damage and senescence caused by metabolic
523 competition. *Nat. communications* **9**, 1–16 (2018).
- 524 **44.** Hsu, C.-L. *et al.* Exploring markers of exhausted CD8 T cells to predict response to immune checkpoint inhibitor
525 therapy for hepatocellular carcinoma. *Liver Cancer* 1–14 (2021).

- 526 **45.** Asp, M. *et al.* A spatiotemporal organ-wide gene expression and cell atlas of the developing human heart. *Cell*
527 **179**, 1647–1660 (2019).
- 528 **46.** Howard, C. M. & Baudino, T. A. Dynamic cell–cell and cell–ECM interactions in the heart. *J. molecular*
529 *cellular cardiology* **70**, 19–26 (2014).
- 530 **47.** Gray, M. O., Long, C. S., Kalinyak, J. E., Li, H.-T. & Karliner, J. S. Angiotensin II stimulates cardiac myocyte
531 hypertrophy via paracrine release of TGF- β 1 and endothelin-1 from fibroblasts. *Cardiovasc. research* **40**,
532 352–363 (1998).
- 533 **48.** Jeunemaitre, X. *et al.* Molecular basis of human hypertension: Role of angiotensinogen. *Cell* **71**, 169–180,
534 [https://doi.org/10.1016/0092-8674\(92\)90275-H](https://doi.org/10.1016/0092-8674(92)90275-H) (1992).
- 535 **49.** Brunnström, H. *et al.* Immunohistochemistry in the differential diagnostics of primary lung cancer: an
536 investigation within the southern swedish lung cancer study. *Am. journal clinical pathology* **140**, 37–46 (2013).
- 537 **50.** Raghavan, P. & Tompson, C. D. Randomized rounding: a technique for provably good algorithms and algorithmic
538 proofs. *Combinatorica* **7**, 365–374 (1987).
- 539 **51.** Datar, M., Immorlica, N., Indyk, P. & Mirrokni, V. S. Locality-sensitive hashing scheme based on p-stable
540 distributions. In *Proceedings of the twentieth annual symposium on Computational geometry*, 253–262 (2004).
- 541 **52.** Stuart, T. *et al.* Comprehensive integration of single-cell data. *Cell* **177**, 1888–1902.e21, [https://doi.org/10.](https://doi.org/10.1016/j.cell.2019.05.031)
542 [1016/j.cell.2019.05.031](https://doi.org/10.1016/j.cell.2019.05.031) (2019).
- 543 **53.** Fan, J., Liao, Y. & Liu, H. An overview of the estimation of large covariance and precision matrices. *The*
544 *Econom. J.* **19**, C1–C32 (2016).
- 545 **54.** Li, Y., Hu, J., Zhang, C., Yu, D.-J. & Zhang, Y. ResPRE: high-accuracy protein contact prediction by coupling
546 precision matrix with deep residual neural networks. *Bioinformatics* **35**, 4647–4655 (2019).
- 547 **55.** DeGroot, M. H. *Optimal statistical decisions*, vol. 82 (John Wiley & Sons, 2005).
- 548 **56.** Spielman, D. A. & Srivastava, N. Graph sparsification by effective resistances. *SIAM J. on Comput.* **40**,
549 1913–1926 (2011).
- 550 **57.** Frieze, A., Kannan, R. & Vempala, S. Fast monte-carlo algorithms for finding low-rank approximations. *J.*
551 *ACM* **51**, 1025–1041, [10.1145/1039488.1039494](https://doi.org/10.1145/1039488.1039494) (2004).
- 552 **58.** Batson, J. D., Spielman, D. A. & Srivastava, N. Twice-ramanujan sparsifiers. In *Proceedings of the Forty-First*
553 *Annual ACM Symposium on Theory of Computing*, STOC '09, 255–262, [10.1145/1536414.1536451](https://doi.org/10.1145/1536414.1536451) (Association
554 for Computing Machinery, New York, NY, USA, 2009).
- 555 **59.** Chang, S. S. & Zadeh, L. A. On fuzzy mapping and control. In *Fuzzy sets, fuzzy logic, and fuzzy systems:*
556 *selected papers by Lotfi A Zadeh*, 180–184 (World Scientific, 1996).

- 557 **60.** Gathigi, S. M., Gichuki, M. N., Otieno, P. A. & Were, H. S. Normality and its variants on fuzzy isotone spaces.
558 *Adv. Pure Math.* **3**, 639–642 (2013).
- 559 **61.** McInnes, L., Healy, J. & Melville, J. Umap: Uniform manifold approximation and projection for dimension
560 reduction. *arXiv preprint arXiv:1802.03426* (2018).
- 561 **62.** Koren, Y. On spectral graph drawing. In *International Computing and Combinatorics Conference*, 496–508
562 (Springer, 2003).
- 563 **63.** Kullback, S. & Leibler, R. A. On information and sufficiency. *The annals mathematical statistics* **22**, 79–86
564 (1951).
- 565 **64.** Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to
566 multiple testing. *J. Royal statistical society: series B (Methodological)* **57**, 289–300 (1995).

567 **Acknowledgements**

568 We would like to express sincere gratitude to Dr Wenji Ma for the suggestions on data collection. We appreciate
569 Ms Wenqian Zhang for the manual cell-type annotation of the heart tissue section. This project was supported by
570 SIRG (CityU SIRG 7020005).

571 **Author Contributions**

572 These authors contributed equally: Jingwan Wang, Shiyong Li

573

574 **Affiliations**

575 **Department of Computer Science, City University of Hong Kong, 83 Tat Chee Ave, Kowloon Tong,**
576 **Hong Kong, China**

577 Jingwan Wang, Shiyong Li, Lingxi Chen & Shuai Cheng Li

578

579 **Contributions**

580 SCL conceived and designed the project. J.W. developed the software. J.W. and SYL performed the analysis and
581 validation experiment. SCL, J.W., SYL, and L.C performed manuscript writing, review, and editing.

582

583 **Corresponding authors**

584 Correspondence to [Shuai Cheng Li](#)

585 **Conflict of Interest**

586 The authors declare no competing interests.

587 Figure Legends

Figure 1. Schematics of STORM. **a-c**, Workflow of the preprocessing module. **a**, The preprocessing module of STORM adopts existing deconvolution software to decompose cell-type mixtures of ST profiles. **b**, The preprocessing module selects a designated amount of cells from the single-cell candidate library, equal to the estimated cell number per cell type in each spot. **c-e**, Workflow of the STORM. **c**, STORM derives the initial cell-cell affinity graph from the single-cell profiles by the LR interactions. **d**, STORM applies partial correlation, spectral graph sparsification, and spatial coordinates refinement on the cell-cell affinity graph to reduce pseudo affinities. **e**, STORM utilizes LFS embedding to embed interactions into a low-dimensional space. **f**, The 2D embedding of the selected single cells reconstructed by STORM. **g**, The determination of dominant ligand-receptor pairs between neighboring cells at single-cell resolution. **h**, The 3D embedding of the HNC data reconstructed by STORM.

Figure 2. Evaluation of the validity and robustness of STORM on simulated datasets. **a**, The expression correlation between each spot and its corresponding single-cell aggregate regarding four cell-number parameters across five repeats of candidate libraries A and B. **b**, The best simulation of each cell-number parameter from libraries A and B, annotated with the statistical significance. Asterisks indicate level of statistical significance: ** - significance 0.01, * - significance 0.05, ns - not significant. **c**, The simulated ST structure (left) and the quasi-structure reconstructed by STORM (right). The color stands for each spot. **d**, The Pearson correlation of the pairwise distance between the reconstructed quasi-structure of STORM and its coupling ST spots. **e**, The standardized gene expression of exemplary genes in ST (left) and reconstructed quasi-structure (right).

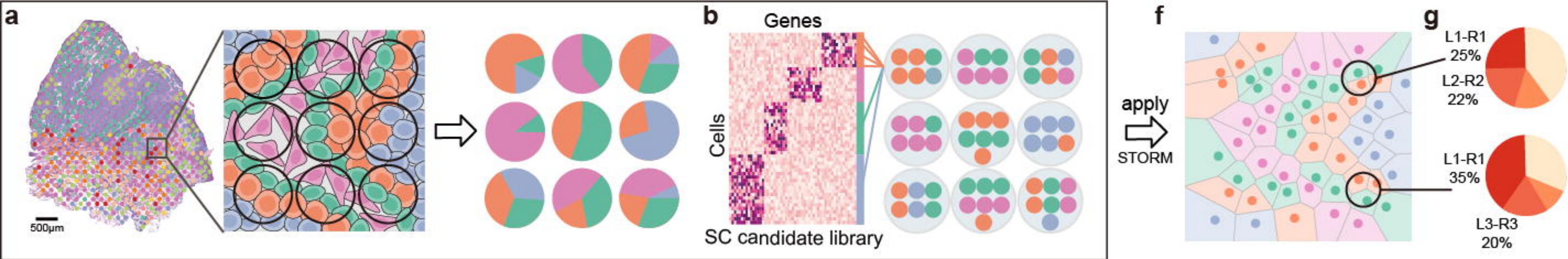
Figure 3. The reconstructed quasi-structure of mouse hippocampus. **a**, The 2D visualization of the ST spots (left) and the reconstructed quasi-structure of the mouse hippocampus (right), colored by cell types. **b**, The standardized gene expression of exemplary genes in ST (top) and reconstructed quasi-structure (bottom). **c**, The cell-type proximity KL divergence for the combinations of two different embedding methods and three sparsification methods. **d**, The pie charts of the LR pair contributions to the interactions of astrocytes with all other cells (top) and with only immune cells (bottom).

Figure 4. Performance of STORM on recapitulating the quasi-structure of SCC dataset. **a**, Spatial (top) and reconstructed quasi-structure (bottom) visualization of SCC, labeled by cell type. **b**, The standardized expression of cell-type marker genes in ST (left) and reconstructed quasi-structure (right). **c**, The cell-type proximity KL divergence of combining two embedding methods and three sparsification methods between ST and reconstructed quasi-structure. **d**, The pie charts of the LR pair contributions to the interaction of T cell and other cells (top), T cell and epithelial cells in particular (bottom).

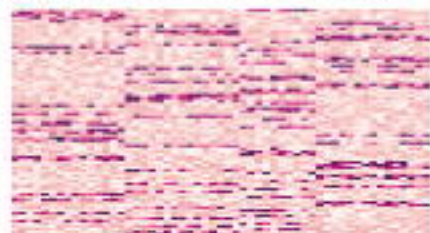
Figure 5. Performance of STORM in rebuilding quasi-structure from PDAC dataset. **a**, ST (top) and reconstructed quasi-structure colored by all cell types (bottom). **b**, The standardized expression of three genes in ST (left) and single-cell aggregates (right). **c**, The cell-type proximity KL divergence for the combination of three sparsification methods and two different embedding methods. **d**, The pie chart of LR pair contributions in *TM4SF1*- and *S100A4*-expressing cells. **e**, The pie chart of LR pair contributions between *TM4SF1*-expressing cells with macrophages (left) and mDC (right).

Figure 6. Application of STORM in restoring the quasi-structure of HCC. **a**, The 3D embedding of the reconstructed quasi-structure of STORM (left) and the prediction of CSOmap (right) on the HCC scRNA-seq data. **b**, Spearman correlation between IHC image-based cell connections (X-axis) and STORM reconstruction (Y-axis). CD8: CD8+ T cells; Tex: exhausted T cell; Treg: Foxp3+ regulatory T cells; M: macrophages; cDC1: CLEC9A+ dendritic cells; O: other cells. **c**, Comparison of cell-type proximity (Spearman) between different embedding and sparsification methods. The green dotted line represents the best Spearman correlation of the CSOmap prediction. **d**, CD8+ T cell and Treg cells in the quasi-structure of STORM colored by cell types (left) and standardized expressions (right). **e**, The pie charts of dominating LR pairs in the interaction of regulatory T cells with CD8+ T cells and exhausted T cells, respectively.

Figure 7. STORM recapitulates the quasi-structure of the developing human heart. **a**, 3D visualization of the reconstructed quasi-structure of developing human heart (left). Ventricular and atrial cardiomyocytes are separately displayed (middle). The tissue section of 6.5 PCW (scale bar: 1 mm), where the ventricular and atrial cardiomyocytes are manually labeled (right). Cell type label is the same as the original data: (0): Capillary endothelium; (1): Ventricular cardiomyocytes; (2): Fibroblast-like (related to cardiac skeleton connective tissue); (3): Epicardium-derived cells; (4): Fibroblast-like (smaller vascular development); (5): Smooth muscle cells / fibroblast-like; (7): Atrial cardiomyocytes; (8): Fibroblast-like (larger vascular development); (9): Epicardial cells; (10): Endothelium / pericytes / adventia; (12): Myoz2-enriched Cardiomyocytes; (14): Cardiac neural crest cells & Schwann progenitor cells. **b**, The normalized Spearman correlation of cell-type proximity in the result of hard-filtering and sparsified graphs embedded by constrained t-SNE (orange) and LFS embedding (blue). **c**, The normalized Spearman correlation between cell-type connections based on spots in the ST section (X-axis) and the quasi-structure reconstructed by STORM (Y-axis), with biases introduced by uneven cell counts among different cell types reduced after normalization. **d**, Mechanism illustration and evaluation of the regulation network between fibroblast and cardiomyocyte. Top-left: schematic diagram of molecular mediation between fibroblast and cardiomyocyte. Bottom-left: standardized expression of above intermediate genes. Right: heatmap of the numbers of neighboring pairs of cells expressing different marker genes.



c single-cell Expression



d Graph Sparsification



Partial Correlation



Spatial Distance



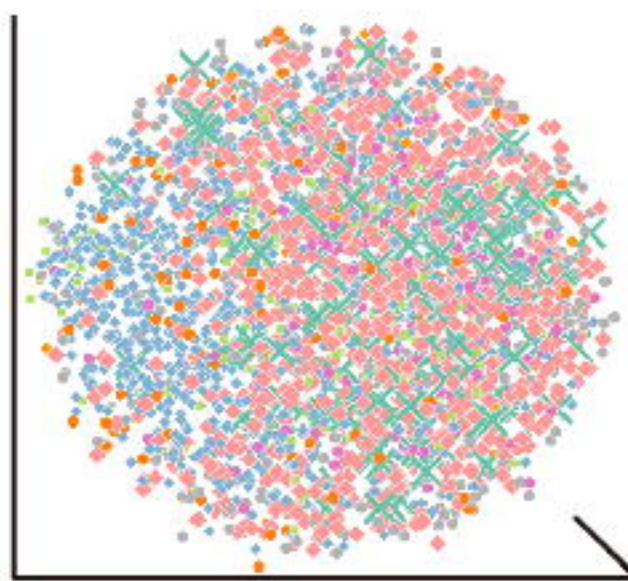
e



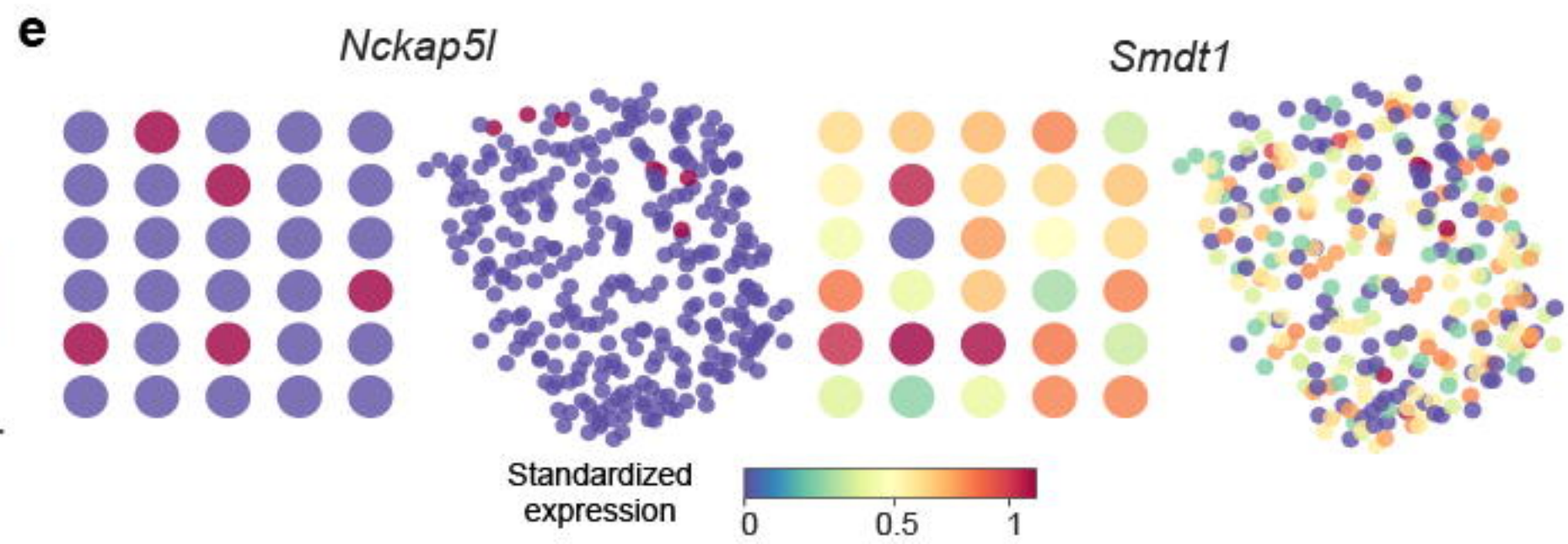
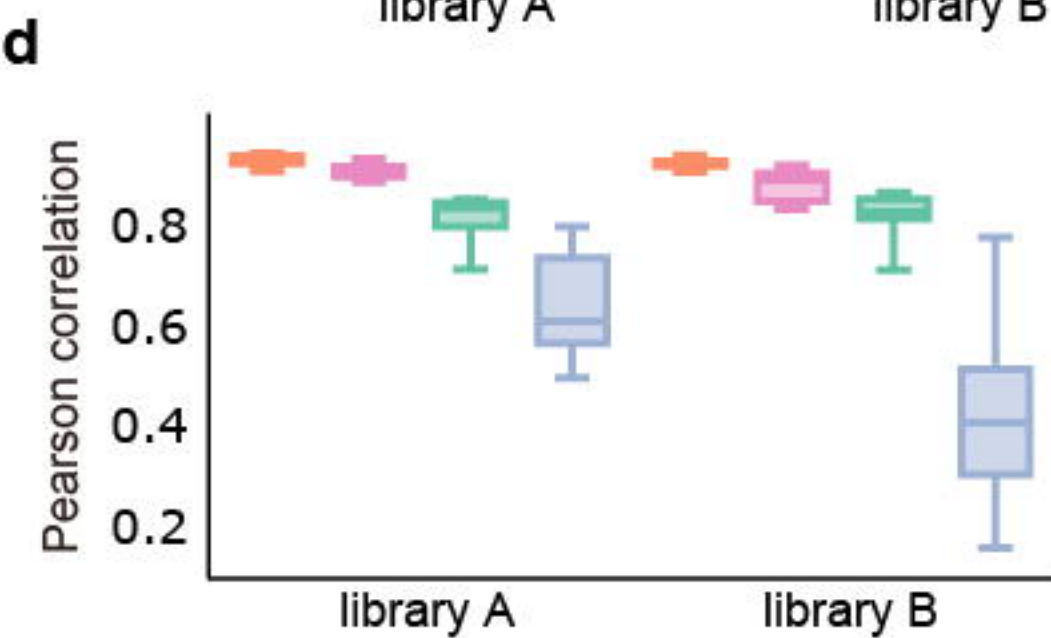
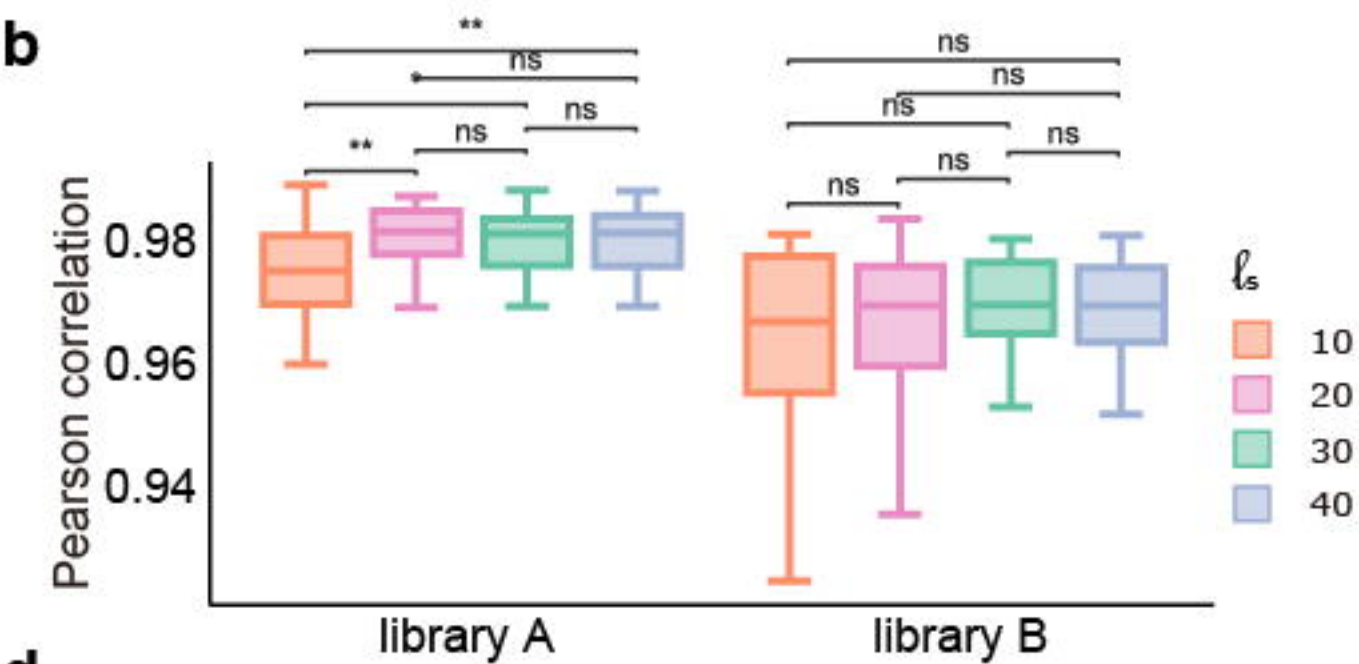
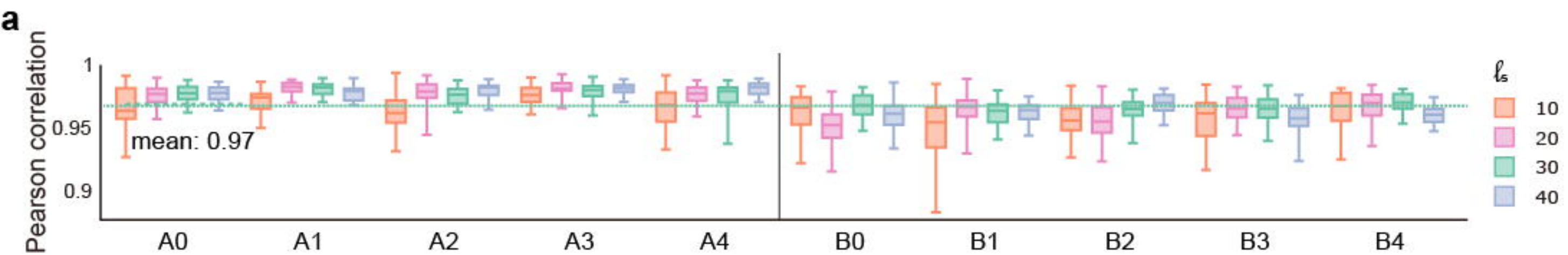
minimize $CE(X, Y)$

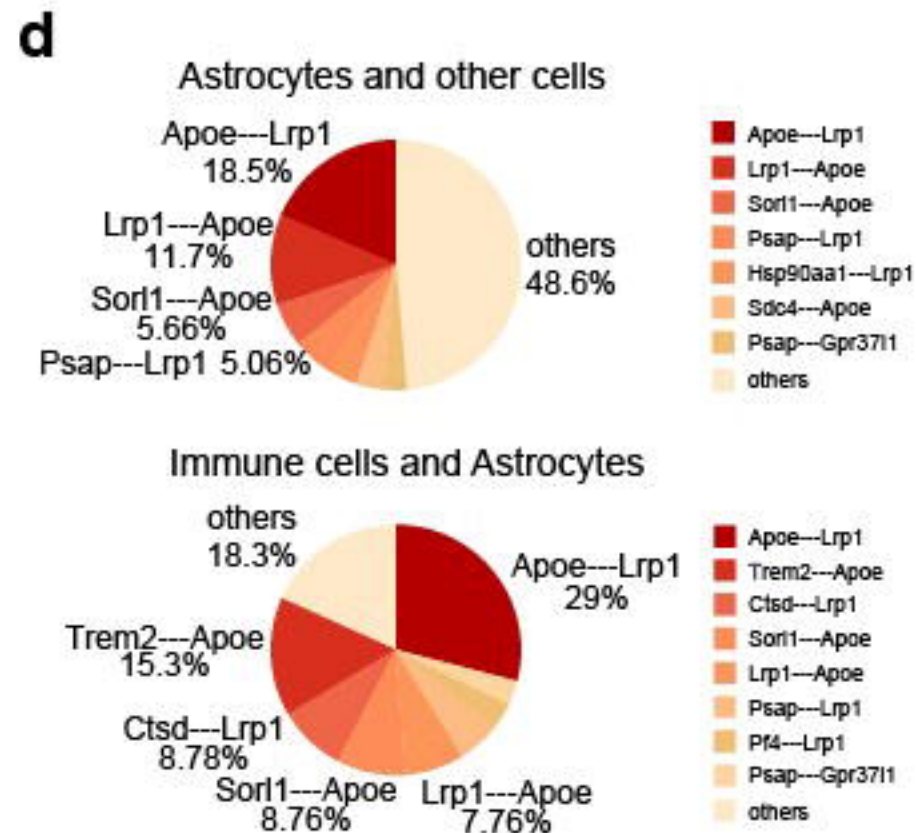
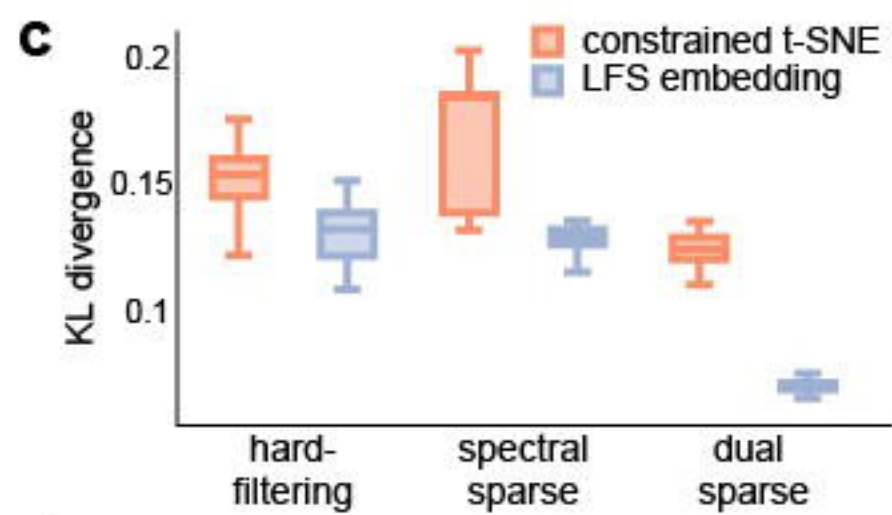
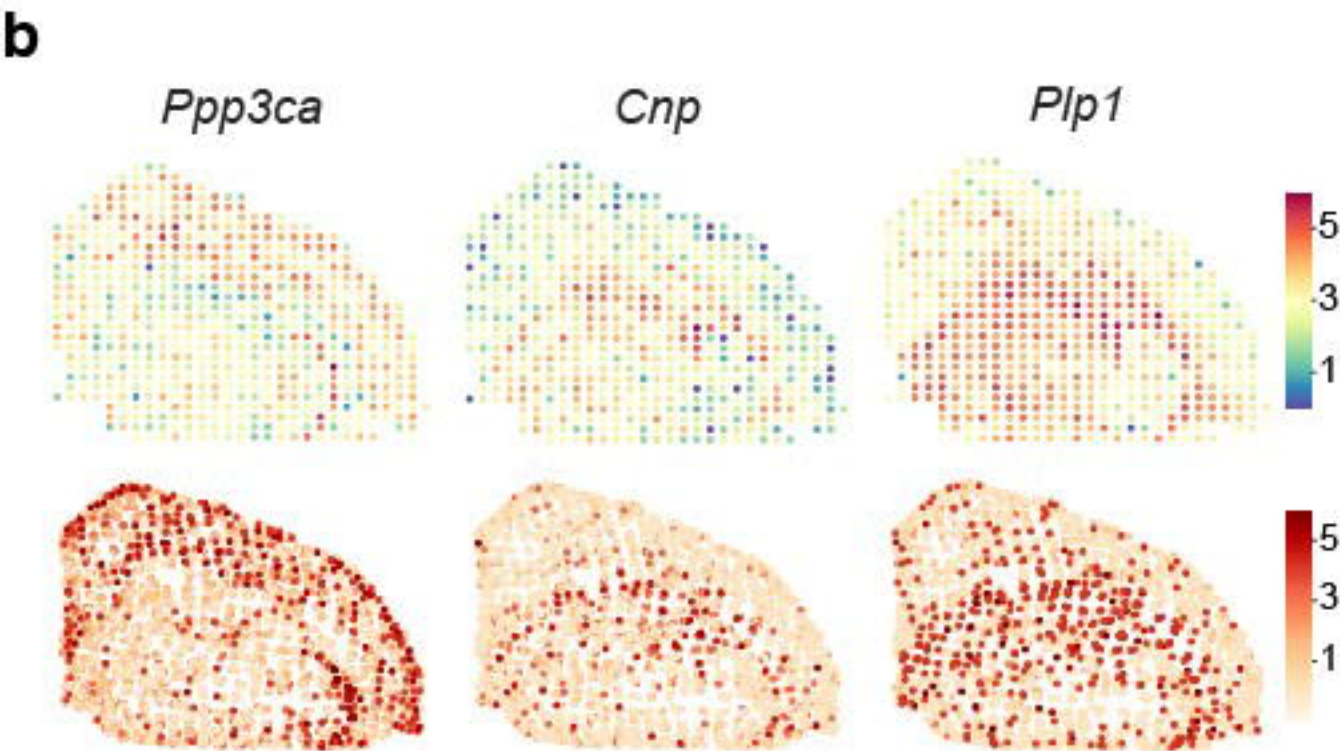


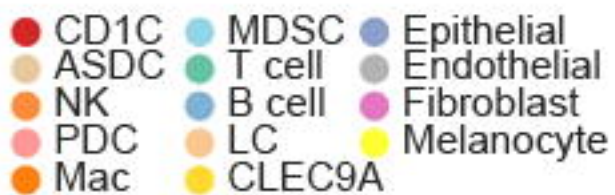
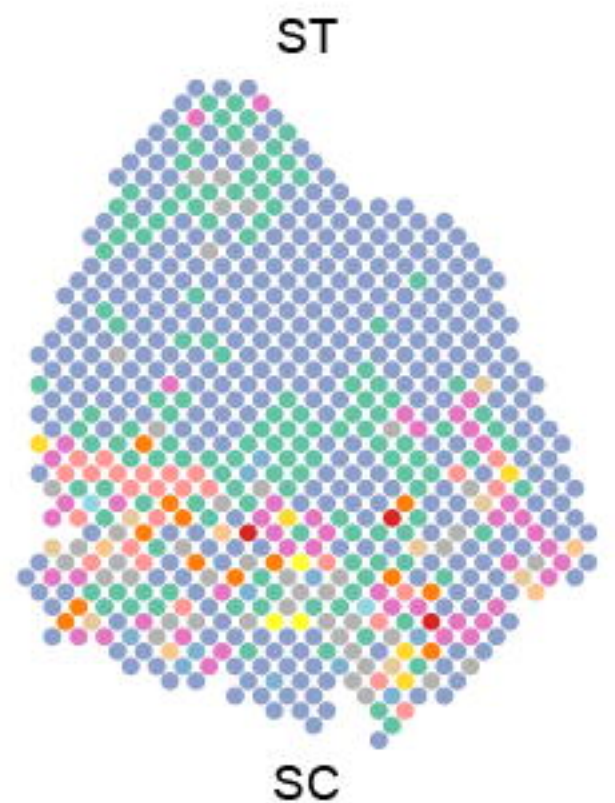
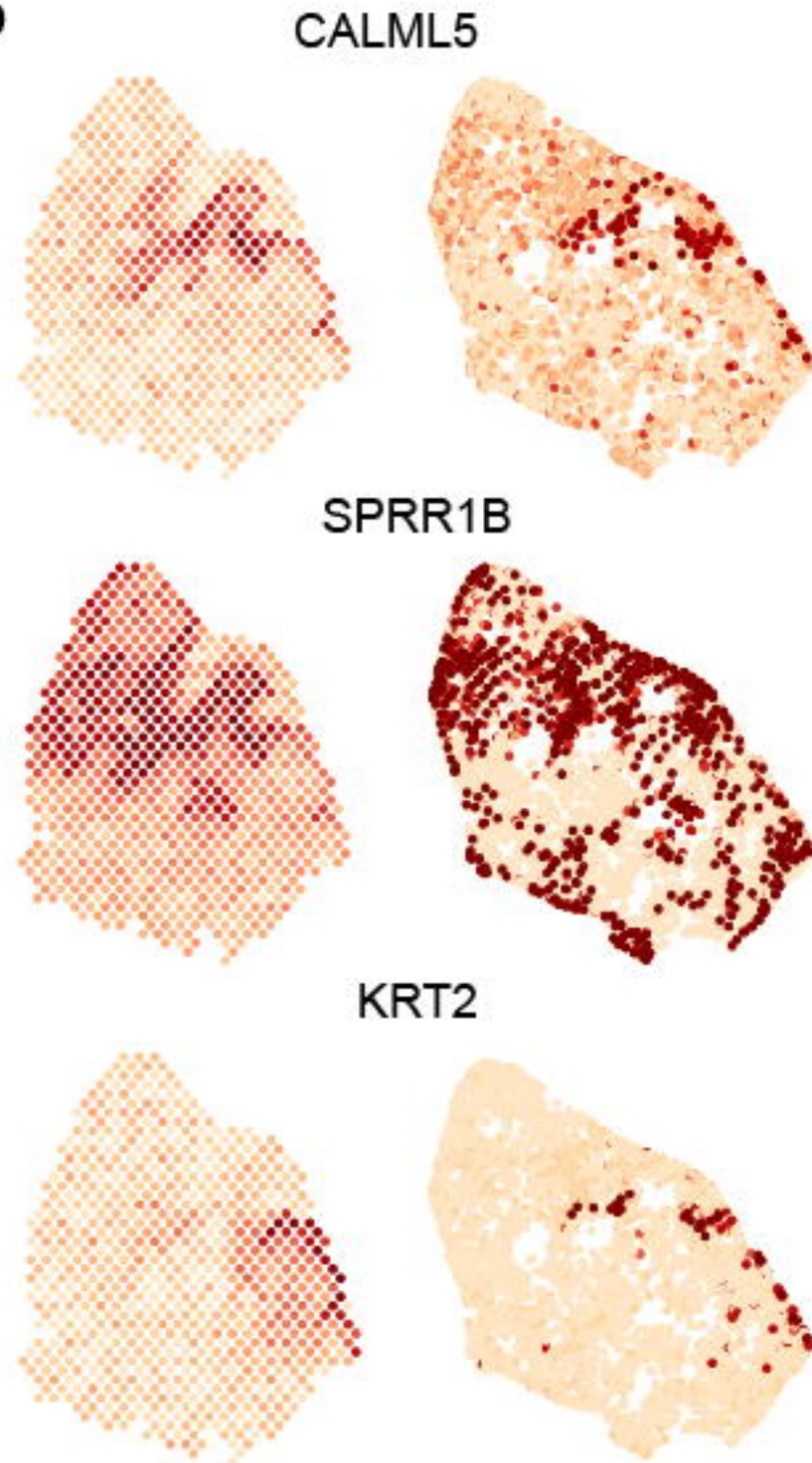
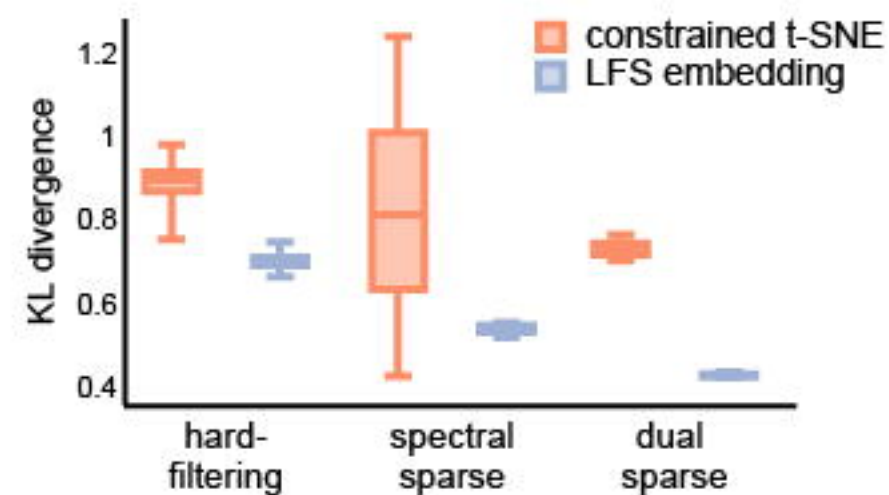
h



× p-EMT
 ● CAF-1
 ● CAF-2
 ● Fibro-others
 ● Endothelial
 ◆ immune cells
 ◆ myocyte
 ◆ malignant





a**b****c****d**