

1 **Unsupervised Spatial Embedded Deep Representation of Spatial**

2 **Transcriptomics**

3

4 Huazhu Fu^{1,*}, Hang Xu^{2,*}, Kelvin Chong², Mengwei Li², Kok Siong Ang², Hong Kai Lee²,
5 Jingjing Ling², Ao Chen³, Ling Shao¹, Longqi Liu³, Jinmiao Chen^{2,†}

6

7 ¹ Inception Institute of Artificial Intelligence, Abu Dhabi, United Arab Emirates

8 ² Singapore Immunology Network (SIgN), Agency for Science, Technology and Research
9 (A*STAR), 8A Biomedical Grove, Immunos Building, 138648, Singapore

10 ³ BGI-ShenZhen, Shenzhen 518103, China.

11

12 *These authors contributed equally to this work.

13 †Corresponding author. Email: chen_jinmiao@immunol.a-star.edu.sg

14

15

16 **Key words:** spatial transcriptomics; graph convolutional network; gene expression; deep
17 learning

18

19

20 Abstract

21 In the past decade, single cell technologies have revolutionized our ability to study cellular
22 heterogeneity. Spatial omics represents the next technological wave, granting spatial context
23 to single cell transcriptomes. Integration analysis of transcripts and spatial information will
24 greatly enable us to dissect tissue organization and inter-cellular communications. Here, we
25 present SEDR, an unsupervised spatial embedded deep representation of both transcript and
26 spatial information. The SEDR pipeline uses a deep autoencoder to construct a gene latent
27 representation in a low-dimensional latent space, which is then simultaneously embedded with
28 the corresponding spatial information through a variational graph autoencoder. SEDR was
29 tested on the 10x Genomics Visium spatial transcriptomics and Stereo-seq datasets,
30 demonstrating its ability to create a better data representation that benefits various follow-up
31 analysis tasks. In benchmarking test, SEDR achieved better clustering accuracy than
32 contemporary methods, and in conjunction with trajectory analysis, it correctly retraced
33 retraces the prenatal development of the human dorsolateral prefrontal cortex. We also found
34 the SEDR representation to be eminently feasible for batch integration. Finally, we used SEDR
35 to characterize the intratumoral heterogeneity of human breast cancer. We identified regions
36 with different immune microenvironments, ranging from pro-inflammatory to immune
37 suppressive areas with infiltrated tumor associated macrophages (TAMs). Analysis suggested
38 a cancer cell dissemination trajectory from cells in pre-metastatic state to invasive carcinoma.

39

40

41

42 Introduction

43 Single-cell omics technologies enable measurements at single cell resolution, leading to
44 discoveries of new subpopulations across various tissues, in both healthy and diseased states.
45 However, tissue dissociation into single cells prior to high throughput omics data acquisition
46 leads to cellular spatial information being lost, hindering our dissection of spatial organization
47 and intercellular interactions of individual cells. While computational tools have been
48 developed to predict cell-cell interactions from ligand and receptor expression, they require
49 validation using immunohistochemistry (IHC) or immunofluorescence (IF). Emerging spatial
50 omics technologies overcome these limitations through simultaneous measurements of
51 gene/protein expression and spatial location of cells. Such spatially-resolved transcriptomes
52 of histological tissues enable the reconstruction of tissue architecture and cell-cell
53 interactions^{1,2,3,4,5,6,7,8,9}. This approach has proven its value in many applications including
54 studies on brain disorders^{2,10}, tumour microenvironment^{3,11}, and embryonic development¹².

55 Among currently available spatial transcriptomics approaches, *in situ* capturing-based
56 technologies such as 10x Genomics Visium and Nanostring GeoMX DSP have gained more
57 popularity owing to their accessibility and ability to profile a large number of mRNA targets
58 within each spot. In principle, a histological section from a tissue sample is permeabilized and
59 the released mRNA is captured by either spatially-arrayed oligos on slide surfaces or by pre-
60 hybridized RNA-target barcodes in manually defined regions of interest (ROI). However, both
61 technologies suffer from limitations in the size of mRNA capture area, where the smallest size
62 is typically ~50 μ m, which is larger than the diameter of a single cell. To overcome this limitation,
63 computational methods have been developed to deconvolute the cell mixture of spatial
64 spot^{13,14,15,16,17,18,19,20}. Recently, improvements in mRNA capture methods have led to smaller
65 subcellular capture areas that are ~1-10 μ m in diameter. These high-resolution spatial
66 transcriptomics methods can obtain spatially-resolved transcriptomes with increased spatial
67 fidelity but without any compromise in the number of genes captured. These methods include

68 Slide-seq⁸, DBiT-seq⁹, with the highest resolution (~1 μ m) thus far obtained by Stereo-seq⁵,
69 PIXEL-seq⁶, and Seq-Scope⁷. These submicrometer-resolution methods usually require voxel
70 binning or cell segmentation to produce a gene-by-cell expression matrix for downstream
71 analysis. Recent technologies have also improved on the size of captured area and thus
72 increased cell throughput, necessitating new computational methods that can handle big
73 spatial data.

74 When analyzing spatial transcriptomics data, combining both gene expression and spatial
75 information to learn a discriminative representation for each cell or spot is crucial. However,
76 established workflows, e.g., Seurat²¹, still employs pipelines for single-cell RNA-seq analysis,
77 which primarily focus on gene expression data and ignore the spatial neighborhood structural
78 relationship. Recently, several new methods have been developed for spatial transcriptomics
79 to overcome this limitation. For example, BayesSpace²² creates a model starts from a Markov
80 Random Filed (MRF) priors which hypothesizes that spots belong to same cell type should be
81 closer to one another and updates models with Bayesian approach. Giotto²³ implements a
82 hidden Markov random filed (HMRF) model to detect domains with coherent patterns by
83 comparing gene expression between cells and their neighbors. SpaGCN²⁴ uses graph
84 convolutional network (GCN) to integrate gene expression, spatial location and histology in
85 spatial transcriptomics data analysis. But the algorithm that SpaGCN integrates histology with
86 spatial location is oversimplified and more evidence should be provided to support its
87 rationality. stLearn²⁵ develops a Spatial Morphological gene Expression (SME) normalization
88 method to normalize spatial omics data. Standard Louvain clustering algorithm is implemented
89 to do unsupervised clustering on SME normalized data. Then stLearn divides broad clusters
90 into sub-clusters according to spatial information if broad clusters spread into several locations.
91 The strategy of stLearn may not make full usage of spatial omics data, because it integrates
92 morphology and spatial information separately at normalization and clustering steps. In
93 general, all stat-of-the-art methods have limitations for properly integration of spatial and
94 morphological information. Moreover, the downstream analyses often require proper low-

95 dimension representation features of the data, which is either neglected or not optimized by
96 state-of-the-art methods.

97 In this work, we developed an unsupervised spatial embedded deep representation (SEDR)
98 method for learning a low-dimensional latent representation of gene expression embedded
99 with spatial information. Our SEDR method consists of two main components, a deep
100 autoencoder network for learning a gene representation, and a variational graph autoencoder
101 network for embedding the spatial information. These two components are optimized jointly to
102 generate a latent representation for spatial transcriptomics data analysis. We applied SEDR
103 on the 10x Genomics Visium spatial transcriptomics dataset and demonstrated its ability to
104 achieve better representation for various follow-up analysis tasks including clustering,
105 visualization, trajectory inference and batch effects correction.

106 Results

107 Overview of SEDR.

108 SEDR learns a gene representation in a low-dimensional latent space with jointly embedded
109 spatial information. As shown in Figure 1, given spatial transcriptomics data, SEDR first learns
110 a nonlinear mapping from the gene expression space to a low-dimensional feature space using
111 a deep autoencoder network. Simultaneously, a variational graph autoencoder is utilized to
112 aggregate the gene representation with the corresponding spatial neighboring relationships to
113 produce a spatial embedding. Then, the gene representation and spatial embedding are
114 concatenated to form the final latent representation used to reconstruct the gene expression.
115 Thereafter, an unsupervised deep clustering method²⁶ is employed to enhance the
116 compactness of learned latent representation. This iterative deep clustering generates a soft
117 clustering by assigning cluster-specific probabilities to each cell, leveraging the inferences
118 between cluster-specific and cell-specific representation learning. Finally, the learned latent
119 representation can be applied towards various analyzing tasks.

120 **Quantitative assessment of SEDR on human dorsolateral prefrontal cortex (DLPFC)**
121 **dataset.**

122 To perform a quantitative comparison between SEDR and other methods, we downloaded the
123 10x Genomics Visium spatial transcriptomics data and the manually annotated layers for LIBD
124 human dorsolateral prefrontal cortex (DLPFC) data². The LIBD data includes 12 slices from
125 the human DLPFC that spans six cortical layers plus white matter. We chose this dataset
126 because the human DLPFC has clear and established morphological boundaries which can
127 serve as the ground truth. We first applied the Seurat standard pipeline²¹ to process and cluster
128 cells using only expression profiles and set the result as the baseline result to benchmark
129 SEDR, in order to investigate the extent to which spatial information improves cell clustering.
130 Moreover, there are some methods that can integrate spatial information and RNA-seq data,
131 including Giotto²³, stLearn²⁵, SpaGCN²⁴, and BayesSpace²². To compare SEDR with these
132 methods, we also employed them to process the same dataset with the recommended default
133 parameters.

134 In slice 151673 (Figure 2A) with 3,639 spots and 33,538 genes, SEDR and BayesSpace had
135 the best performance in terms of both layer borders and ARI. When comparing the results on
136 all 12 DLPFC samples, SEDR had the highest mean ARI (0.426) (Figure 2A bottom right),
137 though the difference between SEDR and BayesSpace (0.418) was not significant (Mann-
138 Whitney U Test²⁷: p-value=0.78). Given the fact that BayesSpace is optimized for clustering,
139 while the objective function of SEDR is to find the best latent representation, comparable
140 clustering performance of SEDR and BayesSpace might indicate that SEDR makes better use
141 of spatial information and gene expression. Besides clustering, BayesSpace does not produce
142 latent representation, in contrast, SEDR derived embedding can be used for not only clustering
143 but also various down-stream analysis tasks such as UMAP visualization, trajectory inference
144 and batch effect correction, and thus provides more flexibility and utilities. Similar to SEDR,
145 SpaGCN also uses GCN to process spatial transcriptomics data. Moreover, it incorporates
146 histology information which is not included in SEDR. However, the clustering performance of

147 SEDR is better than SpaGCN (Mann-Whitney U Test p -value < 0.05). stLearn also integrates
148 histology data, but the performance is likewise poorer. This may indicate that the current
149 approaches utilized by SpaGCN and stLearn to incorporate histological data is not optimal. To
150 make full usage of histology information, we may need to treat it as a separate data modality
151 and use dedicated multi-view algorithms for integration.

152 SEDR generates a set of low dimensional representation features which can be used in
153 various down-stream analyses, such as trajectory inference²⁸. Here, we used Monocle3²⁹ to
154 perform trajectory inference on sample 151673 with the Seurat output (RNA-only) and the
155 SEDR low dimensional representation features. We found that SEDR showed significantly
156 improved performance over Seurat (Figure 2B). In the UMAP plot of SEDR's output, cells
157 belonging to different layers were well-organized, and when we selected white matter (WM)
158 as the root, the pseudo-time reflected the correct "inside-out" developmental ordering of
159 cortical layers (Figure 2B). This demonstrated that compared to RNA-only analyses,
160 incorporating spatial information enabled SEDR to generate a better latent representation that
161 summarized the spatial transcriptomics data. We further confirmed our observations with
162 another trajectory inference method named PArtition-based Graph Abstraction (PAGA)³⁰
163 using the SEDR-derived latent space embedding instead of UMAP coordinates (Figure 2C).
164 The PAGA results showed that the adjacent cortical layers tend to share greater similarity,
165 suggesting spatial adjacency is linked with transcriptomic and even functional similarity.
166 Notably, our trajectory is concordant with the chronological order of cortex development^{31,32,33}.
167 We then compared PAGA graphs generated using Seurat-derived principal components and
168 SEDR embedding. For each of 12 DLPFC slices, we calculated the ratio of weights of edges
169 between adjacent cortical layers over the total sum of weights of all edges. We found
170 significantly higher ratio by SEDR compared to Seurat (Mann-Whitney U test p -value < 0.05)
171 (Figure 2C right).

172 **SEDR corrects for batch effects.**

173 The proliferation of spatial omics application is generating ever increasing volumes of spatially

174 resolved omics data across different labs. However, differences in protocols and technologies
175 complicate comparisons and data integration to produce consensus spatially resolved atlases
176 of tissues. As with scRNA-Seq, removing batch effects in spatial omics dataset is a significant
177 challenge. To date, there are no methods available for batch effects correction of spatial omics.
178 Here, we demonstrate that SEDR is able to learn a joint embedding across multiple batches
179 and project them into a shared latent space. Furthermore, it employs a DEC loss function that
180 enables SEDR to retain biological variations while reducing technical variations. We evaluated
181 the batch correcting performance of SEDR on the DLPFC datasets. We first assessed the
182 batch variations among the 12 datasets and selected 3 sets (151507, 151672, 151673) which
183 exhibit substantial batch effects. The common cortical layers from different batches were
184 separated as shown in the UMAP plot (Figure 3A). We first applied Harmony to remove batch
185 effects based on its superior performance in single-cell RNA-seq data integration³⁴. Harmony
186 was able to mix batches while keeping different layers apart; however, when zoomed into the
187 individual layers, distinct batch-specific sub-clusters were still observable, suggesting that the
188 batch effects were not completely removed (Figure 3B). Next, we applied SEDR on these 3
189 datasets and found that the batch effects were substantially reduced (Figure 3C). Common
190 layers across batches were brought very close and were well-aligned, while different layers
191 were minimally mixed. Further application of Harmony on the SEDR embedding evenly mixed
192 the batches while maintaining separation between layers (Figure 3D). Notably, batch-specific
193 clusters were no longer present within individual layers. Our test showed that by combining
194 SEDR with Harmony, we were able to effectively remove the batch effects present. Among the
195 other spatial omics analysis methods, only stLearn produces a latent space embedding which
196 can be fed to Harmony for batch correction, therefore we benchmarked SEDR against stLearn.
197 As stLearn is unable to jointly project different batches to a shared latent space due to its
198 requirement of histological images as input, we generated a latent space embedding from each
199 dataset and then concatenated them for Harmony integration. The results showed that batches
200 were not well mixed, and the layers were poorly separated (Figure 3E). In conclusion, SEDR
201 combined with Harmony outperforms both Harmony alone and stLearn with Harmony, and this

202 can serve as an effective method for batch correction of spatial omics data.

203 **Dissecting tumor heterogeneity and immune microenvironment using SEDR.**

204 Intratumor heterogeneity in cancer complicates effective treatment formulations and is
205 associated with poor survival prospects³⁵. Spatial transcriptomics is an effective tool for
206 meeting the challenge to dissect and characterize intratumor heterogeneity and tumor-immune
207 crosstalk. Here, we tested SEDR on the 10x Visium spatial transcriptomic data of human breast
208 cancer, which is known for its high intratumoral and intertumoral differences. To aid the
209 interpretation of SEDR results, we performed manual pathology labeling based on the H&E
210 staining. It should be noted that unlike the cerebral cortex that has clear and established
211 morphological boundaries, tumor tissues are highly heterogeneous and encompass complex
212 tumor microenvironments. Manual labeling solely based on tumor morphology is inadequate
213 to characterize such complexity. Based on pathological features, we manually segmented the
214 histological image into 20 regions and grouped them into 4 main morphotypes: Ductal
215 Carcinoma *in Situ*/Lobular Carcinoma *in Situ* (DCIS/LCIS), healthy tissue (Healthy), Invasive
216 Ductal Carcinoma (IDC), and tumor surrounding regions with low features of malignancy
217 (Tumor edge) (Figure 4A). Visually, all five clustering methods agree with the manual
218 annotation at the macroscopic level. Nevertheless, the SEDR clusters presented a smoother
219 segmentation compared to other methods, while Seurat, stLearn and SpaGCN derived clusters
220 appear fragmented and have irregular boundaries. Notably, SEDR found more sub-clusters
221 within the tumor regions, while other methods were prone to divide the healthy regions into
222 sub-clusters, given that all methods were set to generate the same number of clusters. For
223 instance, within DCIS/LCIS_3, SEDR separated an outer “ring” (cluster 7) from the tumor core
224 (cluster 3), and partitioned IDC_2 into 3 sub-clusters. These SEDR sub-clusters suggested
225 transcriptomic heterogeneity within the seemingly homogeneous tumor regions. In addition to
226 clustering analysis, we also employed Seurat3 ‘anchor’-based integration workflow to perform
227 probabilistic transfer of annotations from a reference scRNA-seq data of human breast
228 cancer³⁶ to the spatial data and output, for each spot, a probabilistic classification for each of

229 the scRNA-seq derived classes (Figure 4B, Supplementary Figure 1). The transferred class
230 probabilities were able to delineate the tumor regions and regions where immune cells or
231 fibroblasts were present, which will aid in further dissecting the tumor micro-environment.

232

233 A number of driving forces have been hypothesized for the metastatic transition of tumor cells
234 from a pre-invasive state to invasive carcinoma, including pro-tumor immune
235 microenvironment and reduced cell-cell interactions within the tumor³⁷. Here, we employed
236 PAGA to infer the inter-relatedness between the manually annotated DCIS/LCIS and IDC
237 regions in an attempt to trace the metastatic transition process. The PAGA graph generated
238 using the SEDR embedding showed that DCIS_LCIS_3 was the only DCIS/LCIS region that
239 was likely to spread to its neighboring invasive tumor region IDC_6 (Figure 4C). DEGs
240 between DCIS_LCIS_3 and all other DCIS_LCIS regions and enriched pathways showed that
241 DCIS_LCIS_3 had more immune infiltrates (Supplementary Figure 2A, 2B, 2C), in particular
242 tumor associated macrophages (TAM) (Figure 4B, Supplementary Figure 2D), while the other
243 DCIS_LCIS regions were mainly comprised of epithelial cells that were actively dividing /
244 cycling (Figure 4B) and had up-regulated glycolytic and metabolic processes (Supplementary
245 Figure 2C). TAM infiltration is known to strongly associate with poor survival in solid tumor
246 patients by promoting tumor angiogenesis and inducing tumor migration, invasion and
247 metastasis^{38,39}. We then performed Monocle3 analysis to infer the pseudo-time of the
248 transition from DCIS_LCIS_3 to IDC_6. As DCIS_LCIS_3 and IDC_6 coincide with SEDR
249 clusters 3, 7, and 11 (Figure 4A), we performed Monocle3 on these three clusters and set
250 cluster 3 as the starting point (Figure 4C bottom). We subsequently identified genes that
251 changed expression along Monocle3 pseudotime and revealed sequential waves of gene
252 regulation along the trajectory (Figure 4D). As SEDR cluster 3 and 7 marked the core and
253 outer ring of DCIS/LCIS_3, we identified genes differentially expressed between these two
254 clusters and enriched pathways to further dissect intratumoral heterogeneity (Figure 4E). In
255 cluster 3, we observed the up-regulation of interferon signaling pathways (IFIT1, IFITM1,

256 IFITM3 and TAP1) and NK or neutrophil activities (FCGR3B and TNFSF10) (Figure 4E,
257 Supplementary Figure 2E). In addition, upregulation of RHOB in this region points towards
258 reduced metastatic potential⁴⁰. Cluster 3 represents a region where cancer growth was limited
259 by pro-inflammatory immune responses. On the other hand, in cluster 7, we observed the
260 presence of TAMs (Figure 4B), memory B cells (IGHG1, IGHG3, IGHG4, IGLC2 and IGLC3)
261 and fibroblasts (COL1A1, COL1A2, COL3A1, COL5A1, COL6A1, COL6A2 and FN1) (Figure
262 4E, Supplementary Figure 2E). Upregulated cathepsin activities (CTSB, CTSD and CTSZ)
263 and complement pathway (C1QA, C1S and C4) indicate pro-tumor activities by the TAMs in
264 this region^{41,42,43}. Moreover, upregulated cathepsin activity and metalloprotease inhibitors
265 (TIMP1 and TIMP3) also suggests disturbance in extracellular matrix integrity. Overall, cluster
266 7 represents a region with an immune-suppressed pro-tumor microenvironment and had high
267 potential of cancer metastasis. In summary, SEDR analysis led to the identification of a
268 potentially invasive DCIS region: DCIS/LCIS_3, where the outer ring cluster 7 had TAM
269 infiltration and cancer associated fibroblasts (CAFs) presence, of which both have been
270 reported to facilitate tumor spread^{44,45}. SEDR also enabled the mapping of a molecular path
271 or trajectory from DCIS to IDC. Taken together, SEDR can help dissect intratumoral
272 heterogeneity and understand the relationships between different tumor compartments.

273 **SEDR can handle spatial transcriptomics data with high resolution.**

274 Currently available spatial omics technologies including 10x Visium Spatial Omics, Nanostring
275 GeoMX DSP, SLIDE-seq⁴, and DBIT-seq⁴⁶, do not have single-cell resolution with each
276 capture spot containing 1 to 10 cells. Meanwhile, new emerging methods such as Stereo-seq⁵,
277 PIXEL-Seq⁶ and Seq-Scope⁷ can achieve submicrometer and thus subcellular resolution. With
278 continued advances of spatial omics technologies, spatial resolution and number of cells
279 detected per tissue will significantly improve, producing big datasets with high throughput.
280 Here, we evaluated SEDR's performance on one type of such data (Stereo-seq) of mouse
281 olfactory bulb (Figure 5). Coronal section of the mouse olfactory bulb shows the olfactory nerve
282 layer (ONL), glomerular layer (GL), external plexiform layer (EPL), mitral cell layer (MCL),

283 internal plexiform layer (IPL), granule cell layer (GCL) and rostral migratory stream (RMS)
284 (Figure 5A). We performed unsupervised clustering using Seurat-derived principal
285 components and SEDR-derived embedding to computationally reconstruct the spatial identity
286 of the olfactory bulb profiled with Stereo-seq. Compared to Seurat clusters, SEDR clusters
287 better reflected tissue organization and were more consistent with known anatomical layers
288 (Figure 5B, 5C). We also performed quantitative assessment using local inverse Simpson's
289 index (LISI) and found SEDR produced significantly lower LISI than Seurat showing SEDR
290 clusters were better spatially separated.

291

292 **Discussion**

293 Cell type heterogeneity is a feature of tissue, both healthy and diseased. Capturing this
294 heterogeneity, coupled with their spatial arrangement in the tissue, is crucial when studying
295 the roles of these cells and their cross-talk. Spatial omics technologies represent the state-of-
296 the-art approach to capture omics data with corresponding spatial information from tissue
297 samples. We present SEDR, which leverages on cutting edge machine learning techniques
298 to achieve a better representation of spatial omics data that can be used for clustering and
299 further downstream analyses. SEDR first learns a low dimension latent space representation
300 of the transcriptome information with a deep autoencoder network, which is then aggregated
301 with spatial neighbor information by a variational graph autoencoder to create a spatial
302 embedding. This spatial embedding is then concatenated with the gene expression to be
303 decoded to reconstruct the final gene expression for further analyses. We first demonstrated
304 its efficacy in delineating the different cerebral cortex layers with higher clarity than competing
305 methods, and recapitulated the associated development order by using the joint latent
306 representation with Monocle3.

307 To enhance the analytical power and resolution of spatial omics, we need to integrate multiple
308 datasets from the same tissue. Similar to single-cell transcriptomic data, spatial omics datasets
309 generated in different batches also contain batch-specific systematic variations that present a

310 challenge to batch-effect removal and data integration. In our study, we demonstrated that by
311 combining SEDR and Harmony, we were able to effectively remove batch effects present. In
312 the future, we will integrate Harmony into the SEDR workflow.

313 Spatial omics technologies such as Stereo-seq are able to measure large number of cells per
314 experiment through high spatial resolution and large tissue sizes. In the near future, we expect
315 to see ever increasing throughput from spatial omics experiments, which will result in spatial
316 omics big data that pose significant challenges to data analysis and integration. Computational
317 methods that employ graph neural network require loading the entire graph into GPU memory,
318 which inhibits their applications to very large datasets. We will improve the memory efficiency
319 of SEDR by using GCN min-batch or parallel techniques to construct large-scale graphs for
320 spatial omics data of high throughput and high resolution. Furthermore, technologies with
321 capture spot size smaller than the size of a cell will also require new computational methods
322 that can accurately delineate cells based on capture spots. In the future, we will integrate cell
323 segmentation based on H&E or DAPI staining into SEDR workflow.

324 The current SEDR methodology focuses on gene expression and spatial information, and does
325 not make use of histological images. Contemporary methods such as SpaGCN and stLearn
326 use histological images as input, but in a suboptimal fashion, as demonstrated in our study.
327 SpaGCN utilizes histological image pixels as features by calculating the mean color value from
328 the RGB channels directly. However, the pixel values are easily affected by noise and cannot
329 provide a semantic feature for cell analysis. A more effective approach could be adopting deep
330 CNN model which can learn a high-level representation for histological image. stLearn
331 introduces a deep learning model to extract image features of the spots, and integrates them
332 with the spatial location and gene expression. However, stLearn employs a pre-model trained
333 based on natural images, and does not fine-tune the network towards histological images. In
334 the future, we will incorporate histological images as an additional modality into the SEDR
335 model. We will utilize an image autoencoder network to learn image features, and jointly learn
336 the latent representation by integrating gene expression, image morphology, and spatial

337 information.

338 In summary, SEDR is a promising new approach that builds an integrated representation of
339 cells using both transcriptomic data and spatial coordinates. SEDR derived low dimensional
340 embedding enables more accurate clustering, trajectory inference and batch effect correction.
341 It is able to handle both spatial transcriptomics with capture spot sizes ranging from 50um to
342 less than 1um. Application of SEDR on human breast cancer revealed heterogeneous sub-
343 regions within seemly homogenous tumor regions and shed light on the role of immune
344 microenvironment on tumor invasiveness.

345 **Methods**

346 **Dataset preprocessing.**

347 Our SEDR method takes spatial transcriptomics gene expression and spatial coordinates as
348 inputs. The raw gene expression counts are first normalized using the respective library sizes
349 (by `normalize_total` in Scanpy (v.1.5.0)), with very highly expressed genes excluded from the
350 computation of the normalization factor (size factor) for each cell⁴⁷. Principal component
351 analysis (PCA) is then performed to extract the first 300 principal components to generate the
352 initial gene expression matrix.

353 **Graph construction for spatial transcriptomics data.**

354 To create the graph representing the cell–cell spatial relationships of spatial transcriptomics
355 data, we calculated the Euclidean distance in the image coordinates of all cells and selected
356 the top 10 nearest neighbors of each cell to construction the adjacency matrix. The adjacency
357 matrix, denoted by A , is a symmetric matrix, where $A_{ij} = A_{ji} = 1$ if i and j are neighbors, and
358 0 otherwise.

359 **Deep autoencoder for latent representation learning.**

360 The latent representation of the gene expression is learned through a deep autoencoder. The
361 encoder part, consisting of two fully connected stacked layers, generates the low-dimensional
362 representation $Z_f \in \mathbb{R}^{N \times D_f}$ from the input gene expression matrix $X \in \mathbb{R}^{N \times M}$, while the
363 decoder part with one fully connected layer, reconstructs the expression $X' \in \mathbb{R}^{N \times M}$ from the
364 latent representation $Z \in \mathbb{R}^{N \times D}$, which is obtained by concatenating the low-dimensional
365 representation Z_f and spatial embedding $Z_g \in \mathbb{R}^{M \times D_g}$, where N is the number of cell, M is
366 the number of input genes, and D_f, D_g, D are the dimensions of the learned low-dimensional
367 expression representation of encoder, spatial embedding of GCN, and final latent
368 representation of SEDR with $D = D_f + D_g$. The objective function of the deep autoencoder
369 maximizes the similarity between the input gene and reconstructed expressions, as measured

370 by the mean squared error (MSE) loss function, $\sum(X - X')^2$.

371 **Variational graph autoencoder for spatial embedding.**

372 SEDR utilizes a variational graph autoencoder⁴⁸ (VGAE) to embed the spatial information of
 373 neighborhood cells. With the adjacency matrix A and its degree matrix D , the VGAE learns a
 374 graph embedding Z_g with the formal format as: $g: (A, Z_f) \rightarrow Z_g$, where Z_f is the node/gene
 375 representation from the deep autoencoder. The inference part of VGAE is parameterized by a
 376 two-layer graph convolutional network⁴⁹ (GCN):

377
$$g(Z_g|A, Z_f) = \prod g(z_i|A, Z_f), \text{ with } g(z_i|A, Z_f) = \mathcal{N}(z_i|\mu_i, \text{diag}(\sigma_i^2)),$$

378 where $\mu = GCN_\mu(A, Z_f)$ is the matrix of mean vectors, and $\log\sigma = GCN_\sigma(A, Z_f)$. The two-layer
 379 GCN is defined as $GCN(A, Z_f) = \tilde{A} ReLU(\tilde{A}Z_fW_0)W_1$, with weight W_i and symmetrically
 380 normalized adjacency matrix $\tilde{A} = D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$. The spatial embedding Z_g and reconstructed
 381 adjacency matrix A' are generated as:

382
$$A' = \sigma(Z_g \cdot Z_g^T), \text{ with } Z_g = GCN(A, Z_f).$$

383 The objective function of the VGAE is to minimize the cross-entropy (CE) loss between input
 384 adjacency matrix A and reconstructed adjacency matrix A' , and simultaneously,
 385 minimize Kullback-Leibler (KL) divergence between $g(Z_g|A, Z_f)$ and Gaussian prior $p(Z_g) =$
 386 $\prod_i \mathcal{N}(z_i|0, I)$.

387 **Batch effect correction for spatial transcriptomics**

388 The spatial relationship only exists within single spatial omics, and the cells from different
 389 omics have no direct spatial relation. Let the A^k and Z_f^k denotes the adjacency matrix and deep
 390 gene representation of spatial omics k , we could create a block-diagonal adjacency matrix
 391 A^k and concatenate the deep gene representation in the cell dimension, as:

392
$$A = \begin{bmatrix} A^1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & A^K \end{bmatrix}, Z_f = \begin{bmatrix} Z_f^1 \\ \vdots \\ Z_f^K \end{bmatrix},$$

393 where K is the number of spatial omics. Based on this, we could feed different spatial omics
394 (of potentially different size) as multiple graph instances in the form of one block-diagonal
395 adjacency matrix to the SEDR.

396 For removing batch effects and enhancing the compactness of its latent representation, SEDR
397 employs an unsupervised deep embedded clustering (DEC) method²⁶ to iteratively group the
398 cells into different clusters. To initialize the cluster centers, we employ the KMeans of scikit-
399 learn on the learned latent representations. The number of clusters is pre-defined as a hyper-
400 parameter. With the initialization, the DEC improves the clustering using an unsupervised
401 iterative method with two steps. In the first step, a soft assignment q_{ij} between the cluster
402 center μ_j and latent point z_i is calculated by Student's t-distribution, as:

$$403 \quad q_{ij} = \frac{\left(1 + \|z_i - \mu_j\|^2\right)^{-1}}{\sum_{j'} \left(1 + \|z_i - \mu_{j'}\|^2\right)^{-1}}.$$

404 In the second step, we iteratively refine the clusters by learning from their high confidence
405 assignments with the help of an auxiliary target distribution P based on q_{ij} , as:

$$406 \quad p_{ij} = \frac{q_{ij}^2 / \sum_i q_{ij}}{\sum_{j'} (q_{ij'}^2 / \sum_i q_{ij'})}.$$

407 Based on the soft assignment q_{ij} and auxiliary target distribution p_{ij} , an objective function is
408 defined using the KL divergence:

$$409 \quad KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}.$$

410 The SEDR parameters and cluster centers are then simultaneously optimized by using
411 stochastic gradient descent (SGD) with momentum.

412

413 **Seurat.**

414 Raw mRNA counts were preprocessed to remove low quality genes and sctransformed to
415 remove technical artifacts and normalize the data⁵⁰. We then ran Principal Component

416 Analyses (PCA) to extract the top 30 Principal Components (PCs) and use them to calculate
417 the shared nearest neighbors (SNN). Then the Louvain clustering algorithm was used to
418 identify clusters with the SNN networks. We tried clustering at different resolutions to obtain
419 the same number of clusters as the number of ground truth layers.

420 **SpaGCN, stLearn, BayesSpace, Giotto.**

421 We ran these methods with recommended pipelines and default parameters and set each
422 method to generate the same number of clusters as the number of ground truth layers. stLearn-
423 derived low dimensional embedding was used to for downstream UMAP visualization and
424 harmony batch correction.

425 **Evaluation metric for clustering.**

426 For datasets with cell-type labels (e.g., DLPFC), we employed the adjusted rand index (ARI)
427 to compare the performance of different clustering algorithms. The index calculates the
428 similarity between the clustering labels predicted by algorithm and reference cluster labels as:

$$429 \quad ARI = \frac{RI - E[RI]}{\max(RI) - E[RI]},$$

430 where the unadjusted rand index (RI) is defined as: $RI = (a + b)/C_n^2$, with a as the number of
431 pairs correctly labeled in the same sets, b as the number of pairs correctly labeled as not in
432 the same data set, and C_n^2 as the total number of possible pairs. $E[RI]$ is the expected RI of
433 random labeling. A higher ARI score indicates better performance.

434 **Monocle3.**

435 On DLPFC #151673 slice and breast cancer data, we ran Monocle3 using both Seurat and
436 SEDR outputs. For Seurat, we ran the standard pipeline to get UMAP and used UMAP as input
437 for Monocle3. For SEDR, we first extracted SEDR low dimensional embedding and then used
438 uwot package to calculate UMAP. We then ran Monocle3 on both UMAP using recommended
439 parameters and set white matter (WM) as the start point to generate pseudotime. We then
440 used Moran_I test for detecting significant genes that showed correlation with pseudotime.

441 **Leiden clustering, PAGA trajectory, and UMAP for comparison.**

442 The Leiden clustering, partition-based graph abstraction (PAGA), and uniform manifold
443 approximation and projection (UMAP) of the gene representation and spatial
444 embeddings/principal components (PCs) derived from the SEDR and Seurat were performed
445 using Scanpy (v.1.5.0) package. Briefly, all embeddings or first 30 PCs were directly used to
446 compute a neighborhood graph of observations using n_neighbors of 15, UMAP method to
447 compute the connectivities, and Euclidean method to compute the distance. In order to obtain
448 the same amount of unique Leiden clusters obtained using the SEDR, grid-searching on the
449 Leiden clustering resolution between 0.2 and 2.5 with interval of 0.05/0.01 were performed.
450 Subsequently, the PAGA was performed to quantify the connectivity of Leiden clusters. Lastly,
451 the cluster positions suggested by PAGA were used to initialize the UMAP manifold learning
452 for visualization.

453 **Harmony.**

454 Harmony was used to correct batch effect on low dimensional embeddings. For SEDR, we
455 used latent space embeddings as input. For raw data and stLearn, we used the PCA
456 embeddings as input. We treated different samples as different batches and set all other
457 parameters with default value. For each method, the uncorrected embeddings and batch
458 corrected Harmony embeddings were used to do UMAP analysis.

459 **Prediction of cell type composition of 10x Visium spatial spot.**

460 We downloaded a published scRNA-seq dataset of human breast cancer³⁶ as reference, and
461 ran Seurat to find transfer anchors between the reference and our Visium spatial data. Cell
462 types in the reference are then assigned to the spatial spots by label transferring. We removed
463 cell types that have probability equal to 0 for all spots.

464 **Differential Expression Genes (DEGs) and pathway analyses.**

465 We use Seurat to identify DEGs. Genes with adjusted p-value < 0.05 is used as the input for
466 QIAGEN Ingenuity Pathway Analysis (IPA). For IPA result, the pathway with positive or

467 negative z-score are plotted.

468 **Raw data processing of Stereo-seq data.**

469 Fastq files were generated using MGI DNBSEQ-Tx sequencer. Coordinate identity (CID) and
470 unique molecular identifier (UMI) are contained in the forward reads (CID: 1-25bp, UMI: 26-
471 35bp) while the reverse reads consist of the cDNA sequences. CID sequences on the forward
472 reads were first mapped to the designed coordinates of the *in situ* captured chip, allowing 1
473 base mismatch to correct for sequencing and PCR errors. Reads with UMI containing either N
474 bases or more than 2 bases with quality score lower than 10 were filtered out. CID and UMI
475 associated with each read were appended to each read header. Retained reads were then
476 aligned to the reference genome (mm10) using STAR⁵¹ and mapped reads with MAPQ ≥ 10
477 were counted and annotated to their corresponding genes using an *in-house* script (available
478 at <https://github.com/BGIResearch/handleBam>). UMI with the same CID and the same gene
479 locus were collapsed, allowing 1 mismatch to correct for sequencing and PCR errors. Finally,
480 this information was used to generate a CID-containing expression profile matrix.

481 **Local inverse Simpson's index (LISI).**

482 We first used Seurat and SEDR to generate cell clusters for the stereo-seq data, then used R
483 "lisi" package to calculate LISI with coordinates as X and clustering results of Seurat and SEDR
484 as meta data.

485

486 **Data availability.**

487 (1) The LIBD human dorsolateral prefrontal cortex (DLPFC) Data
488 (<http://spatial.libd.org/spatialLIBD/>);

489

490 **Software availability.**

491 SEDR is written by Python using the PyTorch library. An open-source implementation of SEDR

492 is released on <https://github.com/HzFu/SEDR>

493

494

495

496

497 References

- 498 1. Ståhl, P. L. *et al.* Visualization and analysis of gene expression in tissue sections by
499 spatial transcriptomics. *Science* (2016) doi:10.1126/science.aaf2403.
- 500 2. Maynard, K. R. *et al.* Transcriptome-scale spatial gene expression in the human
501 dorsolateral prefrontal cortex. *Nat. Neurosci.* **24**, 425–436 (2021).
- 502 3. Ji, A. L. *et al.* Multimodal Analysis of Composition and Spatial Architecture in Human
503 Squamous Cell Carcinoma. *Cell* (2020) doi:10.1016/j.cell.2020.08.043.
- 504 4. Rodriques, S. G. *et al.* Slide-seq: A scalable technology for measuring genome-wide
505 expression at high spatial resolution. *Science* (80-.). (2019)
506 doi:10.1126/science.aaw1219.
- 507 5. Chen, A. *et al.* Large field of view-spatially resolved transcriptomics at nanoscale
508 resolution Short title: DNA nanoball stereo-sequencing. *bioRxiv* 2021.01.17.427004
509 (2021).
- 510 6. Fu, X. *et al.* Continuous Polony Gels for Tissue Mapping with High Resolution and
511 RNA Capture Efficiency. *bioRxiv* 2021.03.17.435795 (2021).
- 512 7. Cho, C.-S. *et al.* Seq-Scope: Submicrometer-resolution spatial transcriptomics for
513 single cell and subcellular studies. *bioRxiv* (2021).
- 514 8. Stickels, R. R. *et al.* Highly sensitive spatial transcriptomics at near-cellular resolution
515 with Slide-seqV2. *Nat. Biotechnol.* (2021) doi:10.1038/s41587-020-0739-1.
- 516 9. Liu, Y. *et al.* High-Spatial-Resolution Multi-Omics Sequencing via Deterministic
517 Barcoding in Tissue. *Cell* (2020) doi:10.1016/j.cell.2020.10.026.
- 518 10. Lein, E., Borm, L. E. & Linnarsson, S. The promise of spatial transcriptomics for
519 neuroscience in the era of molecular cell typing. *Science* (2017)
520 doi:10.1126/science.aan6827.
- 521 11. Yoosuf, N., Navarro, J. F., Salmén, F., Ståhl, P. L. & Daub, C. O. Identification and
522 transfer of spatial transcriptomics signatures for cancer diagnosis. *Breast Cancer Res.*
523 (2020) doi:10.1186/s13058-019-1242-9.

- 524 12. van den Brink, S. C. *et al.* Single-cell and spatial transcriptomics reveal somitogenesis
525 in gastruloids. *Nature* (2020) doi:10.1038/s41586-020-2024-3.
- 526 13. Dong, R. & Yuan, G. C. SpatialDWLS: accurate deconvolution of spatial
527 transcriptomic data. *Genome Biol.* **22**, 1–10 (2021).
- 528 14. Song, Q. & Su, J. DSTG: deconvoluting spatial transcriptomics data through graph-
529 based artificial intelligence. *Brief. Bioinform.* (2021) doi:10.1093/bib/bbaa414.
- 530 15. Andersson, A. *et al.* Single-cell and spatial transcriptomics enables probabilistic
531 inference of cell type topography. *Commun. Biol.* **3**, 1–8 (2020).
- 532 16. Biancalani, T. *et al.* Deep learning and alignment of spatially-resolved whole
533 transcriptomes of single cells in the mouse brain with Tangram. *bioRxiv* (2020).
- 534 17. Gayoso, A. *et al.* scvi-tools: a library for deep probabilistic analysis of single-cell omics
535 data. *bioRxiv* (2021).
- 536 18. Lopez, R. *et al.* Multi-resolution deconvolution of spatial transcriptomics data reveals
537 continuous patterns of inflammation. *bioRxiv* (2021).
- 538 19. Elosua-Bayes, M., Nieto, P., Mereu, E., Gut, I. & Heyn, H. SPOTlight: seeded NMF
539 regression to deconvolute spatial transcriptomics spots with single-cell
540 transcriptomes. *Nucleic Acids Res.* (2021) doi:10.1093/nar/gkab043.
- 541 20. Danaher, P. *et al.* Advances in mixed cell deconvolution enable quantification of cell
542 types in spatially-resolved gene expression data. *bioRxiv* (2020).
- 543 21. Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Bioarxiv* (2020).
- 544 22. Zhao, E. *et al.* Bayesspace enables the robust characterization of spatial gene
545 expression architecture in tissue sections at increased resolution. *bioRxiv* (2020)
546 doi:10.1101/2020.09.04.283812.
- 547 23. Dries, R. *et al.* Giotto: a toolbox for integrative analysis and visualization of spatial
548 expression data. *Genome Biol.* **22**, 78 (2021).
- 549 24. Hu, J. *et al.* Integrating gene expression, spatial location and histology to identify
550 spatial 1 domains and spatially variable genes by graph convolutional network 2 3.
551 *bioRxiv* 2020.11.30.405118 (2020).

- 552 25. Pham, D. *et al.* stLearn: Integrating spatial location, tissue morphology and gene
553 expression to find cell types, cell-cell interactions and spatial trajectories within
554 undissociated tissues. *bioRxiv* (2020) doi:10.1101/2020.05.31.125658.
- 555 26. Xie, J., Girshick, R. & Farhadi, A. Unsupervised Deep Embedding for Clustering
556 Analysis. in *ICML* (eds. Balcan, M. F. & Weinberger, K. Q.) vol. 48 478–487 (PMLR,
557 2016).
- 558 27. Mann, H. B. & Whitney, D. R. On a Test of Whether one of Two Random Variables is
559 Stochastically Larger than the Other. *Ann. Math. Stat.* **18**, 50–60 (1947).
- 560 28. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by
561 pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32**, 381–386 (2014).
- 562 29. Cao, J. *et al.* The single-cell transcriptional landscape of mammalian organogenesis.
563 *Nature* **566**, 496–502 (2019).
- 564 30. Wolf, F. A. *et al.* PAGA: graph abstraction reconciles clustering with trajectory
565 inference through a topology preserving map of single cells. *Genome Biol.* **20**, 1–9
566 (2019).
- 567 31. Gilmore, E. G. & Herrup, K. Cortical development: Layers of complexity. *Current*
568 *Biology* (1997) doi:10.1016/s0960-9822(06)00108-4.
- 569 32. Chini, M. & Hanganu-Opatz, I. L. Prefrontal Cortex Development in Health and
570 Disease: Lessons from Rodents and Humans. *Trends in Neurosciences* (2021)
571 doi:10.1016/j.tins.2020.10.017.
- 572 33. Nadarajah, B. & Parnavelas, J. G. Modes of neuronal migration in the developing
573 cerebral cortex. *Nature Reviews Neuroscience* (2002) doi:10.1038/nrn845.
- 574 34. Tran, H. T. N. *et al.* A benchmark of batch-effect correction methods for single-cell
575 RNA sequencing data. *Genome Biol.* (2020) doi:10.1186/s13059-019-1850-9.
- 576 35. Nguyen, P. H. D. *et al.* Intratumoural immune heterogeneity as a hallmark of tumour
577 evolution and progression in hepatocellular carcinoma. *Nat. Commun.* **12**, 1–13
578 (2021).
- 579 36. *et al.* A single-cell RNA expression atlas of normal, preneoplastic and tumorigenic

- 580 states in the human breast. *EMBO J.* **40**, 1–23 (2021).
- 581 37. Friedl, P. & Alexander, S. Cancer invasion and the microenvironment: Plasticity and
582 reciprocity. *Cell* (2011) doi:10.1016/j.cell.2011.11.016.
- 583 38. Kuroda, H. *et al.* Tumor microenvironment in triple-negative breast cancer: the
584 correlation of tumor-associated macrophages and tumor-infiltrating lymphocytes. *Clin.*
585 *Transl. Oncol.* (2021) doi:10.1007/s12094-021-02652-3.
- 586 39. Asiry, S. *et al.* The Cancer Cell Dissemination Machinery as an Immunosuppressive
587 Niche: A New Obstacle Towards the Era of Cancer Immunotherapy. *Front. Immunol.*
588 **12**, 1–19 (2021).
- 589 40. Ju, J. A., Godet, I., DiGiacomo, J. W. & Gilkes, D. M. RhoB is regulated by hypoxia
590 and modulates metastasis in breast cancer. *Cancer Rep.* (2020)
591 doi:10.1002/cnr2.1164.
- 592 41. Olson, O. C. & Joyce, J. A. Cysteine cathepsin proteases: Regulators of cancer
593 progression and therapeutic response. *Nature Reviews Cancer* (2015)
594 doi:10.1038/nrc4027.
- 595 42. Roumenina, L. T. *et al.* Tumor cells hijack macrophage-produced complement C1q to
596 promote tumor growth. *Cancer Immunol. Res.* (2019) doi:10.1158/2326-6066.CIR-18-
597 0891.
- 598 43. Fraser, D., Melzer, E., Camacho, A. & Gomez, M. Macrophage production of innate
599 immune protein C1q is associated with M2 polarization (INM1P.434). *J. Immunol.*
600 (2015).
- 601 44. Monteran, L. & Erez, N. The dark side of fibroblasts: Cancer-associated fibroblasts as
602 mediators of immunosuppression in the tumor microenvironment. *Frontiers in*
603 *Immunology* (2019) doi:10.3389/fimmu.2019.01835.
- 604 45. Lin, Y., Xu, J. & Lan, H. Tumor-associated macrophages in tumor metastasis:
605 Biological roles and clinical therapeutic applications. *Journal of Hematology and*
606 *Oncology* (2019) doi:10.1186/s13045-019-0760-3.
- 607 46. Deng, Y. *et al.* Resource High-Spatial-Resolution Multi-Omics Sequencing via

- 608 Deterministic Barcoding in Tissue II II Resource High-Spatial-Resolution Multi-Omics
609 Sequencing via Deterministic Barcoding in Tissue. *Cell* **183**, 1665-1681.e18 (2020).
- 610 47. Weinreb, C., Wolock, S. & Klein, A. M. SPRING: a kinetic interface for visualizing high
611 dimensional single-cell expression data. *Bioinformatics* **34**, 1246–1248 (2018).
- 612 48. Kipf, T. N. & Welling, M. Variational Graph Auto-Encoders. in *NIPS Workshop on*
613 *Bayesian Deep Learning* (2016).
- 614 49. Kipf, T. N. & Welling, M. Semi-Supervised Classification with Graph Convolutional
615 Networks. in *International Conference on Learning Representations (ICLR)* (2017).
- 616 50. Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell
617 RNA-seq data using regularized negative binomial regression. *Genome Biol.* **20**, 296
618 (2019).
- 619 51. Dobin, A. *et al.* STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics* (2013)
620 doi:10.1093/bioinformatics/bts635.
- 621
- 622

623 **Acknowledgements:**

624 This research was supported by funding from Singapore Immunology Network (SIgN), A *
625 STAR, Singapore.

626

627 **Author contributions:**

628 Huazhu Fu designed and implemented SEDR. Hang Xu, Huazhu Fu, Kelvin Chong, Mengwei
629 Li, Hong Kai Lee and Jingjing Ling performed data analysis. Hang Xu, Huazhu Fu, Mengwei
630 Li generated figures. Jinmiao Chen, Huazhu Fu, Hang Xu, Kok Siong Ang, Kelvin Chong,
631 Jingjing Ling and Ling Shao drafted the manuscript. Ao Chen and Longqi Liu provided Stereo-
632 seq data. Jinmiao Chen conceptualized and supervised the study.

633

634 **Competing interests:**

635 The authors declare no competing interests.

636 **Figure legend:**

637 **Figure 1. Overview of SEDR.** SEDR learns a low-dimensional latent representation of gene
638 expression embedded with spatial information via jointly training a deep autoencoder and a
639 variational graph autoencoder. The low dimensional embedding produced by SEDR can be
640 used for downstream visualization, cell clustering, trajectory inference and batch effect
641 correction.

642 **Figure 2. Quantitative assessment of SEDR on human dorsolateral prefrontal cortex**
643 **(DLPFC) dataset.** A) Ground-truth segmentation of cortical layers; clustering results of Seurat,
644 Giotto, stLearn, SpaGCN, BayesSpace and SEDR on DLPFC slice #151673; Adjusted rand
645 index (ARI) of various cluster sets on 12 DLPFC slices. B) UMAP visualization and Monocle
646 trajectory generated using Seurat-derived PCA embedding (top) and SEDR embedding
647 (bottom); Monocle pseudotimes were visualized on UMAP plot and spatial co-ordinates. C)
648 PAGA graph generated using Seurat-derived PCA embedding and SEDR embedding; SEDR
649 showed a higher percentage of weights of correct PAGA edges compared to Seurat.

650 **Figure 3. Batch effect present in DLPFC datasets and assessment of SEDR's**
651 **performance on batch correction.** A) The slices #151507, #151672 and #151673 showed
652 substantial inter-slice variations before batch effect correction. UMAP plots colored by ground-
653 truth cortical layers (left), slices (right), split by slices and colored by layers (bottom). B)
654 Harmony alone was unable to remove the batch effects present. C) SEDR alone substantially
655 reduced batch effects. D) SEDR combined with Harmony effectively corrected for batch
656 effects. E) stLearn combined with Harmony was unable to correct for batch effects.

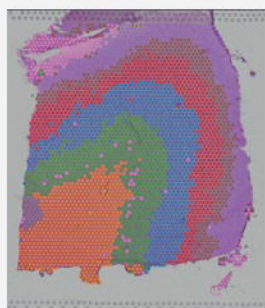
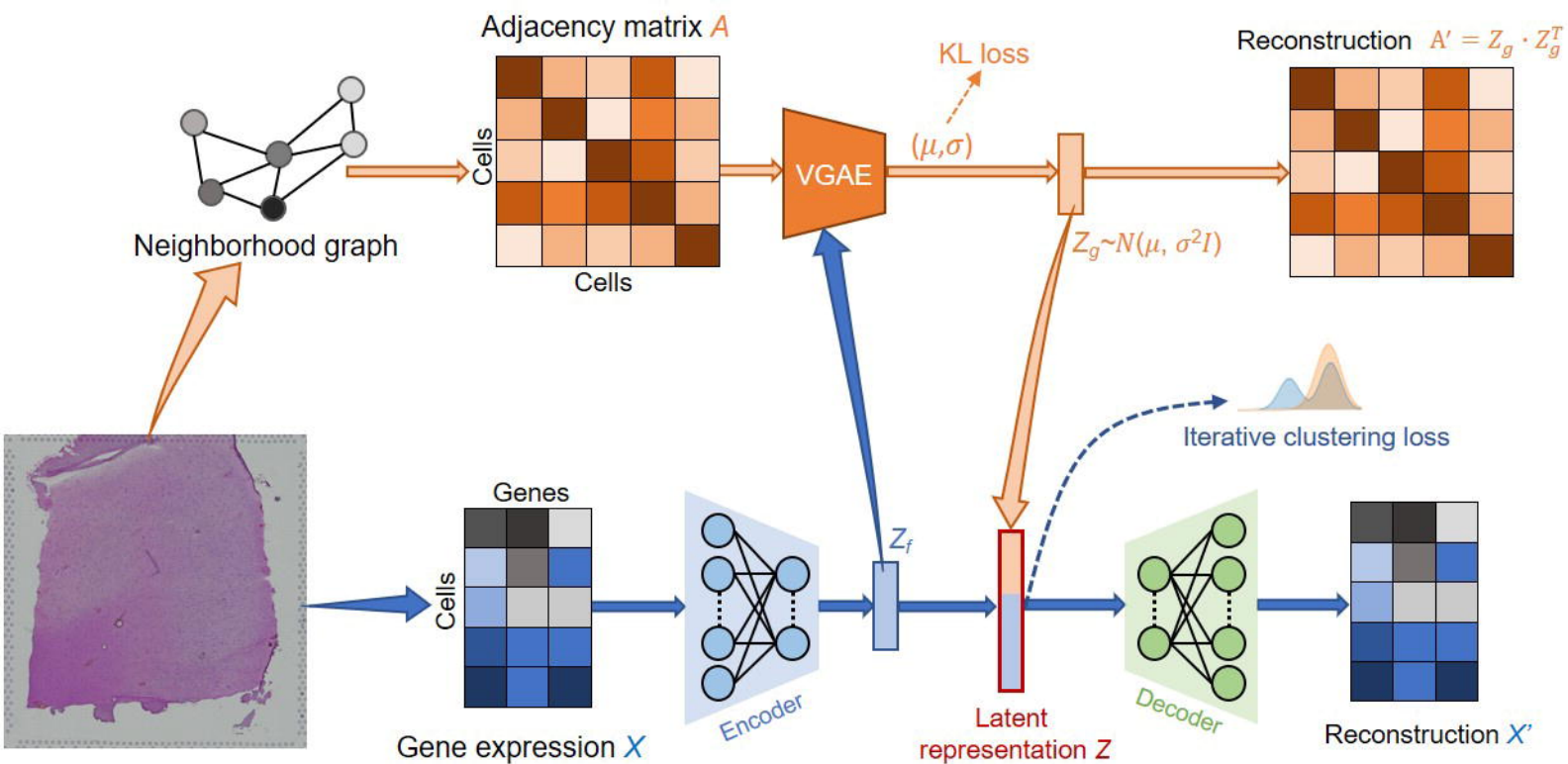
657 **Figure 4. Application of SEDR on 10x Visium spatial transcriptomics data of human**
658 **breast cancer.** A) Manual pathology labeling based on the H&E staining; clustering results of
659 SEDR, Seurat, stLearn, SpaGCN and BayesSpace. B) Seurat3 'anchor'-based integration
660 workflow was used to perform probabilistic transfer of annotations from a reference scRNA-
661 seq data of human breast cancer to the spatial data and output, for each spot, a probabilistic
662 classification for each of the scRNA-seq derived classes. The probabilities of tumor associated
663 macrophage (TAM) and cycling epithelial (C.Epi) were visualized. C) Trajectory analysis
664 results using PAGA (Top) and Monocle3 (Bottom). PAGA graph predicted the inter-relatedness
665 between the manually annotated DCIS/LCIS and IDC regions. Edge width, a measure of
666 connectivity strength, indicates the likelihood of an actual connection being present. Monocle3
667 inferred pseudotimes of spots in SEDR cluster 3, 7 and 11 using Seurat-derived PCA
668 embedding (termed "rna_pseudotime") and SEDR embedding (termed "SEDR_pseudotime").
669 D) Heatmap of genes whose expression changed along Monocle-derived pseudotime. E)
670 Pathways enriched by genes differentially expressed between SEDR cluster 3 and 7. Red bars
671 represent pathways up-regulated in cluster 3.

672 **Figure 5. Application of SEDR on Stereo-seq spatial transcriptomics data of mouse**
673 **olfactory bulb tissue section.** A) Laminar organization of DAPI stained mouse olfactory bulb.
674 B) Unsupervised clustering of the spatial voxels analyzed by Seurat and SEDR. C) Four
675 clusters with the highest number of voxels were selected and visualized. D) Quantitative
676 comparison of Seurat and SEDR clusters using local inverse Simpson's index (LISI).

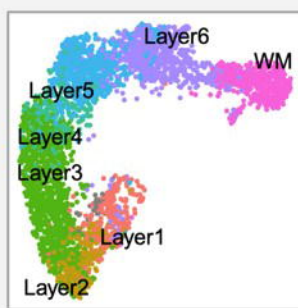
677 **Supplementary:**

678 Figure 1. **Complete deconvolution result for breast cancer sample.**

679 Figure 2. **Differential expression genes and enriched pathways.** A) Position of
680 DCIS_LCIS_3 and other DCIS_LCIS regions. B) Top DEGs between DCIS_LCIS_3 and other
681 DCIS_LCIS regions. C) Enriched pathways of DEGS for DCIS_LCIS_3 vs other DCIS_LCIS
682 regions. D) Percentage of TAM for cluster 3 and cluster 7 of SEDR clustering result. E)
683 Representative DEGs between cluster 3 and cluster 7.



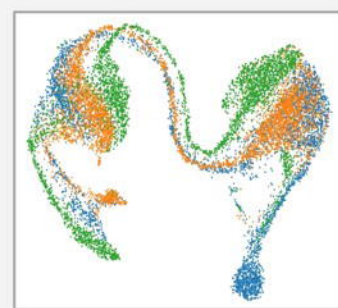
Clustering



Visualization



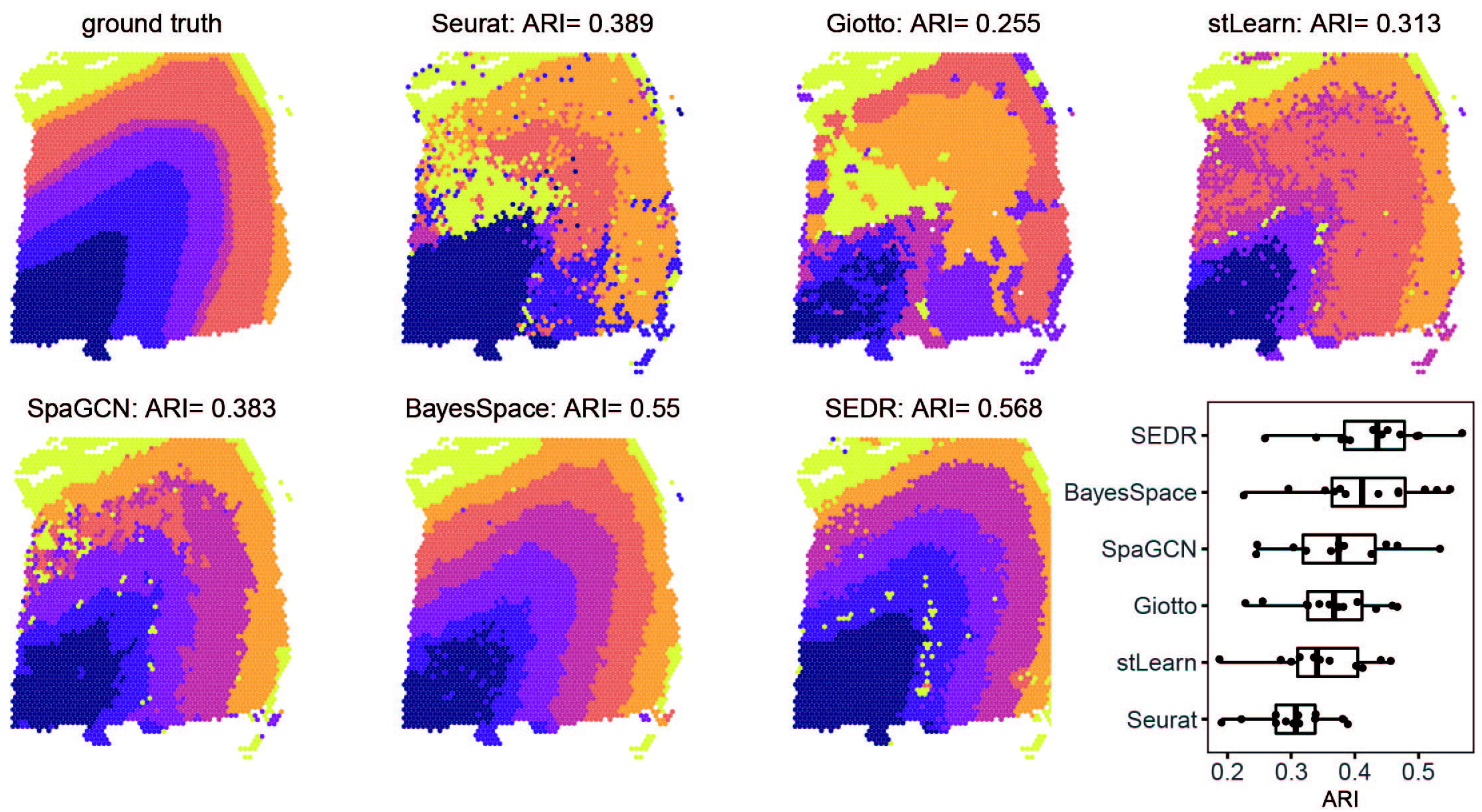
Trajectory inference



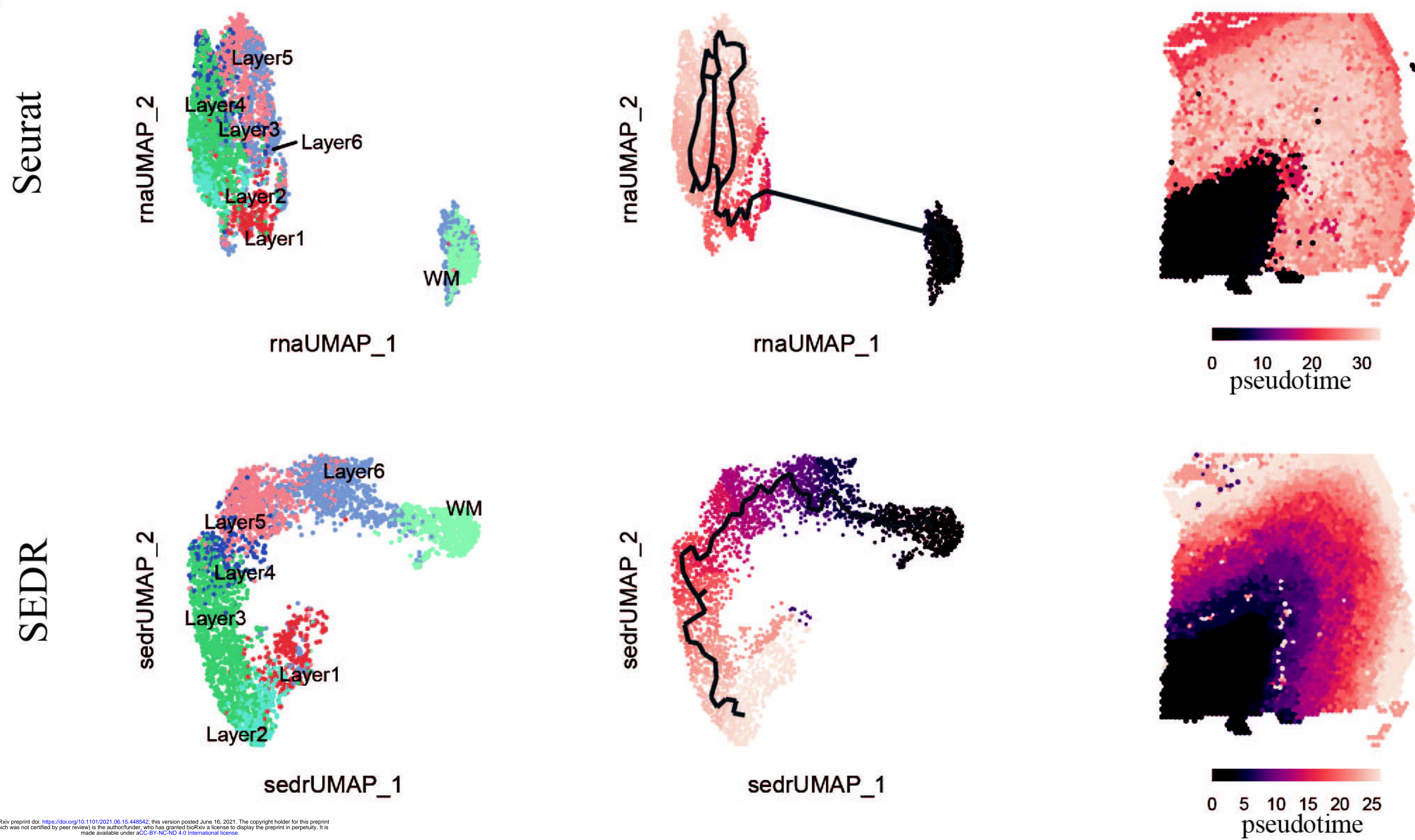
Batch correction

.....

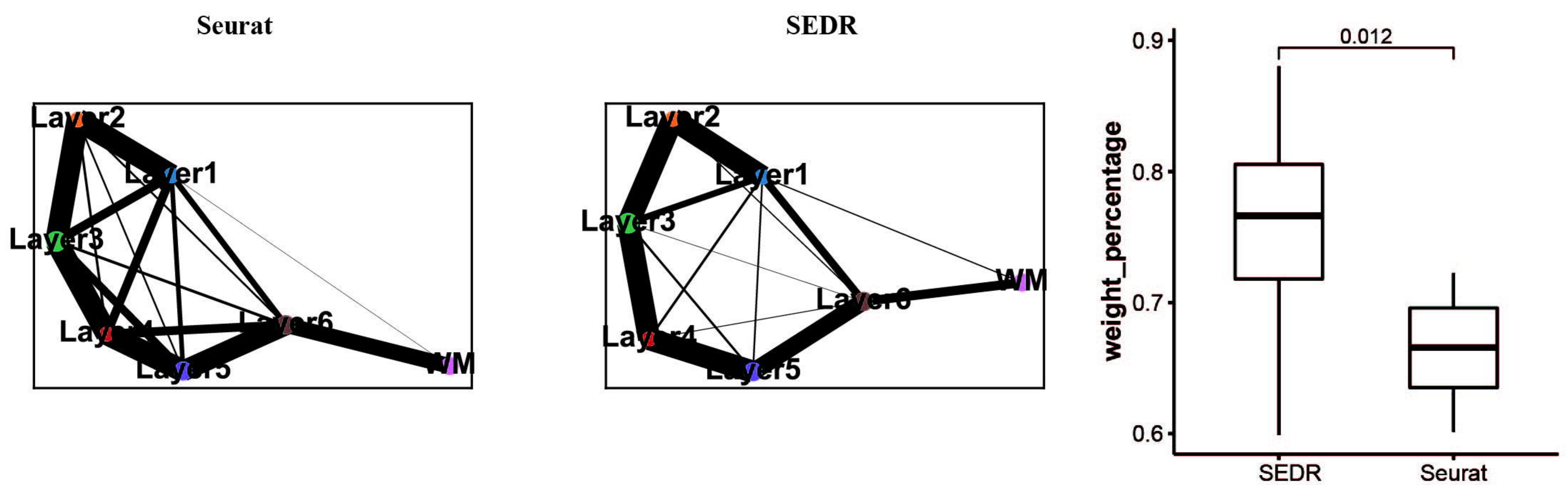
A

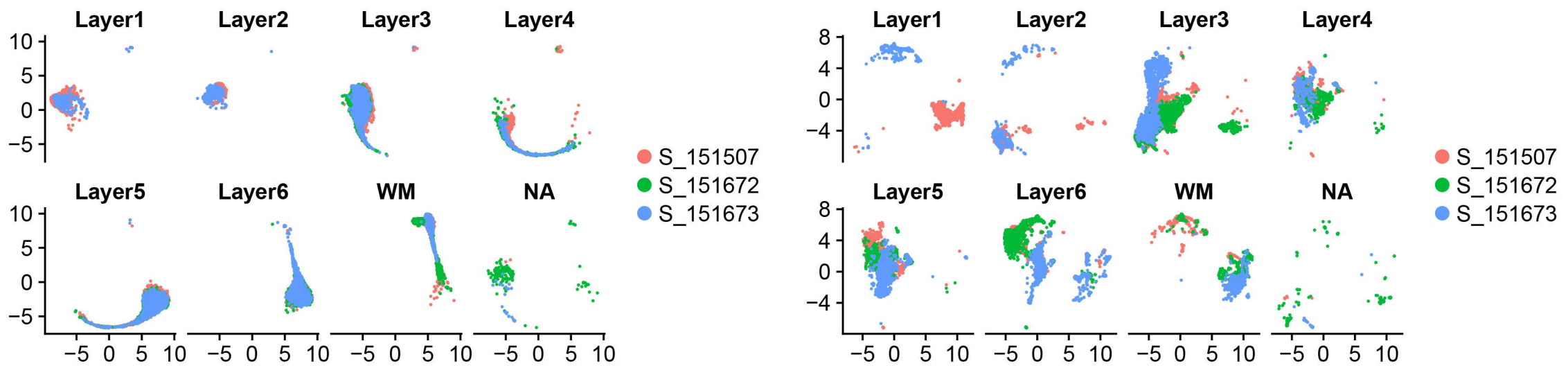
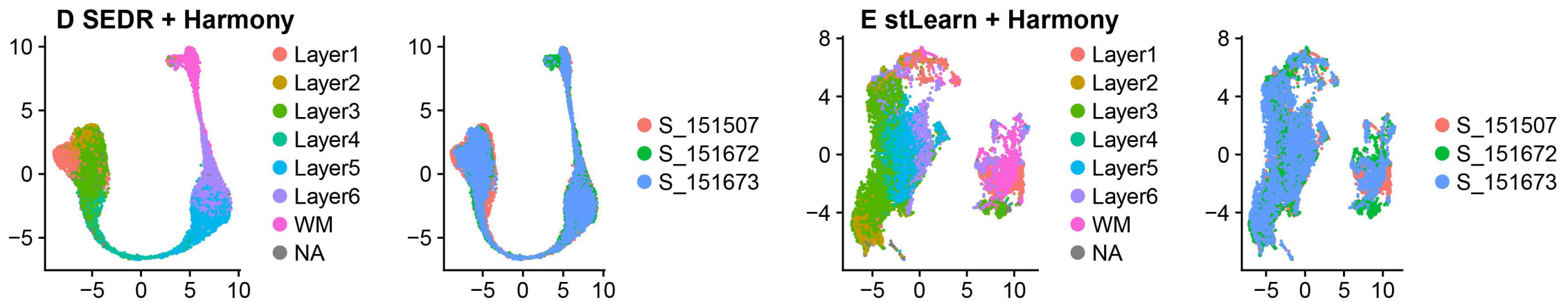
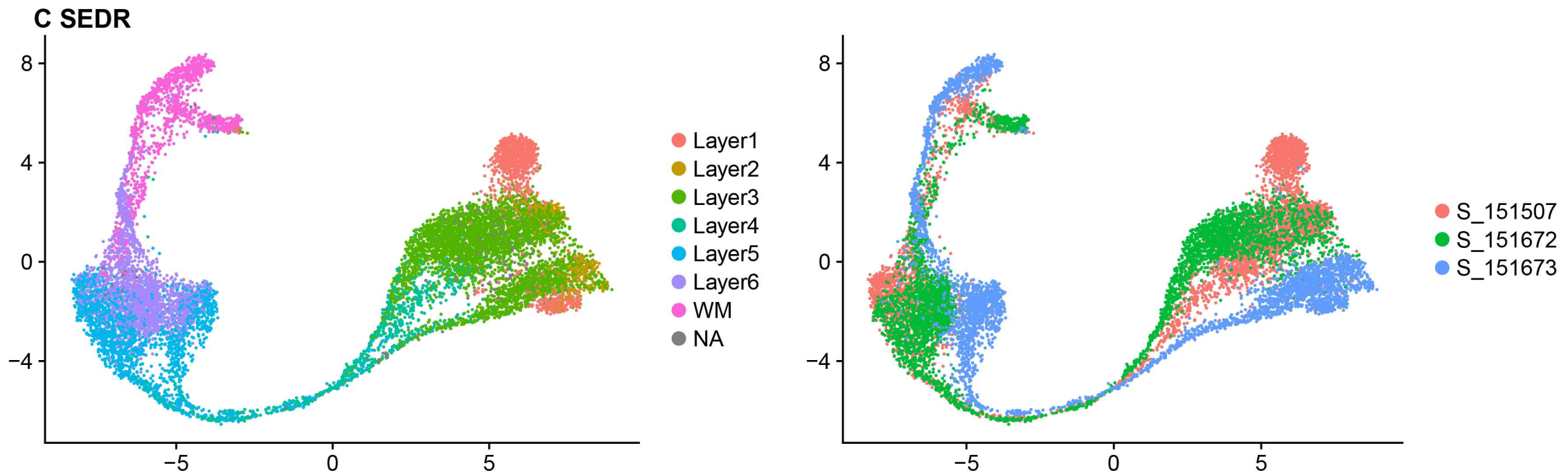
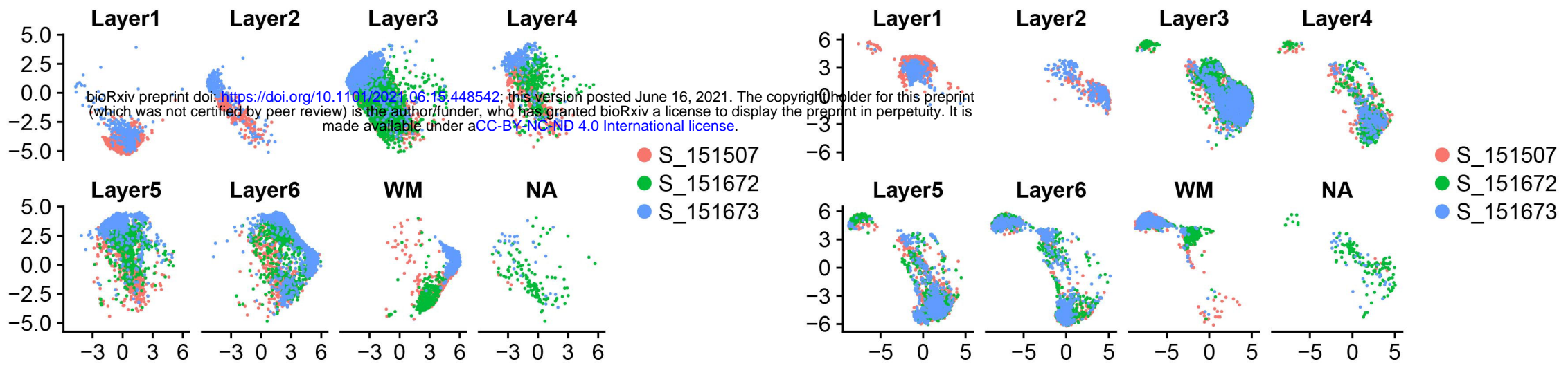
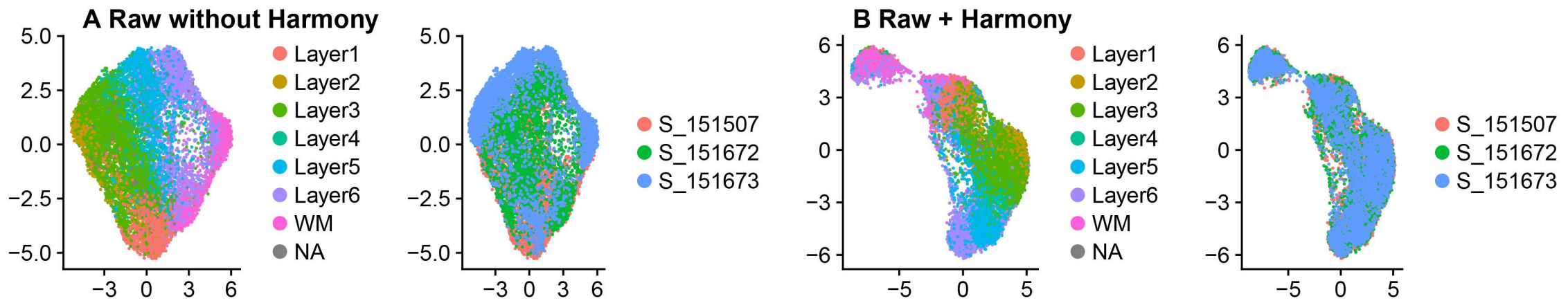


B

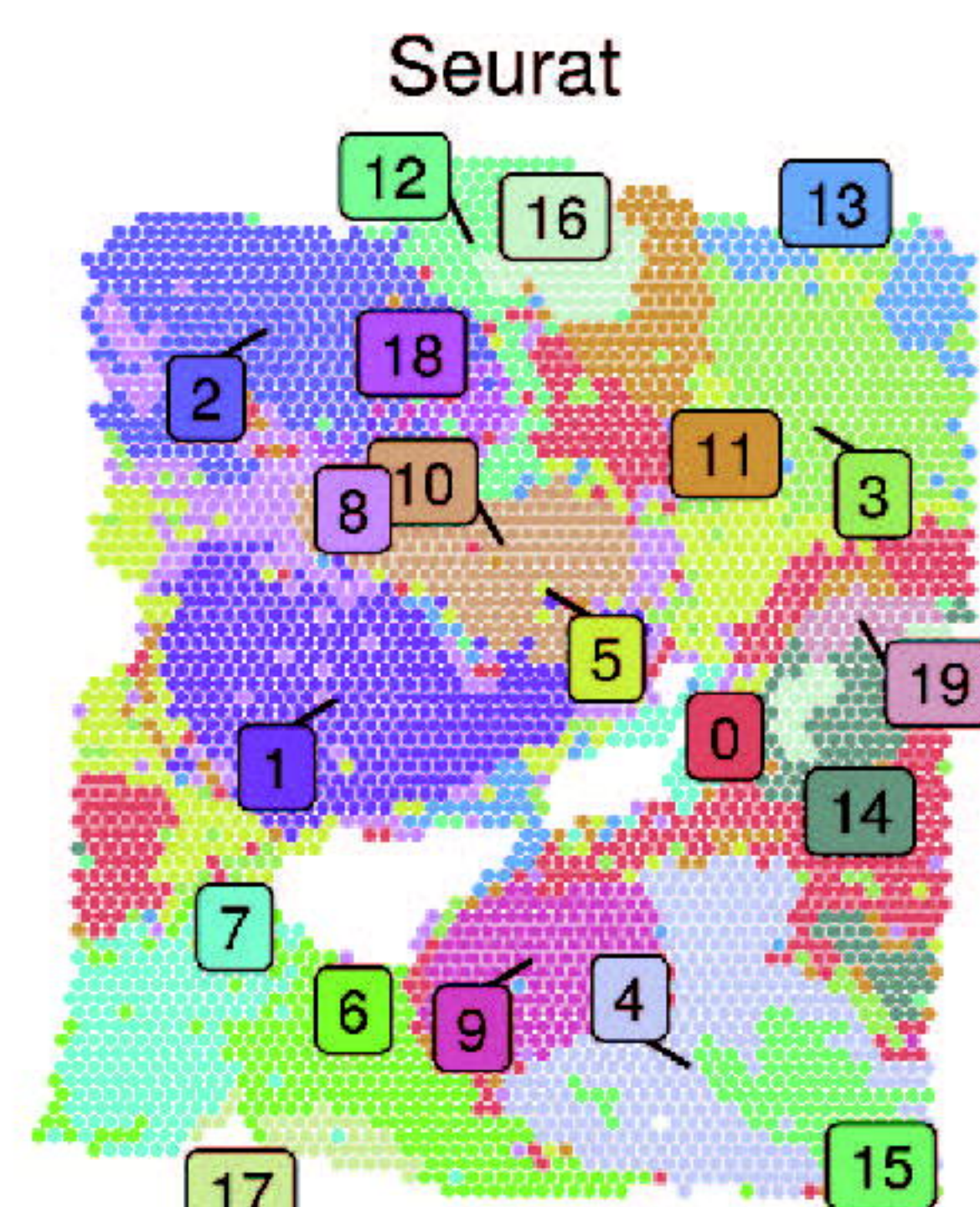
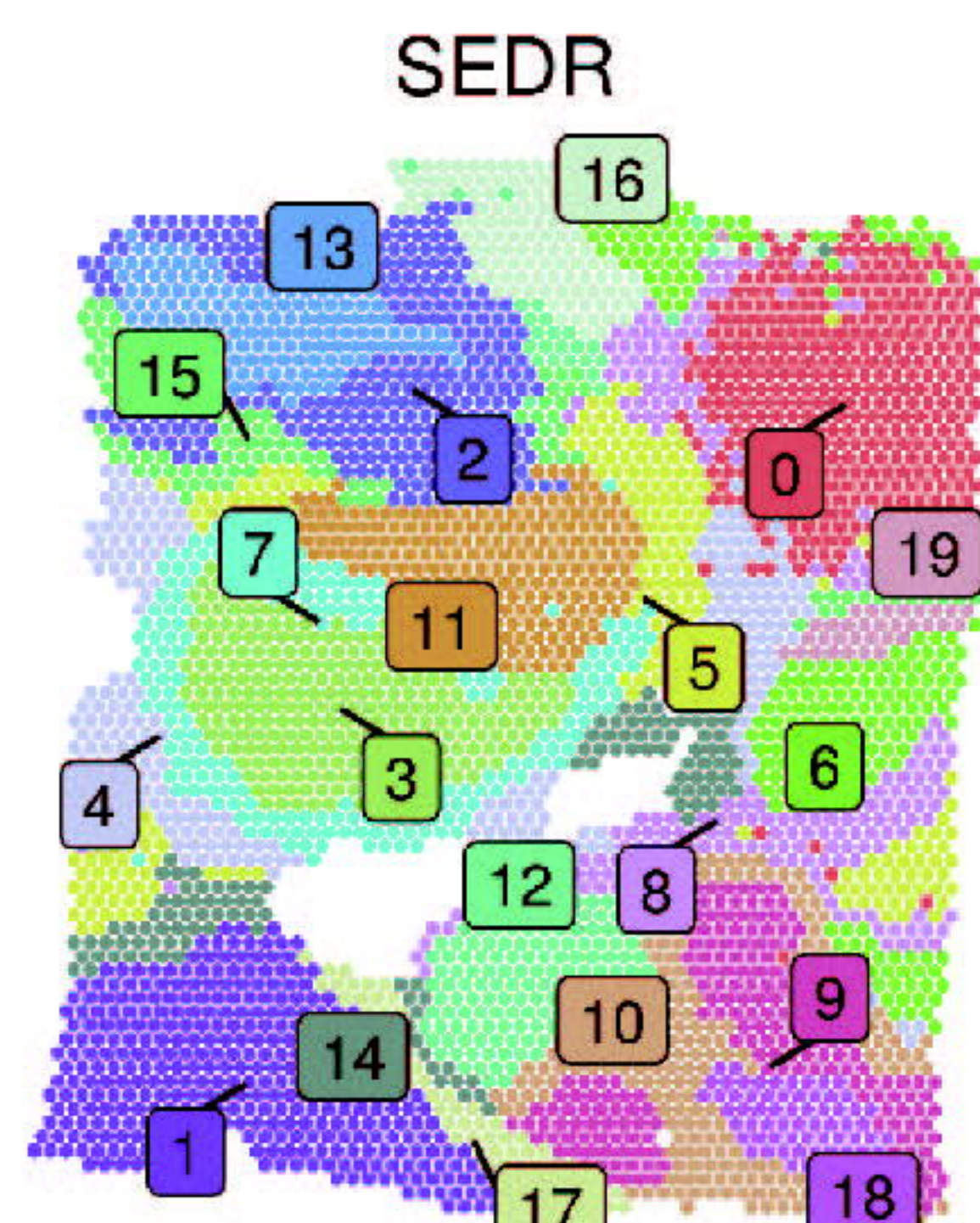
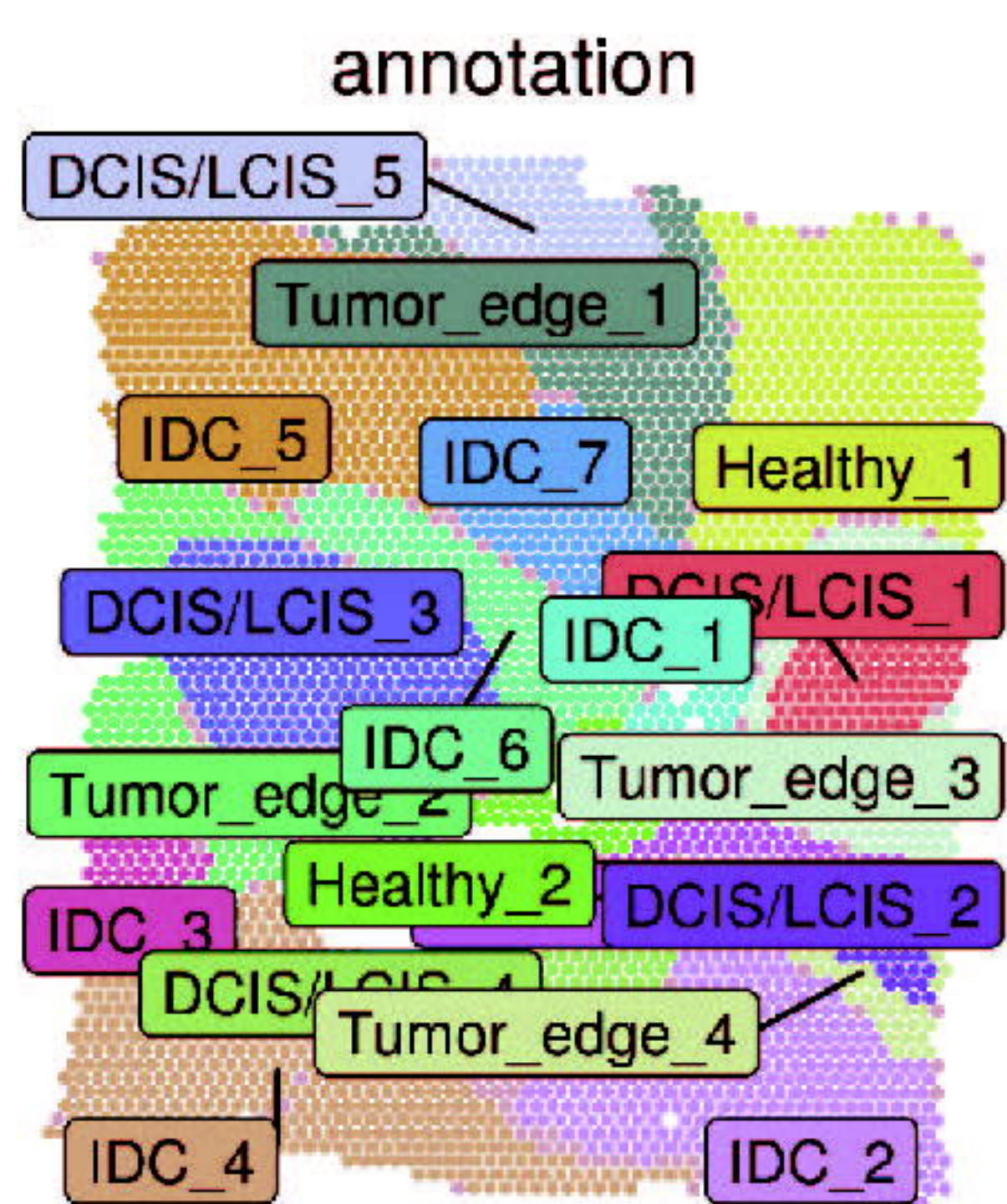


C



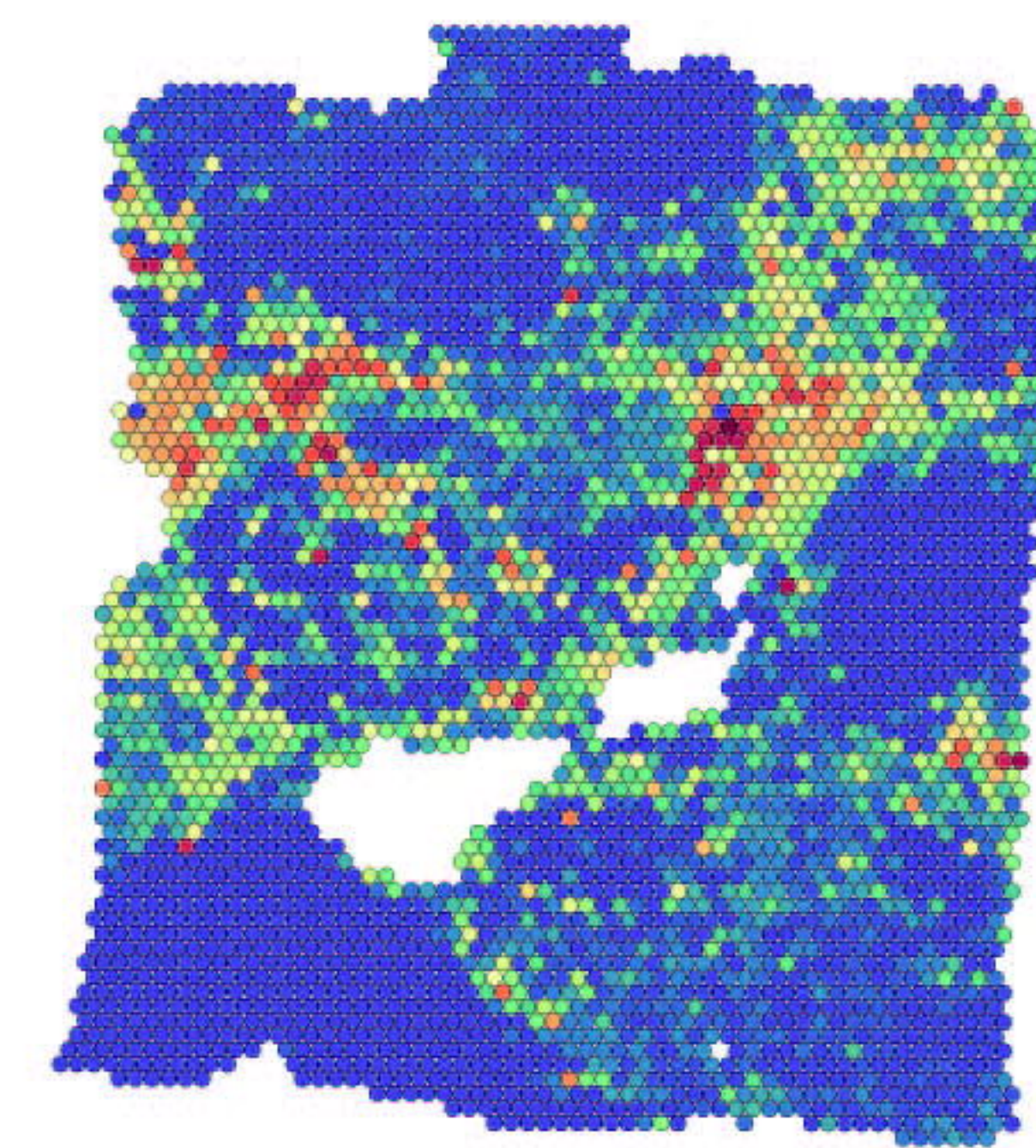


A

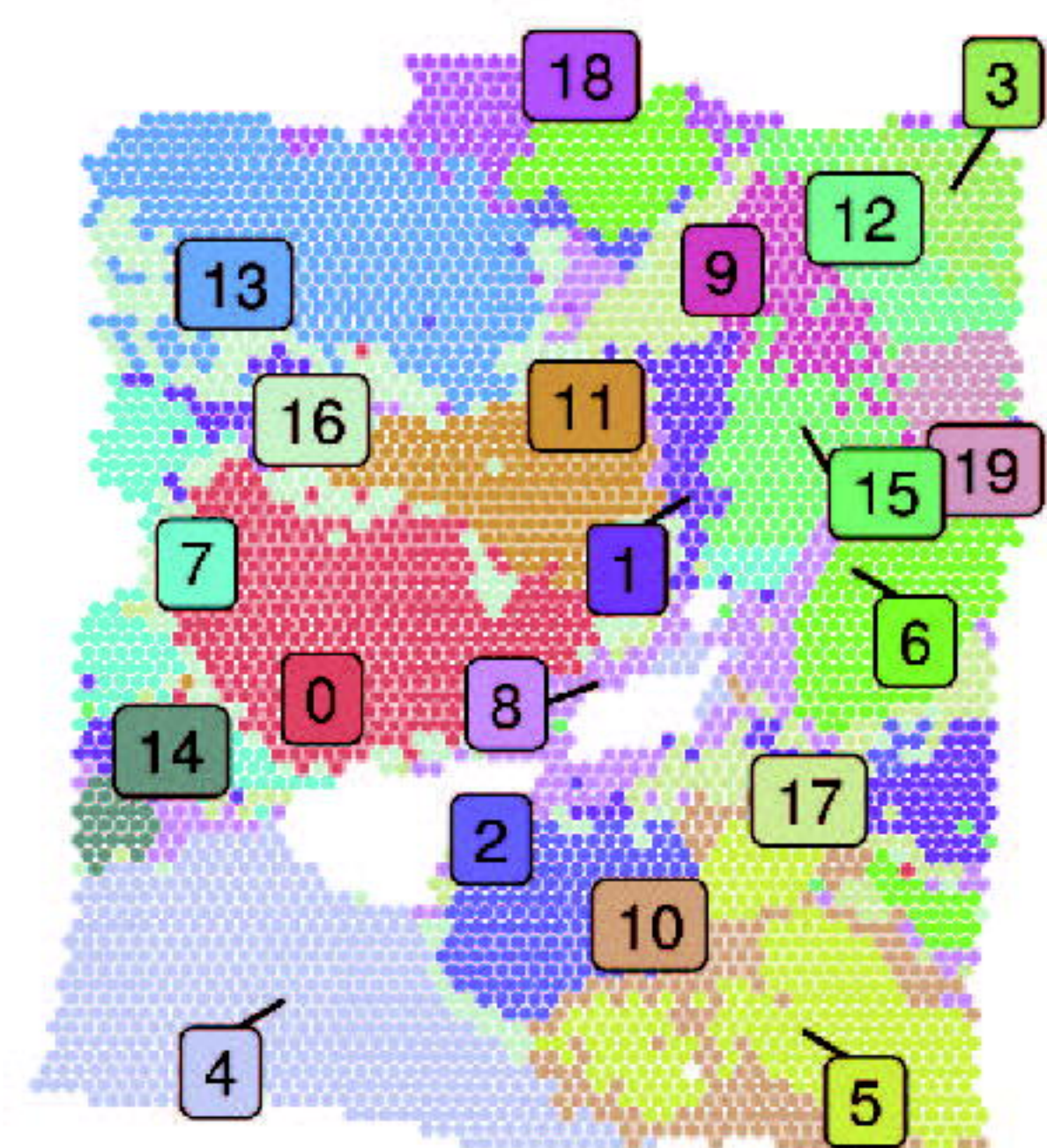


B

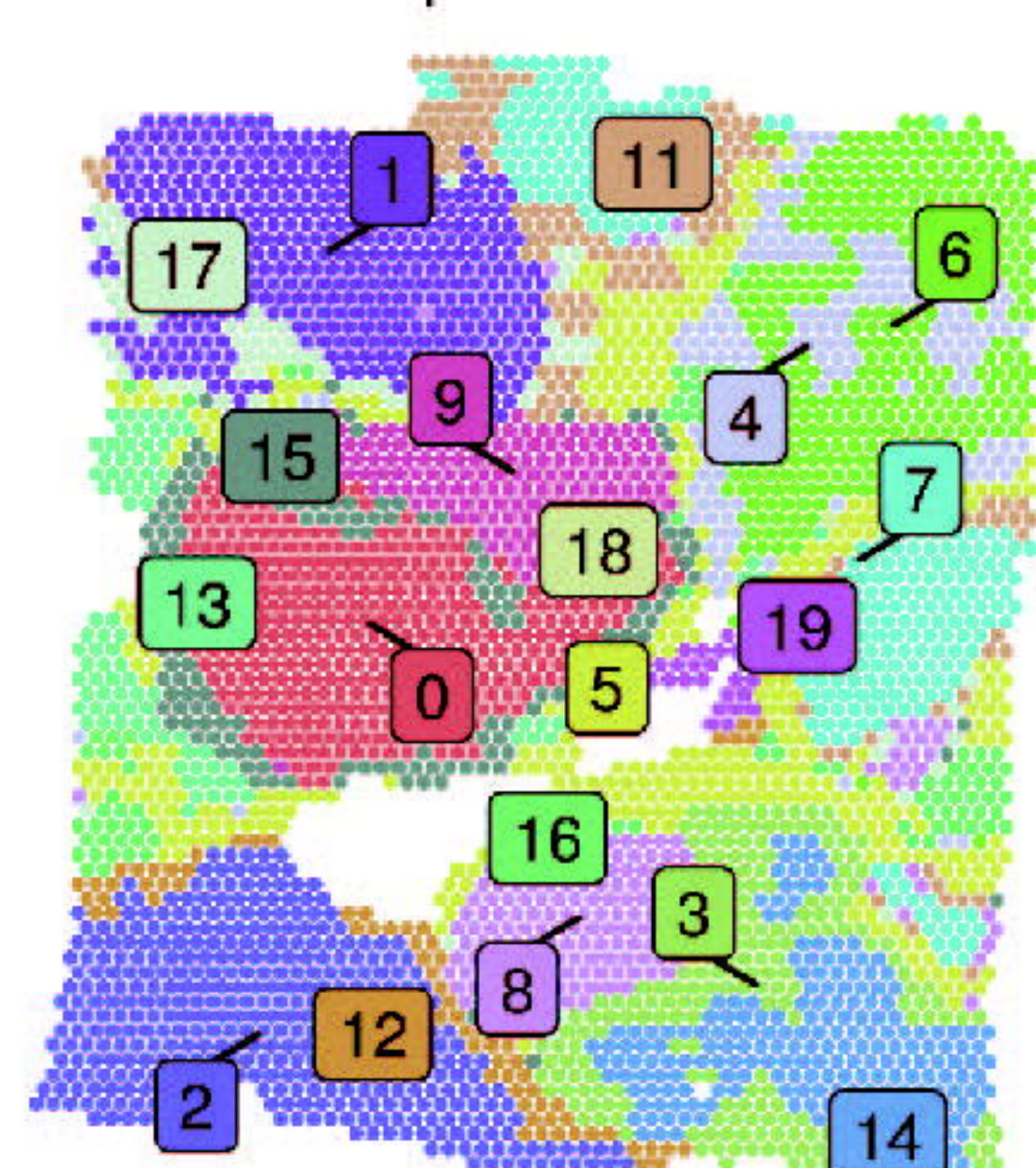
TAM



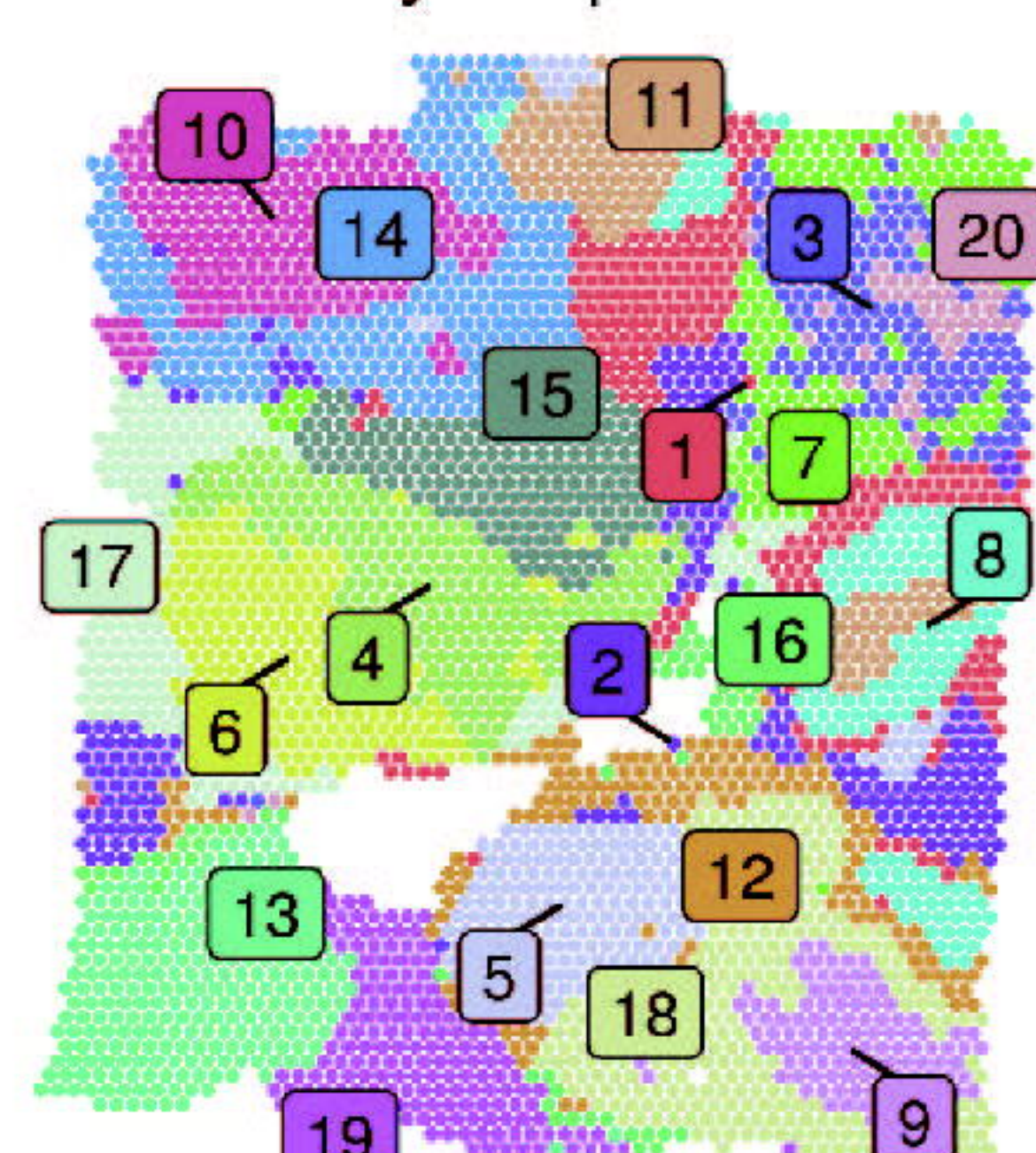
stLearn



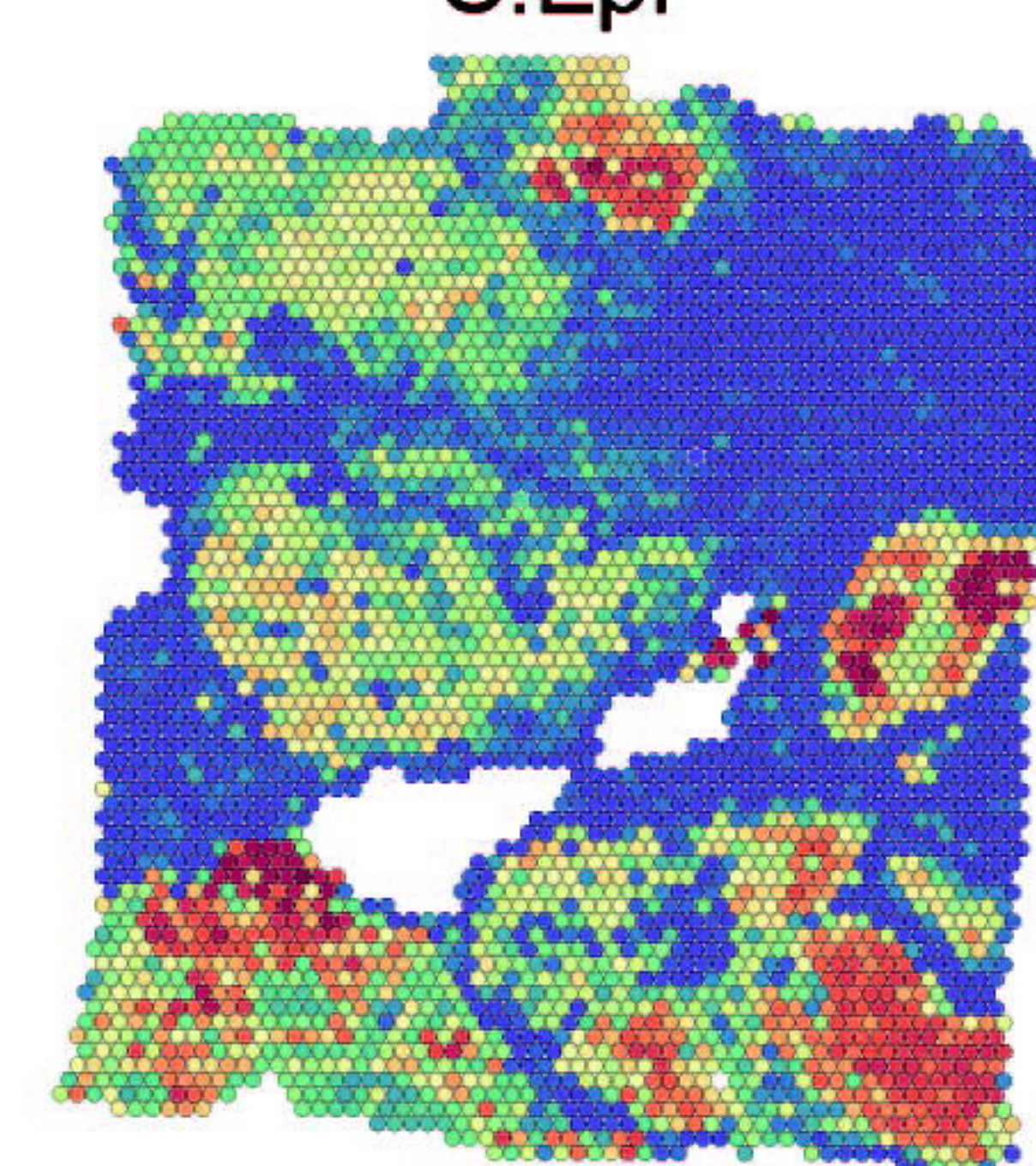
SpaGCN



BayesSpace

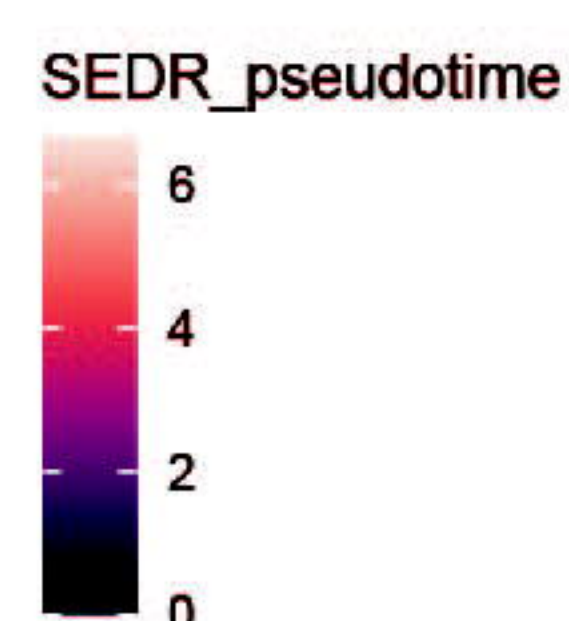
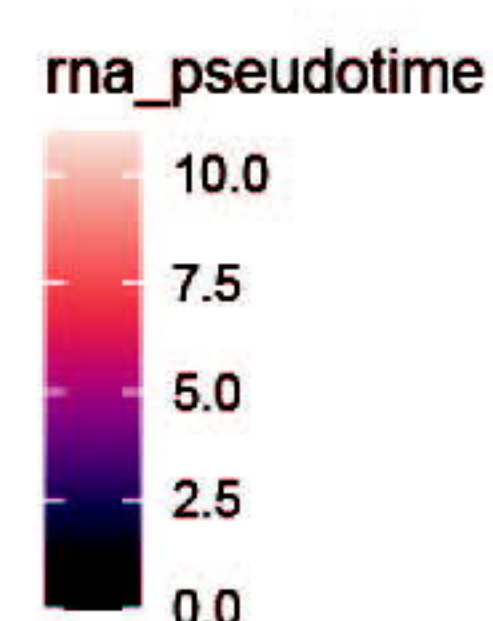
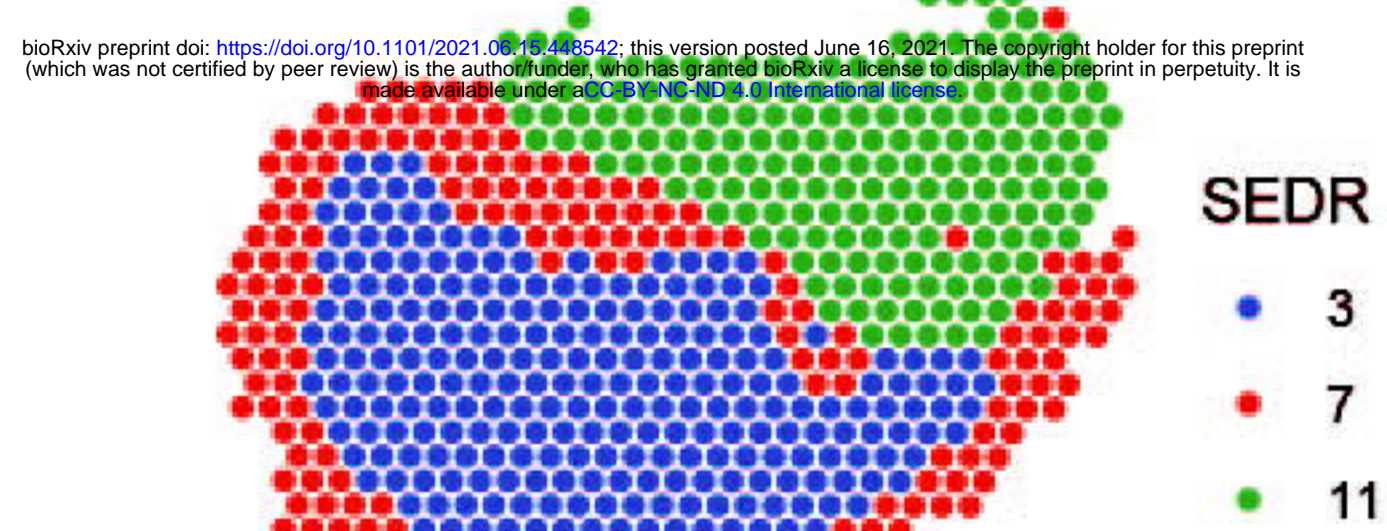
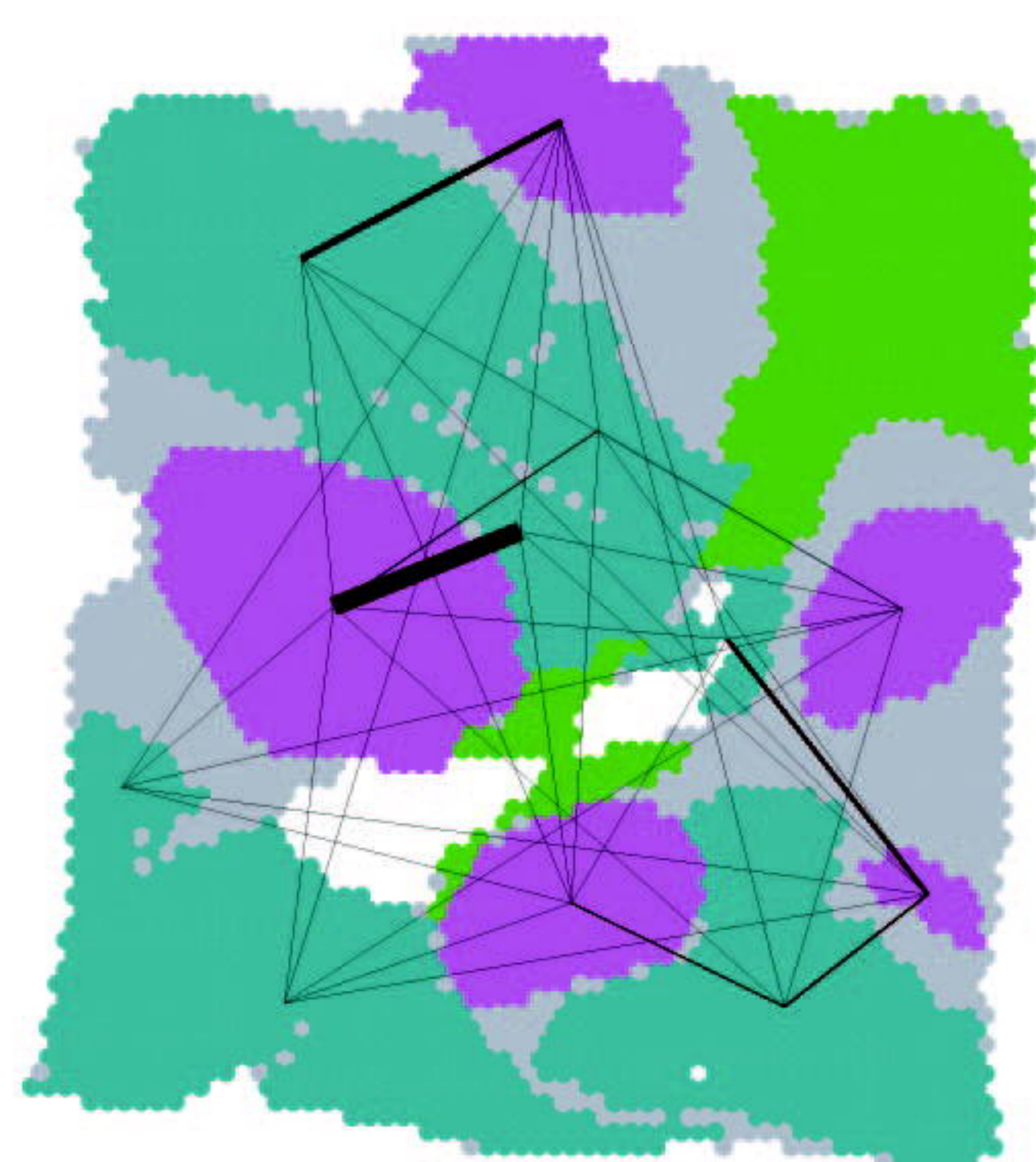


C.Epi

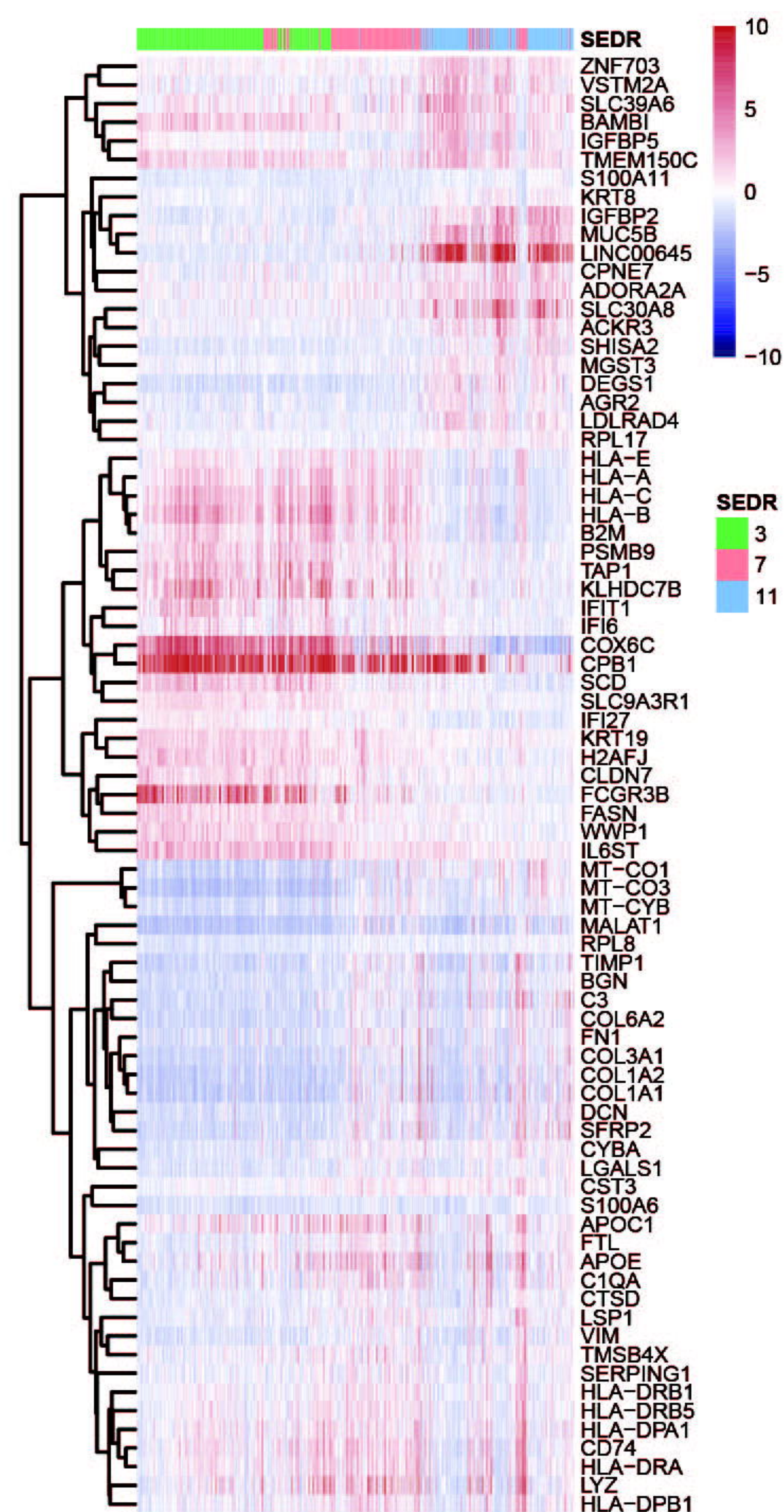


C

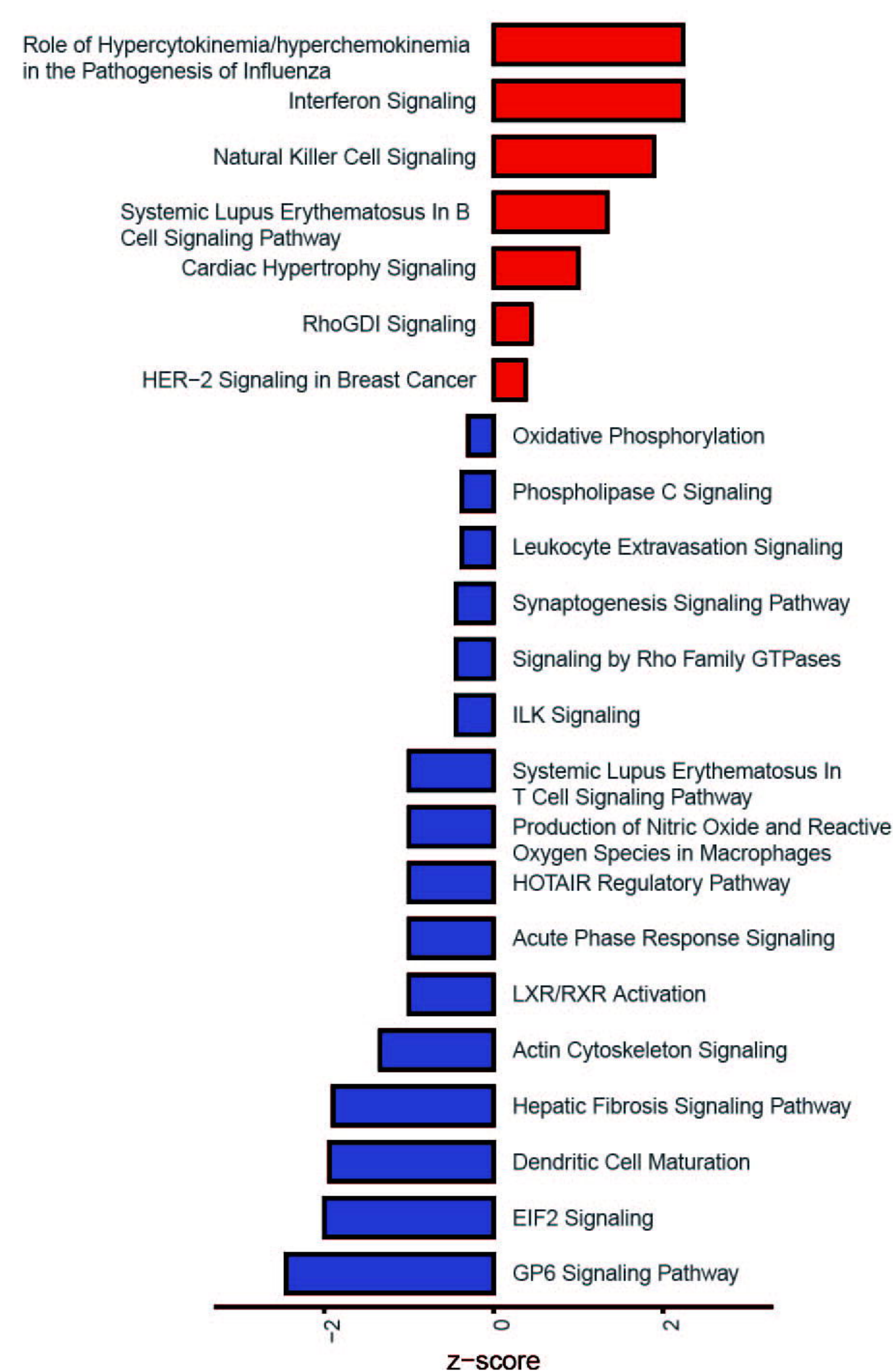
PAGA



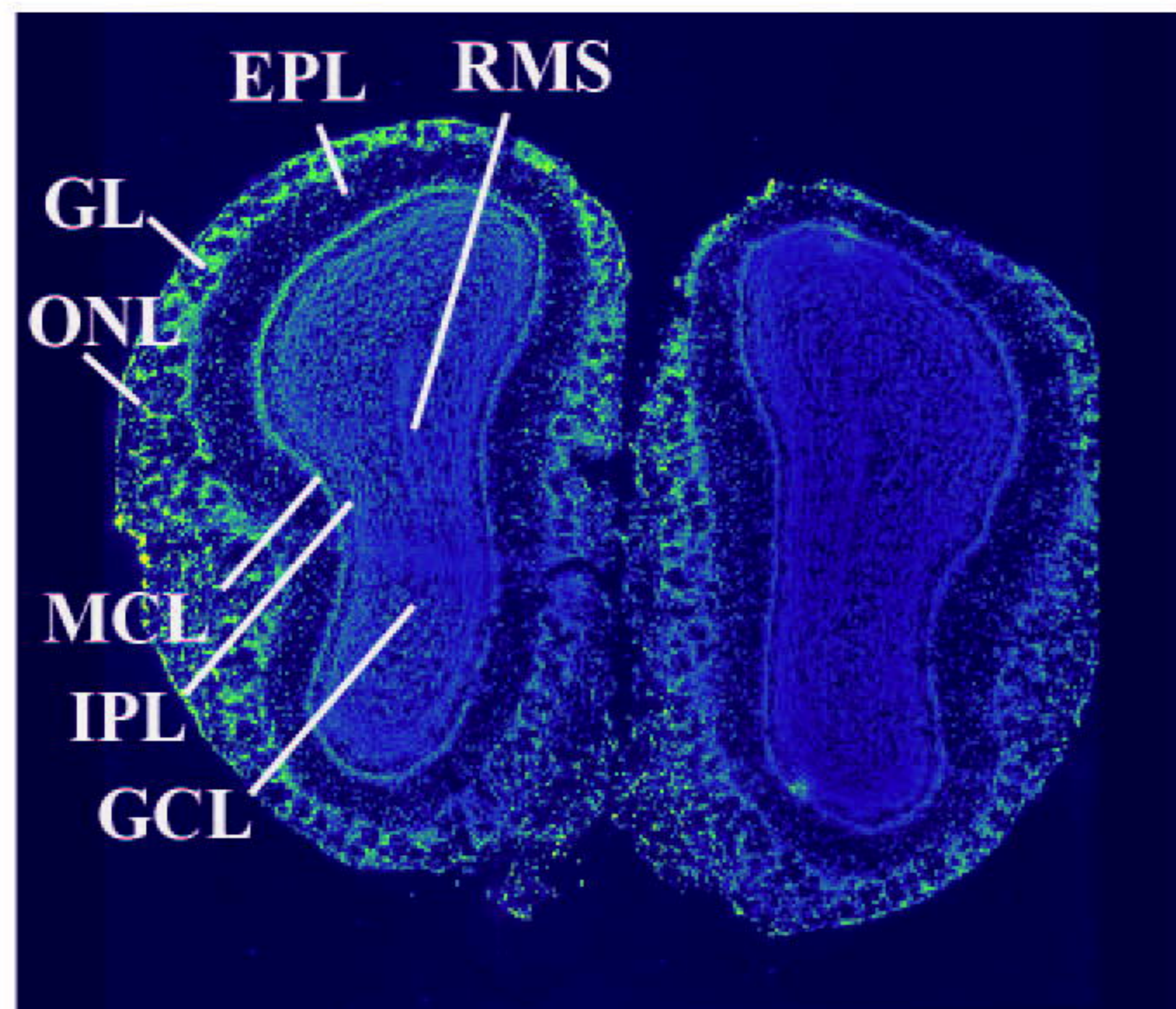
D



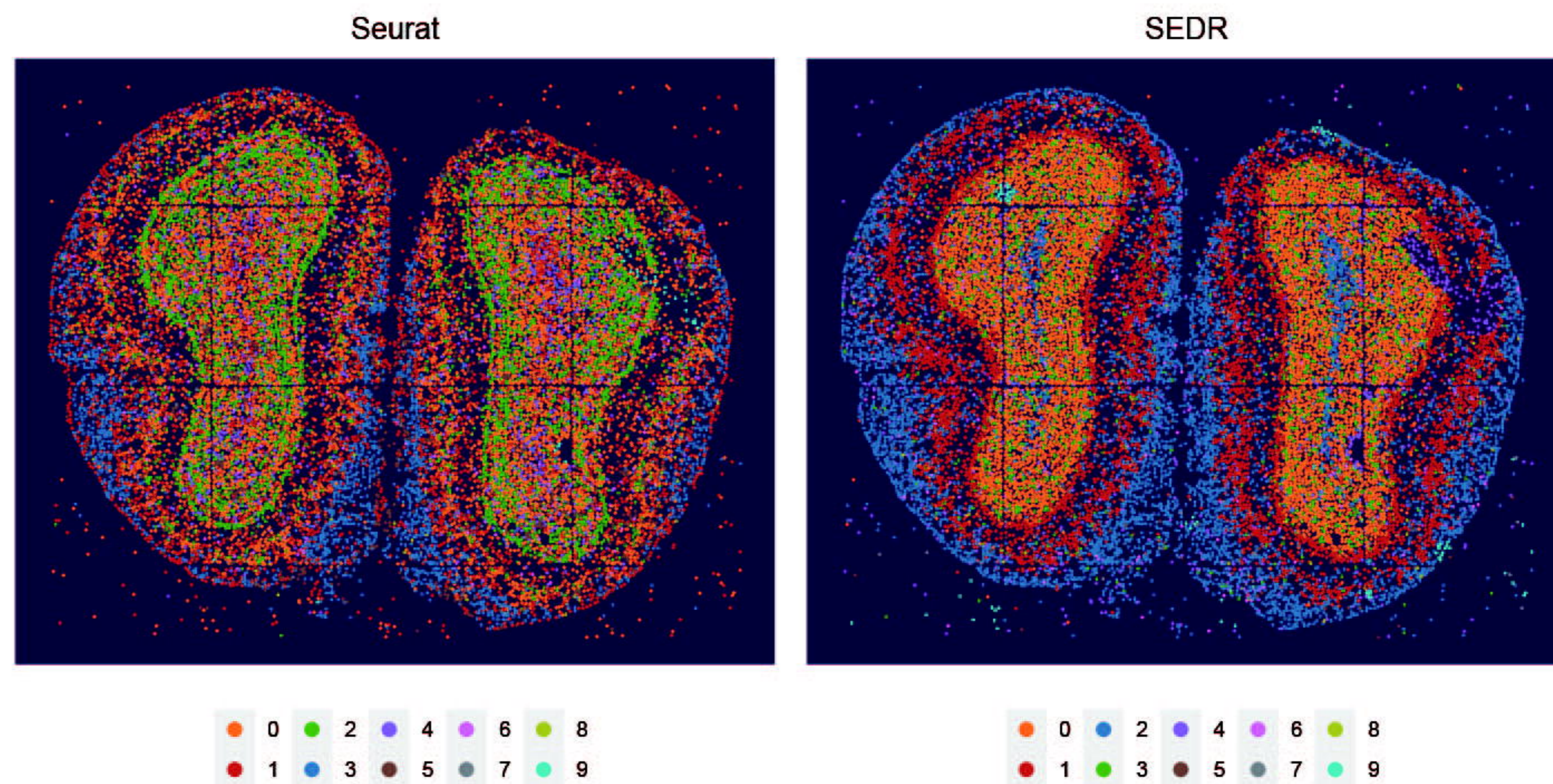
E



A

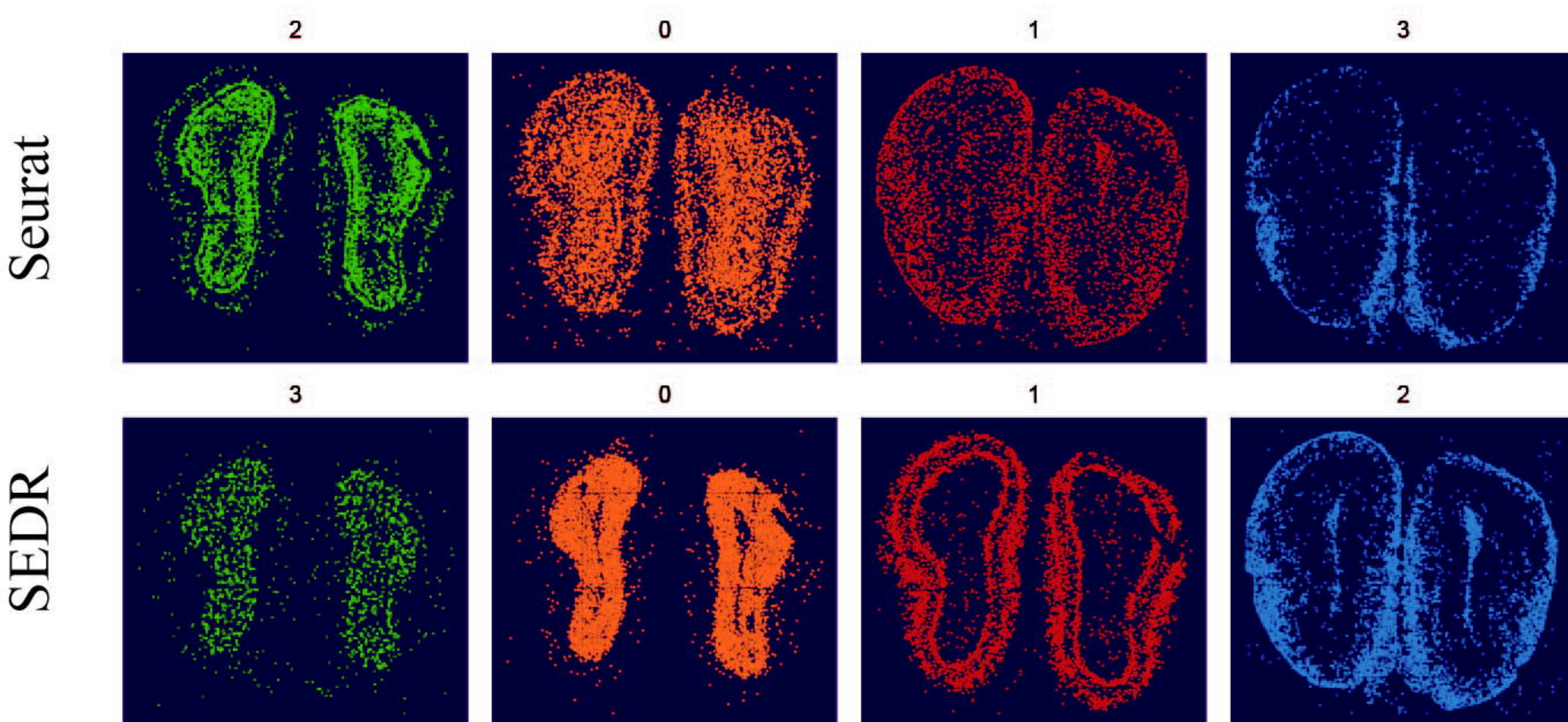


B

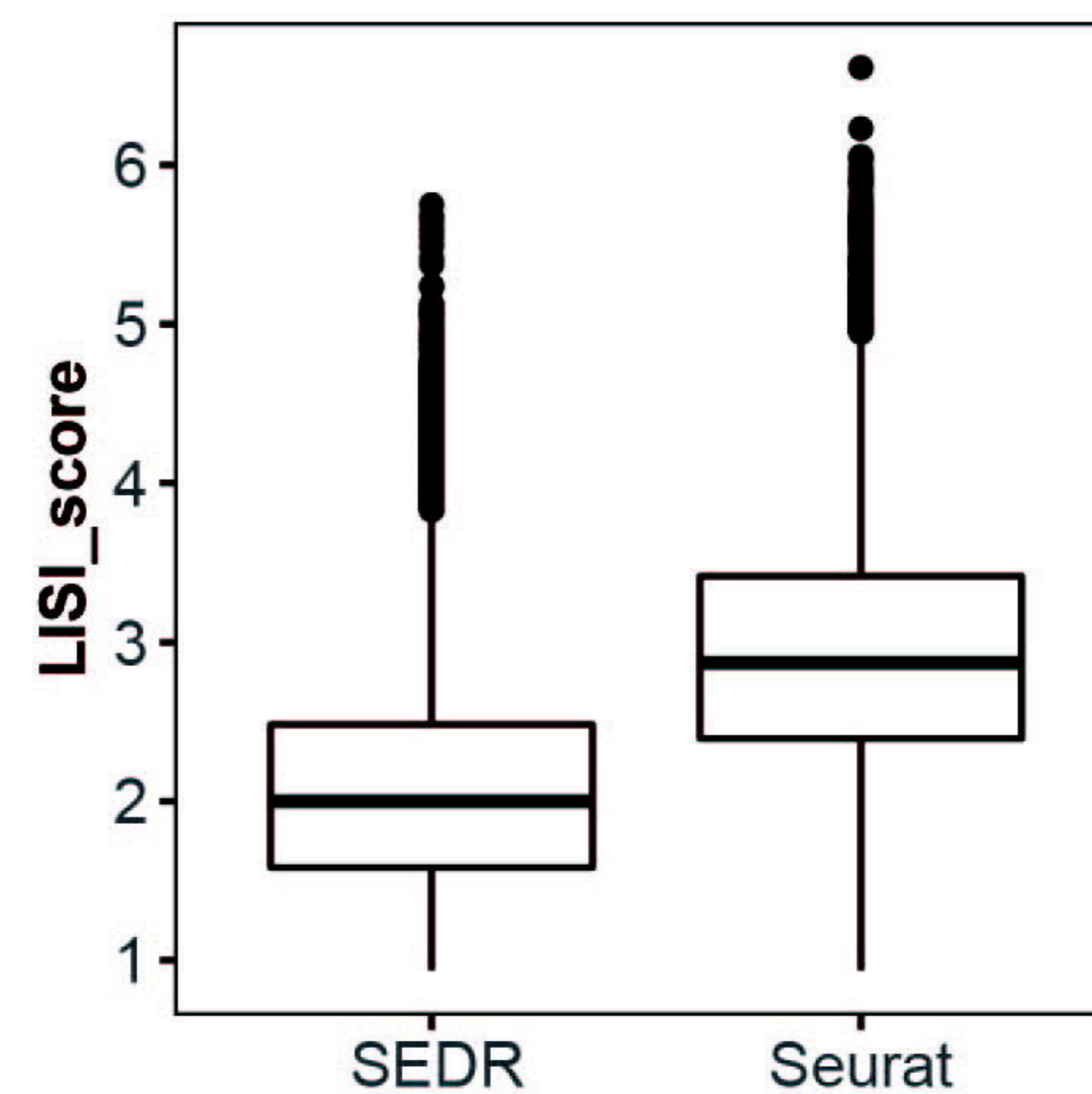


bioRxiv preprint doi: <https://doi.org/10.1101/2021.06.15.448542>; this version posted June 16, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

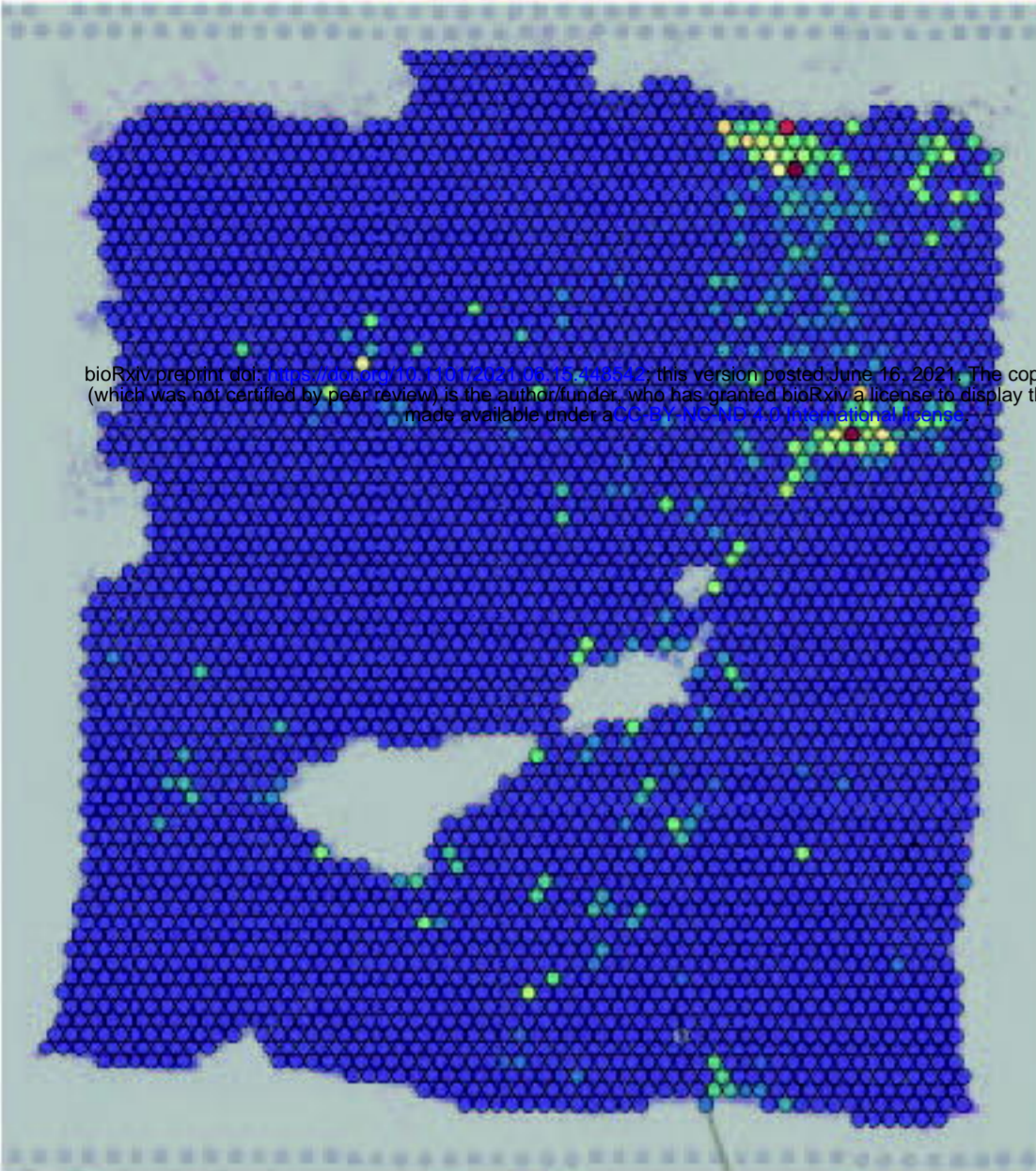
C



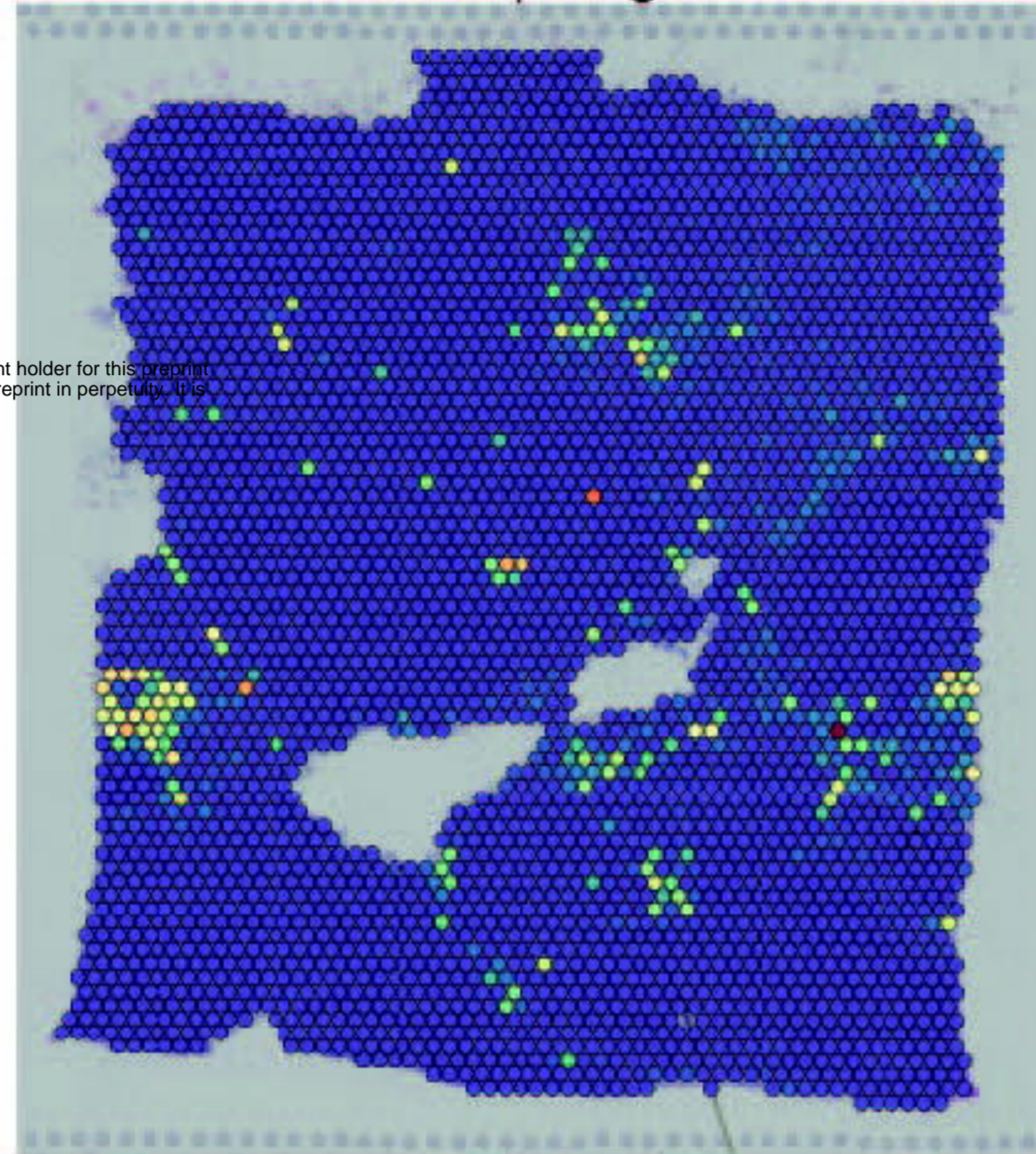
D



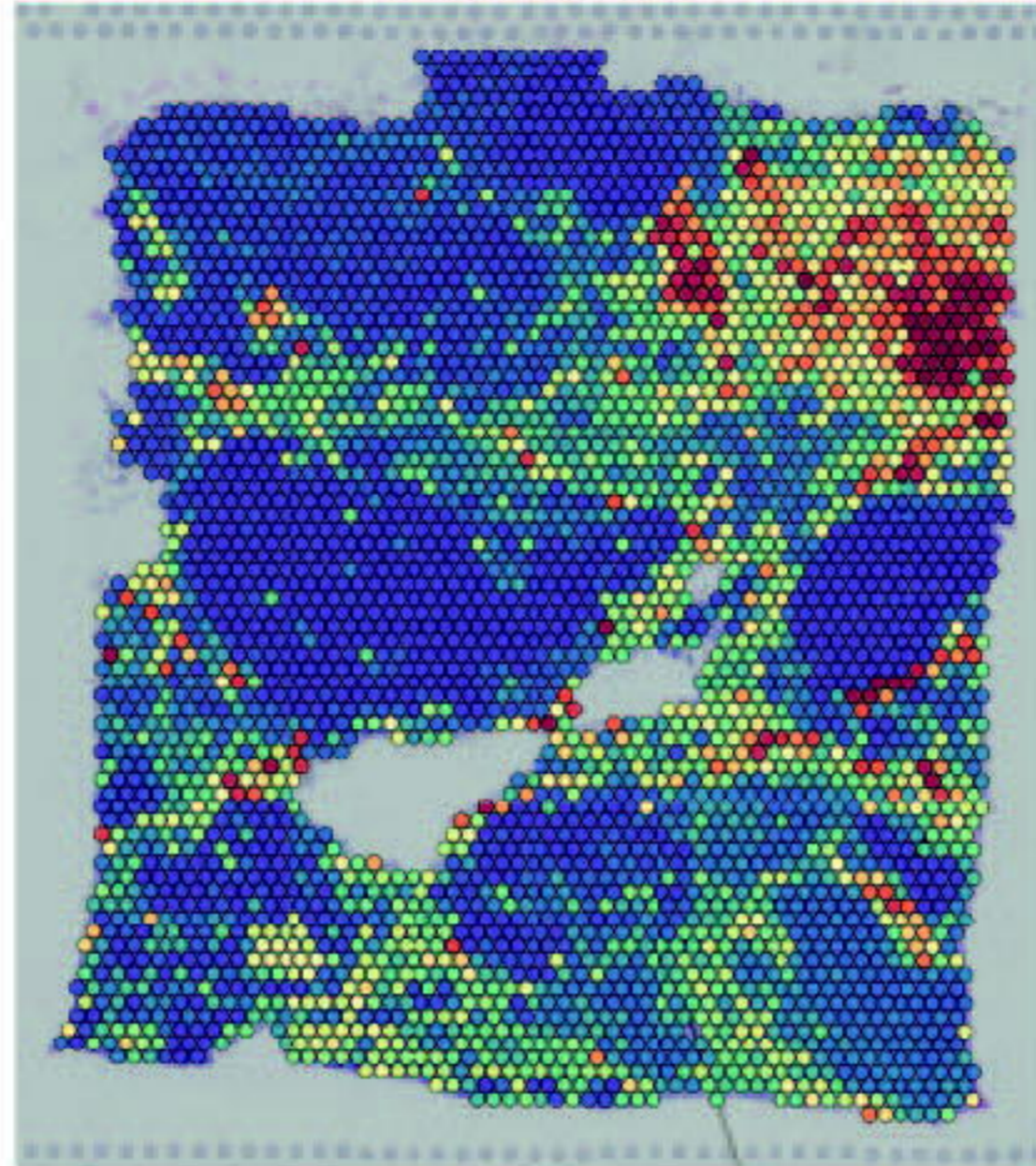
CD4T



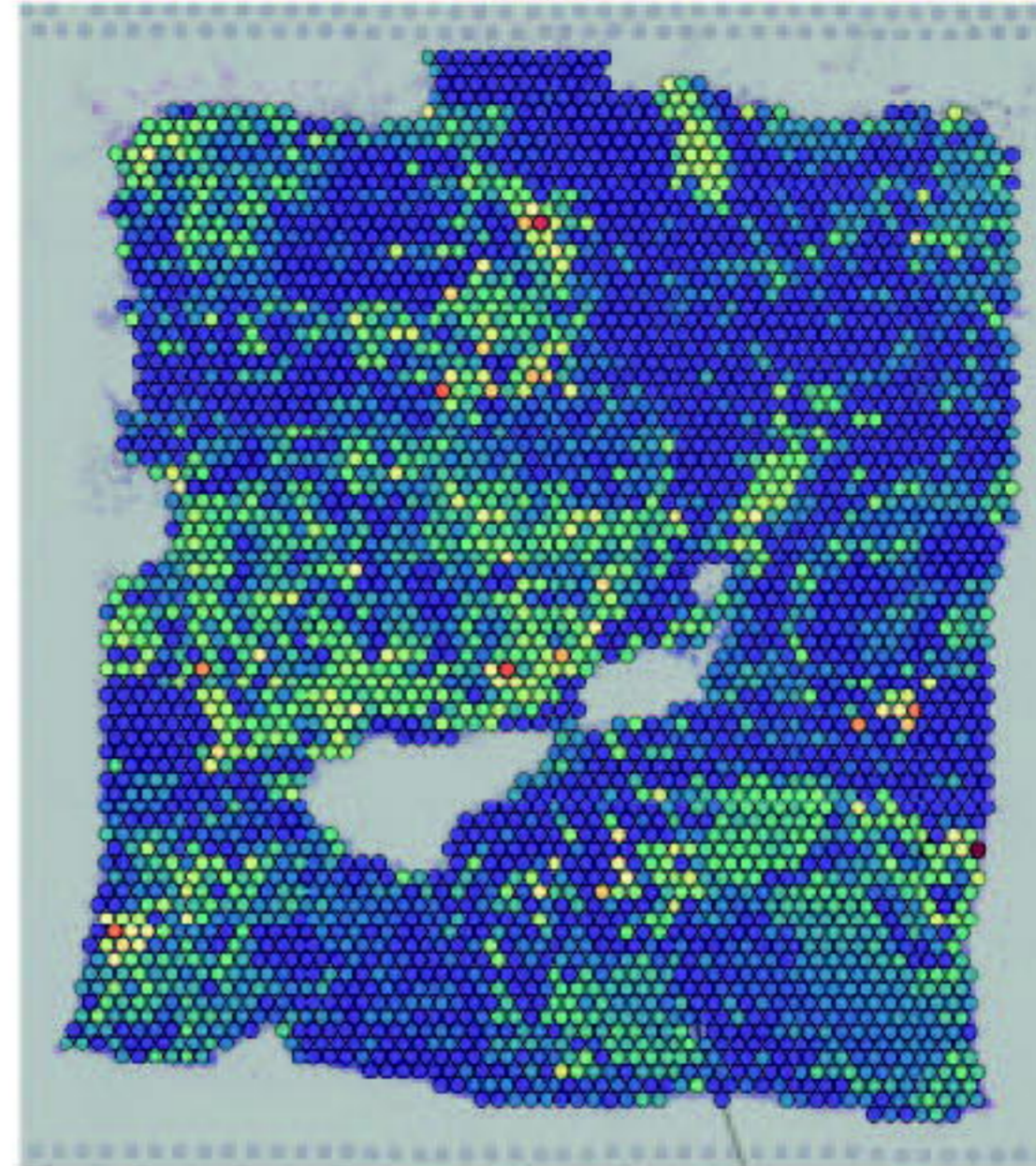
macrophages



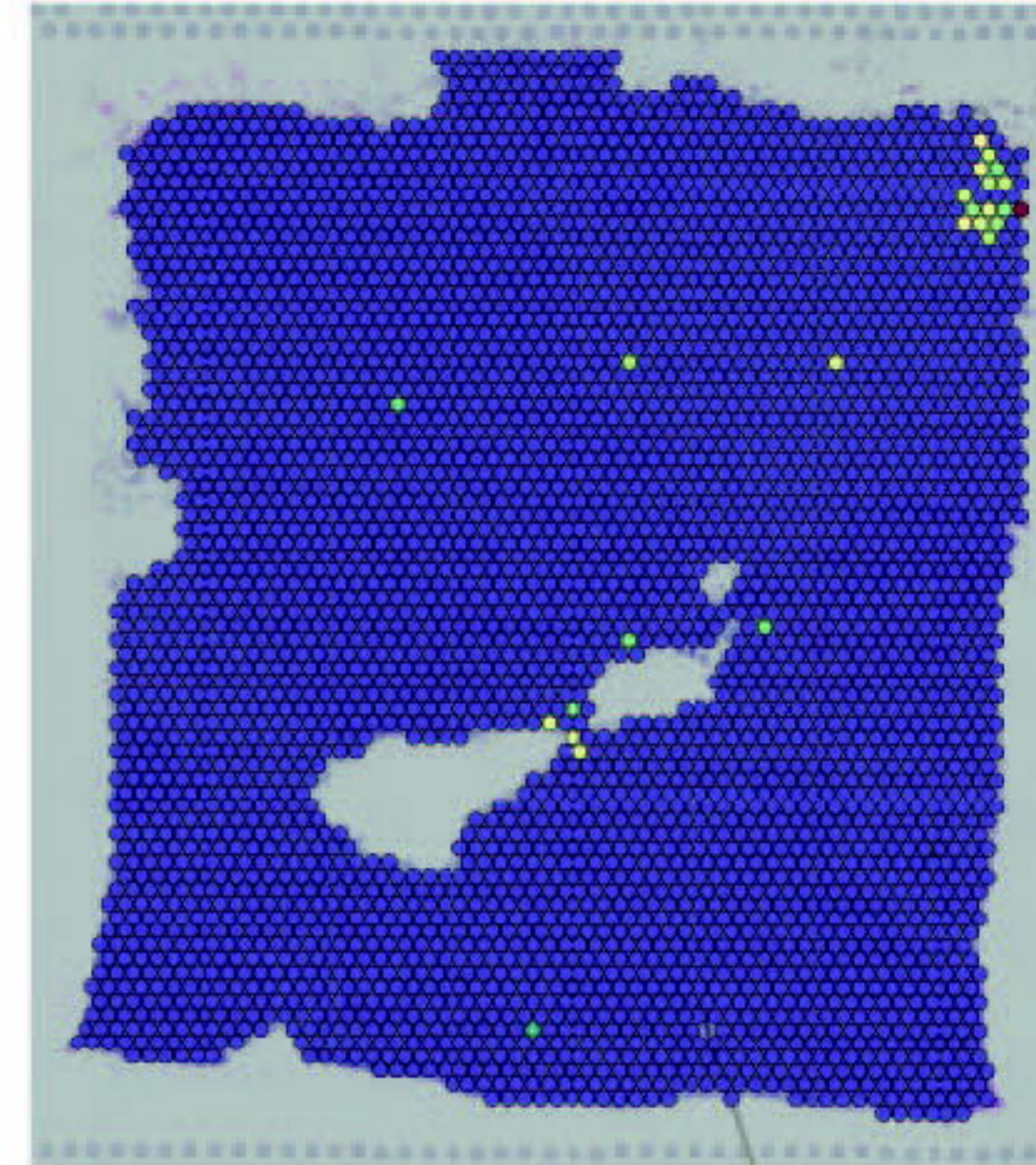
Fibroblast



Endo

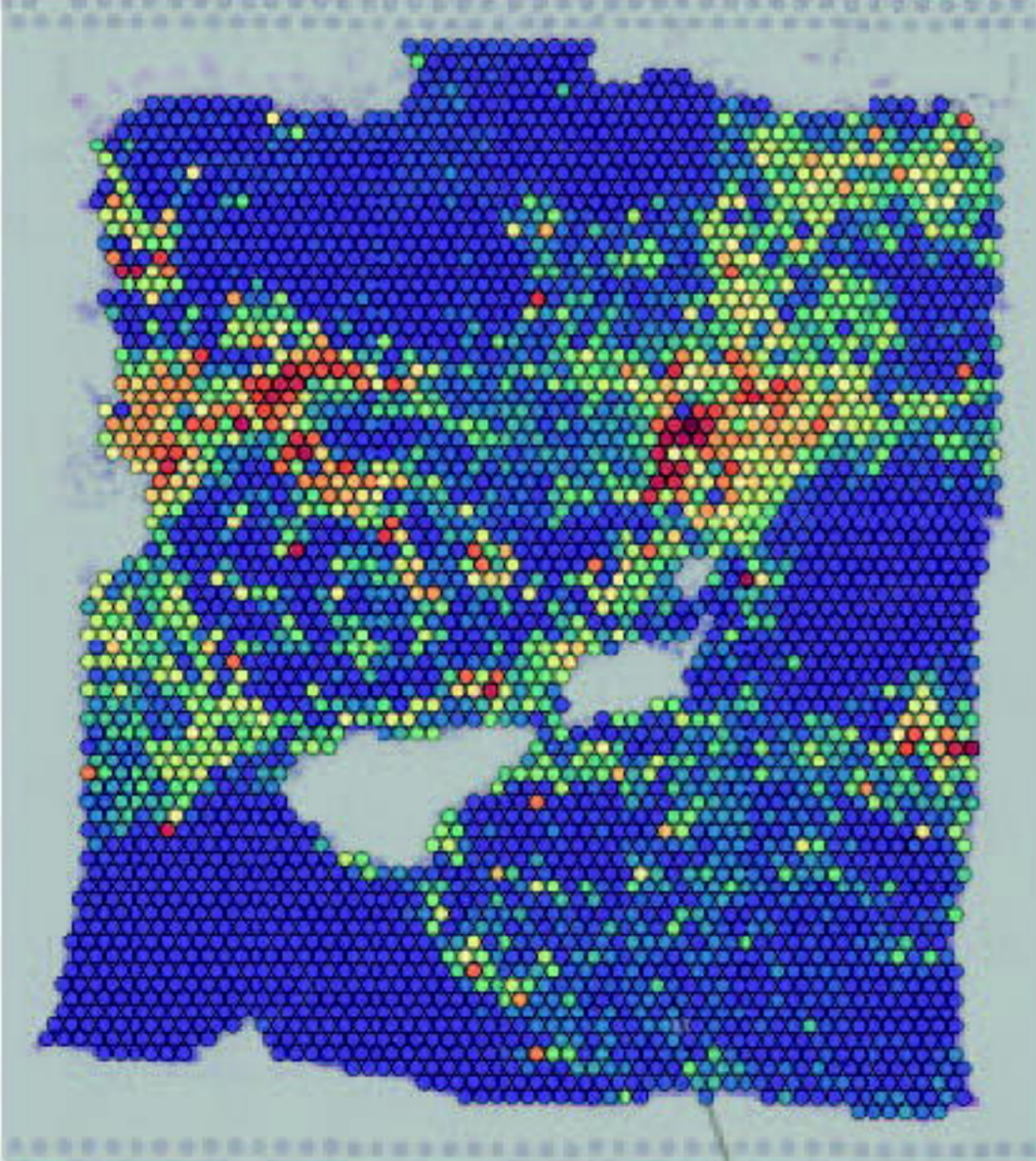


Plasma

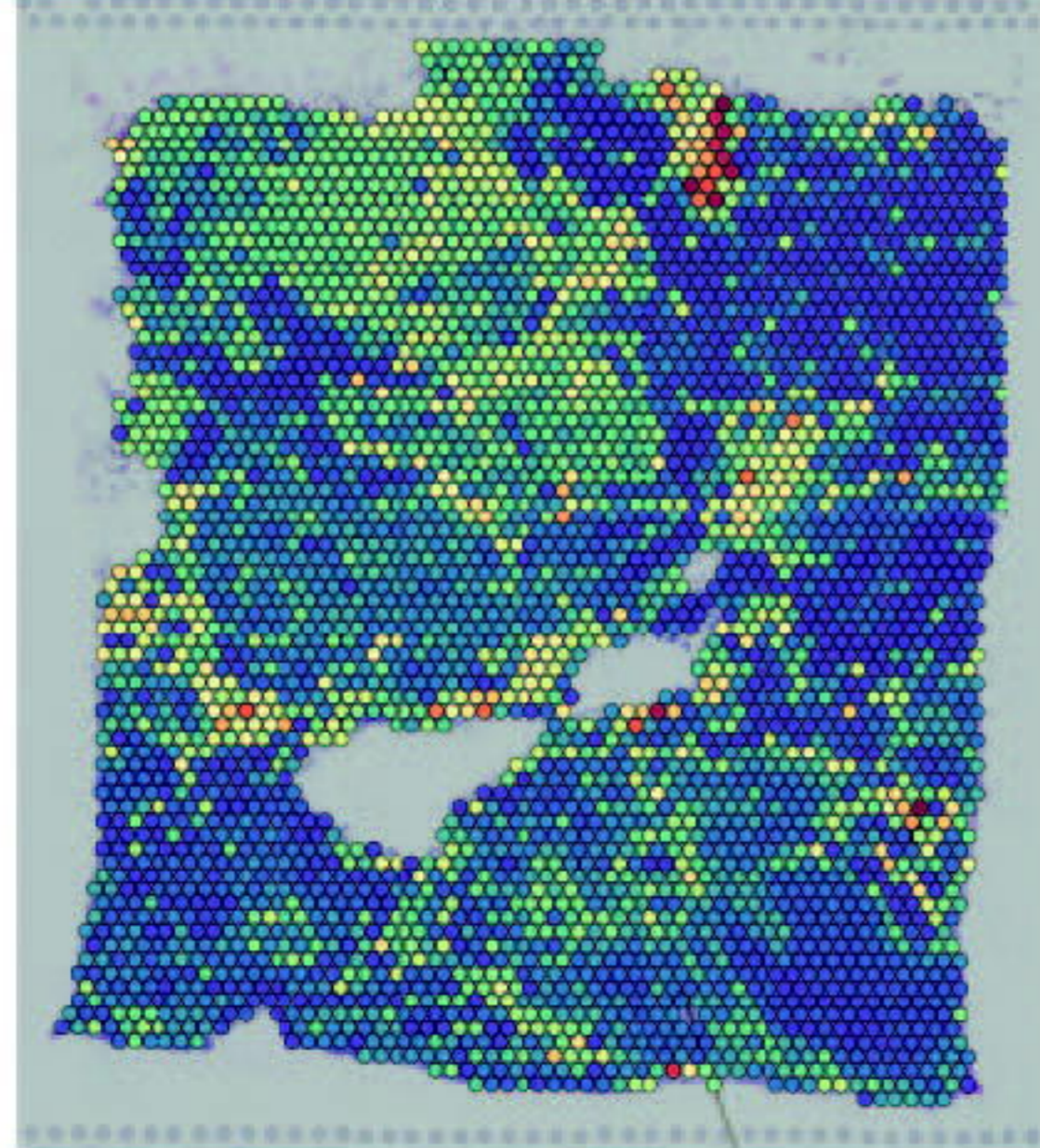


bioRxiv preprint doi: <https://doi.org/10.1101/2021.06.16.202111>; this version posted June 16, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

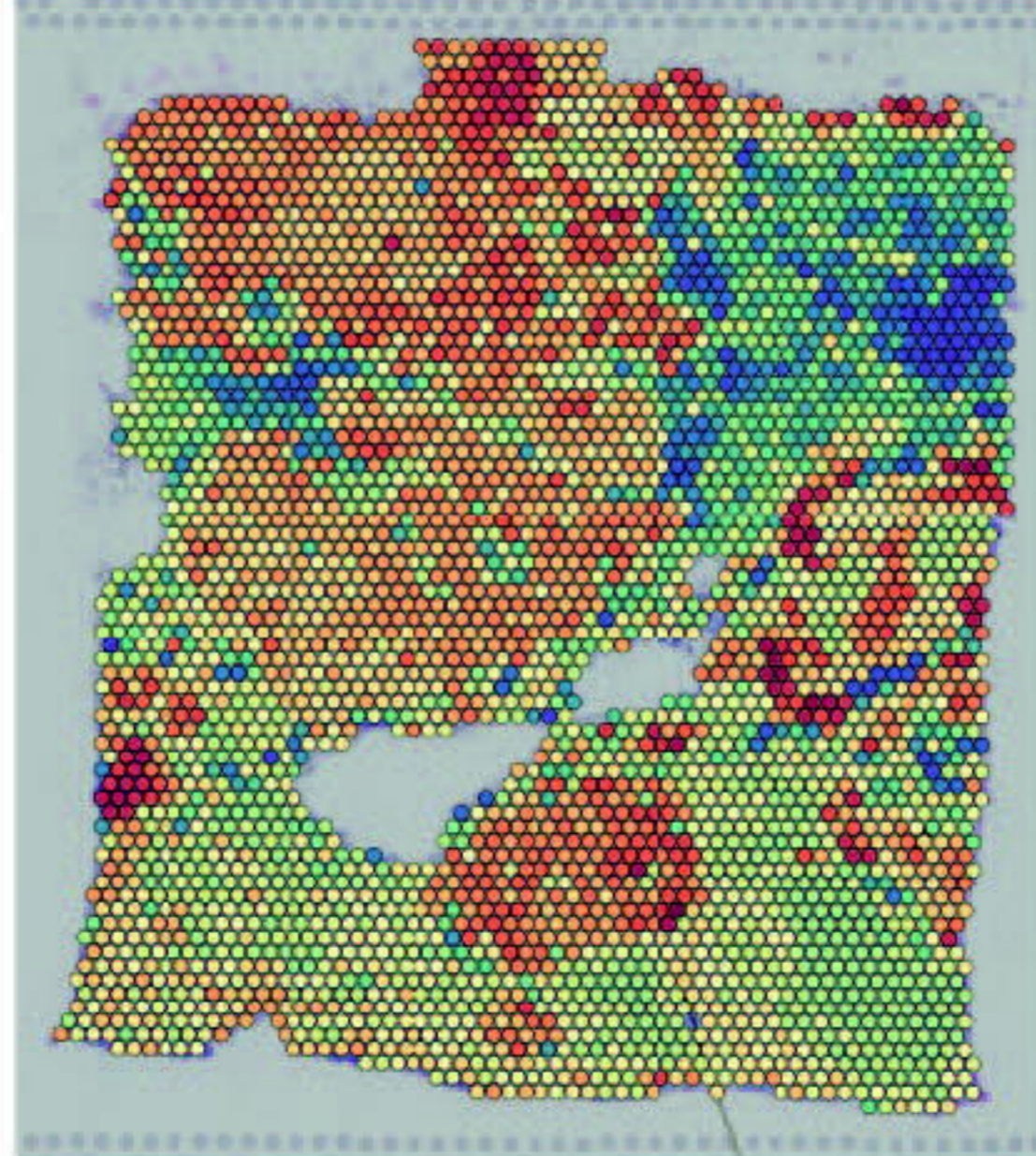
Tumor.associated.macrophages



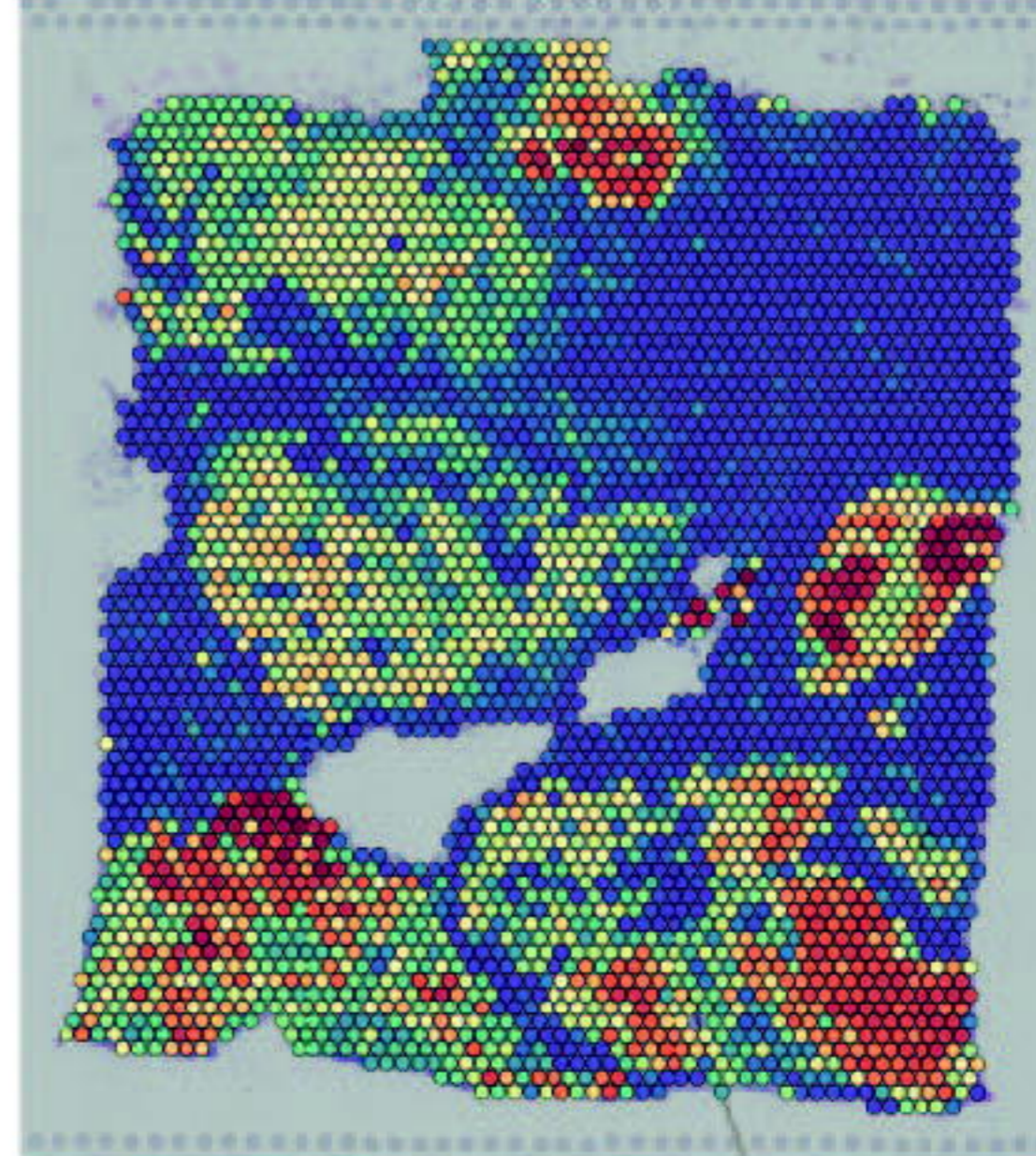
Epi.Basal

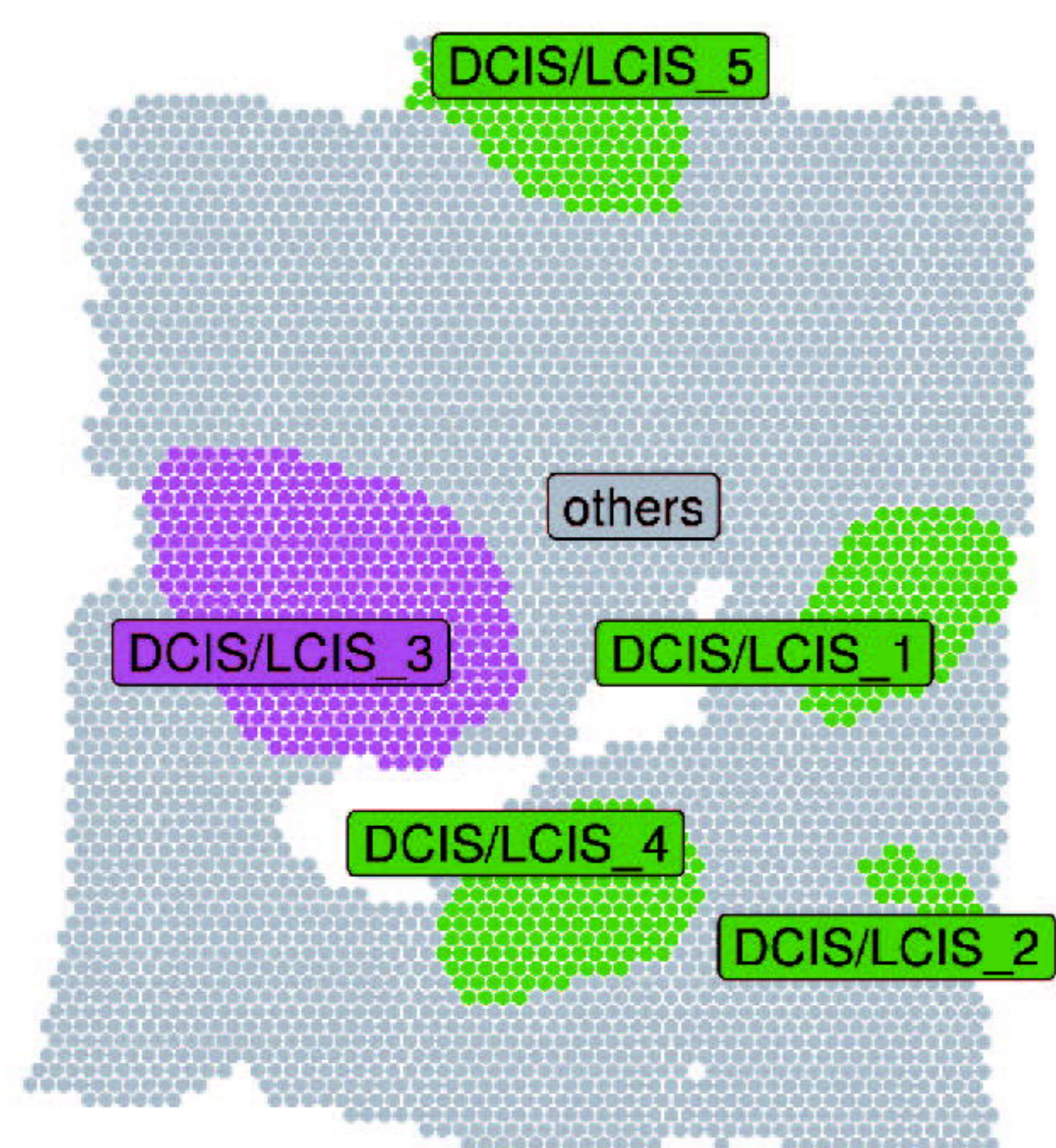
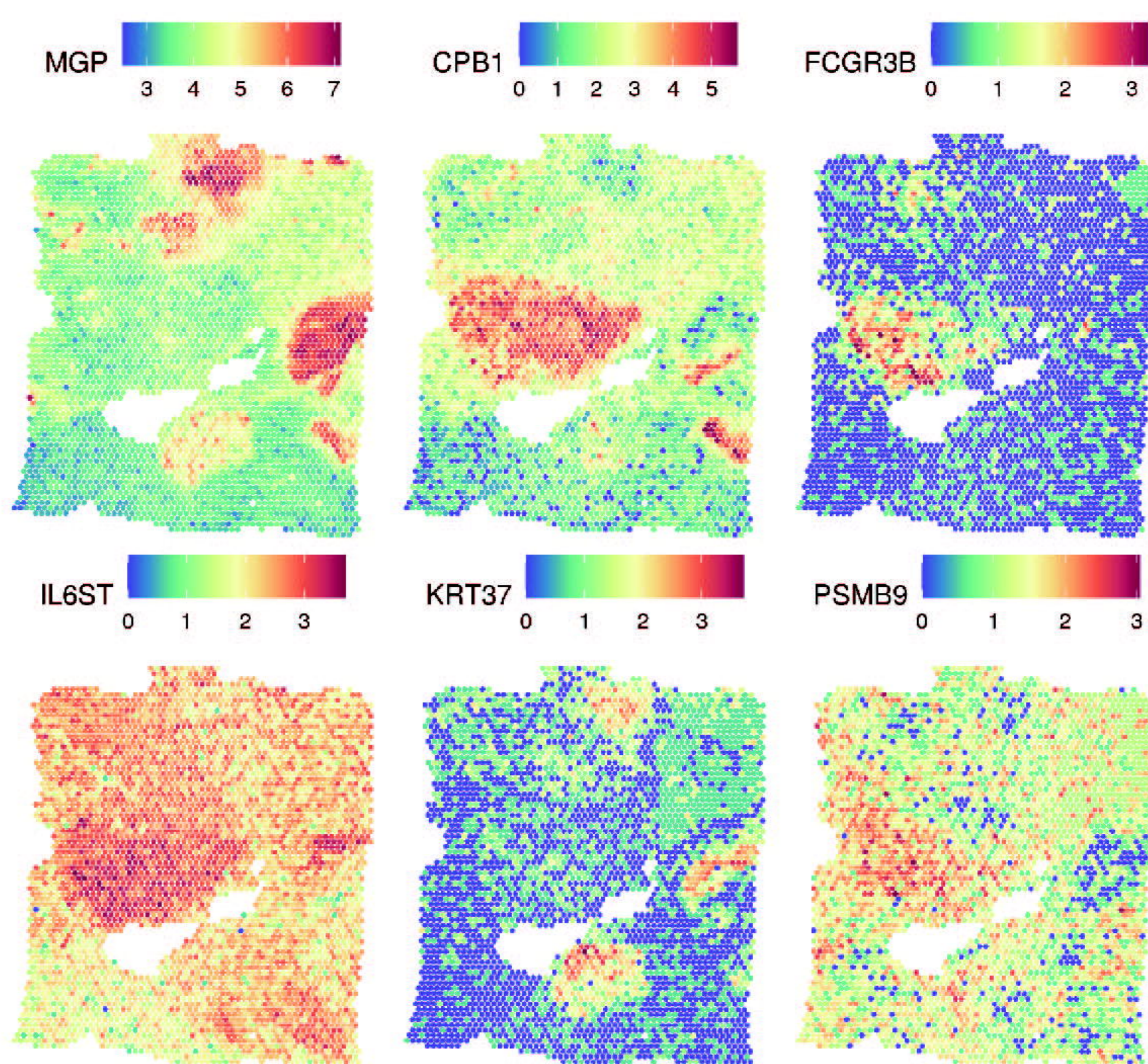
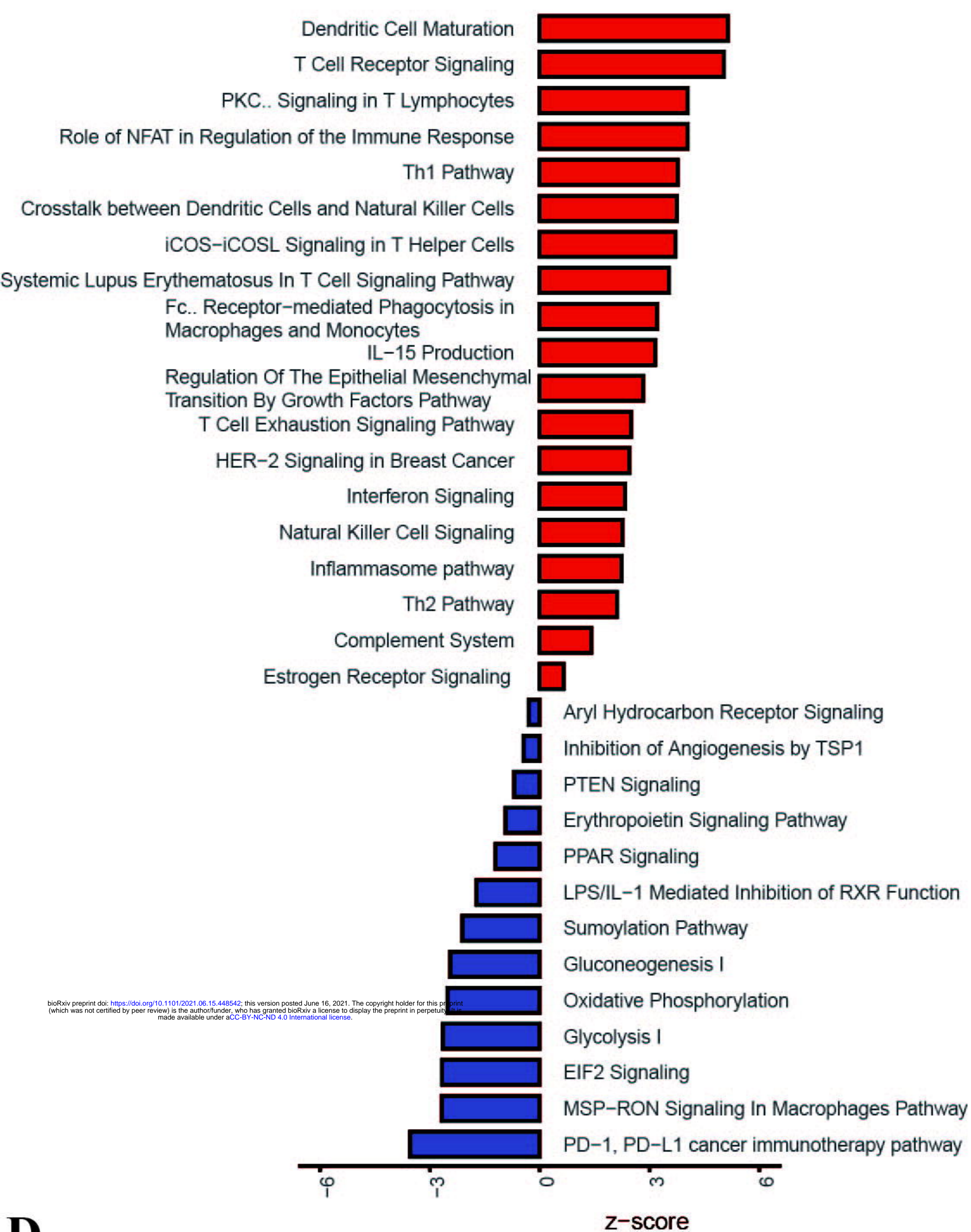
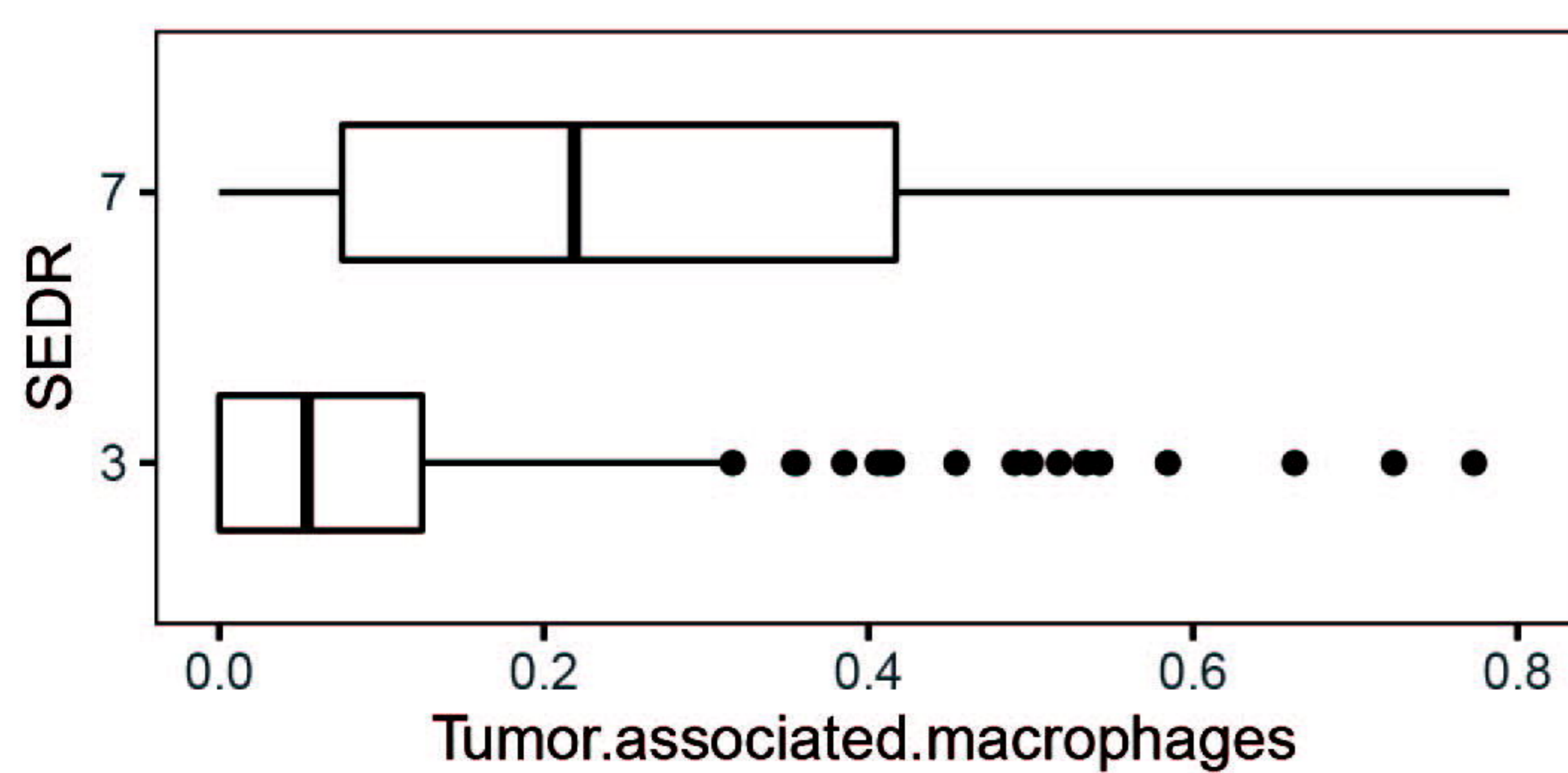


Epi



Cycling.Epi



A**B****C****D****E**