



## Scalable and model-free detection of spatial patterns and colocalization

Qi Liu, Chih-Yuan Hsu and Yu Shyr

*Genome Res.* published online September 9, 2022

Access the most recent version at doi:[10.1101/gr.276851.122](https://doi.org/10.1101/gr.276851.122)

---

<b>P&lt;P</b>	Published online September 9, 2022 in advance of the print journal.
<b>Accepted Manuscript</b>	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
<b>Open Access</b>	Freely available online through the <i>Genome Research</i> Open Access option.
<b>Creative Commons License</b>	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at <a href="http://creativecommons.org/licenses/by-nc/4.0/">http://creativecommons.org/licenses/by-nc/4.0/</a> .
<b>Email Alerting Service</b>	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or <a href="#">click here</a> .

---

---

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

---

To subscribe to *Genome Research* go to:  
<https://genome.cshlp.org/subscriptions>

---

Published by Cold Spring Harbor Laboratory Press

1           **Scalable and model-free detection of spatial patterns and colocalization**

2  
3                           Qi Liu<sup>1,2\*</sup>, Chih-Yuan Hsu<sup>1,2</sup>, Yu Shyr<sup>1,2\*</sup>

4           <sup>1</sup>Department of Biostatistics, Vanderbilt University Medical Center, Nashville, TN 37232, USA

5           <sup>2</sup>Center for Quantitative Sciences, Vanderbilt University Medical Center, Nashville, TN 37232, USA

6           \*Corresponding author: [qi.liu@vanderbilt.edu](mailto:qi.liu@vanderbilt.edu); [yu.shyr@vanderbilt.edu](mailto:yu.shyr@vanderbilt.edu)

7           **Running Title: SpaGene - spatially variable genes detection**

8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19

20 **ABSTRACT**

21 The expeditious growth in spatial omics technologies enable profiling genome-wide molecular events at  
22 molecular and single-cell resolution, highlighting a need for fast and reliable methods to characterize  
23 spatial patterns. We developed SpaGene, a model-free method to discover spatial patterns rapidly in large  
24 scale spatial omics studies. Analyzing simulation and a variety of spatial resolved transcriptomics data  
25 demonstrated that SpaGene is more powerful and scalable than existing methods. Spatial expression  
26 patterns by SpaGene reconstructed unobserved tissue structures. SpaGene also successfully discovered  
27 ligand-receptor interactions through their colocalization.

28

29

30

31

32

33

34

35

36

37

38

## 39 INTRODUCTION

40 Spatial omics technologies map out organizational structures of cells along with their genomics,  
41 transcriptomics, proteomics and epigenomics profiles, providing powerful tools for deciphering  
42 mechanisms of functional and spatial arrangements in normal development and disease pathology  
43 (Larsson et al. 2021; Longo et al. 2021; Marx 2021; Deng et al. 2022; Dhainaut et al. 2022; Ratz et al.  
44 2022; Zhao et al. 2022). The collection of available approaches provides a wide spectrum of throughput  
45 and spatial resolution. Imaging-based approaches generally target pre-selected RNA or proteins at  
46 molecular and single cell resolution, while sequencing-based approaches allow genome-wide profiling  
47 with limited spatial resolution (Lewis et al. 2021; Zhuang 2021). Recent advances in those approaches  
48 move the field rapidly into the direction enabling genome-wide detection with single-cell or subcellular  
49 resolution, presenting a significant computational challenge for scalable and robust methods to derive  
50 biological insights in the spatial context (Atta and Fan 2021).

51 One essential step in spatial omics analysis is to characterize spatial expression patterns and  
52 colocalization. Several methods have been developed to identify spatially variable genes (Edsgard et al.  
53 2018; Svensson et al. 2018; Sun et al. 2020a; Anderson and Lundeberg 2021; Miller et al. 2021; Zhu et al.  
54 2021). Trendsceek uses permutation test to detect significant dependency between the spatial distribution  
55 of points and their expression levels based on marked point processes (Edsgard et al. 2018). Sepal ranks  
56 spatially variable genes by the diffusion time with the rationale that genes with spatial patterns require  
57 more time to reach a homogenous state than those with random spatial distributions (Anderson and  
58 Lundeberg 2021). SpatialDE and SPARK both utilize Gaussian process regression as the underlying data  
59 generative model for spatial covariance structures. SpatialDE decomposes expression variability into  
60 spatial variance and noise, and estimates statistical significance by comparing the likelihoods with and  
61 without a spatial component (Svensson et al. 2018). SPARK extends SpatialDE via generalized linear  
62 spatial error models, with the ability to directly model raw counts and adjust for covariates (Sun et al.  
63 2020a). SPARK-X examines the similarity of expression covariance matrix and distance covariance

64 matrix and tests whether they are more similar than expected by chance (Zhu et al. 2021). The statistical  
65 power of such methods highly depends on spatial covariance models, i.e, how well they match true  
66 underlying expression patterns. Although multiple kernels, including Gaussian, linear and periodic  
67 kernels with different smoothness parameters, are considered to ensure identification of various spatial  
68 patterns, statistical power will be compromised substantially for identifying spatial patterns poorly  
69 modelled by those predefined kernel functions. Furthermore, spatial covariance models are built upon  
70 cellular distances, which would confound true expression variances with those driven by variances in  
71 cellular densities. To take non-uniform cellular densities into consideration, MERINGUE calculates  
72 spatial autocorrelation and cross-correlation based on spatial neighborhood graphs to identify spatially  
73 variable genes and gene interactions (Miller et al. 2021). Above all, even equipped with computationally  
74 efficient algorithms, it would still take days to months for most methods to analyze large-scale spatial  
75 data with genome-wide profiling in tens of thousands of locations (Zhu et al. 2021), resulting in a high  
76 demand for scalable and robust methods for characterizing spatial expression patterns.

77 To address those limitations, we aim to develop a scalable and model-free method for detecting spatial  
78 patterns. Without making assumptions on spatial covariance models and data distributions, the method  
79 will have more degree of freedom and also be more computationally efficient in identifying spatial  
80 patterns than existing methods.

## 81 **RESULTS**

### 82 **Overview of SpaGene**

83 SpaGene is built upon a simple intuition that spatially variable genes have uneven spatial distributions,  
84 meaning that highly expressed cells/spots tend to be more spatially connected than random. Given a set of  
85 spatial locations, SpaGene first builds the spatial network using  $k$ -nearest neighbors. For each gene,  
86 SpaGene then extracts a subnetwork comprising only cell/spots with high expression of the gene from the  
87  $k$ -nearest neighbor graph. SpaGene quantifies the connectivity of the subnetwork by the Earth's Mover  
88 distance between degree distributions of the subnetwork and a fully connected one. Finally, SpaGene

89 compares the observed and the expected distances from random permutations. Genes with significantly  
90 shorter distances than random are identified to be spatially variable (Fig. 1A).

## 91 **Simulation**

92 We first applied SpaGene on two simulation datasets. One simulation was generated from negative  
93 binomial distributions following SPARK-X (Zhu et al. 2021), the other was sampled from real data  
94 following Trendsceek (Edsgard et al. 2018). Cells/spots with higher expression (spiked cells) were  
95 located in one of those five patterns, hotspot, streak, circularity, bi-quarter circularity, and Purkinje layer  
96 in mouse cerebellum (Fig. 1B). The distinctness of the pattern was determined by effect sizes, which were  
97 controlled by the fold change (FC) of expression in spiked cells compared to the background. The pattern  
98 size was determined by the percentage of spiked cells. Higher effect sizes and larger pattern sizes  
99 generated more distinct and bigger patterns, which were easier to be identified. Among the simulated  
100 genes, 500 genes display spatial patterns (details in the Methods section). The area under the curve (AUC)  
101 was used to measure the ability to distinguish between spatially and non-spatially variable genes.

102 We compared SpaGene with SpatialDE and SPARK-X. SpatialDE and SPARK-X both achieved high  
103 computational efficiency and good performance in other studies and SPARK-X is the only method  
104 applicable to data with sample size exceeding 30,000 (Zhu et al. 2021). As expected, effect sizes are the  
105 major factor affecting performance. Larger effect sizes produced more distinct patterns, which were easier  
106 to be distinguished from random spatial distributions and resulted in higher AUC values. For hotspot and  
107 streak patterns, SpaGene, SpatialDE, and SPARK-X successfully distinguished spatially from non-  
108 spatially variable genes when patterns were distinct (AUC=1 at  $FC \geq 5$  for hotspot and AUC=1 at  $FC \geq 8$   
109 for streak patterns). For less distinct patterns, SpaGene performed slightly better than SpatialDE and  
110 SPARK-X for smaller patterns, which obtained AUC of 0.64, 0.52 and 0.55 for SpaGene, SPARK-X and  
111 SpatialDE respectively at  $FC=2$  and  $size=1$  in hotspot patterns, while SPARK-X outperformed SpatialDE  
112 and SpaGene for bigger patterns ( $size > 1$ ) (Fig. 1C). For circularity and bi-quarter circularity patterns,  
113 SpaGene achieved much better performance than SpatialDE and SPARK-X. For the circularity pattern,

114 SpaGene achieved AUC of 0.99 even for the smallest pattern at  $FC=3$  and AUC of 1 at  $FC \geq 5$ . In  
115 comparison, SpatialDE only obtained AUC of 0.73 at  $FC=3$ , and SPARK-X failed to distinguish spatially  
116 from non-spatially variable genes even at  $FC=5$  (AUC=0.5) for the smallest pattern (size=1). SpaGene  
117 and SpatialDE achieved AUC of 1 while SPARK-X only obtained AUC of 0.72 at  $FC=8$  and size=1.  
118 Although the performance of SpatialDE and SPARK-X improved with increasing pattern sizes, SpaGene  
119 was more powerful than SpatialDE and SPARK-X (Fig. 1C). For the bi-quarter circularity pattern,  
120 SPARK-X failed even at the largest effect size for the two small patterns (AUC=0.5 at  $FC=10$ , size=1 or  
121 2), while SpaGene achieved  $AUC \geq 0.9$  and SpatialDE obtained AUC of 0.7-0.83 at  $FC \geq 3$  for any  
122 pattern sizes (Fig. 1C). For the Purkinje layer pattern, SPARK-X failed at any effect sizes (AUC=0.5),  
123 while SpaGene achieved AUC of 0.81 at  $FC=2$ , 0.99 at  $FC=3$  and 1 at  $FC \geq 5$  (Fig. 1C). SpatialDE was  
124 not applied in this setting due to long computational time. To summarize, SpaGene achieved good  
125 performance for all spatial patterns, which obtained  $AUC \geq 0.98$  at  $FC \geq 3$  for relatively big patterns  
126 (size>1) and AUC close to 1 at  $FC \geq 5$  for any pattern sizes. In comparison, SPARK-X seemed to be very  
127 sensitive to pattern shapes, which worked well for hotspot and streak patterns, but not for circularity, bi-  
128 quarter circularity and Purkinje layer patterns even when patterns were strongly distinct from the  
129 background. Furthermore, SpaGene was more robust against pattern sizes than SpatialDE and especially  
130 SPARK-X, which sometimes showed more power to identify indistinct and large patterns than small  
131 distinct patterns. For example, SPARK-X obtained AUC of 0.8 at  $FC=3$  and size=3, but AUC of 0.7 even  
132 at  $FC=8$  and size=1 for circularity patterns. SpatialDE obtained AUC of 0.7 at  $FC=3$  and size=1, but 0.82  
133 at  $FC=2$  and size=5 for bi-quarter circularity patterns. We also simulated scenarios with varying number  
134 of genes and cells/locations (Supplemental Figs. S1-S5). We found that the performance of SpaGene were  
135 less dependent on the number of cells/locations compared to SpatialDE and SPARK-X. The evaluation  
136 on the simulation datasets sampled from real data obtained similar results (Supplemental Figs. S6-S9).  
137 In terms of time complexity, SpaGene and SPARK-X are much more computationally efficient than  
138 SpatialDE. SpatialDE requires several orders of computational time than SpaGene and SPARK-X, and its  
139

140 runtime increases linearly or cubically with the number of genes and the number of cells/locations  
141 (Supplemental Fig. S10A). For example, it takes SpatialDE 4,045 seconds to analyze a data with 10,000  
142 genes and 5,000 cells/location, while it only takes SpaGene and SPARK-X 11 and 22 seconds,  
143 respectively. Additionally, SpaGene and SPARK-X require less memory than SpatialDE. SPARK-X and  
144 SpaGene require 0.5G and 0.6G memory respectively, while SpatialDE demands 1.6G memory to analyze  
145 a data with 10,000 genes and 5,000 locations (Supplemental Fig. S10B).

146

### 147 **Application to MOB by spatial transcriptomics**

148 We applied SpaGene to spatial transcriptomics data from main olfactory bulb (MOB) (Stahl et al. 2016),  
149 involving 16,218 genes measured on 262 spots. The MOB has a roughly concentric arrangement of seven  
150 cell layers (Nagayama et al. 2014). SpaGene identified spatially variable 634 genes (adjusted p-value,  
151  $\text{adj}p < 0.05$ ), including genes known to be located in specific layers. Several examples were shown in Fig.  
152 2A, such as *Pcp4* in Granule cell layer (GCL) ( $\text{adj}p = 3e-6$ ) (Sangameswaran et al. 1989), *Slc17a7* in  
153 Mitral cell layer (MCL) ( $\text{adj}p = 7e-4$ ) (Zhang et al. 2021), *Cck* in Glomerular layer (GL) ( $\text{adj}p = 2e-3$ ) (Sun  
154 et al. 2020b), *Serpine2* in External plexiform layer (EPL) ( $\text{adj}p = 4e-3$ ) (Mansuy et al. 1993), and *Fabp7* in  
155 Olfactory nerve layer (ONL) ( $\text{adj}p = 4e-76$ ) (Young et al. 2013). Based on those identified spatially  
156 variable genes, SpaGene successfully reconstructed the underlying seven-layered MOB structure  
157 (Supplemental Fig. S11). To be noted, SpaGene identified a pattern corresponding to subependymal zone  
158 (SEZ) (pattern 4 in Supplemental Fig. S11). SEZ was unidentifiable by transcriptional profiles-based  
159 clustering, which only discovered five distinct clusters (Supplemental Fig. S12A). SEZ harbors neural  
160 stem cells. *Sp9* is the top gene specifically located in SEZ, which is a transcription factor that regulate  
161 MOB interneuron development (Li et al. 2018).

162 We compared SpaGene with SPARK-X and SpatialDE. Overall, SpaGene and SpatialDE had more  
163 overlapping than SPARK-X (Supplemental Fig. S12B). The original study highlighted 15 genes  
164 differentially expressed in different domains (Stahl et al. 2016). SpaGene detected 12 of 15 genes, while  
165 SPARK-X only found five and SpatialDE identified nine. Since cell clustering based on transcriptional

166 profiles alone uncovered cell types located in MOB layers, genes highly expressed in each layer-specific  
167 cell clusters should be identified to be spatially variable. Using the top 20 markers from each of those cell  
168 clusters as the ground truth, SpaGene achieved a higher true positive rate than SPARK-X and SpatialDE  
169 (Supplemental Fig. S12C). We also calculated scores to measure the enrichment of those top markers in  
170 SpaGene, SPARK-X and SpatialDE. SpaGene obtained high enrichment scores in all layers, suggesting it  
171 successfully identified all layer-specific marker genes as being very significant. In contrast, SPARK-X  
172 obtained high in GCL layers but low scores in other layers. SpatialDE achieved high scores in Mitral cell  
173 layer , but relatively low scores in GCL and EPL layers (Fig. 2B). Moreover, we compiled the top 50  
174 genes with enhanced expression in each layer of the MOB using the “differential search” in Allen Mouse  
175 Brain atlas, which obtained 222 genes in total. Using the 222 genes as the ground truth, SpaGene also  
176 obtained a higher true positive rate than SPARK-X and SpatialDE (Supplemental Fig. S12D) . Finally, we  
177 ranked spatially variable genes by each method and carefully examined those genes identified to be very  
178 significant by one method but insignificant by another method. First, we ranked genes by SpaGene and  
179 listed the top six genes with inconsistent results (Supplemental Fig. S13). *Kif5b*, *Atf5*, *Sorbs1*, *Plekhb1*  
180 and *Mfap3l* were detected to be very significant by SpaGene ( $\text{adj}p < e^{-21}$ ), which were all specifically  
181 expressed in ONL (Supplemental Fig. S11). However, none of them were found by SPARK-X, while *Atf5*,  
182 *Plekhb1* and *Mfap3l* were undiscovered by SpatialDE (Supplemental Fig. S13). Another gene, *Grb2* was  
183 identified by SPARK-X but missed by SpatialDE, showing a very clear GCL pattern (Supplemental Fig.  
184 S13). Then we ranked genes by SPARK-X and checked the top six inconsistent ones (Supplemental Fig.  
185 S14). *Camk2a*, *Psd3*, *Meis2*, *Calm2*, *Arf3* and *Stxbp1* ranked high by SPARK-X, which displayed strong  
186 GCL patterns. All were identified by SpaGene but none by SpatialDE, indicating SpatialDE had limited  
187 power in identifying GCL-specific genes (Supplemental Fig. S14). Finally, we ranked genes by  
188 SpatialDE and examined the top six inconsistent ones (Supplemental Fig. S15). *Spem1*, *Siglec1*, and *Il12a*  
189 only expressed in one or two spots, which were likely to be false signals. *Cck*, *Kif5b*, and *ApoE* exhibited  
190 GL or ONL patterns, which were identified by SpaGene but missed by SPARK-X (Supplemental Fig.  
191 S15). These comparisons demonstrated that SpaGene successfully identified genes with visually distinct

192 patterns, while SPARK-X and SpatialDE missed some genes in certain layers even they showed distinct  
193 patterns.

194

### 195 **Application to mouse preoptic hypothalamus by MERFISH**

196 We applied SpaGene to mouse preoptic hypothalamus data by MERFISH (Moffitt et al. 2018), consisting  
197 of 161 genes measured on 5,665 cells . The 161 genes include 156 pre-selected markers of distinct cell  
198 populations and five blank control genes. Cell clustering based on transcriptional profiles alone identified  
199 multiple cell types, most of which were spatially localized in specific regions, such as mature  
200 oligodendrocyte (OD), ependymal, mural and some inhibitory and excitatory neuron cell types (Fig. 3A).  
201 SpaGene identified those markers from region-specific cell types as top variable genes. Some  
202 representative genes were shown in Fig. 3B, such as *Ntng1* in inhibitory neurons (adjp=5e-108), *Mbp* in  
203 mature OD (adjp=0), *Cd24a* in Ependymal (adjp=0), *Adcyap1* in excitatory neurons (adjp=0), and *Myh11*  
204 in Mural cells (adjp=4e-24).

205 Comparing SpaGene with SPARK-X and SpatialDE, we found their results were highly correlated in  
206 terms of significance (R=0.92 between SpaGene and SpatialDE, R=0.74 between SpaGene and SPARK-  
207 X, and R=0.82 between SPARK-X and SpatialDE) (Fig. 3C). We also compared the number of positive  
208 genes given the number of negative control genes identified (Fig. 3D). The results supported a higher  
209 power of SpaGene. For example, SpaGene detected 149 true positives, while SpatialDE discovered 144  
210 and SPARK-

211 X revealed 128, when one negative control was detected (one false positive). Based on those identified  
212 spatially variable genes, SpaGene successfully reconstructed the underlying spatial organization  
213 (Supplemental Fig. S16).

214

### 215 **Application to mouse cerebellum by Slideseq V2**

216 We applied SpaGene to mouse cerebellum data by Slideseq V2 (Stickels et al. 2021), containing 20,141  
217 genes measured on 11,626 spots. SpaGene identified 619 genes with spatial patterns (adjp<0.05). The

218 cerebellum is made of three layers, molecular, Purkinje and granular layers from outer to inner, and white  
219 matter underneath. SpaGene detected genes, known to be specifically located in three layers and white  
220 matter, to be very significant, such as *Kcnd2* in granular layer (adjp=4e-253) (Varga et al. 2000), *Car8* in  
221 Purkinje layer (adjp=0) (Miterko et al. 2019), *Gad1* in molecular layer (adjp=2e-64) (Kirsch et al. 2012)  
222 and *Mbp* in white matter (adjp=0) (Verity and Campagnoni 1988) (Fig. 4A). Based on those identified  
223 spatially variable genes, SpaGene successfully reconstructed the tightly folded layer structure of  
224 cerebellum. Patterns 1 and 3 corresponded to granular layer, patterns 2, 6 and 8 represented molecular  
225 layer, patterns 4 and 5 stood for Bergmann glia and purkinje neurons in Purkinje layer, and pattern 7  
226 imaged white matter (Supplemental Fig. S17).

227 We compared SpaGene with SPARK-X but not SpatialDE because it would take hours to analyze such  
228 large-scale data. SPARK-X discovered 530 genes, while 230 overlapped with SpaGene (Supplemental  
229 Fig. S18). We examined carefully at those genes detected to be very significant by one method but  
230 insignificant by the other one (Supplemental Fig. S18). Those genes specifically located in Purkinje layer,  
231 such as *Car8*, *Ipr1*, *Pcp2*, and *Pcp4*, were detected as being the most significant by SpaGene (adjp=0)  
232 but undetected by SPARK-X, suggesting SPARK-X had limited power to identify the Purkinje pattern  
233 (Supplemental Fig. S19). In comparison, *Catsperd*, *Ifit3*, and *Ptpst* ranked top by SPARK-X, but  
234 undetected by SpaGene, which didn't seem to have obvious patterns (Supplemental Fig. S20). SpaGene  
235 obtained the significance of *Mog* were just below the cutoff (adjp=0.05), which seemed to be dispersed in  
236 the white matter (Supplemental Fig. S20).

237 Cell clustering based on transcriptional profiles alone found localized cell types, such as molecular layer  
238 neurons, purkinje neurons in the purkinje layer, granule cells in the granule layer (Fig. 4B) . We expected  
239 markers in those spatially-restricted cell types were identified and ranked top by the methods. The  
240 enrichment analysis found that SpaGene obtained high enrichment scores in all three layers, while  
241 SPARK-X got a high score in granular layer, but low scores in other two layers, especially in the Purkinje  
242 layer. This result further demonstrated that SpaGene is more robust to spatial patterns (Fig. 4C).

243 Although there were only 163 common genes between the 619 spatially variable and the top 2000  
244 transcriptionally variable genes, cell clustering derived from these two gene sets were similar  
245 (Supplemental Fig. S21A). Clustering based on the spatially variable genes successfully found those cell  
246 types specifically located in the white matter, molecular, purkinje and granule layers (Supplemental Fig.  
247 S21B). We selected the top 2000 genes by integrating the spatially and transcriptionally variable genes.  
248 Clustering based on the integrative features improved clustering slightly, which showed a higher  
249 percentage of locations expressing cell-type specific marker genes (Supplemental Fig. S21C). The results  
250 suggested that spatially variable genes can serve as a complement to transcriptionally variable genes.

251

### 252 **Application to MOB by HDST**

253 We applied SpaGene to olfactory bulb from high-definition spatial transcriptomics (HDST) (Vickovic et  
254 al. 2019), involving 19,950 genes measured on 181,367 spots. HDST is extremely sparse, where only 21  
255 spots have more than 50 genes detected. In this case, SpaGene used an adaptive strategy to expand the  
256 neighborhood search for genes with high sparsity. SpaGene identified 249 genes as being spatially  
257 variable. The most significant genes included *Ptgds* (adjp=1e-232), *Gphn* (adjp=3e-114) and *Camk1d*  
258 (adjp=3e-61). Although spatial patterns of those genes were not visually distinct due to high sparsity of  
259 the HDST data (Supplemental Fig. S22), there were vague patterns showing *Ptgds* localized in ONL,  
260 *Gphn* in MCL and EPL, and *Camk1d* in GCL (Fig. 5A). Those specific localizations have been reported  
261 before (Rees et al. 2003; Perera et al. 2020) and validated by in situ hybridization in the Allen Brain Atlas  
262 (Fig. 5B).

263 We compared SpaGene with SPARK-X but not SpatialDE because it would take months to analyze such  
264 large-scale data. SPARK-X detected 133 genes, which overlapped significantly with SpaGene (90 in  
265 common). Among the 40 genes most associated with each MOB layer (top five genes in eight patterns in  
266 Supplemental Fig. S11), SpaGene found 12 genes (*Ptgds*, *Fabp7*, *Gad1*, *Vtn*, *Kctd12*, *Kif5b*, *Apod*, *Pcp4*,  
267 *Gpsm1*, *Slc1a2*, *Nrgn*, and *Map1b*), while SPARK-X only detected six (*Ptgds*, *Fabp7*, *Kctd12*, *Kif5b*,  
268 *Apod*, and *Pcp4*).

269

**270 Identification of spatially colocalized ligand-receptor pairs**

271 We extended SpaGene to identify cell-cell communications mediated by colocalized ligand and receptor  
272 pairs. SpaGene found 35 ligand-receptor interactions from the MOB data by spatial transcriptomics. The  
273 two most significant ligand-receptor pairs were IGFBP5-CAV1 (adjp=3e-31) and APOE-LRP6 (adjp=2e-  
274 18), both happening between ONL and GL. *ApoE* is known to be enriched in ONL and GL and also  
275 identified to be very significant by SpaGene (adjp=1e-50). Most spots with high *ApoE* expression were  
276 surrounded with spots with high *Lrp6* expression (Fig. 6A), suggesting potential interactions between  
277 them. However, a number of spots with high *Lrp6* expression were not adjacent to those with high *ApoE*  
278 expression, indicating other ligands might colocalize with *Lrp6* as well. APOE-LRP6 mediates Wnt  
279 signaling, which is important for the regulation of synaptic integrity and cognition (Zhao et al. 2018). The  
280 identification of APOE-LRP6 between ONL and GL layers might be suggestive of the potential  
281 regulation of Wnt signaling in the establishment of periphery–CNS olfactory connections.

282 SpaGene found 13 ligand-receptor interactions from the mouse cerebellum data by Slideseq V2. The most  
283 significant pair was PSAP-GPR37L1 (adjp=1e-27) (Fig. 6B). *Gpr37l1* was known to be strongly  
284 expressed in Purkinje layer and also identified by SpaGene (adjp=8e-130). *Psap*, in contrast, was not as  
285 specifically localized as *Gpr37l1* (adjp=6e-8). PSAP-GPR37L1 protects neural cells from cellular damage  
286 (Li et al. 2017). The identification of PSAP-GPR37L1 between Purkinje layer and surrounding layers  
287 further supports its important role in brain function. Additionally, PTN-PTPRZ1, identified as the only  
288 interaction by MERINGUE (Miller et al. 2021), ranked the top four by SpaGene (adjp=2e-7).

289

**290 DISCUSSION**

291 Recent advances in spatial omics technologies increase the demand for scalable and robust methods to  
292 characterize spatially variable patterns. Here, we developed SpaGene, a fast and model-free method to  
293 identify spatially variable genes. SpaGene has been extensively evaluated on seven datasets generated  
294 from a variety of spatial technologies, ranging from low to high throughput and spatial resolution.

295 Additional analyses on breast cancer from spatial transcriptomics, mouse brain from 10X Visium, and  
296 olfactory bulb from Slide-seqV2 were shown in Supplemental Figs. S23-S33. The results consistently  
297 demonstrated that SpaGene successfully identified known spatially variable genes and also markers in  
298 spatially-restricted cell clusters. Simple factor analysis on those identified genes reconstructed underlying  
299 tissue structures, further demonstrating the ability of SpaGene to characterizing spatial patterns.

300 SpaGene builds upon a simple intuition that spatially variable genes show uneven spatial distributions. As  
301 a model-free and distribution-free method, SpaGene is more robust to pattern shapes, data distribution and  
302 sparsity, non-uniform cellular densities, and the number of spatial locations than existing approaches. The  
303 power of SpatialDE, SPARK and SPARK-X highly depend on spatial covariance models, that is, how  
304 well those predefined kernel functions match the true underlying spatial patterns. Moreover, SpatialDE  
305 and SPARK use parametric modeling based on the assumption of spatial data following Gaussian or  
306 Poisson distributions. Therefore, their performance would be compromised significantly for those genes  
307 whose expression misalign the model defined by those kernel functions and whose distribution violate  
308 Gaussian or Poisson distributions. SpaGene, in contrast, is a model-free and distribution-free method.  
309 Without any assumption, SpaGene is able to identify any spatial patterns and applied on any spatial omics  
310 data, such as identification of spatially localized clones and histone markers in spatial genomics and  
311 epigenomics data. The significance from SpaGene reflects the distinctness of spatial patterns rather than  
312 the extent of match to the defined model. SpaGene uses neighborhood graphs to represent spatial  
313 connections, making it more robust to non-uniform cellular densities common in tissues. Furthermore,  
314 SpaGene is highly computationally efficient in terms of runtime and memory requirement. It only took  
315 SpaGene seconds to minutes to analyze large-scale spatial transcriptomics data (Supplemental Fig. S10C),  
316 which required hours, days or even months for most methods (Zhu et al. 2021).

317 SpaGene uses equal weights by default. Its power can be further improved if we adjust the weight  
318 parameter to assign unequal weights to different degrees (Supplemental Figs. S34A and S34B). Since  
319 clustered connections are more informative than scattered ones in defining spatial patterns, putting more  
320 weights on higher degrees strengthen the ability of SpaGene to distinguish visually distinct patterns from

321 vague ones. For example, *Nppa* displayed a more distinct expression pattern than *Smim36* (*Gm45716*).  
322 *Nppa* is locally expressed in a specific region, whereas *Smim36* is expressed everywhere. SpaGene with  
323 unequal weights successfully ranked *Nppa* much more statistically significant (adjp=e-35) to be spatially  
324 variable than *Smim36* (adjp=e-8). SpaGene with equal weights, however, ranked the opposite  
325 (Supplemental Fig. S34C). Another example on *Wfdc2* and *Zfp235* was given in Supplemental Fig. S34D.  
326 In general, the performance of SpaGene is insensitive to the parameter  $k$  to build the nearest neighbor  
327 graph. The results were highly correlated across four different  $k$ -values (4, 8, 24, and 48) on three large-  
328 scale spatial transcriptomics data (Supplemental Fig. S35). For very sparse data, SpaGene provides an  
329 option to tune  $k$ -values automatically based on the expression sparsity of each gene. Moreover, SpaGene  
330 can incorporate the cell type information to find spatially variable genes within the same cell type. For  
331 example, SpaGene identified *Aldoc* as the most spatially variable genes within the Purkinje layer  
332 (adjp=4e-90) (the function `SpaGene_CT` was provided in the package), which has been demonstrated to  
333 show a regional enrichment pattern that was consistent with the known paths of parasagittal stripes across  
334 individual lobules (Kozareva et al. 2021). Furthermore, SpaGene was easily extended to find colocalized  
335 gene pairs. It successfully identified *Psap-Gpr37L1* and *Ptn-Ptprz1* in mouse cerebellum, and *Fnl1-Cd44*  
336 in invasive breast cancer regions (Supplemental Fig. S33). The default neighborhood search regions could  
337 be further adjusted to identify those long-distance interactions. Finally, potential extensions of SpaGene  
338 to find common and specific spatial patterns across multiple samples would further expands its  
339 application. SpaGene provides two functions `FindPattern_Multi` and `PlotPattern_Multi` to detect and  
340 visualize common and different patterns across samples. An example on two mouse brain datasets from  
341 anterior and posterior regions was given in the GitHub.

342 Although SpaGene is powerful in characterizing localized and co-localized patterns, it has some  
343 limitations. SpaGene binarizes gene expression into high and low, which increases the speed but loses the  
344 quantitative information of expression abundances. The binarization might underpower its performance  
345 on the identification of patterns with a gradient. SpaGene is able to identify long-distance interactions

346 with a large  $k$ -value. However, it lacks the ability of modeling the diffusivity properties of ligands and  
 347 receptors and their activity range.

348

## 349 **METHODS**

### 350 **Identification of spatially variable genes**

351 Spatially variable genes show uneven spatial distribution of expression, where cells/spots with high  
 352 expression are more likely to be spatially connected than random. SpaGene constructs the  $k$ -nearest  
 353 neighbor graph based on spatial locations. For each gene, SpaGene extracts a subnetwork comprising  
 354 only cells/spots with high expression of the gene from the  $k$ -nearest neighbor graph. SpaGene quantifies  
 355 the connectivity of the subnetwork using the Earth's mover's distance between degree distributions of the  
 356 subnetwork and a fully connected one. The degree distribution is more powerful and flexible than the  
 357 total number of connections (Ren et al. 2020) to define spatial connectivity. The reason is that sparsely  
 358 scattered connections are less informative and important than clustered ones in defining spatial patterns.  
 359 For example, it is hard to shape a spatial pattern from a number of scattered connections. The utilization  
 360 of degree distribution allows to assign different weights to different degrees rather than treating them  
 361 equally.

362 Earth mover's distance ( $EMD^g$ ) quantifies the distance from the observed degree distribution of the  
 363 subnetwork of the gene  $g$  to a distribution from a fully connected network (Equation 1). Therefore, shorter  
 364 EMD distances indicate higher spatial connectivity. The degree distribution  $p_i^g$  is defined to be the  
 365 fraction of cells/spots with degree of  $i$  in the subnetwork for the gene  $g$ ,  $w_i$  is the weight assigned to the  
 366 degree of  $i$ , and  $k$  is the number of nearest neighbors to build the spatial network. Since clustered  
 367 connections are more important than scattered ones in defining spatial patterns, at least equal or more  
 368 weights should be assigned to higher degrees, that is,  $w_i \leq w_j$ , if  $i \leq j$ .  $EMD$  with equal weights  
 369 ( $w_i = 1, i = 0, 1, \dots, 2 * k$ ) is reduced to the average number of non-connections.

$$370 \quad EMD^g = \sum_{i=0}^{2*k} w_i p_i^g (2 * k - i) \quad (1)$$

371 To generate the null distribution of  $EMD$ , the same number of cells/spots is randomly sampled and the  
 372 spatial connection of those cells/spots is quantified as  $EMD'$ . The mean and the standard deviation of  
 373  $EMD'$  are estimated after random permutations (default: 500). The observed  $EMD$  is compared to the null  
 374 distribution of  $EMD'$  to evaluate its significance. The Benjamini-Hochberg procedure is used to adjust p-  
 375 values for FDR control.

$$p(x < EMD^g) = p(z < \frac{EMD^g - mean(EMD')}{Sd(EMD')})$$

376

### 377 **Identification of spatial patterns**

378 Non-negative matrix factorization is applied on spatially variable genes detected by SpaGene to identify  
 379 distinct spatial patterns. NMF is implemented by the RcppML R package (DeBruine et al. 2021). It is  
 380 challenging to choose the optimal number of NMF factors. Although several approaches have been  
 381 proposed (Brunet et al. 2004; Frigyesi and Hoglund 2008; Hutchins et al. 2008), the computation is very  
 382 lengthy and results from different approaches are inconsistent. Therefore, selecting the number of ranks  
 383 based on the prior knowledge of the tissue structure is recommended. For example, the number of ranks  
 384 of eight to 12 is recommended for ST MOB data with a roughly arrangement of seven layers. The  
 385 Spearman's correlation between expression of spatially variable genes and cells/spots factor matrix from  
 386 NMF is used to find the most representative genes in each pattern.

387

### 388 **Adaptive strategy to tune neighborhood search regions**

389 SpaGene uses an adaptive strategy to expand neighborhood search regions in very sparse datasets, where  
 390 a single  $k$ -value to build the nearest neighbor graph will not work well for all genes. To improve  
 391 sensitivity, SpaGene increases the  $k$ -value for genes with high sparsity. SpaGene groups genes into  
 392 different bins ( $b_j, j = 1, 2, \dots, J$ ) based on the number of cells/spots with detected expression, where  
 393 different bins  $b_j$  correspond to different  $k$ -values. In this way, SpaGene chooses the  $k$ -value automatically  
 394 based on the sparsity level of the gene.

395  $J = \text{round}(\log_2(n_{max}/n_{min}))+1$

396  $b_1 = (+\infty, n_{max}]$ ,  $b_j = [n_{max} * 2^{-(j-1)}, n_{max} * 2^{-(j-2)})$ ,  $j=2,3\dots J$

$k_{j+1} = k_j + 8 * j$ ,  $k_1 = 8$

397 where  $J$  is the number of bins, determined by the maximum and the minimum number of cells/spots with  
 398 detected expression *that* users set ( $n_{max}$  and  $n_{min}$ ).  $b_j$  is the bin  $j$  that one gene is assigned to by the  
 399 number of cells/spots with the gene expression detected and  $k_j$  is the corresponding  $k$ -value for the bin  $j$ .  
 400 For example, if one gene has the number of cells/spots with detected expression greater than  $n_{max}$ , this  
 401 gene is grouped into  $b_1$  with  $k_1=8$ .

402

#### 403 **Identification of ligand-receptor interactions**

404 SpaGene is extended to identify ligand-receptor interactions. For each ligand-receptor pair, SpaGene  
 405 estimates the spatial connectivity of the subnetwork comprising only connections between cells/spots  
 406 with both high expression of the ligand and the receptor. SpaGene uses the Earth's mover's distance  
 407 based on the degree distribution of the subnetwork to quantify its spatial connectivity.

408

#### 409 **Enrichment analysis of cell type-specific marker genes**

410 Cell clustering based on transcriptional profiles alone discovers cell types localized in specific spatial  
 411 regions. Therefore, marker genes in those spatially-restricted cell types should be identified as spatially  
 412 variable genes. The gene set is built from the top markers based on the fold change between the  
 413 expression in the cell type compared to others. Top 20 are selected for ST MOB, while top 50 are chosen  
 414 for other datasets. The results from SpaGene, SpatialDE and SPARK-X are ranked from the most to the  
 415 least significant. Unweighted gene set enrichment analysis (Subramanian et al. 2005) is implemented to  
 416 evaluate the enrichment of the gene set in the high ranking of pre-ranked gene lists of SpaGene,  
 417 SpatialDE and SPARK-X.

418

## 419 **Simulation designs**

420 We followed simulation designs of SPARK-X and Trendsceek. Briefly we generated two datasets with  
421 five spatial expression patterns, local hotspot, streak, circularity, bi-quarter circularity and mouse purkinje  
422 layer. For the first four patterns, spatial locations of cells were generated by a random-point-pattern  
423 Poisson process. The spatial locations of the pattern of mouse purkinje layer was obtained from Slideseq  
424 V2 mouse cerebellum data. The expression values were either generated from negative binomial  
425 distributions following SPARK-X or bootstrap-sampled from spatial transcriptomics MOB data  
426 following Trendsceek. Simulation datasets varied on a number of parameters: 1) the number of genes  
427 varied from 1000, 3000, and 10,000, among of which 500 genes are spatially variable; 2) the number of  
428 cells varied from 300, 1000, 2000 and 5000 except for the purkinje layer pattern; 3) the fold change of  
429 expression in the spatial region compared to those in the background; For the negative binomial  
430 distribution, the fold change varied from 2, 3,5, 8 to 10. For the resampled real dataset, the expression of  
431 spiked cells were generated from 65%, 70%, 80% to 90% quantile of the expression distribution; 4) the  
432 number of spiked cells except for the purkinje layer pattern. For the hotspot and the streak patterns, the  
433 percentage of spiked cells varied from 5%, 10%, 20% to 30%. For the circularity and bi-quarter  
434 circularity patterns, the width of circularity varied between 0.05, 0.075, 0.1, 0.125 and 0.15.

435

## 436 **Spatial transcriptomics datasets**

437 SpaGene was applied on seven spatial transcriptomics datasets, covering a variety of platforms with low  
438 and high throughput and spatial resolution. Two spatial transcriptomics data from mouse olfactory bulb  
439 and human breast cancer contained genome-wide expression profiles on only hundreds of spots (low  
440 spatial resolution) (Stahl et al. 2016). MERFISH on the mouse preoptic region of the hypothalamus  
441 targeted only 160 genes at single cell resolution (Moffitt et al. 2018). 10X Visium on the mouse brain  
442 comprised of whole transcriptomics on thousands of spots with a spatial resolution of 55  $\mu\text{m}$ , which can  
443 be downloaded from 10x Genomics website (<https://support.10xgenomics.com/spatial-gene-expression/datasets>). Two Slideseq V2 from mouse cerebellum and olfactory bulb contained whole

445 transcriptomics on tens of thousands of spots with a spatial resolution of 10  $\mu\text{m}$  (Stickels et al. 2021).  
446 HDST from mouse olfactory bulb measured whole transcriptomics on hundreds of thousands of spots  
447 with a spatial resolution of 2 $\mu\text{m}$  (Vickovic et al. 2019).

448

#### 449 **Software availability**

450 SpaGene, an R package (R Core Team 2021), is freely available at the GitHub repository  
451 <https://github.com/liuqivandy/SpaGene> . Source codes and seven transcriptomics data are also available  
452 as Supplemental Code. Vignettes on seven spatial transcriptomics data with raw data, codes and results,  
453 including spatial variable genes identification, pattern identification and visualization, co-localized  
454 ligand-receptor pairs identification and visualization, are also available at the GitHub.

455

#### 456 **COMPETING INTEREST STATEMENT**

457 The authors declare no competing interests.

458

#### 459 **ACKNOWLEDGMENTS**

460 This work is supported by National Cancer Institute grants (U2C CA233291 and U54 CA217450),  
461 National Institutes of Health (P01 AI139449), and Cancer Center Support Grant (P30CA068485).

462

463

#### 464 **FIGURE LEGENDS**

465

466 **Figure 1. Schematic of SpaGene and simulation results.** A) Schematic of SpaGene; B) Visualization of  
467 five spatial patterns; C) AUC plots of SpaGene (red), SpatialDE (gray) and SPARK-X (blue) in simulated  
468 datasets with different effect sizes (x axis) and pattern sizes (point shapes) and 10,000 genes and 1,000  
469 cells/locations. Simulated data were generated from negative binomial distributions.

470 **Figure 2. Application of SpaGene to spatial transcriptomics of main olfactory bulb data (MOB).** A)  
471 Visualization of five known spatially variable genes located in specific MOB layers (high expression in  
472 red, and low in blue), with adjusted p-values from SpaGene; B) Enrichment scores of markers in location-  
473 restricted cell types by SpaGene, SpatialDE and SPARK-X.

474 **Figure 3. Application of SpaGene to MERFISH of mouse preoptic hypothalamus data.** A) Cell  
475 clustering based on transcriptional profiles alone; B) Visualization of five spatial variable genes (high

476 expression in red and low in blue) with adjusted p-values from SpaGene; C) Pairwise correlation of  
 477 results from SpaGene, SpatialDE and SPARK-X; D) Power plot shows the number of genes with spatial  
 478 expression pattern (y axis) identified by SpaGene, SpatialDE and SPARK-X versus the number of blank  
 479 control genes identified at the same threshold.

480  
 481 **Figure 4. Application of SpaGene to Slideseq V2 of mouse cerebellum data.** A) Visualization of four  
 482 known spatially variable genes located in specific cerebellum layers (high expression in red, and low in  
 483 blue), with adjusted p-values from SpaGene; B) Cell clustering based on transcriptional profiles alone; C)  
 484 Enrichment scores of markers in location-restricted cell types by SpaGene and SPARK-X.

485 **Figure 5. Application of SpaGene to HDST of MOB data.** Visualization of three spatially variable  
 486 genes. A) gene-expression levels from HDST (high in red, low in blue), with adjusted p-values from  
 487 SpaGene; B) in situ hybridization results for the three genes obtained from the Allen Brain Atlas.

488 **Figure 6. Extension of SpaGene to identify ligand-receptor interactions.** A) Visualization of IGFBP5-  
 489 CAV1 and APOE-LRP6 interactions for ST MOB data, with adjusted p-values from SpaGene. B)  
 490 Visualization of the PSAP-GPR37L1 interaction for Slideseq V2 mouse cerebellum data, with the  
 491 adjusted p-value from SpaGene. Left is the relative expression of the ligand and the receptor, right is the  
 492 interaction strength.

493

494

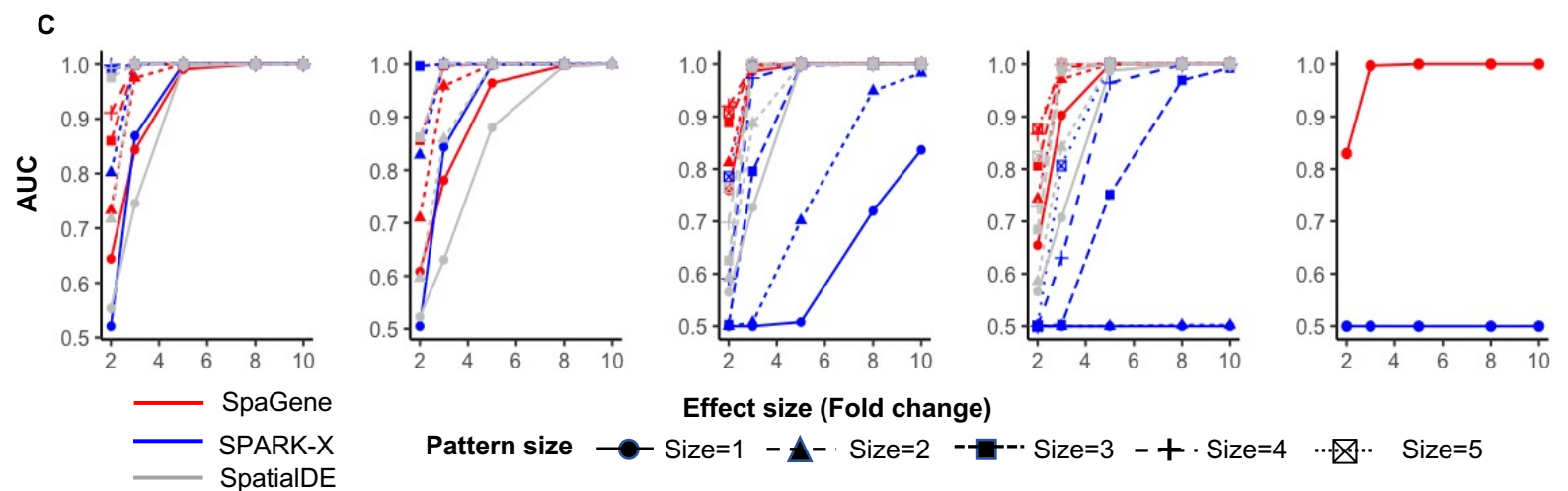
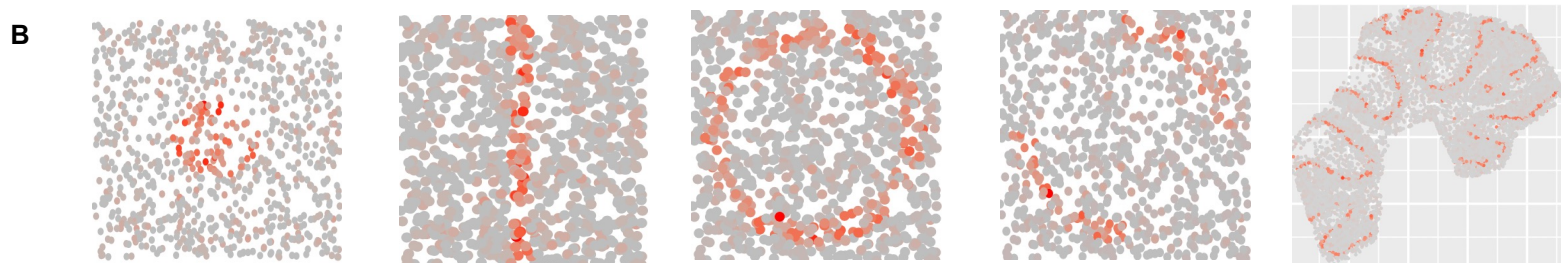
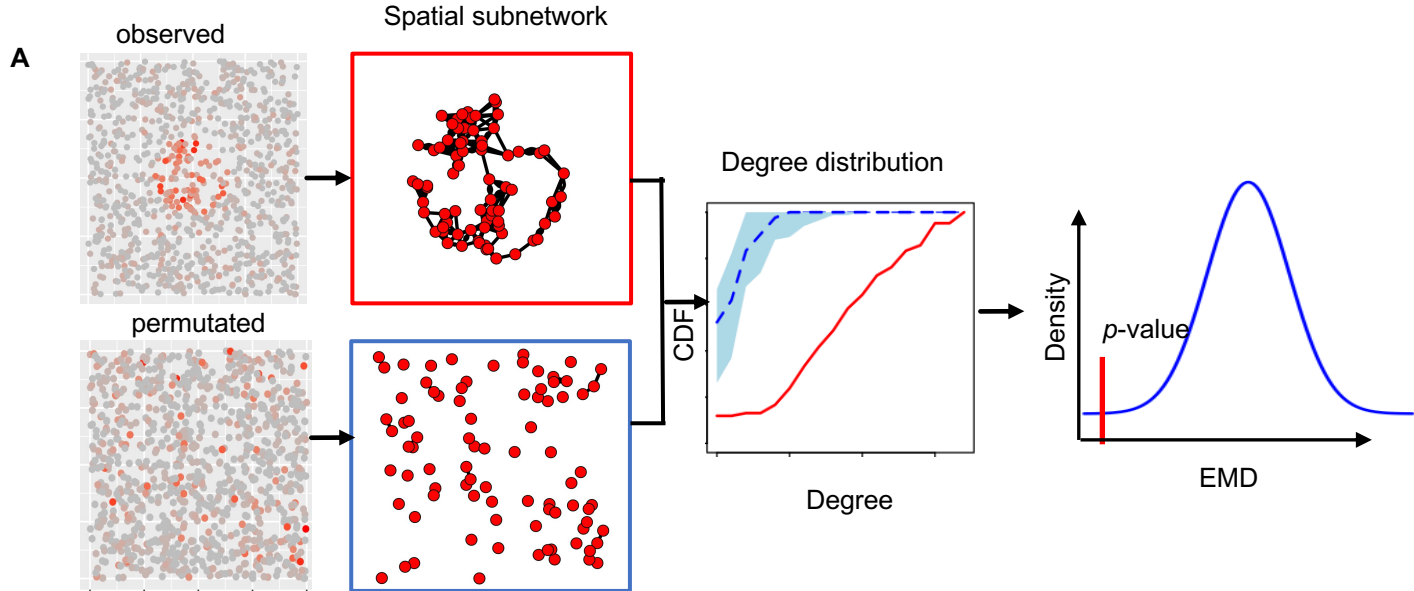
## 495 REFERENCES

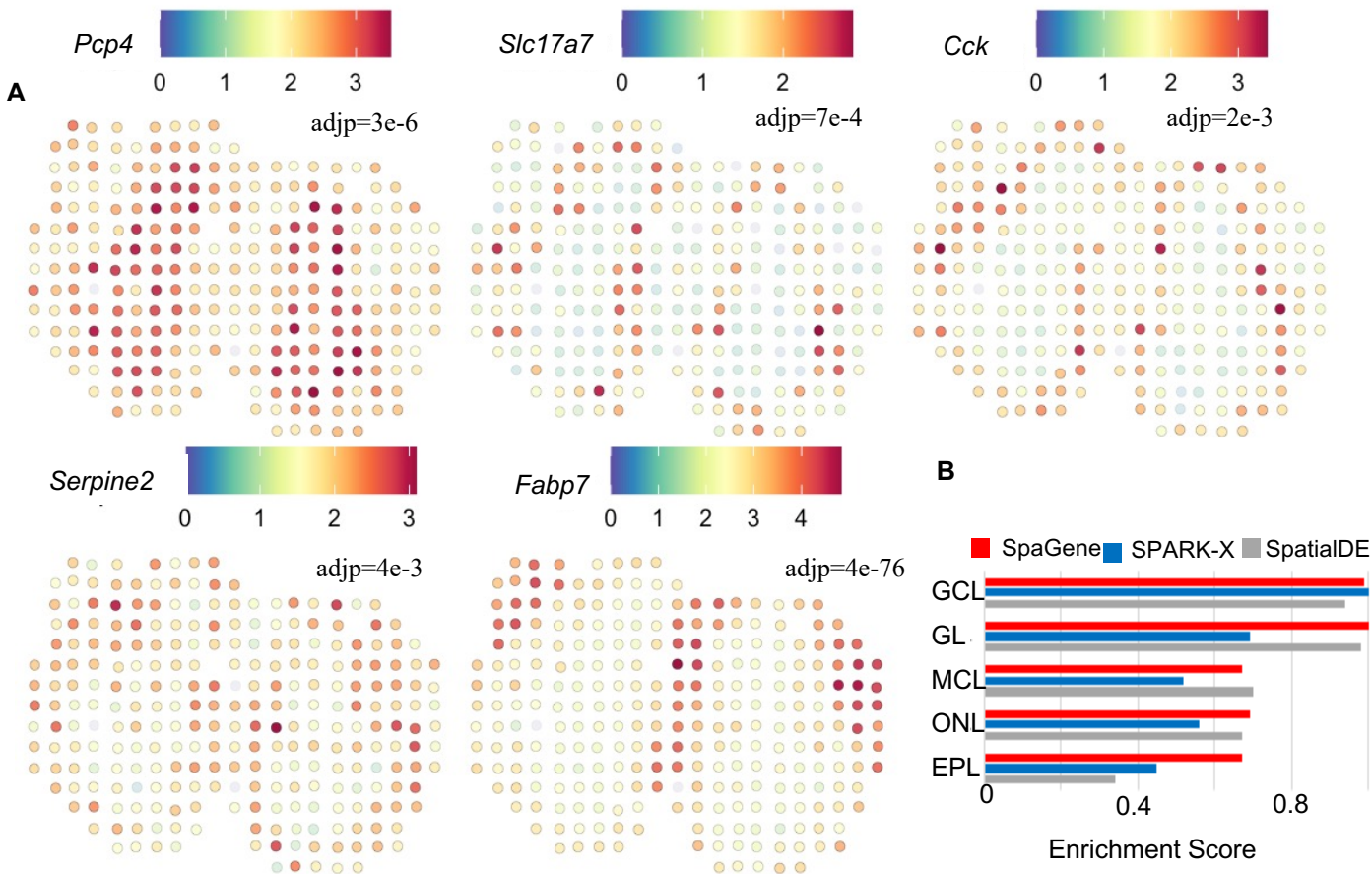
- 496 Anderson A, Lundeberg J. 2021. sepal: Identifying Transcript Profiles with Spatial Patterns by  
 497 Diffusion-based Modeling. *Bioinformatics* doi:10.1093/bioinformatics/btab164.  
 498 Atta L, Fan J. 2021. Computational challenges and opportunities in spatially resolved  
 499 transcriptomic data analysis. *Nat Commun* **12**: 5283.  
 500 Brunet JP, Tamayo P, Golub TR, Mesirov JP. 2004. Metagenes and molecular pattern discovery  
 501 using matrix factorization. *Proc Natl Acad Sci U S A* **101**: 4164-4169.  
 502 DeBruine ZJ, Melcher K, Triche TJ. 2021. Fast and robust non-negative matrix factorization for  
 503 single-cell experiments. *Biorxiv* doi:<https://doi.org/10.1101/2021.09.01.458620>.  
 504 Deng Y, Bartosovic M, Kukanja P, Zhang D, Liu Y, Su G, Enniful A, Bai Z, Castelo-Branco G, Fan R.  
 505 2022. Spatial-CUT&Tag: Spatially resolved chromatin modification profiling at the  
 506 cellular level. *Science* **375**: 681-686.  
 507 Dhainaut M, Rose SA, Akturk G, Wroblewska A, Nielsen SR, Park ES, Backup M, Roudko V, Pia L,  
 508 Sweeney R et al. 2022. Spatial CRISPR genomics identifies regulators of the tumor  
 509 microenvironment. *Cell* doi:10.1016/j.cell.2022.02.015.  
 510 Edsgard D, Johnsson P, Sandberg R. 2018. Identification of spatial expression trends in single-  
 511 cell gene expression data. *Nat Methods* **15**: 339-342.  
 512 Frigyesi A, Hognlund M. 2008. Non-negative matrix factorization for the analysis of complex gene  
 513 expression data: identification of clinically relevant tumor subtypes. *Cancer Inform* **6**:  
 514 275-292.

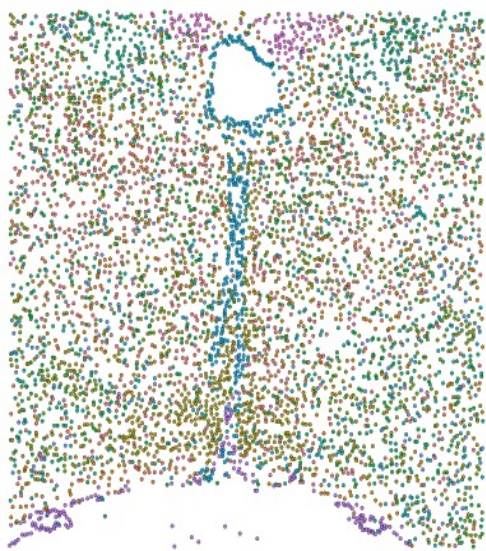
- 515 Hutchins LN, Murphy SM, Singh P, Graber JH. 2008. Position-dependent motif characterization  
516 using non-negative matrix factorization. *Bioinformatics* **24**: 2684-2690.
- 517 Kirsch L, Liscovitch N, Chechik G. 2012. Localizing genes to cerebellar layers by classifying ISH  
518 images. *PLoS Comput Biol* **8**: e1002790.
- 519 Kozareva V, Martin C, Osorno T, Rudolph S, Guo C, Vanderburg C, Nadaf N, Regev A, Regehr WG,  
520 Macosko E. 2021. A transcriptomic atlas of mouse cerebellar cortex comprehensively  
521 defines cell types. *Nature* **598**: 214-219.
- 522 Larsson L, Frisen J, Lundeberg J. 2021. Spatially resolved transcriptomics adds a new dimension  
523 to genomics. *Nat Methods* **18**: 15-18.
- 524 Lewis SM, Asselin-Labat ML, Nguyen Q, Berthelet J, Tan X, Wimmer VC, Merino D, Rogers KL,  
525 Naik SH. 2021. Spatial omics and multiplexed imaging to explore cancer biology. *Nat*  
526 *Methods* **18**: 997-1012.
- 527 Li J, Wang C, Zhang Z, Wen Y, An L, Liang Q, Xu Z, Wei S, Li W, Guo T et al. 2018. Transcription  
528 Factors Sp8 and Sp9 Coordinately Regulate Olfactory Bulb Interneuron Development.  
529 *Cereb Cortex* **28**: 3278-3294.
- 530 Li X, Nabeka H, Saito S, Shimokawa T, Khan MSI, Yamamiya K, Shan F, Gao H, Li C, Matsuda S.  
531 2017. Expression of prosaposin and its receptors in the rat cerebellum after kainic acid  
532 injection. *IBRO Rep* **2**: 31-40.
- 533 Longo SK, Guo MG, Ji AL, Khavari PA. 2021. Integrating single-cell and spatial transcriptomics to  
534 elucidate intercellular tissue dynamics. *Nat Rev Genet* **22**: 627-644.
- 535 Mansuy IM, van der Putten H, Schmid P, Meins M, Botteri FM, Monard D. 1993. Variable and  
536 multiple expression of Protease Nexin-1 during mouse organogenesis and nervous  
537 system development. *Development* **119**: 1119-1134.
- 538 Marx V. 2021. Method of the Year: spatially resolved transcriptomics. *Nat Methods* **18**: 9-14.
- 539 Miller BF, Bambah-Mukku D, Dulac C, Zhuang X, Fan J. 2021. Characterizing spatial gene  
540 expression heterogeneity in spatially resolved single-cell transcriptomic data with  
541 nonuniform cellular densities. *Genome Res* **31**: 1843-1855.
- 542 Miterko LN, White JJ, Lin T, Brown AM, O'Donovan KJ, Sillitoe RV. 2019. Persistent motor  
543 dysfunction despite homeostatic rescue of cerebellar morphogenesis in the Car8  
544 waddles mutant mouse. *Neural Dev* **14**: 6.
- 545 Moffitt JR, Bambah-Mukku D, Eichhorn SW, Vaughn E, Shekhar K, Perez JD, Rubinstein ND, Hao  
546 J, Regev A, Dulac C et al. 2018. Molecular, spatial, and functional single-cell profiling of  
547 the hypothalamic preoptic region. *Science* **362**.
- 548 Nagayama S, Homma R, Imamura F. 2014. Neuronal organization of olfactory bulb circuits.  
549 *Front Neural Circuits* **8**: 98.
- 550 Perera SN, Williams RM, Lyne R, Stubbs O, Buehler DP, Sauka-Spengler T, Noda M, Micklem G,  
551 Southard-Smith EM, Baker CVH. 2020. Insights into olfactory ensheathing cell  
552 development from a laser-microdissection and transcriptome-profiling approach. *Glia*  
553 **68**: 2550-2584.
- 554 R Core Team. 2021. R: A Language and Environment for Statistical Computing. In *R Foundation*  
555 *for Statistical Computing*, Vienna, Austria.
- 556 Ratz M, von Berlin L, Larsson L, Martin M, Westholm JO, La Manno G, Lundeberg J, Frisen J.  
557 2022. Clonal relations in the mouse brain revealed by single-cell and spatial  
558 transcriptomics. *Nat Neurosci* **25**: 285-294.

- 559 Rees MI, Harvey K, Ward H, White JH, Evans L, Duguid IC, Hsu CC, Coleman SL, Miller J, Baer K et  
560 al. 2003. Isoform heterogeneity of the human gephyrin gene (GPHN), binding domains  
561 to the glycine receptor, and mutation analysis in hyperekplexia. *J Biol Chem* **278**: 24688-  
562 24696.
- 563 Ren X, Zhong G, Zhang Q, Zhang L, Sun Y, Zhang Z. 2020. Reconstruction of cell spatial  
564 organization from single-cell RNA sequencing data based on ligand-receptor mediated  
565 self-assembly. *Cell Res* **30**: 763-778.
- 566 Sangameswaran L, Hempstead J, Morgan JI. 1989. Molecular cloning of a neuron-specific  
567 transcript and its regulation during normal and aberrant cerebellar development. *Proc*  
568 *Natl Acad Sci U S A* **86**: 5651-5655.
- 569 Stahl PL, Salmen F, Vickovic S, Lundmark A, Navarro JF, Magnusson J, Giacomello S, Asp M,  
570 Westholm JO, Huss M et al. 2016. Visualization and analysis of gene expression in tissue  
571 sections by spatial transcriptomics. *Science* **353**: 78-82.
- 572 Stickels RR, Murray E, Kumar P, Li J, Marshall JL, Di Bella DJ, Arlotta P, Macosko EZ, Chen F. 2021.  
573 Highly sensitive spatial transcriptomics at near-cellular resolution with Slide-seqV2. *Nat*  
574 *Biotechnol* **39**: 313-319.
- 575 Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A,  
576 Pomeroy SL, Golub TR, Lander ES et al. 2005. Gene set enrichment analysis: a  
577 knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl*  
578 *Acad Sci U S A* **102**: 15545-15550.
- 579 Sun S, Zhu J, Zhou X. 2020a. Statistical analysis of spatial expression patterns for spatially  
580 resolved transcriptomic studies. *Nat Methods* **17**: 193-200.
- 581 Sun X, Liu X, Starr ER, Liu S. 2020b. CCKergic Tufted Cells Differentially Drive Two Anatomically  
582 Segregated Inhibitory Circuits in the Mouse Olfactory Bulb. *J Neurosci* **40**: 6189-6206.
- 583 Svensson V, Teichmann SA, Stegle O. 2018. SpatialDE: identification of spatially variable genes.  
584 *Nat Methods* **15**: 343-346.
- 585 Varga AW, Anderson AE, Adams JP, Vogel H, Sweatt JD. 2000. Input-specific immunolocalization  
586 of differentially phosphorylated Kv4.2 in the mouse brain. *Learn Mem* **7**: 321-332.
- 587 Verity AN, Campagnoni AT. 1988. Regional expression of myelin protein genes in the developing  
588 mouse brain: in situ hybridization studies. *J Neurosci Res* **21**: 238-248.
- 589 Vickovic S, Eraslan G, Salmen F, Klughammer J, Stenbeck L, Schapiro D, Aijo T, Bonneau R,  
590 Bergenstrahle L, Navarro JF et al. 2019. High-definition spatial transcriptomics for in situ  
591 tissue profiling. *Nat Methods* **16**: 987-990.
- 592 Young JK, Heinbockel T, Gondre-Lewis MC. 2013. Astrocyte fatty acid binding protein-7 is a  
593 marker for neurogenic niches in the rat hippocampus. *Hippocampus* **23**: 1476-1483.
- 594 Zhang L, Hernandez VS, Gerfen CR, Jiang SZ, Zavala L, Barrio RA, Eiden LE. 2021. Behavioral role  
595 of PACAP signaling reflects its selective distribution in glutamatergic and GABAergic  
596 neuronal subpopulations. *Elife* **10**.
- 597 Zhao N, Liu CC, Qiao W, Bu G. 2018. Apolipoprotein E, Receptors, and Modulation of  
598 Alzheimer's Disease. *Biol Psychiatry* **83**: 347-357.
- 599 Zhao T, Chiang ZD, Morriss JW, LaFave LM, Murray EM, Del Priore I, Meli K, Lareau CA, Nadaf  
600 NM, Li J et al. 2022. Spatial genomics enables multi-modal study of clonal heterogeneity  
601 in tissues. *Nature* **601**: 85-91.

- 602 Zhu J, Sun S, Zhou X. 2021. SPARK-X: non-parametric modeling enables scalable and robust  
603 detection of spatial expression patterns for large spatial transcriptomic studies. *Genome*  
604 *Biol* **22**: 184.
- 605 Zhuang X. 2021. Spatially resolved single-cell genomics and transcriptomics by imaging. *Nat*  
606 *Methods* **18**: 18-22.
- 607

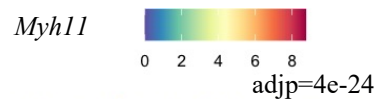
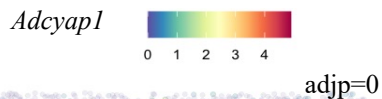
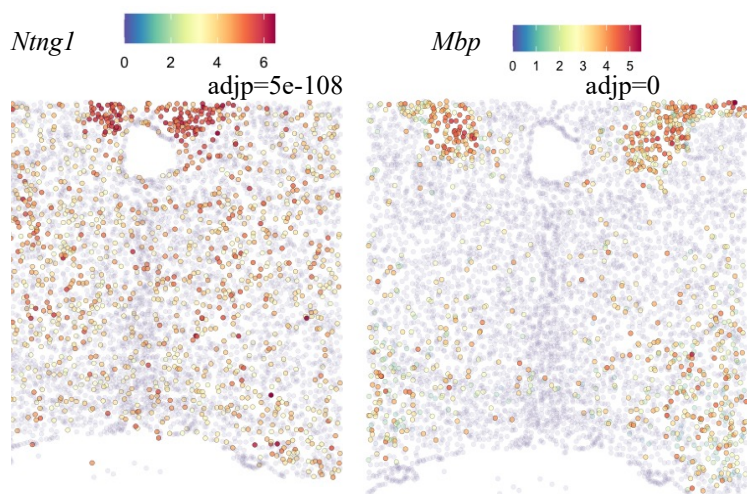
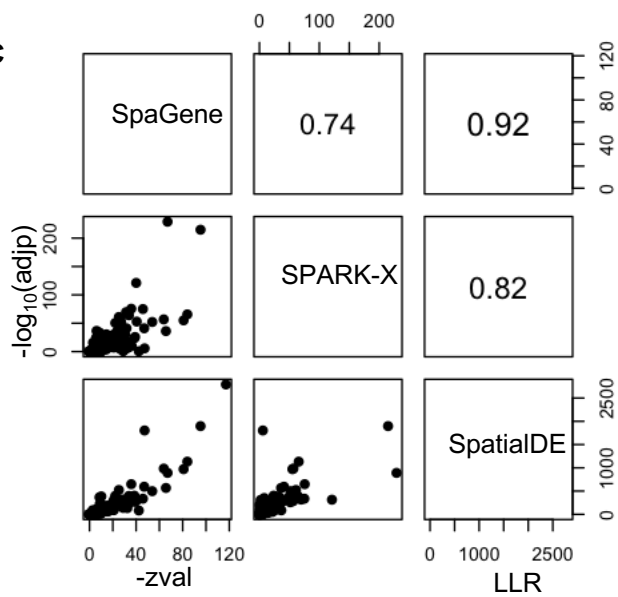
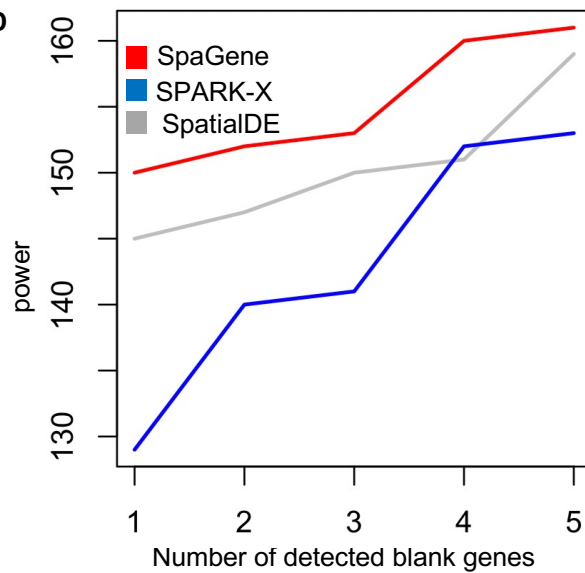


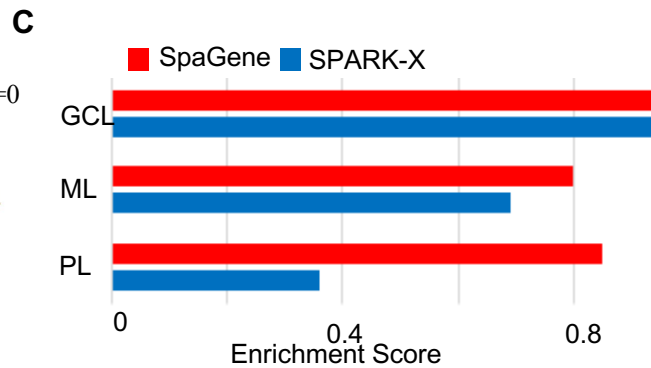
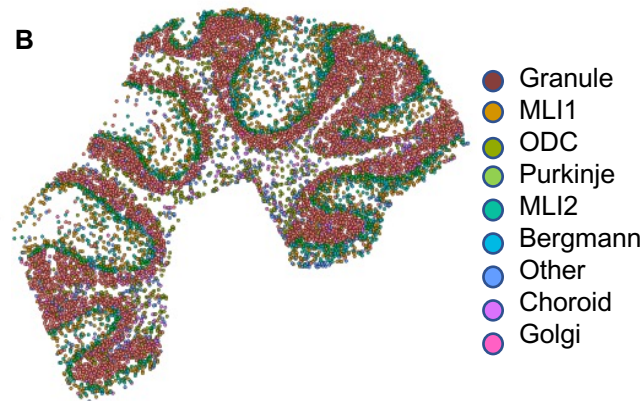
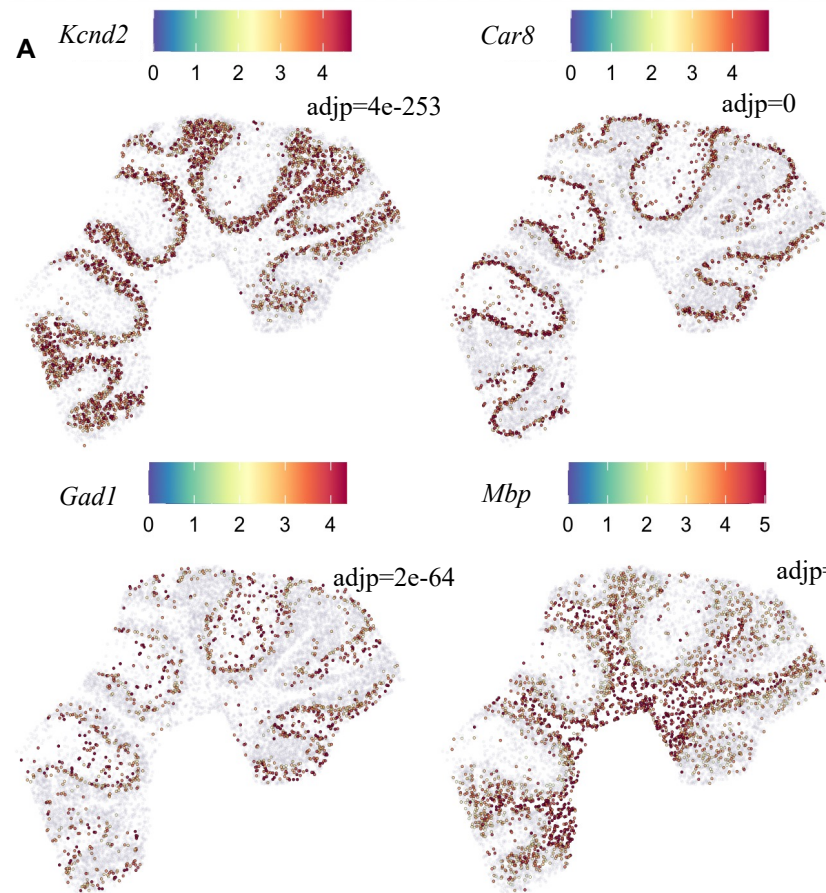


**A**

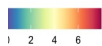
Cell type

- Inhibitory 1
- Astrocyte
- Excitatory 1
- Inhibitory 2
- Inhibitory 3
- OD
- Endothelial
- Ependymal
- Ambiguous
- Mural
- Excitatory 2
- OD immature

**B****C****D**



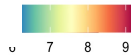
*Ptgds*



*Gphn*

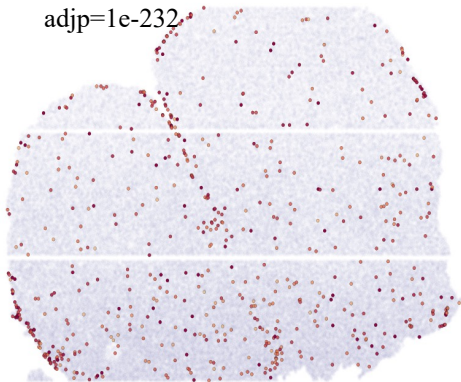


*Camk1d*

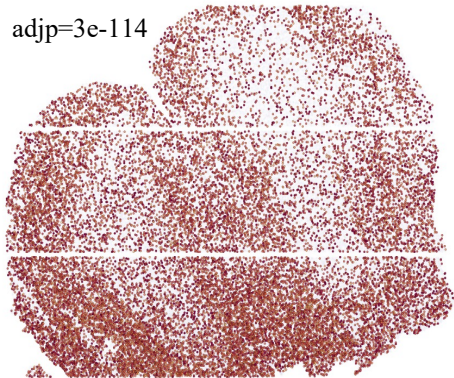


**A**

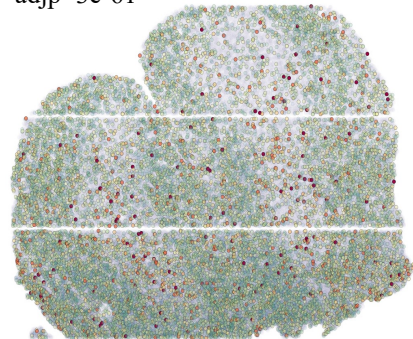
adjp=1e-232



adjp=3e-114



adjp=3e-61



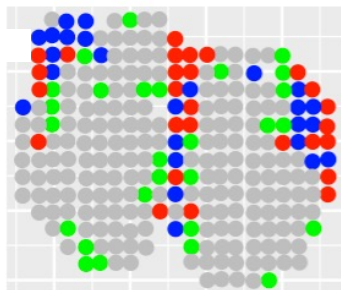
**B**



IGFBP5-CAV1

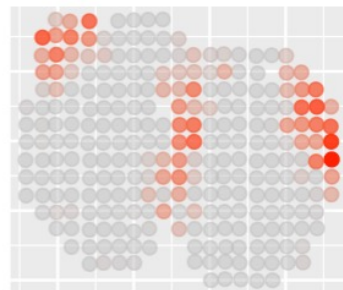
adjp=3e-31

A



Exp

- Both low
- Ligand high
- Receptor High
- Both High



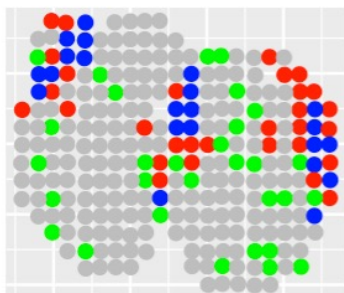
LR

5  
4  
3  
2  
1  
0

APOE-LRP6

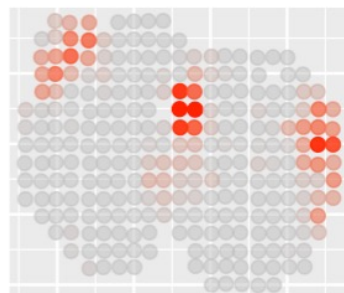
adjp=2e-18

B



Exp

- Both low
- Ligand high
- Receptor High
- Both High



LR

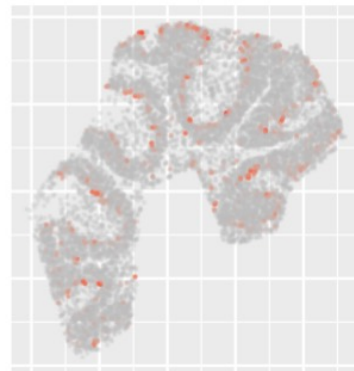
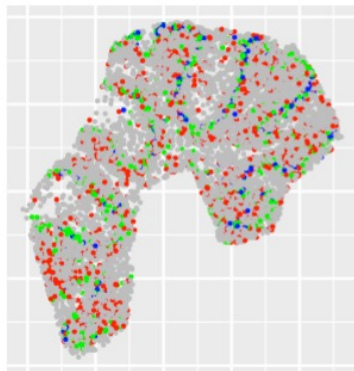
4  
3  
2  
1  
0

PSAP-GPR37L1

adjp=1e-27

Exp

- Both low
- Ligand high
- Receptor High
- Both High



LR

15  
10  
5  
0