

1 ELLA: Modeling Subcellular Spatial Variation of Gene Expression 2 within Cells in High-Resolution Spatial Transcriptomics

3

4

5 Jade Xiaoqing Wang^{1,2}, and Xiang Zhou^{1,2,#}

6

7 1. Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA

8 2. Center for Statistical Genetics, University of Michigan, Ann Arbor, MI 48109, USA

9 #: correspondence to XZ (xzhousph@umich.edu)

11 **Abstract**

12 Spatial transcriptomic technologies are becoming increasingly high-resolution, enabling precise
13 measurement of gene expression at the subcellular level. Here, we introduce a computational
14 method called subcellular expression localization analysis (ELLA), for modeling the subcellular
15 localization of mRNAs and detecting genes that display spatial variation within cells in high-
16 resolution spatial transcriptomics. ELLA creates a unified cellular coordinate system to anchor
17 diverse cell shapes and morphologies, utilizes a nonhomogeneous Poisson process to model spatial
18 count data, leverages an expression gradient function to characterize subcellular expression
19 patterns, and produces effective control of type I error and high statistical power. We illustrate the
20 benefits of ELLA through comprehensive simulations and applications to four spatial
21 transcriptomics datasets from distinct technologies, where ELLA not only identifies genes with
22 distinct subcellular localization patterns but also associates these patterns with unique mRNA
23 characteristics. Specifically, ELLA shows that genes enriched in the nucleus exhibit an abundance
24 of long noncoding RNAs or protein-coding mRNAs, often characterized by longer gene lengths.
25 Conversely, genes containing signal recognition peptides, encoding ribosomal proteins, or
26 involved in membrane related activities tend to enrich in the cytoplasm or near the cellular
27 membrane. Furthermore, ELLA reveals dynamic subcellular localization patterns during the cell
28 cycle, with certain genes showing decreased nuclear enrichment in the G1 phase while others
29 maintain their enrichment patterns throughout the cell cycle. Overall, ELLA represents a calibrated,
30 powerful, robust, scalable, and versatile tool for modeling subcellular spatial expression variation
31 across diverse high-resolution spatial transcriptomic platforms.

32

33 **Keywords:** spatially resolved transcriptomics, subcellular mRNA localization, spatial variable
34 genes, spatial variation, gene expression, ELLA, nonhomogeneous Poisson process.

35 **Introduction**

36 Spatial transcriptomics is a collection of new genomics technologies designed to measure gene
37 expression within tissues while preserving spatial localization information. Recent technological
38 advancements have substantially improved the spatial resolution of spatial transcriptomics,
39 facilitating expression measurements at cellular and subcellular levels. Specifically, *in situ* RNA-
40 sequencing techniques, such as ISS [1], FISSEQ [2], STARmap [3], and Ex-seq [4], achieve a
41 spatial resolution under 1 μm , which is much smaller than the size of a typical cell. Recent high-
42 throughput sequencing-based techniques, such as Seq-Scope [5], VisiumHD [6], Open-ST [7], and
43 Stereo-seq [8], offer spatial resolutions in the range of 0.5-2 μm . *In situ* imaging techniques, such
44 as MERFISH [9], SeqFISH+ [10], MERSCOPE [11], CosMx [12], and 10X Xenium [13], provide
45 spatial resolutions as fine as 0.1-0.2 μm (Fig. S1). Together, these high-resolution spatial
46 transcriptomics technologies have enabled expression measurement at subcellular resolution,
47 providing unprecedented opportunities to interrogate the intracellular localization and distribution
48 of mRNAs within cells.

49
50 The intracellular localization and distribution of mRNAs are vital for cellular functions. They
51 ensure the targeted delivery of mRNAs and facilitate localized protein synthesis, enabling precise
52 regulation of gene expression within specific subcellular compartments. The spatial localization
53 of mRNAs empowers cells to respond rapidly to local cues and signals, adapting effectively to
54 changing environments and supporting specialized cellular functions [14]. For example, the
55 localization of mRNAs encoding for β -actin at the leading edges of fibroblasts or the lamellipodia
56 of myoblasts ensures localized protein synthesis of actin, supporting proper cell polarity and
57 motility [15]. In addition, the spatial localization of mRNA contributes to cellular organization and
58 differentiation, aiding in the establishment and maintenance of distinct cellular identities and
59 functions, influencing asymmetric cell division and cell fate determination across various
60 organisms. For example, the spatially localized expression of *Oskar* at the posterior end of the
61 embryo is essential for the development and assembly of the germ plasm in *Drosophila*, facilitating
62 germ cell formation [16]. As another example, *Ash* mRNA localizes to the bud tip in *S. cerevisiae*
63 to establish asymmetry of HO endonuclease gene expression, which is important for mating type
64 switching [17]. Given the importance of proper mRNA spatial localization, their misplacement
65 often leads to detrimental effects and has been associated with multiple diseases [16]. For example,
66 disruptions in axonal mRNA transport and localization contribute to neurodegeneration in
67 Huntington's disease [18]. Therefore, understanding the spatial localization and distribution of
68 mRNA within cells is crucial for unraveling the complexity of cellular structure and function, as
69 well as for elucidating the cellular mechanisms underlying disease etiology.

70
71 Despite its importance, however, characterizing the subcellular spatial organization of mRNAs in
72 high-resolution spatial transcriptomics turns out to be a computationally challenging task. Only
73 two methods have been developed for this purpose, each with its own limitations. Specifically,
74 Bento [19] employs pre-trained random forest classifiers to categorize each gene into five pre-
75 defined subcellular RNA localization patterns, while SPRAWL [20] relies on four metrics to
76 identify four pre-specified subcellular patterns. However, both methods are limited to imaging-
77 based spatial transcriptomics data, failing to leverage the vast amount of high-resolution spatial
78 transcriptomics obtained from recent sequencing-based technologies. Additionally, they are
79 constrained to detect genes with pre-defined localization patterns, thus limiting the discovery of
80 any new spatial localization patterns and suffering from low statistical power. Besides these major

81 limitations, Bento requires nuclear boundary information, which may not be readily available in
82 some spatial transcriptomics datasets. In addition, Bento is only applicable to analyzing a single
83 cell and lacks the ability to borrow the spatial localization pattern shared across multiple cells.
84 Conversely, SPRAWL is only applicable to analyzing multiple cells, not a single cell, and is unable
85 to distinguish between enrichment and depletion in the pre-specified localization patterns due to
86 the nature of its two-sided tests.
87

88 To address the above limitations, we present subcellular Expression LocaLization Analysis
89 (ELLA), a computational method for modeling the subcellular localization of mRNAs and
90 detecting genes that display spatial variation within cells across a range of high-resolution spatial
91 transcriptomics technologies. ELLA comes with three unique features. First, it creates a unified
92 cellular coordinate system, which allows for anchoring diverse cell shapes and morphologies
93 regardless of the spatial transcriptomics technology, thus enabling the joint modeling of cells with
94 distinct shapes. Second, it constructs a novel nonhomogeneous Poisson process model to directly
95 and explicitly model the spatial occurrence of expression measurements within cells, thus
96 facilitating powerful and effective modeling of both spatial count data from sequencing-based
97 techniques and spatial binary data from imaging-based techniques. Finally, it devises an expression
98 intensity function to model the subcellular spatial distribution of mRNAs along the cellular radius,
99 which, when paired with a range of kernel functions, is capable of capturing a wide variety of
100 expression distribution patterns within cells without the need to restrict to pre-defined localization
101 patterns. As a result, ELLA can be applied to an arbitrary number of cells and detect a wide variety
102 of subcellular localization patterns across diverse spatial transcriptomic techniques, all with
103 effective type I error control and high statistical power. With a computationally efficient algorithm,
104 ELLA is also scalable to tens of thousands of genes across tens of thousands of cells.
105

106 We illustrate the benefits of ELLA through comprehensive simulations and applications to four
107 spatial transcriptomics datasets. In the real data applications, ELLA not only identifies genes with
108 distinct subcellular localization patterns but also reveals that these patterns are associated with
109 unique mRNA characteristics. Specifically, genes enriched in the nucleus show an abundance of
110 long noncoding RNAs (lncRNAs) and protein-coding mRNAs, often characterized by longer gene
111 lengths. Conversely, genes containing signal recognition peptides, encoding ribosomal proteins,
112 or involved in membrane related activities such as synaptic transmission and G protein coupled
113 receptor activities, tend to enrich in the cytoplasm or near the cellular membrane. Moreover, genes
114 exhibit dynamic subcellular localization during the cell cycle, with some showing decreased
115 nuclear enrichment in the G1 phase, while others maintain their patterns of enrichment regardless
116 of cell cycle phases.
117
118

119 Results

120 Method overview

121 ELLA is described in [Methods](#), with its technical details provided in [Supplementary Notes](#) and
122 method schematic displayed in [Fig. 1a](#). Briefly, ELLA is a statistical method for modeling the
123 subcellular localization of mRNAs and detecting spatially variable genes with subcellular spatial
124 expression patterns in high-resolution spatial transcriptomics ([Fig. S1](#)). ELLA examines one gene
125 at a time, creates a unified cellular coordinate system through defining a cellular radius in each
126 cell that points from the center of the nucleus towards the cellular boundary, relies on a
127 nonhomogeneous Poisson process (NHPP) to capture the spatial distribution of expression
128 measurements within cells, devises an expression intensity function and computes a P value to
129 capture any subcellular expression patterns observed along the cellular radius. ELLA is capable of
130 borrowing information across cells through a joint likelihood framework to substantially improve
131 detection power, while taking advantage of multiple intensity kernel functions to capture the
132 distinct subcellular expression patterns that may be encountered in various biological settings to
133 ensure robust performance. In addition, ELLA relies on a fast binning algorithm for approximate
134 position computation and utilizes Adam optimization for scalable inference. As a result, ELLA is
135 computationally efficient and is easily scalable to tens of thousands of genes measured in tens of
136 thousands of cells. ELLA is implemented in Python, freely accessible from
137 <https://xiangzhou.github.io/software/>.

138

139 Simulations

140 We performed comprehensive simulations on imaging-based spatial transcriptomics to evaluate
141 the performance of ELLA and compared it with three methods. The three methods include
142 SPRAWL [20], Bento [19], and Wilcox, where Wilcox denotes a modified Wilcoxon rank sum
143 test [21] that uses expression measurements normalized by the area of subcellular regions to
144 examine the difference in expression between nuclear and cytoplasmic areas. All methods examine
145 one gene at a time and all methods except Bento produce a P value for each gene; Bento outputs
146 five prediction probabilities for five pre-specified cellular localization patterns, which cannot be
147 converted to a P value. Among these methods, ELLA can analyze either one or multiple cells;
148 SPRAWL and Wilcox can only analyze multiple cells; and Bento can only analyze one cell.
149 Therefore, we compared ELLA with SPRAWL and Wilcox in all our main simulations on multiple
150 cells while compared ELLA with Bento in additional simulations on only one cell. Unlike ELLA
151 and SPRAWL, both Bento and Wilcox require nuclear boundary information in addition to cell
152 boundary information ([Tab. S1](#)). We provide the actual nuclear boundary information to Bento
153 and Wilcox, although this information may not be readily available in certain sequencing-based
154 techniques such as Seq-Scope [5] and Stereo-seq [8] and may not be accurately inferred in other
155 techniques.

156

157 Simulation details are provided in [Methods](#). Briefly, we sampled n different embryonic fibroblast
158 cells from seqFISH+ data ([Fig. S2](#)) and simulated expression counts for 1,000 genes to be spatially
159 distributed within these cells. We examined type I error control of different methods in null
160 simulations, where the simulated gene expression counts are randomly distributed spatially within
161 each cell without any specific subcellular spatial expression patterns ([Fig. 1b, S3](#)). We also
162 examined the power of different methods in alternative simulations, where the simulated gene
163 expression counts are enriched in specific subcellular regions within the cells, exhibiting either
164 symmetric (consisting of eleven distinct symmetric patterns; [Fig. 1c, S4-5](#)) or asymmetric patterns

165 (three distinct asymmetric patterns, [Fig. 1d, S6](#)). In the simulations, we first created a baseline
166 setting and then varied the number of cells (n), the gene expression level (m), and in the alternative
167 settings, the strength of the subcellular expression patterns (s ; [Methods](#)), one at a time on top of
168 the baseline setting, to create additional settings. In total, we examined 13 null and 40 alternative
169 settings, with 1,000 replicates per setting.
170

171 In the null simulations, the P values from ELLA are well calibrated across settings, and so are the
172 P values from SPRAWL, although SPRAWL failed to produce P values for the radial and punctate
173 metrics in $m=1$ (one count per cell) settings ([Fig 1b, S7-8](#)). The inability of SPRAWL to produce
174 P values in these settings arises from its inability to make use of cells with less than two counts of
175 the gene, which constitutes a large fraction of gene-cell combinations in the real data (e.g. 57.8%
176 in the MERFISH data) [22]. Wilcox yielded inflated P values, especially in settings where the gene
177 expression level is low, or where the number of cells is large ([Fig 1b, S7-8](#)). The P value inflation
178 observed in Wilcox suggests that the simple normalization procedure and the non-parametric
179 Wilcoxon test are not sufficient to control for variance heterogeneity and subsequently type I error
180 ([Fig. S9, Tab. S2](#)).
181

182 In the alternative simulations, because some methods failed to control for type I error, we evaluated
183 power based on a fixed false-discovery rate (FDR) to ensure a fair comparison across methods
184 ([Methods](#)). We first examined the eleven subcellular expression patterns in the symmetric pattern
185 category, including two patterns with nucleus enrichment, two patterns with nuclear edge
186 enrichment, five patterns with cytoplasmic enrichment, and two patterns with membrane
187 enrichment. Based on an FDR threshold of 0.05, ELLA achieves consistently higher power
188 (average=0.63, range=0.41-0.80) than the other methods (SPRAWL: average=0.04, range=0.00-
189 0.09; Wilcox: average=0.04, range=0.00-0.15) in detecting each of the eleven patterns ([Fig. 1c](#)).
190 For SPRAWL, its radial and punctate metrics tend to exhibit very low power in detecting any of
191 the patterns (average=0.01, range=0.00-0.03), presumably because these metrics are not well
192 suited for detecting symmetric patterns. The peripheral and central metrics of SPRAWL have low
193 power for detecting the cytoplasmic enrichment patterns (average=0.01, range=0.00-0.04) but
194 have slightly higher powers for detecting the membrane and nuclear enrichment patterns
195 (average=0.05, range=0.01-0.09), as one might expect. Also as expected, the power of ELLA,
196 SPRAWL, and Wilcox all improves with increasing number of cells, increasing expression level,
197 and increasing pattern strength across all eleven patterns, although the power of ELLA improves
198 much faster compared to the other two methods ([Fig. S10](#)). For example, at an FDR of 0.05, the
199 power of ELLA in detecting the first nucleus pattern is 0.01 with 10 cells but increases to 1.00
200 with 300 cells, while the power SPRAWL's central metric only increases from 0.01 to 0.42 and
201 the power of SPRAWL's peripheral metric only increases from 0.00 to 0.66. The exceptions are
202 Wilcox and SPRAWL's radial metric, whose power for detecting nucleus patterns remains below
203 0.05 and barely improves as the number of cells increases.
204

205 ELLA is also more powerful than the other methods in detecting two of the three asymmetric
206 subcellular expression patterns. These include the radial-cyto and punctate-cyto patterns, where
207 gene expression is enriched in either a circular sector or a small subcellular disc in the cytoplasm
208 ([Fig. 1d](#)). Specifically, for the radial-cyto pattern, ELLA achieved a power of 0.55 while Wilcox
209 achieved a power of 0.00. For SPRAWL, its peripheral, central, radial, and punctate metrics
210 achieved a power of 0.10, 0.00, 0.16, and 0.23, respectively. For the punctate-cyto pattern, ELLA

211 achieved a power of 0.65 while Wilcox had zero power. For SPRAWL, its peripheral, central,
212 radial, and punctate metrics achieved a power of 0.17, 0.00, 0.16, and 0.39, respectively ([Fig. 1d](#)).
213 Certainly, because ELLA models expression patterns along the cellular radius, it is not powered
214 to detect radial-unif asymmetric pattern, where gene expression is enriched in a circular sector of
215 the cell completely uniformly ([Fig. 1d](#)), a scenario unlikely in practical biological applications.
216

217 Importantly, ELLA not only achieves high power in detecting genes with various subcellular
218 expression patterns but also accurately estimates these patterns ([Fig. S11-12](#)). Specifically, the
219 average KL-divergences achieved by ELLA for estimating the two pattern categories are 0.12 and
220 0.29, respectively ([Tab. S3](#)). To further summarize the observed subcellular pattern, ELLA
221 computes a subcellular pattern score for each gene. This score represents the relative position of
222 subcellular expression enrichment, with zero indicating enrichment in the cell nucleus and one
223 indicating enrichment on the cell membrane; [Methods](#)). The majority of the pattern scores (77%)
224 are within 0.1 of the truth across the three pattern categories, underscoring the accuracy of ELLA
225 ([Fig. S13-14](#), [Tab. S4](#)).
226

227 We performed additional simulations with only one cell in order to compare ELLA with Bento
228 ([Fig. 1e](#)). Bento is capable of detecting five pre-specified patterns including enrichment in nucleus,
229 nuclear edge, cytoplasm, cell boundary, and none. To favor the comparison towards Bento, we
230 focused on comparing ELLA with Bento under five symmetric patterns that Bento specifically
231 models, where gene expression is enriched in nucleus (including 2 patterns), nuclear edge (1),
232 cytoplasm (1), or cellular boundary (1) under a relatively high expression level ($m=30$) and a high
233 pattern strength ($s=9$) ([Fig. S15](#)). Because Bento cannot produce P values, we used the prediction
234 probabilities output from Bento to rank genes, with which we measured powers based on FDR
235 ([Methods](#)). We are able to compute FDR for Bento in simulations only because we know the truth,
236 which is certainly unknown for any real data applications. In the simulations, ELLA achieves high
237 power ([Fig. 1e](#), average=0.86, range=0.43-0.99) and accuracy ([Fig. S16-17](#), [Tab. S5](#)) across all
238 five patterns, consistently outperforming Bento (average=0.10, range=0.00-0.75).
239

240 Seq-Scope mouse liver data

241 We applied ELLA to analyze four published datasets obtained using different high-resolution
242 spatial transcriptomics technologies ([Methods](#)). The four datasets include a liver data by Seq-
243 Scope [5], an embryo data by Stereo-seq [8], an NIH/3T3 embryonic fibroblast cell line data by
244 seqFISH+ [10], and a brain data by MERFISH [22].
245

246 We first analyzed the Seq-Scope mouse liver data ([Fig. 2a, S18-30](#)), which contains 497 to 1,349
247 genes measured on 870 cells from four cell types, with 82 to 276 cells per cell type ([Fig. S31-32](#)).
248 The four cell types include periportal hepatocyte (PP; $n=276$) and pericentral hepatocyte (PC;
249 $n=276$) in normal mice, and PP ($n=236$) and PC ($n=82$) cells in early-onset liver failure mice (TD,
250 [23]; [Fig. S33](#)). We were only able to apply ELLA to the data as SPRAWL and Bento are not
251 applicable to sequencing-based data and the nuclear boundary information required for Wilcox
252 and Bento was not available.
253

254 At an FDR of 5%, ELLA identified 84, 123, 98, and 40 genes that display subcellular expression
255 patterns in normal PP, PC cells and TD PP, PC cells, respectively. 77 of these genes, including
256 one transcription factor (*Mlxpl*), were detected in two or more cell types. Based on their

257 subcellular spatial expression patterns, we clustered the detected genes into five distinct pattern
258 clusters ([Fig. 2b, Method](#)): 101 genes (29%) display a nuclear expression pattern (clusters 1), 34
259 (10%) genes display a nuclear edge expression pattern (cluster 2), and 210 genes (61%) display
260 one of the three cytoplasmic expression patterns near the cellular membrane (cluster 3-5). Example
261 cells from the five clusters are shown in [Fig. 2c](#).

262 The detected genes from ELLA allow us to comprehensively investigate the properties of genes
263 that display distinct enrichment patterns within cells. For nuclear genes (clusters 1) with
264 subcellular enrichment near the nuclear center, we found them to have significantly higher snRNA
265 expression in a similar cell type from a separate study [24] (clusters 1 vs clusters 2-5 fold
266 enrichment=15.31, Mann-Whitney U test P value = 5e-22; cluster 1 vs all the remaining genes -
267 consisting of cluster 2-5 genes and nonsignificant genes, fold enrichment=4.80, P value=1e-11;
268 [Fig. 2d](#)) with significantly higher unsplice/splice ratio supporting their nuclear enrichment (cluster
269 1 vs clusters 2-5 fold enrichment=2.52, Mann-Whitney U test P value=8e-21; cluster 1 vs all the
270 other genes fold enrichment=1.08, P value=0.58; [Fig. 2e](#)). For genes with subcellular enrichment
271 in the cytoplasm (cluster 4-5), we found them to frequently encode a signal recognition peptide
272 (SRPs; proportion=40.57%) as compared to the genes in the nuclear cluster 1 (proportion=16.83%;
273 Fisher's exact test P value=2e-5) or the remaining genes (proportion=14.95%; P value=9e-21; [Fig.](#)
274 [2g](#)). SRPs are short sequence segments located at the N-termini of newly synthesized proteins that
275 are translated at the endoplasmic reticulum (ER) and sorted towards the secretory pathway [25].
276 Importantly, we found that nuclear genes (cluster 1) have significantly longer gene lengths
277 compared to genes in the other clusters or the remaining genes, both in terms of the average isoform
278 length (Mann-Whitney U test P value=6e-4 and 0.02), the longest isoform length (P value=1e-6
279 and 4e-3), and the total length across exons (P value=2e-7 and 6e-3; [Fig. 2f](#)). These new findings
280 suggest that long genes may require additional time to be transcribed and exported [26] and their
281 enrichment in the nucleus may serve as a reservoir so that they can be quickly exported to the
282 cytoplasm for translation in response to stimuli [27]. Notably, this novel discovery is consistently
283 observed across the other three datasets we analyzed, as presented in later sections.
284

285 We explored additional biological insights by focusing on the normal PC cell type, which has the
286 largest number of genes with subcellular spatial expression patterns, to carefully examine the 123
287 genes detected by ELLA ([Fig. S19](#)). Among the 34 nuclear (cluster 1) genes ([Fig. S34, Tab. S6](#)),
288 three of them (*Malat1*, *Neat1* and *Gm13775*) are long non-coding RNAs that are previously known
289 to be localized to the nucleus [28]. 30 of them are protein encoding genes including two previously
290 known nuclear-enriched mRNAs *Chd9* and *Ppara*, dovetailing recent findings that retention of
291 mRNAs in the nucleus may help buffer noise in the stochastic mRNA production process [21].
292 Seven of them (*Malat1*, *Neat*, *n-R5-8s1*, *Gm24601*, *Mlxipl*, *Mafb*, and *Echdc2*) were also found
293 among the top 10 nuclear-enriched genes identified in the original Seq-Scope study, which
294 explicitly searched for genes enriched within 10 μ m from the nuclear center [5]. Among the seven
295 genes, four encode transcription factors or proteins with transcription factor activity. For example,
296 *Mlxipl*, one of these genes, is a transcription factor retained in the nuclear speckles in the liver [29].
297 Finally, all nine significant mitochondrial genes were detected as cytoplasmic localized (cluster 3;
298 [Fig. S35a, Tab. S7](#)) and all three significant PC cell type marker genes were detected as
299 cytoplasmic or membrane localized (clusters 3 and 5; [Fig. S35b, Tab. S7](#)).
300

301

302 **Stereo-seq mouse embryo data**

303 Next, we analyzed the Stereo-seq mouse embryo data, focusing on two major cell types localized
304 in the cardiothoracic region on slice E1S3 ([Fig. 3a, S36-37](#)): precursor muscle cells, or myoblasts
305 (596 cells with 2,008 genes); and mature muscle cells, or cardiomyocytes (553 cells with 1,743
306 genes; [Fig. S38](#)). We were only able to apply ELLA to the data as SPRAWL and Bento are not
307 applicable to sequencing-based data, and the nuclear boundary information required for Wilcox
308 and Bento was not available in this data.
309

310 At an FDR of 5%, ELLA identified 264 and 304 genes to be spatially variable within myoblasts
311 and cardiomyocytes, respectively ([Fig. S39](#)). 89 genes were detected in both cell types including
312 12 transcription factors. Based on their subcellular spatial expression patterns, we clustered the
313 detected genes into five distinct clusters ([Methods, Fig. 3b](#)): 56 genes (10%) display a nuclear
314 expression pattern (clusters 1), 346 genes (61%) display one of the two nuclear edge expression
315 patterns (cluster 2-3), and 166 genes (29%) display one of the two cytoplasmic expression patterns
316 (cluster 4-5). Example cells from the five clusters are shown in [Fig. 3c](#).
317

318 The genes detected by ELLA again allow us to comprehensively investigate the properties of genes
319 that display distinct enrichment patterns within cells. For nuclear genes (clusters 1-3) with
320 subcellular enrichment near the nuclear center, we found them to have significantly higher
321 unsplice/splice ratio (clusters 1-3 vs clusters 4-5, fold enrichment=2.42, P value=1e-24; clusters
322 1-3 vs all the remaining genes, fold enrichment=3.36, P value=5e-96; [Fig. 3d](#)), which is also
323 negatively correlated with the expression pattern score (Pearson correlation=-0.464, P value=1e-
324 152). In addition, genes in clusters 1-3 are enriched with transcription factors (proportion=14.43%)
325 as compared to the other clusters (clusters 4-5, proportion=5.42%, Fisher's exact test P value=2e-
326 3) or the remaining genes (proportion=5.20%, P value=2e-10; [Fig. 3f](#)). For the genes with
327 subcellular enrichment in the cytoplasm (clusters 4-5), we found them to contain a significantly
328 higher proportion of ribosomal protein (RP) genes (clusters 4-5, 7.23% vs clusters 1-3, 0%,
329 Fisher's exact test P value=3e-7; cluster 3-4 vs all the remaining genes, 4.38%, P value=0.086; [Fig.](#)
330 [3g, Methods](#)), supporting their localized synthesis. Importantly, nuclear genes (clusters 1-3) also
331 tend to have longer gene lengths compared to genes in the other clusters or the remaining genes,
332 in terms of the average isoform length (P value=2e-12 and 2e-60), the median isoform length (P
333 value=2e-8 and 4e-34), the longest isoform length (P value=1e-16 and 5e-84), and the total length
334 across exons (P value=1e-14 and 1e-87; [Fig. 3e](#)). Finally, in terms of 3'UTR length
335 ([Supplementary Notes 1, Fig. S40](#)), 19 genes display significant variation across five expression
336 pattern clusters ([Fig. S41](#)), 21 genes display significant correlation with expression pattern strength
337 ([Fig. S42](#)), and 18 genes display significant correlation with expression pattern score ([Fig. S43](#)).
338

339 We further investigated the shared and distinct features of the genes detected by ELLA in both
340 myoblasts and cardiomyocytes to reveal additional biological insights ([Fig. S44](#)). Both cell types
341 exhibit a similar proportion of genes across the five expression pattern clusters, with common
342 genes displaying similar estimated expression intensities ([Fig. S45-46](#)). Among the detected genes,
343 12 transcription factors are detected in both cell types (16 unique in myoblasts and 27 unique in
344 cardiomyocytes; [Fig. S47a](#)). These transcription factors are enriched in GO gene sets related to
345 regulation of transcription, development, and various regulatory categories ([Fig. S47b-d](#)). In
346 addition, among the detected genes, 7 long noncoding genes are detected in both cell types (5
347 unique in myoblasts and 0 unique in cardiomyocytes; [Fig. S48](#)), including five (*Xist, Meg3,*

348 *Kcnq1ot1*, *Malat1*, and *Airn*) localized near the nuclear center (cluster 1-3). One mitochondrial
349 gene (*mt-Nd1*) is detected in both cell types (0 unique in myoblasts and one unique in
350 cardiomyocytes; [Fig. S49a](#)) and it displays a cytoplasmic localized pattern (cluster 4; [Fig. S49b](#)).
351 Two cardiomyocyte cell type marker genes (*Acta1* and *Myh3*) are detected as localized within
352 cytoplasmic region (clusters 4) in the cardiomyocyte cell type ([Fig. S50](#)).
353

354 **SeqFish+ mouse embryonic fibroblast data**

355 Next, we analyzed the NIH/3T3 mouse embryonic fibroblast cell line data generated by seqFISH+
356 [10], which contains 2,747 genes measured on 171 embryonic fibroblast cells ([Fig. 4a, S51](#)). We
357 were unable to apply SPRAWL due to its heavy computational burden but were able to apply
358 Bento as this data contains nucleus segmentation information.
359

360 At an FDR of 5%, ELLA identified 2,744 genes to display subcellular spatial expression patterns,
361 with 244 being transcription factors. The subcellular expression patterns of the detected genes can
362 be clustered into five distinct clusters ([Fig. 4b, Methods](#)): 32 genes (1%) display a nuclear
363 expression pattern (cluster 1), 1,073 genes (39%) display one of the two nuclear edge expression
364 patterns (clusters 2-3), and 1,639 genes (60%) display one of the two cytoplasmic expression
365 patterns (clusters 4-5). The identified genes included 57 out of 60 genes with subcellular
366 localization patterns detected through an *ad hoc* procedure in the seqFISH+ original study. The
367 localization categorization of the 57 genes closely aligns with the pattern reported in the original
368 study but with finer details: for example, 20 genes detected as enriched generally in the nuclear
369 and perinuclear regions in the original study were clustered here as either cluster 2 (8 genes) or
370 cluster 3 (12 genes) genes ([Fig. S52](#)). Example cells from the five clusters are shown in [Fig. 4c](#).
371

372 Because Bento is only applicable to individual cells, we randomly selected 20 cells ([Fig. S53](#)) and
373 applied both ELLA and Bento to analyze one cell at a time on 356-1,213 (mean=808) genes with
374 more than 10 counts. Across cells, Bento classified 38.2% genes to one of the four compartmental
375 patterns, 21.5% genes to a pattern called “none”, and the remaining 40.39% genes to either none
376 of these five patterns or multiple patterns ([Fig. S54a](#)). Certainly, Bento is unable to produce P
377 values nor quantifications of statistical significance for any of the genes. ELLA was able to allocate
378 all genes to five identified patterns, with 13.4% genes achieving statistical significance (5% FDR;
379 [Fig. S54b-c](#)). For genes detected by ELLA and classified by Bento to patterns other than none,
380 their expression pattern classifications are largely consistent with each other, although ELLA
381 offers more detailed results ([Fig. S55](#)). For example, 92.6% of the “nuclear” patterned genes
382 detected by Bento were also identified as nuclear genes by ELLA, and these genes were classified
383 by ELLA into two separate clusters (66.0% genes in cluster 1 with nuclear pattern and 26.5% genes
384 in cluster 2 with nuclear edge pattern).
385

386 The genes detected by ELLA allow us to comprehensively investigate the properties of genes that
387 display distinct enrichment patterns within cells. For genes with subcellular enrichment near the
388 nuclear center (clusters 1-3), we found them to have significantly longer gene lengths compared
389 to genes in the other clusters (clusters 4-5) or the remaining genes, in terms of the average isoform
390 length (P value=8e-19 and 1e-18), the median isoform length (P value=6e-14 and 7e-14), the
391 longest isoform length (P value=6e-11 and 7e-11), and the total length across exons (P value=1e-
392 4 and 1e-4; [Fig. 4d](#)). These four types of gene lengths are also significantly negatively correlated
393 with the ELLA pattern scores (Pearson correlation ranges from -0.17 to -0.07; P values range from

394 1e-19 to 2e-4). Genes with enrichment near the nuclear center (clusters 1-3) are also enriched with
395 transcription factors (proportion=12.56%) as compared to the other clusters (clusters 1 and 4-5,
396 proportion=7.73%, P value=3e-4) or the remaining genes (proportion=7.72%, P value=2e-4; Fig.
397 4e).

398
399 ELLA also provides a unique opportunity for us to explore whether cell cycle may influence the
400 subcellular spatial localization of gene expression, as the data is collected from cultured cells that
401 undergo continuous cell division. To do so, we first clustered fibroblast cells into three distinct
402 cell-cycle phases, including G1 (n=36, 21%), S (n=83, 49%), and G2M (n=52, 30%). We then
403 applied ELLA to analyze each cell phase separately and detected 728, 2,368, and 1,726 genes with
404 subcellular spatial expression patterns, respectively (Fig. 4f). We found that genes significant in
405 the G1 phase are less likely to be enriched close to the nuclear center and display larger pattern
406 scores compared to the genes in the S and G2M phases, regardless which cluster the genes belong
407 to (pattern score fold enrichment in G1 vs S and G2M=2.33, 1.31, 1.12, 1.16, and 1.02, for the five
408 clusters, respectively; one side Mann-Whitney U test P value=0.37, 0.01, 1e-4, 4e-67, 1e-3; Fig.
409 4g), suggesting that DNA replication during the S phase enhances nuclear enrichment in S and
410 G2M phases. Among the detected genes, 723 are shared across three cell cycles, including 6 (1%),
411 105 (15%), 123 (17%), 442 (61%), and 47 (7%) genes for each of the five clusters, respectively.
412 Among the shared genes, a subset of genes in clusters 2-5 display decreasing pattern scores through
413 G1, S, and G2M phases, corresponding to increasing enrichment towards the nucleus. In contrast,
414 all six cluster 1 genes retain the nuclear pattern across all cell cycle phases (Fig. 4g, Methods),
415 suggesting that some genes are capable of retaining their nuclear enrichment throughout the cell
416 cycle phases [21].
417

418 **MERFISH mouse brain data**

419 Lastly, we analyzed the adult mouse brain data generated by MERFISH [22] (Fig. 5a, S56). We
420 focused on four major cell types residing in midbrain: excitatory neurons (EX, n=577), inhibitory
421 neurons (IN, n=525), astrocytes (Astr, n=480), and oligodendrocytes (Olig, n=948) with 557-878
422 genes per cell type (Fig. S57). Besides ELLA, we were able to also apply SPRAWL to the data,
423 but unable to apply Wilcox and Bento as the nuclear boundary information required for these two
424 methods were not available in this data.
425

426 At an FDR of 5%, ELLA identified 235, 250, 169, and 147 (total=801, total distinct=502) genes
427 that display subcellular spatial expression patterns in EX, IN, Astr, and Olig cells, respectively
428 (Fig. S57). 227 of these genes, including 36 transcription factors, were detected in two or more
429 cell types. The subcellular spatial expression patterns of the detected genes can be clustered into
430 four distinct pattern clusters (Fig. 5b, Method): 337 genes (42%) display a nuclear expression
431 pattern (cluster 1), 125 (16%) genes display a nuclear edge expression pattern (cluster 2), and 339
432 genes (42%) display one of the two cytoplasmic expression patterns (clusters 3-4). Example cells
433 from the four clusters are shown in Fig. 5c. Compared to the number of genes (801) detected by
434 ELLA, the peripheral, central, radial, and punctate metrics of SPRAWL detected 572, 305, 138,
435 and 238 genes, respectively, with 434 distinct genes in total, the majority of which (345; 79.49%)
436 are overlapped with ELLA (Fig. S58). Note that SPRAWL radial and punctate metrics excluded
437 57.8% of the unqualified gene-cell pairs that have less than two counts of a gene in a cell, which
438 likely leads to their lower power as well as their failure in producing P values for a small percentage
439 of genes across cell types (2.3%, 272 genes).

440
441 The genes detected by ELLA allow us to comprehensively investigate the properties of genes that
442 display distinct enrichment patterns within cells. For genes with subcellular enrichment near the
443 nuclear center (clusters 1-2), we found them to have significantly higher snRNA expression in the
444 same cell types from a separate study (clusters 1-2 vs clusters 3-4, fold enrichment=1.25, P value
445 = 1e-15; cluster 1-2 vs all remaining genes fold enrichment=1.29, P value=2e-31; [Fig. 5d](#); [30]).
446 We also found them to have significantly longer gene lengths compared to genes in the other
447 clusters or the remaining genes, in terms of the average isoform length (P value=0.083 and 0.021),
448 the longest isoform length (P value=6e-5 and 1e-6), and the total length across exons (P value=2e-
449 6 and 2e-9; [Fig. 5e](#)). In addition, the cluster 4 genes contain a lower proportion of transcription
450 factors (proportion=11.04%) as compared to the other clusters (cluster 1-3, proportion=18.08%, P
451 value=0.041) or the remaining genes (proportion=17.16%, P value=0.598; [Fig. 5f](#)). Gene sets
452 enriched with the cluster 1-2 genes are related to various functions including transcription
453 regulation ([Fig. S59](#)), while gene sets enriched with cluster 3-4 genes are particularly related to
454 dendrites and synaptic transmission and signaling ([Fig. 5g-h](#)).
455
456 We investigated the shared and distinct features of the genes detected by ELLA in the two neuronal
457 cell types, excitatory and inhibitory neurons to reveal additional biological insights. Excitatory
458 neurons contain a slightly higher proportion of nuclear localized genes (cluster 1) and a lower
459 proportion of cell membrane localized genes (cluster 4) compared to inhibitory neurons ([Fig. S60](#)).
460 A fraction of the detected genes (Jaccard index=31.8%) are shared between the two neuronal types,
461 with 38, 3, 8, and 11 shared genes detected across clusters 1-4 and with similar estimated
462 expression patterns ([Fig. S61-62](#)). In addition, the majority of the detected transcription factors
463 (126) are shared between the two neuronal types, while 20 are uniquely detected in excitatory
464 neurons and 9 uniquely detected in inhibitory neurons ([Fig. S63](#)). The 126 shared transcription
465 factors are enriched in 112 gene sets related to various transcription regulations and neuron
466 differentiation ([Fig. S64-65](#)). Four out of eight long noncoding genes are detected in both cell types
467 ([Fig. S66a](#)). Three out of four of the common long noncoding genes are localized in the nucleus in
468 both cell types (cluster 1: *A830036E02Rik*, *B020031H02Rik*, and *Dlx6os1*), and one gene (*Rmst*)
469 is localized in nuclear edge in the excitatory neurons (cluster 2; [Fig. S66b](#)). Most cell type marker
470 genes detected by ELLA belong to clusters 3-4 with cytoplasmic or membrane localization patterns
471 except for one gene (*Cux2* being detected as cluster 2, nuclear edge localized, in excitatory neurons;
472 [Fig. S67](#)).
473
474
475

476 **Discussion**

477 We have presented ELLA, a statistical method for modeling and detecting spatially variable genes
478 within cells that display various subcellular spatial expression patterns in high-resolution spatial
479 transcriptomic studies. ELLA models the spatial distribution of gene expression measurements
480 along the cellular radius using a nonhomogeneous Poisson process, leverages multiple kernel
481 functions to detect a variety of subcellular spatial expression patterns, and is capable of analyzing
482 a large number of genes and cells. We have illustrated the benefits of ELLA through simulations
483 and real data applications.

484

485 We have primarily focused on utilizing ELLA to capture the spatial variation of gene expression
486 along the cellular radius within cells, which is inherently one-dimensional and rotation invariant.
487 Detecting rotation-invariant and radial symmetric patterns enables information sharing across
488 multiple cells, thereby enhancing statistical power. In addition, rotation-invariant patterns facilitate
489 results interpretation, as the detected genes can be naturally categorized into cellular compartments,
490 including nucleus, nucleus membrane, and cellular membrane. The framework of ELLA, however,
491 is general and can be extended to two- or three-dimensional cellular space, enabling modeling of
492 2D cellular space with kernels defined on a unit circle or 3D cellular space with kernels defined
493 on a unit ball. Use of different kernels on higher dimensional space may further enhance the power
494 of ELLA. For example, radial kernel functions may be particularly effective in detecting genes
495 with radial patterns in 2D cellular space -- a pattern that, although unlikely to be biological, the
496 one-dimensional version of ELLA is ill equipped to detect, as shown in the simulations. Such
497 extensions, however, necessitate careful consideration, as additional modeling features, such as
498 rotation invariance, may need to be incorporated into the kernel structure to effectively utilize
499 information from multiple cells.

500

501 ELLA leverages nuclear center and cellular boundary information extracted from the spatial
502 transcriptomics data or its accompanying histology image data to register and segment cells
503 through multiple pre-processing steps. These pre-processing steps can vary substantially across
504 different spatial transcriptomics technologies. For example, the accompanying H&E and nucleic
505 acid staining images in Seq-Scope and Stereo-seq need to be registered with the spatial
506 transcriptomics data to obtain the cellular boundary information, while the DAPI images in
507 imaging-based datasets have already aligned with the spatial transcriptomics data without the need
508 for further registration. Similarly, the nucleus center in imaging-based datasets is determined as
509 the geometric center of the nuclear segmentation, while in sequencing-based datasets is determined
510 based on the enrichment of unspliced sequencing read counts. Importantly, ELLA provides
511 accompanying scripts tailored to distinct spatial transcriptomics platforms to streamline these pre-
512 processing steps. In addition to the nuclear center and cellular boundary information, additional
513 data such as nuclear boundary information can also be integrated into ELLA as needed. In such
514 cases, the registration step of ELLA can be extended to register cells based on the nuclear center,
515 nuclear boundary, as well as cellular boundary. Furthermore, the modeling framework of ELLA
516 can be extended to accommodate this additional information, as well as other subcellular
517 compartmentalization information, such as annotations on nucleolus or ER membrane, as
518 technologies continue to improve in the future. Investigating the effectiveness of ELLA in the
519 context of additional feature information represents an important avenue for future research.

520

521 Methods

522 ELLA overview

523 Subcellular resolution spatial transcriptomics and data preprocessing

524 We consider a high-resolution spatial transcriptomics study that collects gene expression
525 measurements at subcellular level. For sequencing-based techniques such as Seq-Scope [5] and
526 Stereo-seq [8], the expression of a gene is measured on a set of pre-specified spatial locations and
527 is represented as the number of read counts mapped to the mRNA transcripts of the gene. For
528 imaging-based techniques such as seqFISH+ [10] and MERFISH [9], the expression of a gene is
529 measured as the presence of a hybridization signal of its mRNA across spatial locations.
530 Regardless of the techniques, we assume that G genes are measured on S spatial locations, where
531 these locations have known two-dimensional x and y spatial coordinates that are recorded during
532 the experiment. For a gene g , its raw expression measurement at each location is represented either
533 as a count (for sequencing-based techniques) or as a binary label (for imaging-based techniques)
534 depending on the spatial transcriptomic technique.

535

536 To facilitate joint modeling across cells, we create a unified cellular coordinate system to anchor
537 diverse cell shapes and morphologies. To do so, for the high-resolution spatial transcriptomics data,
538 we first follow standard data preprocessing procedures to segment the tissue into cells. We cluster
539 these cells into different cell types based on marker gene expression. For each cell in turn, we
540 obtain the center of its nucleus and assign the spatial coordinates to all expression measurement
541 locations within the cell. For each measured location inside the cell, we calculate two distances:
542 its distance to the nuclear center d_1 , and its distance to the cell boundary d_2 in the opposite
543 direction from the nuclear center (Fig. S68a). With these two distances, we further calculate the
544 relative position of the measured location inside the cell as the ratio between the nuclear distance
545 and the summation of the two distances $d_1' = d_1/(d_1 + d_2)$. The relative position ranges between
546 0 and 1 and allows us to create a unified coordinate system across cells, enabling the joint modeling
547 of multiple cells regardless of their sizes and shapes (Fig. S68b). Importantly, we compute the
548 cellular distances for each measured location efficiently using a binning-based numerical
549 approximation approach. Specifically, we first divided each cell from the center of nucleus into
550 100 circular sectors of equal angle measure. In each sector v , we denote r_v as the maximum
551 distance between the center of the nucleus and the cellular boundary in the sector using cell
552 segmentation boundary or mask. For each expression measurement location within the sector, we
553 obtain its distance from the center of the nucleus and normalize it by r_v to obtain its relative
554 position. This binning-based approximation approach speeds up computation through eliminating
555 the requirement of computing the distance of each measurement location to the cell boundary,
556 facilitating parallel computation across cells and sectors.

557

558 ELLA model for detecting genes with subcellular spatial expression patterns

559 With the expression measurements and their relative positions within each cell, we aim to identify
560 spatially variable genes that display subcellular spatial expression patterns along the cellular radius
561 that points from the center of the nucleus towards the cellular boundary. The genes with subcellular
562 spatial expression patterns are often localized in certain cellular compartments such as nucleus,
563 cytoplasm, Golgi apparatus, or cell membrane and may display distinct enrichment associated with
564 such compartmentalization. To identify those genes, we examine one gene at a time and jointly
565 model its expression measurements within n cells that belong to a given cell type. For the i th cell

566 ($i = 1, \dots, n$), we assume that the gene is measured on m_i spatial locations. For the j th measured
567 location ($j = 1, \dots, m_i$), we denote the measured gene expression value as y_{ij} , which is either a
568 count or a binary value. We denote the relative position of j th measured location as $r_{ij} \in [0, 1]$,
569 where 0 corresponds to the center of the nucleus and 1 corresponds to the cellular boundary.
570

571 We model the subcellular spatial localization of gene expression within each cell using a one-
572 dimensional nonhomogeneous Poisson process (NHPP) model. Specifically, we assume that the
573 gene expression counts summed across all relative positions within a given interval $[a, b] \subset [0, 1]$
574 on the cellular radius follow a Poisson distribution, with the rate parameter being the integration
575 of an underlying NHPP density function in the interval $[a, b]$, where the NHPP density function
576 may vary with respect to the relative position within the cell. Mathematically, the model is
577 expressed as:

578
$$\sum_{r_{ij} \in [a, b]} y_{ij}(r_{ij}) \sim Poi\left(\int_a^b \lambda_i^*(r) dr\right),$$

579 where Poi denotes a Poisson distribution and $\lambda_i^*(r)$ is the unknown NHPP density function
580 depending on the relative position r . We assume that the NHPP density function $\lambda_i^*(r)$ is expressed
581 as a product of three terms

582
$$\lambda_i^*(r) = c_i s(r) \lambda(r), \quad (2)$$

583 where c_i , the total read depth for the i th cell, is used for normalization purpose and is calculated
584 as the summation of the total read counts across all genes within the cell; $s(r) = 2\pi r$ is another
585 normalization term to capture the fact that the cellular area corresponds to a particular radius r is
586 proportional to the size of the radius (that is, the area between r and $r + \Delta r$ is $\pi(r + \Delta r)^2 -$
587 $\pi r^2 = 2\pi r + \pi(\Delta r)^2 \rightarrow 2\pi r$ as $\Delta r \rightarrow 0$); and $\lambda(r)$ is the key term of interest, the subcellular
588 spatial expression intensity function that captures the subcellular spatial expression pattern along
589 the cellular radius.

590 With the above NHPP model, we can write down the joint likelihood of the subcellular gene
591 expression across n cells as:

593
$$L = \prod_{i=1}^n \left\{ e^{-\Lambda_i^*} \prod_{j=1}^{m_i} \lambda_i^*(r_{ij})^{y_{ij}} \right\}$$

594 with $\Lambda_i^* = \int_0^1 \lambda_i^*(r) dr$. Note that we have assumed that the subcellular spatial expression intensity
595 function $\lambda(r)$ is shared across cells, allowing us to borrow information across cells to enhance the
596 detection of subcellular spatial expression patterns.

597 The intensity function $\lambda(r)$ is key for modeling the subcellular spatial expression pattern of the
598 given gene. In particular, if a gene does not display subcellular spatial expression pattern and is
600 instead uniformly distributed within the cells, then $\lambda(r)$ is expected to be a constant that is
601 invariant to the relative position r . In contrast, if a gene displays subcellular spatial expression
602 pattern, then $\lambda(r)$ is expected to vary as a function of the relative position r .

603 Therefore, in the above NHPP model, identifying genes that display subcellular spatial expression
604 pattern within cells is equivalent to testing whether $\lambda(r)$ is a constant or not. The statistical power
605 of such hypothesis test will inevitably vary depending on how the specified expression intensity

607 function $\lambda(r)$ matches the true underlying subcellular spatial expression pattern displayed by the
608 gene of focus. For example, an intensity function enriched near zero will be particularly useful for
609 detecting subcellular expression patterns that are also enriched in the nuclear, while an intensity
610 function enriched near one will be particularly useful for detecting subcellular expression patterns
611 that are also enriched near the cellular membrane. However, the true underlying subcellular spatial
612 pattern for any gene is unfortunately unknown and may vary across genes. To ensure robust
613 identification of subcellular spatial expression genes across various spatial patterns, we consider
614 using a total of $k=22$ different kernel functions $\varphi_1(r), \dots, \varphi_k(r)$ inside the intensity function $\lambda(r)$
615 to capture a wide variety of possible subcellular spatial expression patterns (Fig. S68c). In
616 particular, each function is a Beta probability density function defined on the interval [0,1],
617 characterized by one of the 22 sets of shape parameters (Tab. S8) with a mode centering on 0, 0.1,
618 0.2, ..., or 1. Note that, while we use these 22 kernel functions as default kernels in the present
619 study, our method and software implementation can easily incorporate various numbers or types
620 of intensity kernels as desired by the user.
621

622 For each kernel $l = 1, \dots, k$ in turn, we model the intensity function in the form of $\lambda(r) = \alpha_l +$
623 $\beta_l \varphi_l(r)$, where α_l is the nonnegative intercept parameter and β_l is the nonnegative scaling
624 parameter for the l th kernel function. With the functional form of $\lambda(r)$, we can test the null
625 hypothesis $H_0: \beta_l = 0$, that $\lambda(r)$ is a constant. Rejecting the null hypothesis allows us to detect
626 genes that display subcellular spatial expression patterns captured by the particular kernel. We
627 perform inference and hypothesis test for each kernel in turn using a likelihood ratio test. In
628 particular, we first maximize the log likelihood both under the null and under the alternative using
629 the Adam optimizer in PyTorch [31]. Afterwards, we obtain the corresponding P value
630 asymptotically based on an equal mixture of two chi-square distributions with degrees of freedom
631 being zero and one [32]. Afterwards, we combine the k different P values calculated using
632 different kernels into a single P value using the Cauchy combination rule [33, 34]. Specifically,
633 we convert each of the k P values into a Cauchy statistic, aggregate the k Cauchy statistics through
634 summation, and convert the summation back to a single P value based on the standard Cauchy
635 distribution. The Cauchy rule takes advantage of the fact that a combination of Cauchy random
636 variables also follows a Cauchy distribution regardless of whether these random variables are
637 correlated or not. Therefore, the Cauchy combination rule allows us to effectively combine
638 multiple potentially correlated P values into a single P value for every gene. Finally, we control
639 FDR across genes using the Benjamini-Yekutieli procedure, which is effective for arbitrary
640 dependency among test statistics. We used an FDR cutoff of 0.05 for declaring significance.
641

642 Estimation of the subcellular spatial expression pattern with ELLA

643 While the primary focus of ELLA is on hypothesis testing, it can also be used to estimate the
644 subcellular spatial expression pattern for the detected genes. Specifically, for gene g , we can first
645 obtain the k estimated intensity functions for each of the k kernel functions as

$$646 \hat{\lambda}_l(r) = \hat{\alpha}_l + \hat{\beta}_l \varphi_l(r), \quad l = 1, \dots, k.$$

647 where $\hat{\alpha}_l$ and $\hat{\beta}_l$ are the estimates for the corresponding parameters. Because each of the k
648 estimated intensity functions captures a particular aspect of the overall subcellular spatial
649 expression intensity function $\lambda(r)$, we estimate $\lambda(r)$ with a weighted combination of the estimated
650 intensity functions in the form of

651

$$\hat{\lambda}(r) = \sum_{l=1}^k w_l \hat{\lambda}_l(r),$$

652 where w_l is the weight for the l th intensity function with $\sum_{l=1}^k w_l = 1$. The weights can be derived
653 based on Bayesian model averaging [35]. In particular, we denote the model with l th kernel
654 function as M_l and denote the data as D . The posterior distribution for $\lambda(r)$ is in the form of:
655 $P(\lambda(r)|D) = \sum_{l=1}^k P(\lambda(r)|M_l, D)P(M_l|D)$, with the posterior mean estimate being $\hat{\lambda}(r) =$
656 $\mathbb{E}[P(\lambda(r)|D)] = \sum_{l=1}^k \mathbb{E}[P(\lambda(r)|M_l, D)]P(M_l|D) = \sum_{l=1}^k \hat{\lambda}_l(r)P(M_l|D)$. Therefore, the weights
657 are in the form

658

$$w_l = P(M_l|D) = \frac{P(D|M_l)P(M_l)}{\sum_{j=1}^k P(D|M_j)P(M_j)} = \frac{P(D|M_l)}{\sum_{j=1}^k P(D|M_j)},$$

659
660 where the last equation holds due to the equal prior assumption on each model, with $P(M_j) = 1/k$
661 ($j = 1, \dots, k$). We approximate $P(D|M_l)$ with the maximized likelihood estimates to obtain the
662 weights and subsequently $\hat{\lambda}(r)$ (Supplementary Notes 2).

663
664 ELLA is implemented in python, with an underlying PyTorch Adam for efficient CPU or GPU
665 computation. The software ELLA, together with all analysis code used in the present study, are
666 freely available at <https://xiangzhou.github.io/software/>.

667
668 **Compared methods**
669 We compared ELLA with three methods: (1) SPRAWL [20], (2) Bento [19], and (3) Wilcox. For
670 both SPRAWL and Bento, we followed the tutorial on their corresponding GitHub pages and used
671 the recommended default parameter settings.

672
673 SPRAWL takes RNA location information from subcellular multiplexed imaging datasets as
674 inputs and does not explicitly require nuclear boundary or nuclear center information. SPRAWL
675 examines one gene at a time and uses four localization metrics to capture four different types of
676 subcellular spatial enrichment patterns that include peripheral, central, radial, and punctate.
677 Specifically, the peripheral metric is used to identify peripheral/anti-peripheral patterns where the
678 expression enrichment is either proximal or distal from the cell membrane. The central metric is
679 used to identify central/anti-central patterns where the expression enrichment is either proximal or
680 distal from the cell centroid. The radial metric is used to identify radial/anti-radial patterns where
681 a gene is either aggregated or depleted in a sector of the cell. The punctate metric is used to identify
682 punctate/anti-punctate patterns where a gene displays either self-colocalizing/self-aggregating or
683 self-repulsion inside the cell. Because the radial and punctate metrics can only be computed for
684 cells with no less than two expression counts, we had to filter out cells with less than two counts
685 when analyzing a given gene for these two metrics. For each gene and each metric in turn,
686 SPRAWL computes a score for every cell and averages them across cells in a particular cell type
687 to obtain the per-cell-type score. SPRAWL then converted the per-cell-type score to a P value
688 based on a standard normal distribution and used the Benjamini Hochberg (BH) procedure for
689 FDR control. We used an FDR threshold of 0.05 to obtain significant genes.

690
691 Bento takes RNA location information from subcellular multiplexed imaging datasets as inputs
692 and requires nuclear and cell boundaries as additional information. For each gene-cell pair in turn,

693 Bento computes 13 spatial summary statistics and uses its RNAforest function, which consists of
694 five independent pretrained binary random forest classifiers, to produce five binary labels that
695 classify gene expression pattern into one of the five patterns including nuclear, nuclear edge,
696 cytoplasmic, cell edge, and none. For each gene in the cell, we obtained the classification
697 probability p_c for each pattern c and used $1 - p_c$ to rank genes for the pattern, which allowed us
698 to measure powers based on FDR in the simulations. However, due to its use of classification
699 probability, it is not feasible to obtain FDR control in any real datasets with Bento.
700

701 Wilcox, a Wilcoxon rank sum test-based approach developed in the present study, detects genes
702 that are differentially expressed between two subcellular regions: the nucleus and the cytoplasm.
703 We focus on these two subcellular regions because we can extract the nuclear boundary and cell
704 boundary in many spatial transcriptomics studies. To detect those genes, for each cell in turn, we
705 first extracted the gene expression counts within the nucleus as well as the gene expression counts
706 in the cytoplasm. We then normalized the two counts by the corresponding cellular areas for the
707 two subcellular regions. Afterwards, we performed Wilcoxon rank sum test across cells to detect
708 genes that are differentially expressed between the nucleus and the cytoplasm.
709

710 **Simulations**

711 We performed comprehensive simulations based on imaging data to evaluate the performance of
712 ELLA and compare it with other methods. We did not perform simulations based on sequencing
713 data as neither SPRAWL nor Bento can be applied to analyze these data. For simulations, we first
714 extracted the cell boundaries of the embryonic fibroblast cells from the seqFISH+ data, calculated
715 the minimal and maximal radius of each cell, obtained a list of 90 reasonably shaped cells with the
716 ratio of minimal and maximal radius ≥ 0.3 , and extracted their nuclear centers and boundaries. We
717 then sampled with replacement n cells from these cells. For each cell in turn, we applied the same
718 binning strategy used in ELLA preprocess to divide the cell from the center of nucleus into 100
719 circular sectors of equal angle measure. In each sector v , we denote r_v as the maximum distance
720 between the center of nucleus and the cellular boundary in the sector. We calculated the
721 approximate area of the sector v as $\pi r_v^2 / 100$. We also denote $\theta_{v,min}$ and $\theta_{v,max}$ as the minimal and
722 maximum angle measurement of the sector, respectively. For the alternative simulations, we
723 further divided each circular sector into 25 annulus sectors with equal distances.
724

725 With the above preparations, we simulated gene expression for 1,000 genes, where each gene is
726 expressed as a binary count on m subcellular localizations in each cell as imaging data. In the null
727 simulations, none of these genes display cellular spatial expression patterns. In the alternative
728 simulations, 800 genes are null while 200 genes display different types of subcellular expression
729 patterns. Specifically, in the null simulations, we first randomly sampled the number of measured
730 locations inside each sector (m_v). We set m_v to be proportional to the area of the sector using the
731 function “np.random.choice” with the constraint $\sum_v m_v = m$. For each of the m_v locations in
732 sector v , we obtained two independent random variables, u_1 and u_2 , from a uniform distribution
733 $U(0,1)$, and converted them into the radius (r) and angle (θ) coordinates for the location, where
734 $r = r_v \sqrt{u_1}$ and $\theta = \theta_{v,min} + u_2(\theta_{v,max} - \theta_{v,min})$. The radius and angle coordinates are further
735 converted to the x and y coordinates in the form of $x = r \cos(\theta)$ and $y = r \sin(\theta)$.
736

737 In the alternative simulations, we simulated gene expression to exhibit subcellular expression
738 patterns from three pattern categories: symmetric, radial, and punctate. For the symmetric pattern

739 category, we considered eleven different expression patterns, including two patterns with nucleus
740 enrichment, two patterns with nuclear edge enrichment, five patterns with cytoplasmic enrichment,
741 and two patterns with membrane enrichment. For each pattern, we first randomly sampled the
742 number of measured locations inside each sector (m_v). We set m_v to be proportional to the area of
743 the sector using the function “np.random.choice” with the constraint $\sum_v m_v = m$. We then
744 constructed the expression intensity function $\lambda^{\text{true}}(r)$ in the form of $\lambda^{\text{true}}(r) = \alpha + \beta\varphi(r)$, where
745 $\varphi(r)$ is set to be one of the eleven beta probability density functions described earlier (upper panel
746 in Fig. S68c). Each beta probability density function is characterized by one of the eleven sets of
747 shape parameters (Set 1 in Tab. S8), with a mode centering on 0, 0.1, 0.2, ..., or 1. With $\lambda^{\text{true}}(r)$,
748 we define the pattern strength s as $(\max \lambda^{\text{true}}(r) - \min \lambda^{\text{true}}(r)) / \min \lambda^{\text{true}}(r)$. We also compute
749 $\lambda_i^{\text{true}}(r) = 2\pi r \lambda^{\text{true}}(r)$ and further $p_q = \int_{r_{q,\min}}^{r_{q,\max}} \lambda_i^{\text{true}}(r) dr$ [36], which represents the
750 probability of observing an expression measurement in the q -th annulus sector. Afterwards, we
751 simulated the number of expression measurement locations in each annulus sector,
752 $m_{v1}, \dots, m_{v20} \sim \text{Multinomial}(m_v, p_1, \dots, p_{20})$, with the total number of measured locations in the
753 sector being $m_v = \sum_q m_{vq}$. We then applied the same strategy described in the above paragraph
754 to simulate the x and y coordinates for each of the m_{vq} locations within each annular sector q .
755

756 In the symmetric pattern, we created different simulation settings by varying the number of cells
757 (n), expression level (m), the subcellular expression patterns, and pattern strength (s). To do so,
758 for each pattern, we first create a baseline simulation setting where we set the number of cells to
759 be $n=100$, the expression level to be $m=5$, and in the case of alternative simulations, the pattern
760 strength to be moderate ($s=0.6$). We then varied the cell number ($n=10, 20, 50, 100, 200, 300$ or
761 500), expression level ($m=1, 2, 10, 20, 50, 100$), and pattern strength (s ranges from 0.1 to 1.0
762 with increments of 0.1), one parameter at a time on top of the baseline settings for each of the 11
763 symmetric patterns to create 22 simulations settings. The detailed parameters for each simulation
764 setting are listed in Tab. S9. We performed 10 simulation replicates in each setting.
765

766 For the radial patterns, we first consider a radial-unif setting where gene expression is enriched in
767 one sector of the cell with the expression counts within the sector being randomly distributed. For
768 each cell and each gene in turn, we randomly selected a sector with a central angle $\pi/2$. We
769 sampled the number of measuring locations in the sector, m_1 , from a binomial distribution
770 $\text{Bin}(m, 0.5)$. We also sampled the number of measuring locations in the complementary sector
771 with a central angle of $3/2\pi$, m_2 , to be $m - m_1$. Afterwards, we randomly sampled the x and y
772 coordinates for each measurement location in the same way as described in the null simulations.
773 Therefore, the gene expression is enriched in one sector of the cell with a fold enrichment of 3.0.
774 Next, we consider a radial-cyto setting where the gene expression is not only enriched within the
775 sector but are also further enriched in the cytoplasm. To do so, on top of the radial-uniform setting,
776 we used the intensity function described in the symmetric pattern #7 to simulate the x and y
777 coordinates for the measurement locations in the selected circular sector that has a central angle of
778 $\pi/2$. In addition, we randomly sampled the x and y coordinates in the complementary circular
779 sector with a central angle of $3/2\pi$ for the measurement locations in the same way as described in
780 the null simulations. Therefore, the average gene expression inside the sector is also 3.0 times
781 higher than that in the remaining parts of the cell, while the expression within the sector is enriched
782 in the cytoplasmic region due to symmetric pattern #7 with a fold enrichment of approximately
783 5.1.

784
785 For the punctate pattern, we consider a punctate-cyto setting where gene expression is enriched in
786 a small subcellular disc in the cytoplasm. To do so, we set the radius coordinate for the center of
787 the punctate disc to be 0.8 and randomly sampled the corresponding angle coordinate θ from a
788 uniform distribution $U(0, 2\pi)$. We then converted the radius and angle coordinates to the location
789 coordinates (x_c, y_c) . Afterwards, we set the radius of the punctate disc to be 1/10 of the average
790 cell diameter, which consists of 30 pixels for seqFISH+ cells. We sampled the number of
791 measurement locations within the punctate disc, m_1 , from a binomial distribution $\text{Bin}(m, 0.2)$. We
792 randomly sampled the x and y coordinates for the m_1 locations inside the punctate disc as well as
793 those for the remaining $m - m_1$ locations in the entire cell including the punctate disc using the
794 same strategy in the null simulations. The expression in the punctate disc is on average 5.03 times
795 higher than that in the remaining parts of the cell. For radial and punctate patterns, we also
796 performed 10 simulations replicates for each of the three settings.

797 In summary, we used the simulations introduced above to systematically evaluate and compare the
798 performance of different methods in terms of type I error control and statistical power. Specifically,
799 type I error control was assessed by generating QQ plots of the log-transformed p-values (\log_{10})
800 under the null hypothesis, which allowed us to examine how closely the observed distribution of
801 P values matched the expected uniform distribution under the null. Power was assessed by
802 calculating the proportion of genes detected exhibiting any subcellular expression pattern,
803 calculated as the fraction of true positive genes detected at an FDR threshold of 0.05, across various
804 simulation scenarios, including different pattern types (symmetric, radial, and punctate) and
805 varying parameters such as the number of cells, expression levels, and pattern strength.

806 **Analyzed datasets**

807 We examined four public high-resolution spatial transcriptomics datasets described below.

808 Seq-Scope mouse liver data

809 Seq-Scope is a spatial barcoding technology with a spatial resolution comparable to an optical
810 microscope. It is based on a solid-phase amplification of randomly barcoded single-molecule
811 oligonucleotides using an Illumine sequencing platform. These RNA-capturing barcoded clusters
812 represent the pixels of Seq-Scope and are $\sim 0.5\text{-}0.8 \mu\text{m}$ apart from each other with an average
813 distance of $0.6 \mu\text{m}$, capturing 848 UMI on average per $10 \mu\text{m}$ diameter bin.

814
815 We downloaded the mouse liver data from the Seq-Scope resources website [37]. The data contains
816 5.88 ± 4.22 (mean \pm sd) number of genes per pixel, with a total of 32,976 genes measured across
817 $\sim 2 \times 10^7$ locations. The Seq-Scope mouse liver data contains 10 tiles sequenced on one MiSeq
818 flow cell with each tile being a 1mm-wide circular imaging area. Among these 10 tiles, six of them
819 are from a normal mouse fragmented frozen liver section and four of them are from an early-onset
820 liver failure mouse model section (TD; [23]). The tiles cover liver portal-central tissue zonation
821 and contain two main cell types: hepatocytes and non-parenchymal cells (NPC) such as
822 macrophages, hepatic stellate cells, endothelial cells, and red blood cells. Our analyses focus on
823 the hepatocytes which can be further divided into periportal (PP) and pericentral (PC) cells. The
824 two Seq-Scope tissue sections (normal and TD) each comes with multiple H&E staining images,
825 including high-resolution images (10X) covering a portion of the normal and TD tile areas and

827 low-resolution images (4X) covering nearly all the normal and TD tile areas. We used the low
828 resolution (4X) images to ensure high coverage of the tiles.
829

830 The Seq-Scope mouse liver data consists of two data modalities, namely the spatial transcriptomics
831 data and the accompanying H&E staining images. For the spatial transcriptomics data, we obtained
832 the unspliced and spliced gene expression counts on each measured location using STARsolo from
833 the raw fastq files. For the H&E staining images, we concatenated all the images from the normal
834 tissue section or the TD tissue section, segmented individual cells on the concatenated image using
835 Cellpose ([38]; [Fig. S18-23](#)), and obtained cells that overlapped with the tile areas. On each tile,
836 we plotted the unspliced expression reads to visualize cell nucleus and plotted the total UMI counts
837 to visualize the tissue boundaries ([Fig. S24-26](#)). These nucleus and tissue boundary information
838 were used to manually align each spatial transcriptomics tile to the concatenated normal or TD
839 H&E images ([Fig. S27-28](#)). After modality alignment, we assigned each spatial location to a cell
840 based on the aligned cell segmentation results ([Fig. S29-30](#)). For each cell in turn, we used
841 numpy.argmax function in python to declare its nuclear center, which is defined to be the location
842 within 200 units (~2 μ m) from the cell boundary where the maximum of unspliced read counts
843 density is observed. In each tile, we filtered out cells with a low-quality nuclear center where the
844 unspliced read count density values at the nuclear center is below the 95% quantile value across
845 locations or where the spliced read count density at the nuclear center is above the 95% quantile
846 value across locations. In addition, we obtained cell type marker genes for each of the three cell
847 types (PP, PC, and NPC; [Tab. S10](#); [5]) and obtained the total counts of cell type marker genes
848 for each cell. Note that the NPC cells, such as macrophages, hepatic stellate cells, endothelial, and
849 red blood cells, are relatively rare across the tiles and are hard to segment due to their small sizes
850 on the H&E-based images. Therefore, following the original Seq-Scope study, we removed NPC
851 cells that are characterized by NPC marker gene counts above the 95% quantile across all cells.
852 Afterwards, we normalized the PP and PC marker gene counts for the remaining cells first across
853 genes to have zero mean and unit standard deviation and then across cells to have zero mean and
854 unit standard deviation. We then summed the normalized PP and PC marker genes separately in
855 each cell to obtain a PP score and a PC score per cell. We annotated a cell as a PP cell if its PP
856 score is greater than the PC score and annotated a cell as a PC cell otherwise. Such annotations
857 largely align with Seq-Scope's original cell type annotations ([Fig. S32](#)). We removed cells with
858 extreme sizes, including extremely large cells with x or y coordinate range (max-min) exceeding
859 the 95% quantile value across cells within the cell type or extremely small cells with x or y
860 coordinate range below the 5% quantile value. After quality control, we obtained 276 normal PP
861 cells, 276 normal PC cells, 236 TD PP cells, and 82 TD PC cells. Genes expressed in more than
862 50 cells and with more than 3 counts in at least 5 cells were retained, leading to 497 to 1,349 genes
863 per cell type.
864

865 Stereo-seq mouse embryo data

866 Stereo-seq combined DNA nanoball (DNB)-patterned arrays and *in situ* RNA capture to enhance
867 the spatial resolution of omics-sequencing. Standard DNB chips have spots with approximately
868 0.22 μ m diameter and a center-to-center distance of 0.5 or 0.715 μ m, providing up to 400 spots
869 per 100 μ m² for tissue RNA capture. Stereo-seq captured UMI counts range on average from 69
870 per 2 μ m diameter bin (for bin3, 3 \times 3 DNB) to 1,450 per 10 μ m diameter bin (for bin 14, 14 \times 14
871 DNB, equivalent to ~one medium size cell).
872

873 We downloaded the raw sequencing data on slice E1S3 of the Stereo-seq mouse embryo data from
874 CNGB Nucleotide Sequence Archive [39]. We downloaded the processed gene expression (bin1)
875 data and the accompanying nucleic acid staining image from MOSTA [40]. Slice E1S3 is a profiled
876 sagittal frozen tissue section with 10 μ m thickness from a C57BL/6 mouse embryo on day E16.5.
877 It covers all major tissues and organs including Epidermis, Meninges, Cartilage, Jaw and tooth,
878 Choroid plexus, Kidney, GI tract, Spinal cord, Muscle, Heart, Bone, Cartilage primordium, Brain,
879 Adrenal gland, Connective tissue, Thymus, Blood vessel, Liver, Olfactory epithelium, Lung,
880 Pancreas, and Mucosal epithelium. The nucleic acid staining image of the slice was stained using
881 BM purple and was imaged using a Ti-7 Nikon Eclipse microscope. We considered 25 cell types
882 along with cell type marker genes from the Stereo-seq study (Tab. S11). The 25 cell types include
883 Cardiomyocyte, Chondrocyte, Choroid plexus, Dorsal midbrain neuron, Ganglion, Endothelial cell,
884 Keratinocyte, Epithelial cell, Erythrocyte, Facial fibroblast, Fibroblast, Forebrain neuron,
885 Forebrain radial glia, Hepatocyte, Immune cell, Limb fibroblast, Macrophage, Meninges cell, Mid-
886 /hindbrain and spinal cord neuron, Myoblast, Olfactory epithelial cell, Radial glia, Smooth muscle
887 cell, Spinal cord neuron, and Diencephalon neuron. We processed the Stereo-seq data in the same
888 way as we did for the Seq-Scope data except for the modality alignment step which is omitted here
889 as the Stereo-seq slice was accompanied by nucleic acid staining that has already been aligned
890 with the slices (Fig. S36). The processed data contains cell label of each location, cell center, cell
891 boundary, cell type, and read depth of each cell (Fig. S37). We annotated cell types based on 75
892 cell type marker genes provided by the original study, resulting in an average of 3,689 cells
893 (median=3,968, min=782, max=5,314) per cell type (Fig. S38). We focused on a cardiothoracic
894 region on slice E1S3 (Fig. 3a) and two major cell types: precursor muscle cells, or myoblasts, and
895 mature muscle cells, or cardiomyocytes. Similar quality control steps were conducted as described
896 in the Seq-Scope data preprocessing. We retained genes expressed in more than 30 cells.
897

898 SeqFISH+ mouse fibroblast data

899 SeqFISH+ performs super-resolution imaging and multiplexing of 10,000 genes in a single cell
900 using sequential hybridizations and imaging with a standard confocal microscope. We obtained
901 the seqFISH+ NIH/3T3 fibroblast data preprocessed by Bento from [41]. The raw seqFISH+ data
902 consists of two modalities: the spatial transcriptomics measurements and an accompanying DAPI
903 staining image. The spatial transcriptomics modality of the data contains 3,726 genes with at least
904 10 counts expressed in at least one cell and 179 cells with nuclear segmentation results, with a
905 resolution of 103nm. The downloaded seqFISH+ data comes with cell segmentation boundaries
906 and nuclear segmentation boundaries, each represented by a set of points densely scattered along
907 the boundaries. With the nucleus segmentation information, we computed the nuclear center of
908 each cell as the k-means center of all nucleus boundary points. We computed the average nuclear
909 radius of each cell by averaging the distance of all nuclear boundary points to the nuclear center.
910 We computed the average cell radius of each cell by averaging the distance of all cell boundary
911 points to the nuclear center. Afterwards, we computed the nucleus-cell ratio of each cell by
912 dividing the average nuclear radius with the average cell radius. We excluded eight cells that have
913 a nuclear-cell ratio beyond two standard deviations from the mean (Fig. S51a). We focused on the
914 remaining 171 cells for analysis. These cells have an average nuclear-cell ratio of 0.46 (Fig. S51b).
915 We retained genes expressed in more than 50 cells and with more than 3 counts in at least 5 cells,
916 resulting in 2,747 genes for analysis.
917

918 MERFISH adult mouse brain data
919 The mouse brain MERFISH dataset contains over 200 adult mouse brain slices from 4 mice and
920 covers a panel of ~1,100 selected genes with around 8 million cells. The dataset consists of two
921 data modalities, namely the spatial transcriptomics data and the accompanying DAPI and polyA
922 staining images. We focused on one coronal slice of mouse 2 from the
923 220501_wb3_co2_15_5z18R_merfish5 experiment and obtained the preprocessed data from [42].
924 The obtained data was measured on a coronal tissue slice with 10 μm thickness and contains five
925 1.5- μm -thick optical z-stacks, with 1,147 genes measured on ~100,000 cells. The data also
926 includes cell segmentation information in the form of sets of points densely scattered along the
927 boundaries for each z-stack (0-4), along with cell centroid information shared across z-stacks (Fig.
928 S56). For each measured transcript, we calculated its relative position to nuclear center based on
929 the cell segmentation on the z-stack that it belongs to as well as the shared cell centroid. We
930 exclude cells whose centroid is outside or too close to ($< 0.5 \mu\text{m}$) its segmentation boundaries on
931 the baseline stack ($z=0$). In addition, we measured the variability of cell segmentation boundaries
932 on each non-baseline stack ($z>0$) versus that on the baseline stack ($z=0$) by KL divergence. We
933 excluded cells whose cell segmentation boundaries are highly variable across z-stacks based on a
934 KL divergence threshold of 0.5. We obtained cell type marker genes (Tab. S12) from the Stereo-
935 seq study for four cell types that include excitatory neurons (EX), inhibitory neurons (IN),
936 astrocytes (Astr), and oligodendrites (Olig). We then carried out the same cell typing procedure as
937 described in the Seq-Scope and Stereo-seq datasets above. We focused on four major cell types
938 residing in the midbrain: excitatory neurons (EX, $n=577$), inhibitory neurons (IN, $n=525$),
939 astrocytes (Astr, $n=480$), and oligodendrocytes (Olig, $n=948$), with 557-878 genes per cell type.
940 Similar quality control steps were conducted as described in the Seq-Scope data preprocessing.
941 After quality control, we retained 480-948 cells per cell type. We retained genes expressed in more
942 than 50 cells, resulting in 557-878 genes per cell type for analysis.
943

944 **Real data analysis details**

945

946 **Subcellular expression pattern score**

947 After obtaining the estimated subcellular expression intensity function $\hat{\lambda}(r)$, we computed a
948 subcellular expression pattern score r^* , defined as the relative position corresponding to the
949 mode/peak of the estimated expression intensity function: $r^* = \underset{r \in [0,1]}{\operatorname{argmax}} \hat{\lambda}(r)$. Therefore, r^*

950 ranges from zero to one, with a value close to zero indicating expression enrichment in the center
951 of the cell nucleus and a value close to one indicating expression enrichment on the cell boundary.
952

953 **Gene clustering based on the estimated expression pattern**

954 We clustered genes into different spatial pattern categories based on their estimated intensity
955 functions. To do so, for each detected gene, we evaluated its estimated intensity function $\hat{\lambda}(r)$ at
956 21 equidistant points, ranging from $r=0$ to $r=1$ with increments of 0.05. Additionally, we
957 calculated the difference between consecutive functional values to obtain 20 differences. We then
958 pooled the 21 functional values and 20 differences for each gene and used them as input for k-
959 means clustering. We determined the optimal number of gene clusters using the Elbow method
960 [43].

961

962 **Transcription factor analysis**

963 To examine the subcellular localization of transcription factors, we obtained a list of 1,358 mouse
964 transcription factors from FANTOM5 SSTAR [26]. For all datasets, we examined the proportions
965 of transcription factors that are measured in the datasets between pairs of gene clusters with
966 Fisher's exact tests.
967

968 Computing the unspliced-spliced ratio

969 In the sequencing-based datasets (Seq-Scope and Stereo-seq), for each gene in turn, we calculated
970 the unspliced-spliced ratio for each cell by dividing the total unspliced counts (plus a pseudo count
971 of one) by the total spliced counts (plus a pseudo count of one). We then computed the average
972 value of this ratio across cells. We applied Mann-Whitney U tests to test the unspliced-spliced
973 ratios between pairs of gene clusters across cell types.
974

975 snRNA-seq analysis

976 We examined the genes detected in the Seq-Scope dataset using a matched single-nucleus RNA
977 sequencing (snRNA-seq) dataset. The snRNA-seq data was collected on mouse hepatocytes and
978 was downloaded from the BRAIN Initiated Cell Census Network (BICCN) consortium 2021 [24].
979 For each gene in turn, we defined its sn-sc ratio as the average gene counts per nucleus in the
980 snRNA-seq data divided by the average gene counts per cell in the Seq-Scope data. We applied
981 Mann-Whitney U tests to test the sn-sc ratios between pairs of clusters across cell types. We also
982 examined the genes detected in the MERFISH adult mouse brain dataset using a matched snRNA-
983 seq dataset [30] that was collected on adult brain sections and calculated sn-sc ratios for the
984 corresponding four cell types (EX, IN, Astr, and Olig) in the same way. We filtered out lowly
985 expressed nonsignificant genes (average sc counts <1.5) due to the sparsity of the data.
986

987 Gene length analysis

988 We performed gene length analysis in the four datasets. To do so, we excluded mitochondrial
989 genes and genes that have not been mapped to a chromosome as their gene length information is
990 unavailable. We extracted four types of gene length measurements using GTF tools [44] from the
991 same reference genome (mm10.gtf) that were used for alignment. The four measurements include
992 (i) mean, (ii) median, (iii) longest single isoform, and (iv) total length across exons, all in the unit
993 of base pairs. We then applied Mann-Whitney U tests to test the difference between pairs of gene
994 clusters for each measurement.
995

996 SRP and RP analysis

997 In the Seq-Scope data, for each gene in turn, we used DeepSig [45] with Gencode [46] to predict
998 whether the corresponding protein contains SRP. To do so, we downloaded protein sequences in
999 the form of protein coding transcripts fasta files from Gencode release M28, used DeepSig to
1000 analyze the protein sequence, and referred to the genes corresponding to proteins with SRPs as
1001 SRP-coded genes. For genes with multiple protein isoforms, we used the longest isoform for SRP
1002 prediction. We examined the proportions of SRP-coded genes between pairs of gene clusters with
1003 Fisher's exact tests. In the Stereo-seq data, we identified a list of ribosomal protein (RP) genes
1004 whose gene ID starts with RPS or RPL. These are genes of the nuclear genome that encode the
1005 protein subunits of the ribosome. These genes are expected to be enriched in the cytoplasm as
1006 ribosomal subunits are exported from the nucleus to the cytoplasm after their assembly in the
1007 nucleolus. We examined the proportions of RP genes between pairs of gene clusters with Fisher's
1008 exact tests.
1009

1009

1010 Cell by cell analysis in the seqFISH+ data

1011 We randomly picked 20 cells from the seqFISH+ data. For each cell in turn, we kept genes with
1012 more than 10 counts in this analysis. We applied Bento following its instruction to classify each
1013 gene in each cell into five binary labels, corresponding to “nuclear”, “nuclear edge”, “cytoplasmic”,
1014 “cell edge”, and “none” patterns. We also applied ELLA to analyze each cell separately. We
1015 collected the estimated expression intensities $\hat{\lambda}(r)$ of all genes and carried out k-means clustering
1016 ([Methods](#)) to obtain their pattern cluster labels.

1017

1018 Cell cycle-based analysis in the seqFISH+ data

1019 In the seqFISH+ data, we computed the single cell gene counts and used Seurat to classify the
1020 fibroblasts into three cell subclusters corresponding to G1 (n=36, 21%), S (n=83, 49%), and G2M
1021 (n=52, 30%) cell cycle phases. We kept genes that are expressed in at least 30 cells and that have
1022 more than 3 counts in at least 5 cells, resulting in 756, 2475, and 1776 genes for G1, S, and G2M
1023 cells respectively. We then applied ELLA to analyze one gene at a time for each cell cycle
1024 subcluster. Afterwards, we retrieved ELLA genes pattern cluster labels (1-5) obtained using all
1025 fibroblasts, calculated pattern scores for genes obtained in the cell cycle specific ELLA analysis,
1026 examined these patterns cores across gene clusters, and carried out one side Mann-Whitney U test
1027 to compare the pattern scores between G1 and S/G2M subclusters ([Fig. 4g](#), upper panel). We also
1028 focused on 723 genes commonly detected across the three cell-cycle subclusters and identified a
1029 set of genes in each gene pattern cluster with decreasing pattern scores from G1 to S and from S
1030 to G2M. Specifically, in each pattern cluster, we computed the increase in pattern scores for
1031 each gene from G1 to S (denoted as s_1) and another increase in pattern score from S to G2M
1032 (denoted as s_2). We identified a set of genes with increasing pattern scores from G1 to S to G2M,
1033 characterized by $s_1 \leq 0$, $s_2 \leq 0$, and $s_1 + s_2 < 0$. We applied Mann-Whitney U tests to compare
1034 the pattern scores between the identified genes and the remaining genes at G1, S, and G2M phases
1035 across gene pattern clusters ([Fig. 4g](#), lower panel).

1036

1037 **Data availability**

1038 The original public data used in this study can be accessed through the following links: (normal
1039 and diseased) mouse liver data by Seq-Scope available at <https://lee.lab.medicine.umich.edu/seq->
1040 scope; mouse embryo E1S3 data by Stereo-seq available at
1041 <https://db.cngb.org/search/project/CNP0001543/>; mouse embryonic fibroblast data by seqFISH+
1042 at https://figshare.com/articles/dataset/Bento_spatial_AnnData_formatted_datasets/15109236/2;
1043 and adult mouse brain data by MERFISH available at
1044 <https://download.brainimagerlibrary.org/29/3c/293cc39ceea87f6d/>. Details about the data we used
1045 in this study are provided in the [Analyzed datasets](#).

1046

1047 **Code availability**

1048 The ELLA software program and source code have been deposited at
1049 <https://xiangzhou.github.io/software/> and <https://github.com/jadexq/ELLA>. All scripts used to
1050 reproduce all the analyses are also available on the websites.
1051

1052 **Acknowledgements**

1053 This study was supported by the National Institutes of Health (NIH) grants R01GM126553,
1054 R01HG011883 and R01GM144960, all to X.Z.
1055

1056 **Ethics declarations**

1057 Competing interests:
1058 The authors declare that they have no competing interests.

1059 **Figure legends**
1060

1061 **Fig. 1 | Schematic of ELLA and simulation results.** **a.** ELLA is a method designed for modeling
1062 the subcellular localization of mRNAs and detecting genes that display spatial variation within
1063 cells in high-resolution spatial transcriptomics. ELLA takes as inputs the high-resolution spatial
1064 gene expression data along with nuclear center and cell segmentation information. It first performs
1065 data pre-processing including normalization and registration to create a unified cellular coordinate
1066 system to anchor diverse cell shapes and morphologies through defining a cellular radius in each
1067 cell that points from the center of the nucleus towards the cellular boundary. It then fits a
1068 nonhomogeneous Poisson process model for each gene to capture its spatial distribution within
1069 cells, computes a P value to capture any subcellular expression pattern observed along the cellular
1070 radius, and estimates such pattern in the form of estimated pattern expression intensity and pattern
1071 score. ELLA is capable of borrowing information across cells through a joint likelihood framework
1072 to substantially improve detection power, while taking advantage of multiple intensity kernel
1073 functions to capture the distinct subcellular expression patterns that may be encountered in various
1074 biological settings to ensure robust performance. **b.** Quantile-quantile plots of the expected and
1075 observed -log₁₀ P values for different methods in the baseline null simulation, where gene
1076 expression is randomly distributed spatially within cells. ELLA was compared to SPRAWL, Bento,
1077 and Wilcox. **c.** Radar plots show the power of different methods in the alternative simulations with
1078 multiple cells across eleven symmetric subcellular expression patterns, where gene expression is
1079 enriched in specific subcellular regions within cells. ELLA was compared to SPRAWL and
1080 Wilcox and power was evaluated based on 5% FDR. **d.** Radar plots of the power of different
1081 methods in the alternative simulations with multiple cells across three asymmetric subcellular
1082 expression patterns, where gene expressions exhibit distinct asymmetric patterns. ELLA was
1083 compared to SPRAWL and Wilcox and power was evaluated based on 5% FDR. **e.** Radar plots
1084 show the power of different methods in the additional alternative simulations with one cell across
1085 five symmetric subcellular expression patterns, where gene expression is enriched in specific
1086 subcellular regions within cells. ELLA was compared to Bento and power was evaluated based on
1087 5% FDR.
1088

1089 **Fig. 2 | Seq-Scope mouse liver data analysis.** **a.** Data snapshot for tissue tile 2107. Left panel
1090 displays the full tile with expression from four gene sets (PC marker genes, blue dots; PP marker
1091 genes, green dots; [Tab. S10](#); mitochondria genes, black dots; nuclear genes, red dots; [Tab. S13](#))
1092 along with cell segmentation boundaries. Middle panel zooms into a subregion and displays
1093 unspliced expression densities (colored green), nuclear centers (crosses), gene expression from the
1094 four lists, along with cell segmentation boundaries. Right panel displays the H&E staining image
1095 of the subregion. **b.** Estimated spatial expression pattern for genes in each of the five gene pattern
1096 clusters identified by ELLA. Upper panel shows the number and proportion of genes across five
1097 pattern clusters. Middle panel displays the estimated expression intensities for genes across
1098 clusters. Lower panel displays the estimated pattern score for genes across clusters. **c.** Example
1099 genes and cells for the five pattern clusters. One example gene is shown for each pattern cluster.
1100 Upper panel shows the gene name, ELLA P value, and the estimated expression intensities
1101 overlaid on the density heat map. Each row of the density heat map visualizes the number of
1102 counts standardized by area in 10 randomly selected cells across relative positions with intervals
1103 of 0.1. Lower panel displays the expressions of the corresponding genes within five selected cells,
1104 overlaid with cell boundaries and aligned nuclear centers (crosses). **d.** Bar plot shows the average

1105 sn/sc RNA ratio in the form of snRNA expression level normalized by scRNA expression level
1106 across genes in pattern clusters 1 (red), 2-5 (green), and non-cluster 1 genes (i.e. clusters 2-5 plus
1107 the nonsignificant genes; grey). Nuclear enriched genes (cluster 1) tend to exhibit higher relative
1108 snRNA expression levels. **e**. Bar plot shows the average unspliced/spliced expression ratio across
1109 genes in pattern clusters 1, 2-5, and non-cluster 1 genes. Nuclear enriched genes (cluster 1) tend
1110 to exhibit higher unspliced/spliced expression ratios. **f**. Bar plot displays the average gene length,
1111 measured by four metrics (x-axis), across genes in pattern clusters 1, 2-5, and non-cluster 1 genes.
1112 Nuclear enriched genes (cluster 1) tend to exhibit longer gene lengths. **g**. Bar plot displays the
1113 proportions of SRP-coded genes for genes in pattern clusters 1, 2-5, and non-cluster 1 genes.
1114 Cytoplasmic enriched genes (clusters 2-5) frequently encode SRP. Statistical significance for pair-
1115 wise comparisons (*: <0.05; **: <0.01; ***: <0.001) is based on Mann-Whitney U test (**d-f**) or
1116 Fisher's exact test (**g**). The error bars represent the 25th and 75th percentiles, and data points
1117 beyond this range are not included (**d-f**).

1118

1119 **Fig. 3 | Stereo-seq mouse embryo data analysis.** **a**. Data snapshot for tissue slice E1S3. Left
1120 panel displays the nucleic acid staining image of the slice. Middle panel displays expression from
1121 three gene sets (Erythrocyte marker genes, grey dots; Myoblast marker genes, blue dots;
1122 Cardiomyocyte marker genes, green dots; [Tab. S11](#)). Right panel zooms into a subregion and
1123 displays expression from three gene sets (Myoblast marker genes, blue dots; Cardiomyocyte
1124 marker genes, green dots; *Malat1*, red dots), nuclear centers (crosses) overlayed on the nucleic
1125 acid staining image. **b**. Estimated spatial expression pattern for genes in each of the five gene
1126 pattern clusters identified by ELLA. Upper panel shows the number and proportion of genes across
1127 five pattern clusters. Middle panel displays the estimated expression intensity for genes across
1128 clusters. Lower panel displays the estimated pattern score for genes across clusters. **c**. Example
1129 genes and cells for the five pattern clusters. One example gene is shown for each pattern cluster.
1130 Upper panel shows the gene name, ELLA P value, and the estimated expression intensities
1131 overlayed on the density heat map. Each row of the density heat map visualizes the number of
1132 counts standardized by area in 10 randomly selected cells across relative positions with intervals
1133 of 0.1. Lower panel displays the expressions of the corresponding genes within five selected cells,
1134 overlayed with cell boundaries and aligned nuclear centers (crosses). **d**. Bar plot shows the average
1135 unspliced/spliced expression ratio across genes in pattern clusters 1-3 (red), 4-5 (green), and non-
1136 cluster 1-3 genes (i.e. clusters 4-5 plus the nonsignificant genes; grey). Genes enriched close to
1137 nuclear center (cluster 1-3) tend to exhibit higher unspliced/spliced expression ratios. **e**. Bar plot
1138 displays average gene length, measured by four metrics (x-axis), across genes in pattern clusters
1139 1-3, 4-5, and non-cluster 1-3 genes. Genes enriched close to nuclear center (cluster 1-3) tend to
1140 exhibit longer gene lengths. **f**. Bar plot displays the proportions of transcription factors (TFs) for
1141 genes in pattern clusters 1-3, 4-5, and non-cluster 1-3 genes. Genes enriched close to nuclear center
1142 (cluster 1-3) contain a higher proportion of TFs. **g**. Bar plot displays the proportions of ribosomal
1143 protein (RP) genes for genes in pattern clusters 1-3, 4-5, and non-cluster 1-3 genes. Cytoplasmic
1144 enriched genes (clusters 4-5) contain a higher proportion of RP genes. Statistical significance for
1145 pair-wise comparisons (*: <0.05; **: <0.01; ***: <0.001) is based on Mann-Whitney U test (**d-e**)
1146 or Fisher's exact test (**f-g**). The error bars represent the 25th and 75th percentiles, and data points
1147 beyond this range are not included (**d-e**).

1148

1149 **Fig. 4 | SeqFISH+ mouse embryonic fibroblast data analysis.** **a**. Data snapshot for the
1150 embryonic fibroblast cell line. Left panel displays all transcripts (grey dots) measured across 17

1151 batches. Middle panel zooms into batch 14 and displays all transcripts (grey dots), nuclear centers
1152 (crosses), along with nuclear segmentation boundaries (grey dashed line) and cell segmentation
1153 boundaries (grey solid line). Right panel displays batch 14 with expression from three gene sets
1154 (nuclear or nuclear edge genes, red dots; cytoplasmic genes, green dots; protrusion genes, blue
1155 dots; [Tab. S14](#)) and nuclear centers (crosses) overlayed with nuclear segmentation boundaries
1156 (grey dashed line) and cell segmentation boundaries (grey solid line). **b.** Estimated spatial
1157 expression pattern for genes in each of the five gene pattern clusters identified by ELLA. Upper
1158 panel shows the number and proportion of genes across five pattern clusters. Middle panel displays
1159 the estimated expression intensities for genes across clusters. Lower panel displays the estimated
1160 pattern score for genes across clusters. **c.** Example genes and cells for the five pattern clusters. One
1161 example gene is shown for each pattern cluster. Upper panel shows the gene name, ELLA P value,
1162 and the estimated expression intensity overlayed on the density heat map. Each row of the density
1163 heat map visualizes the number of counts standardized by area in 10 randomly selected cells across
1164 relative positions with intervals of 0.1. Lower panel displays the expressions of the corresponding
1165 genes within one selected cell, overlayed with cell boundary and nuclear center (cross). **d.** Bar plot
1166 displays average gene length, measured by four metrics (x-axis), across genes in pattern clusters
1167 1-3 (red), 4-5 (green), and non-cluster 1-3 genes (i.e. clusters 4-5 plus the nonsignificant genes;
1168 grey). Genes enriched close to nuclear center (clusters 1-3) tend to exhibit longer gene lengths. **e.**
1169 Bar plot displays the proportions of transcription factors (TFs) for genes in pattern clusters 1-3, 4-
1170 5, and non-cluster 1-3 genes. Genes enriched close to nuclear center (clusters 1-3) contain a higher
1171 proportion of TFs. **f.** Estimated spatial expression pattern of genes across cell cycle phases (from
1172 left to right: G1, S, and G2M). Upper panels display the estimated expression intensities of genes
1173 across cell cycle phases. Middle panels display the estimated pattern score of genes across cell
1174 cycle phases. Lower Panel display the distribution of the estimated pattern score of genes across
1175 cell cycle phases. **g.** Upper panel displays Violin plots of the estimated pattern scores across cell
1176 cycle phases for genes in different pattern clusters. Genes significant in the G1 phase are less likely
1177 to be enriched close to nuclear center and display larger pattern scores compared to the genes in
1178 the S and G2M phases based on one side Mann-Whitney U test. Lower panel displays line plots
1179 of patterns score trajectories with respect to cell cycle phases for genes commonly detected in all
1180 cell cycle phases across pattern clusters. A subset of genes in clusters 2-5 display decreasing
1181 pattern scores through G1, S, and G2M phases (colored lines) while all six cluster 1 genes retain
1182 the nuclear pattern across all cell cycle phases (grey lines). Statistical significance for pair-wise
1183 comparisons (*: <0.05; **: <0.01; ***: <0.001) is based on Mann-Whitney U test (**d**) or Fisher's
1184 exact test (**e**). The error bars represent the 25th and 75th percentiles, and data points beyond this
1185 range are not included (**d**).
1186

1187 **Fig. 5 | MERFISH mouse brain data analysis.** **a.** Data snapshot for tissue slice hemi-brain region.
1188 Left panel shows the DAPI staining image. Middle panel displays expression for four gene sets
1189 (EX marker genes, blue dots; IN marker genes, green dots; Astr marker genes, red dots; Olig
1190 marker genes, orange dots; [Tab. S12](#)). Right panel zooms into a subregion and displays expression
1191 for the four gene sets along with cell centroids (crosses) overlayed with cell segmentation
1192 boundaries across five z stacks. **b.** Estimated spatial expression pattern for genes in each of the
1193 four gene pattern clusters identified by ELLA. Upper panel shows the number and proportion of
1194 genes across four pattern clusters. Middle panel displays the estimated expression intensity for
1195 genes across clusters. Lower panel displays the estimated pattern score for genes across clusters.
1196 **c.** Example genes and cells for the four pattern clusters. One example gene is shown for each

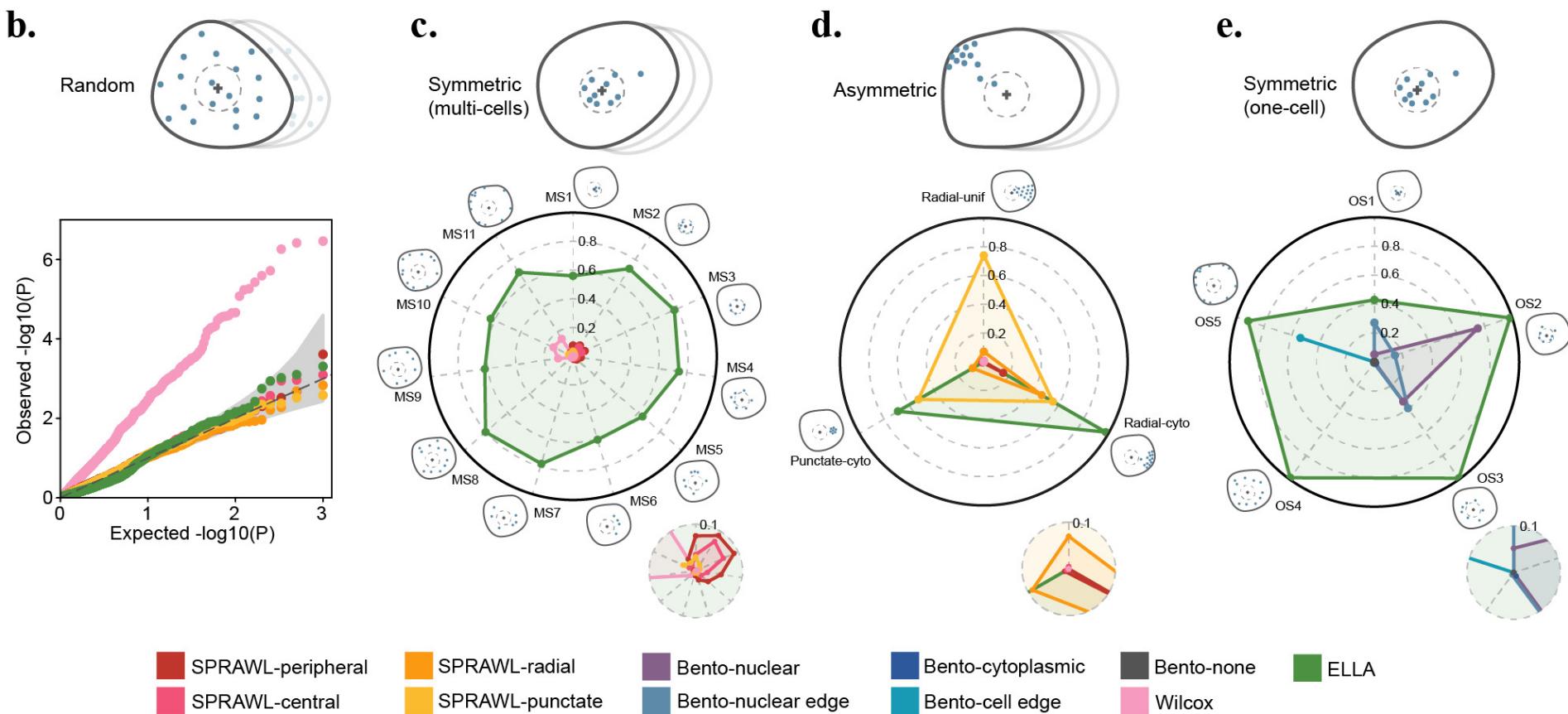
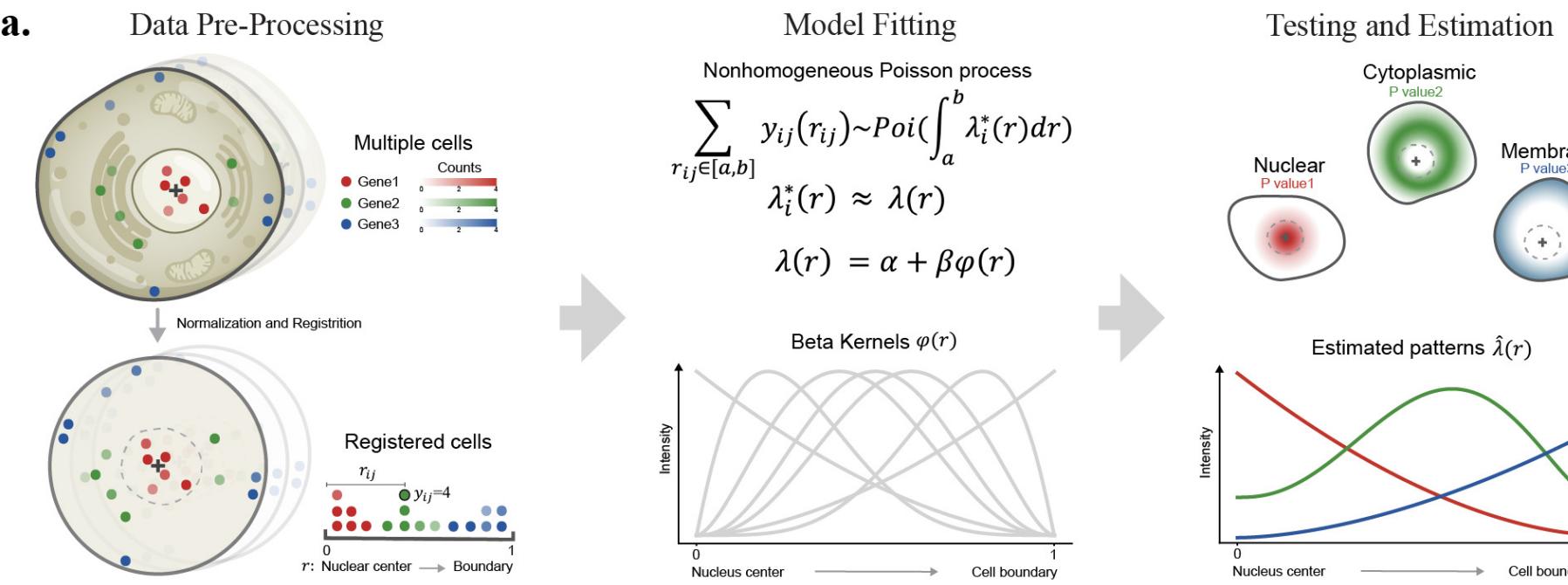
1197 pattern cluster. Upper panel shows the gene name, ELLA P value, and the estimated expression
1198 intensities overlayed on the density heat map. Each row of the density heat map visualizes the
1199 number of counts standardized by area in 10 randomly selected cells across relative positions with
1200 intervals of 0.1. Lower panel displays the expression of the corresponding genes within five
1201 selected cells, overlayed with cell boundaries and aligned nuclear centers (crosses). **d.** Bar plot
1202 shows the average sn/sc RNA ratio in the form of snRNA expression level normalized by scRNA
1203 expression level across genes in pattern clusters 1-2 (red), 3-4 (green), and non-cluster 1-2 genes
1204 (i.e. clusters 3-4 plus the nonsignificant genes; grey). Genes enriched close to nuclear center
1205 (clusters 1-2) tend to exhibit higher relative snRNA expression levels. **e.** Bar plot displays average
1206 gene length, measured by four metrics (x-axis), across genes in pattern clusters 1-2, 3-4, and non-
1207 cluster 1-2 genes. Genes enriched close to nuclear center (clusters 1-2) tend to exhibit longer gene
1208 lengths. **f.** Bar plot displays the proportions of transcription factors (TFs) for genes in pattern
1209 clusters 1-3 (orange), 4 (blue), and non-cluster 4 genes (i.e. clusters 1-3 plus the nonsignificant
1210 genes; grey). Genes enriched close to cell boundary (cluster 4) contain a lower proportion of TFs.
1211 **g** and **h**. Stem plots show the -log10 P values of the top 10 enriched gene sets in GSEA analysis
1212 for genes in pattern cluster 3 and 4 respectively. Gene sets enriched with cluster 3 or 4 genes are
1213 related to dendrites and synaptic transmission and signaling. Statistical significance for pair-wise
1214 comparisons (*: <0.05; **: <0.01; ***: <0.001) is based on Mann-Whitney U test (**d-e**) or Fisher's
1215 exact test (**f**). The error bars represent the 25th and 75th percentiles, and data points beyond this
1216 range are not included (**d-e**).
1217

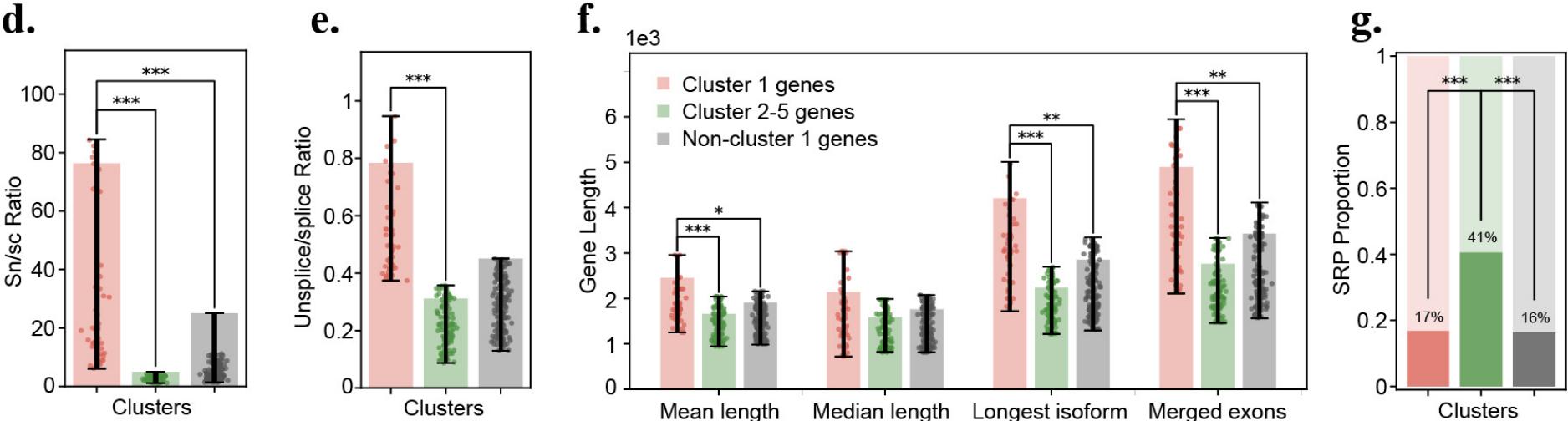
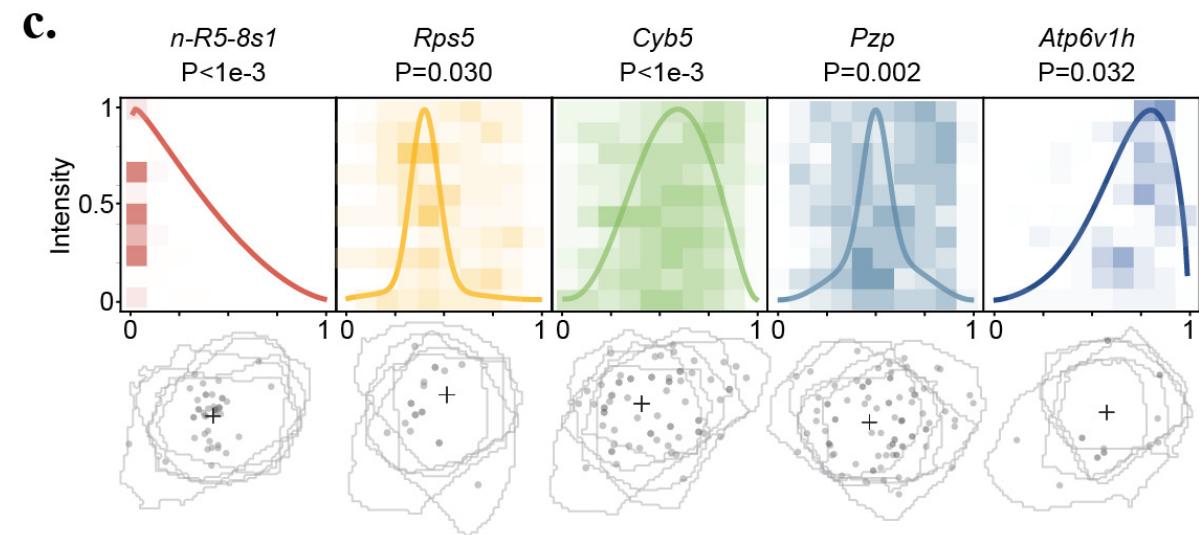
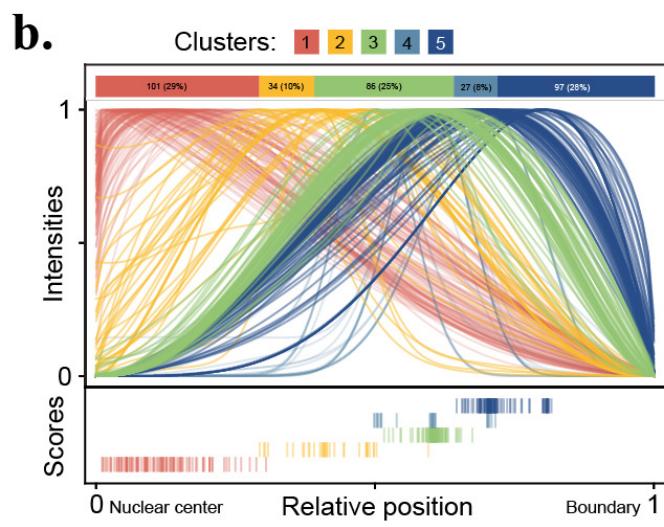
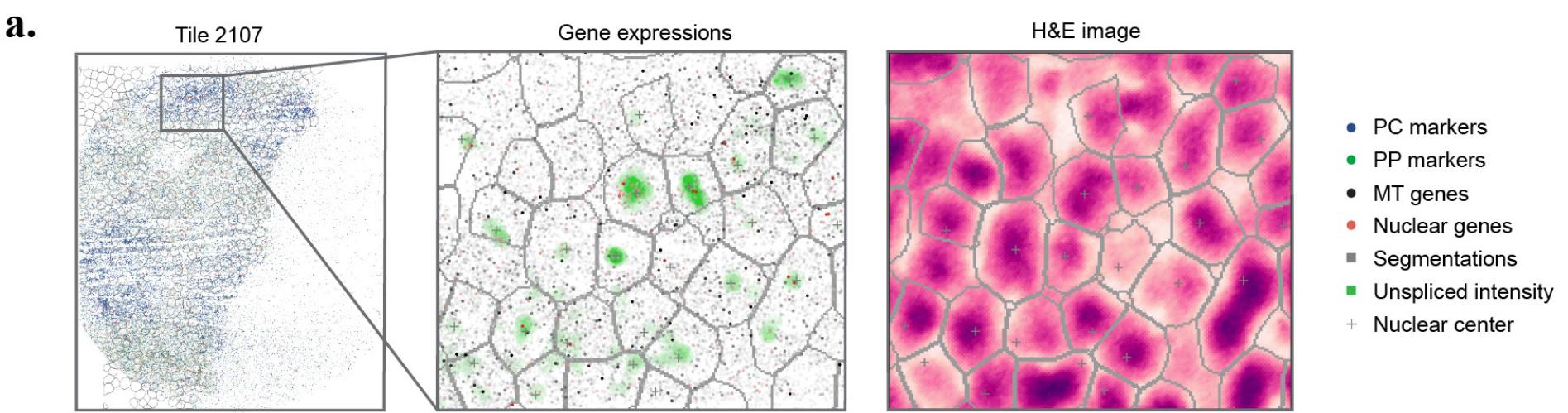
1218 References

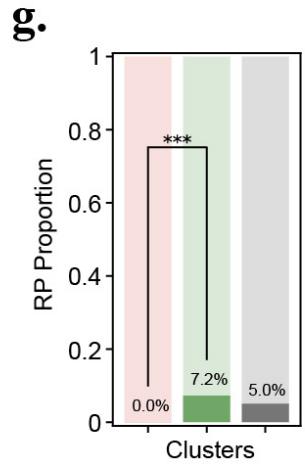
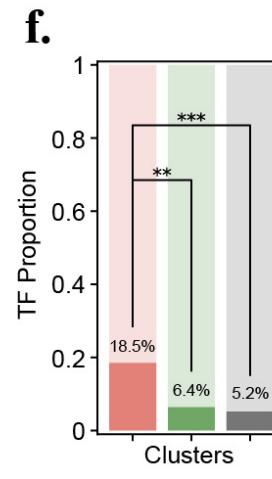
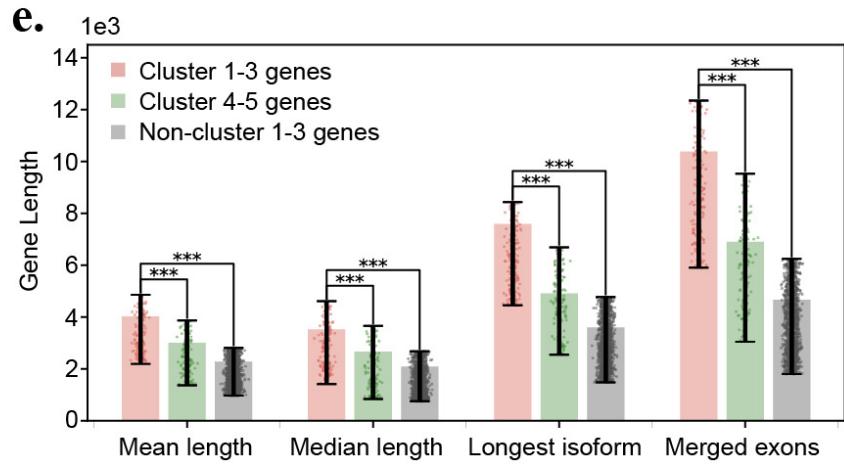
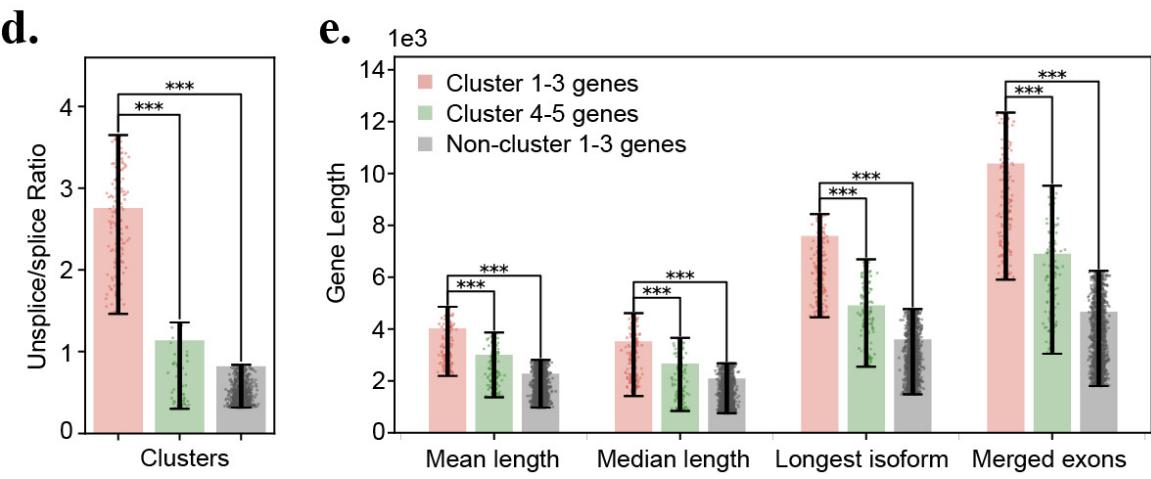
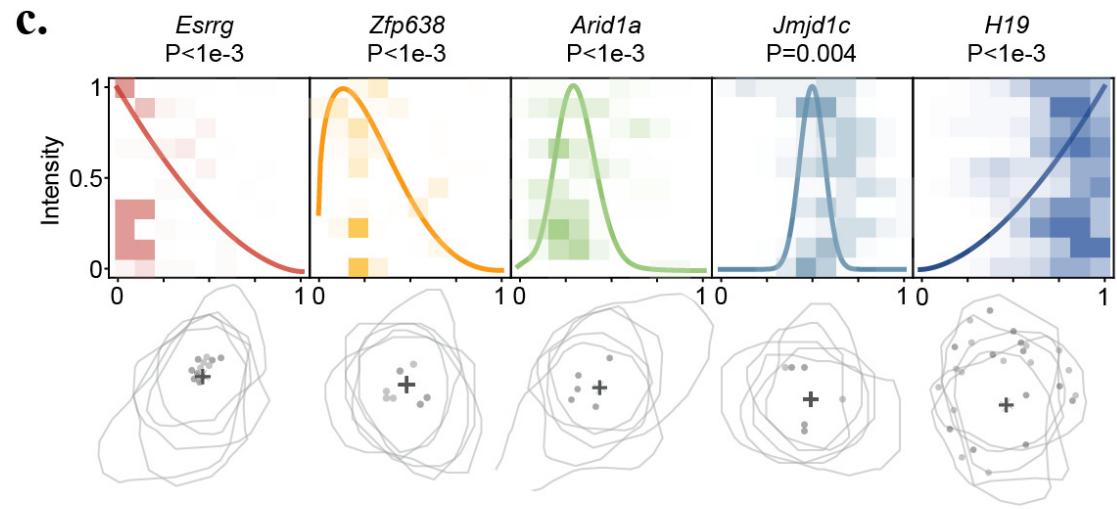
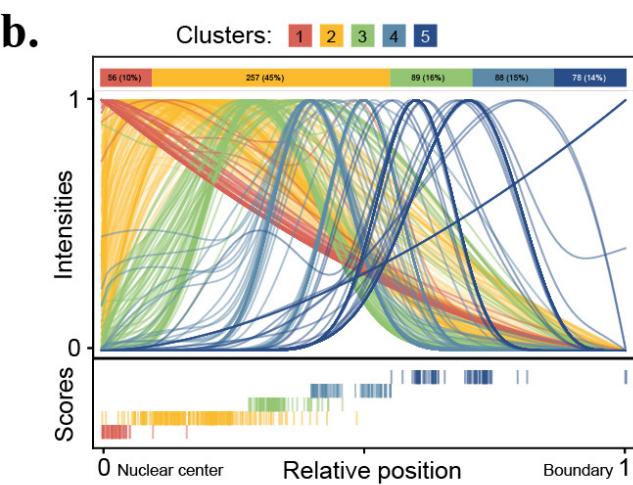
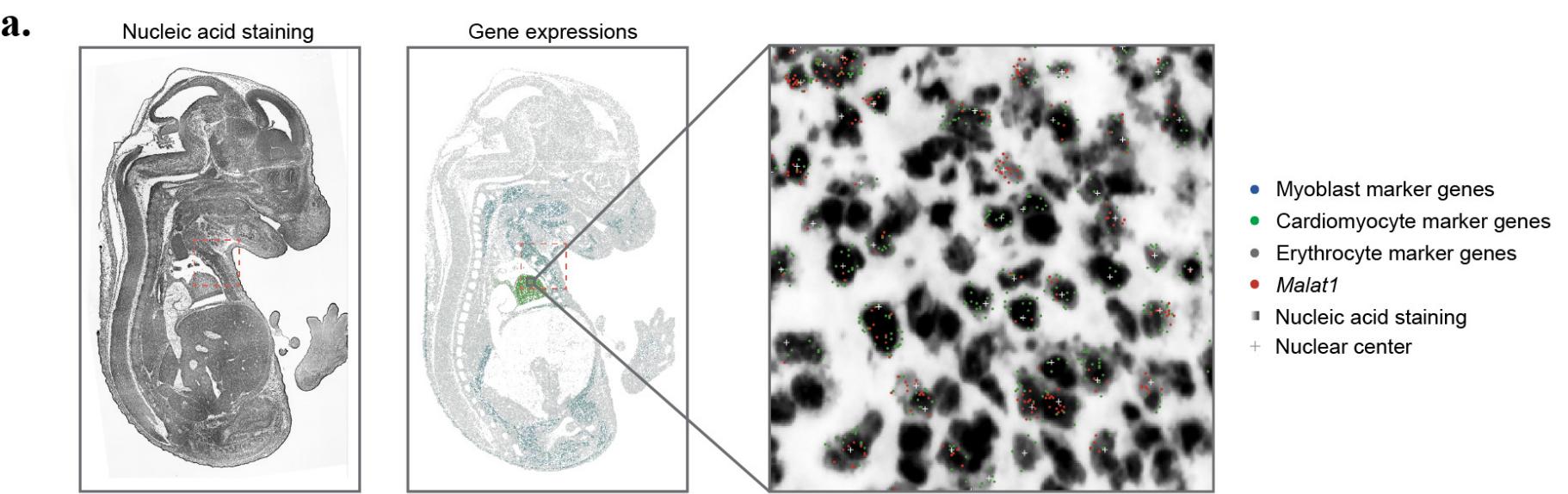
- 1219 1. Ke, R., et al., *In situ sequencing for RNA analysis in preserved tissue and cells*. Nature
1220 methods, 2013. **10**(9): p. 857-860.
- 1221 2. Lee, J.H., et al., *Fluorescent in situ sequencing (FISSEQ) of RNA for gene expression*
1222 *profiling in intact cells and tissues*. Nature protocols, 2015. **10**(3): p. 442-458.
- 1223 3. Wang, X., et al., *Three-dimensional intact-tissue sequencing of single-cell transcriptional*
1224 *states*. Science, 2018. **361**(6400): p. eaat5691.
- 1225 4. Alon, S., et al., *Expansion sequencing: Spatially precise in situ transcriptomics in intact*
1226 *biological systems*. Science, 2021. **371**(6528): p. eaax2656.
- 1227 5. Cho, C.-S., et al., *Microscopic examination of spatial transcriptome using Seq-Scope*.
1228 Cell, 2021. **184**(13): p. 3559-3572. e22.
- 1229 6. Genomics, x. *Visium HD Spatial Gene Expression: High-resolution spatial discovery at*
1230 *single-cell scale*. 2023; Available from: <https://www.10xgenomics.com/products/visium-hd>.
- 1232 7. Schott, M., et al., *Open-ST: High-resolution spatial transcriptomics in 3D*. Cell, 2024.
1233 **187**(15): p. 3953-3972. e26.
- 1234 8. Chen, A., et al., *Spatiotemporal transcriptomic atlas of mouse organogenesis using DNA*
1235 *nanoball-patterned arrays*. Cell, 2022. **185**(10): p. 1777-1792. e21.
- 1236 9. Chen, K.H., et al., *Spatially resolved, highly multiplexed RNA profiling in single cells*.
1237 Science, 2015. **348**(6233): p. aaa6090.
- 1238 10. Eng, C.-H.L., et al., *Transcriptome-scale super-resolved imaging in tissues by RNA*
1239 *seqFISH+*. Nature, 2019. **568**(7751): p. 235-239.
- 1240 11. Vizgen. *Combine single-cell and spatial transcriptomics analysis with MERSCOPE*
1241 *spatial imaging*. . 2024 [cited 2024 13/05]; Available from:
1242 <https://vizgen.com/products/>.
- 1243 12. NanoString. *CosMX Spatial Molecular Imager*. 2024 [cited 2024 13/05]; Available from:
1244 <https://nanostring.com/products/cosmx-spatial-molecular-imager/>.
- 1245 13. genomics, X. *Xenium In Situ*. 2024 [cited 2024 13/05]; Available from:
1246 [https://www.10xgenomics.com/platforms/xenium?utm_medium=search&utm_source=go
gle&utm_content=website-page&utm_campaign=7011P000001Pw8ZQAS&gad_source=1](https://www.10xgenomics.com/platforms/xenium?utm_medium=search&utm_source=google&utm_content=website-page&utm_campaign=7011P000001Pw8ZQAS&gad_source=1).
- 1249 14. Buxbaum, A.R., G. Haimovich, and R.H. Singer, *In the right place at the right time: visualizing and understanding mRNA localization*. Nature reviews Molecular cell
1250 biology, 2015. **16**(2): p. 95-109.
- 1252 15. Lawrence, J.B. and R.H. Singer, *Intracellular localization of messenger RNAs for*
1253 *cytoskeletal proteins*. Cell, 1986. **45**(3): p. 407-415.
- 1254 16. Taliaferro, J.M., E.T. Wang, and C.B. Burge, *Genomic analysis of RNA localization*.
1255 RNA biology, 2014. **11**(8): p. 1040-1050.
- 1256 17. Martin, K.C. and A. Ephrussi, *mRNA localization: gene expression in the spatial*
1257 *dimension*. Cell, 2009. **136**(4): p. 719-730.
- 1258 18. Romo, L., E.S. Mohn, and N. Aronin, *A fresh look at Huntington mRNA processing in*
1259 *Huntington's disease*. Journal of Huntington's Disease, 2018. **7**(2): p. 101-108.
- 1260 19. Mah, C.K., et al., *Bento: a toolkit for subcellular analysis of spatial transcriptomics data*.
1261 Genome Biology, 2024. **25**(1): p. 82.
- 1262 20. Bierman, R., et al., *Statistical analysis supports pervasive RNA subcellular localization*
1263 *and alternative 3'UTR regulation*. eLife, 2023. **12**.

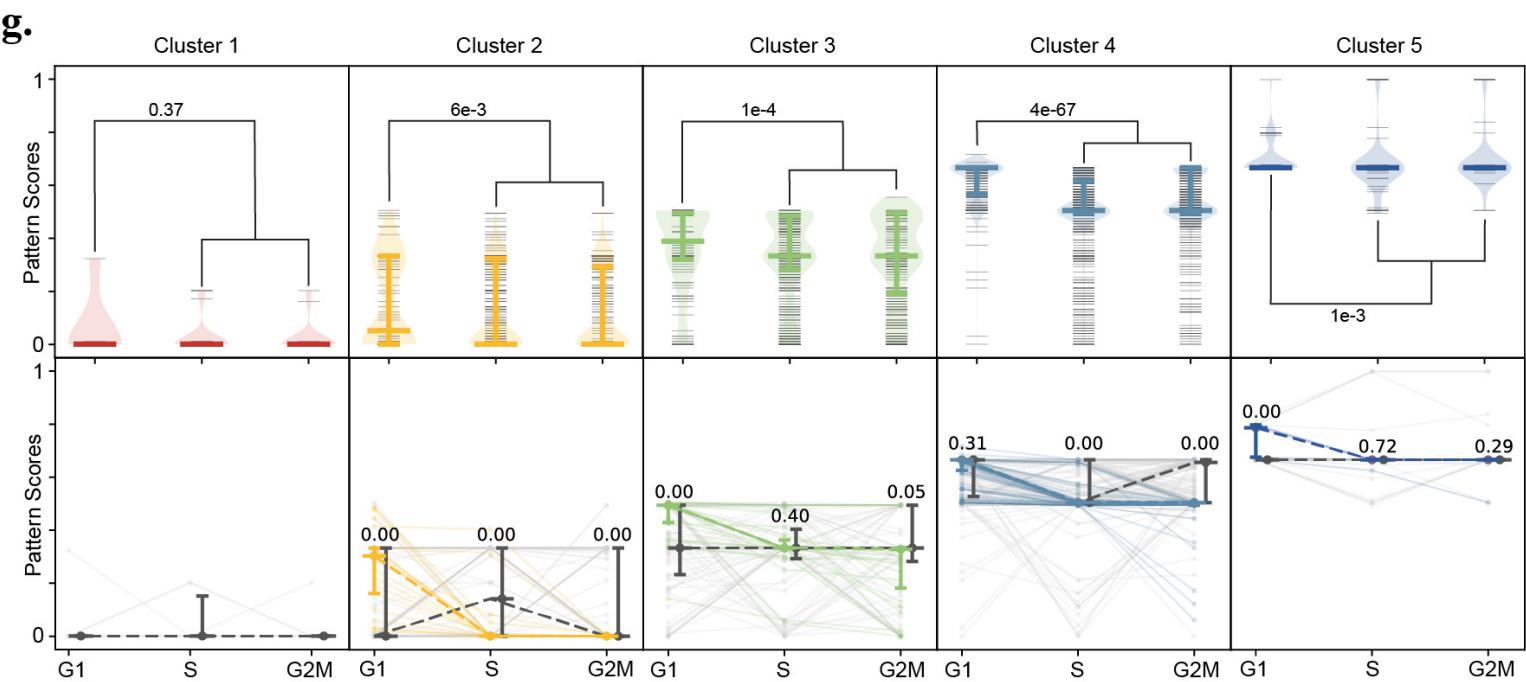
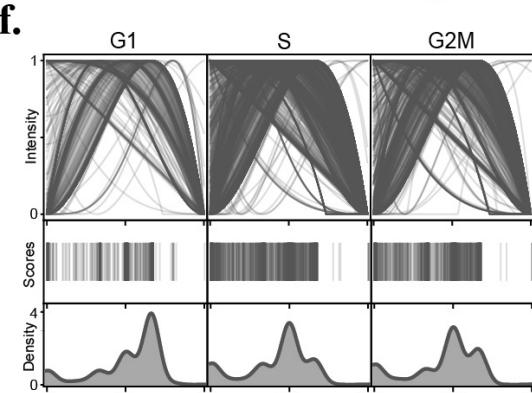
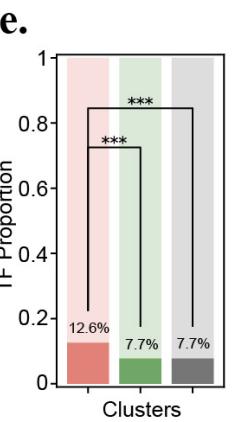
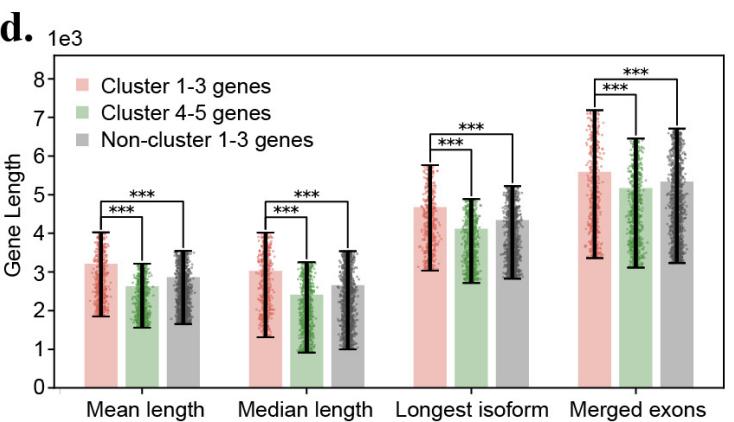
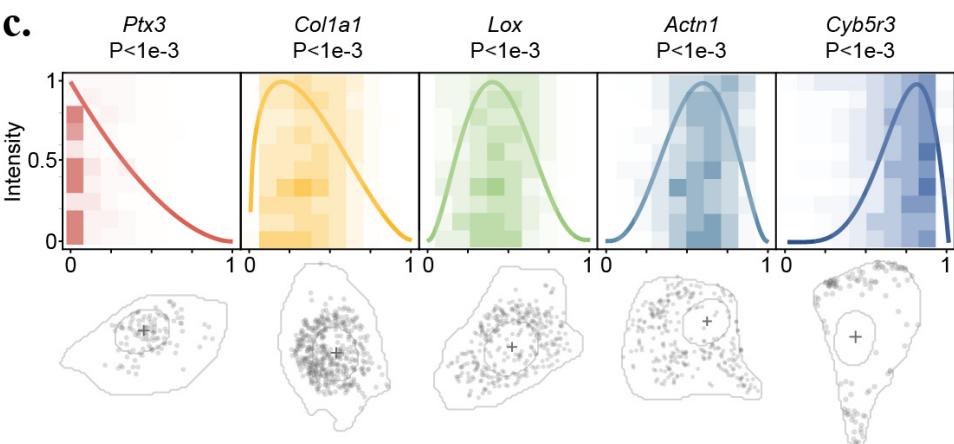
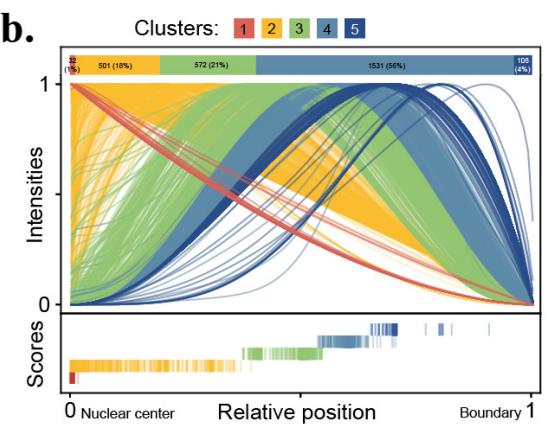
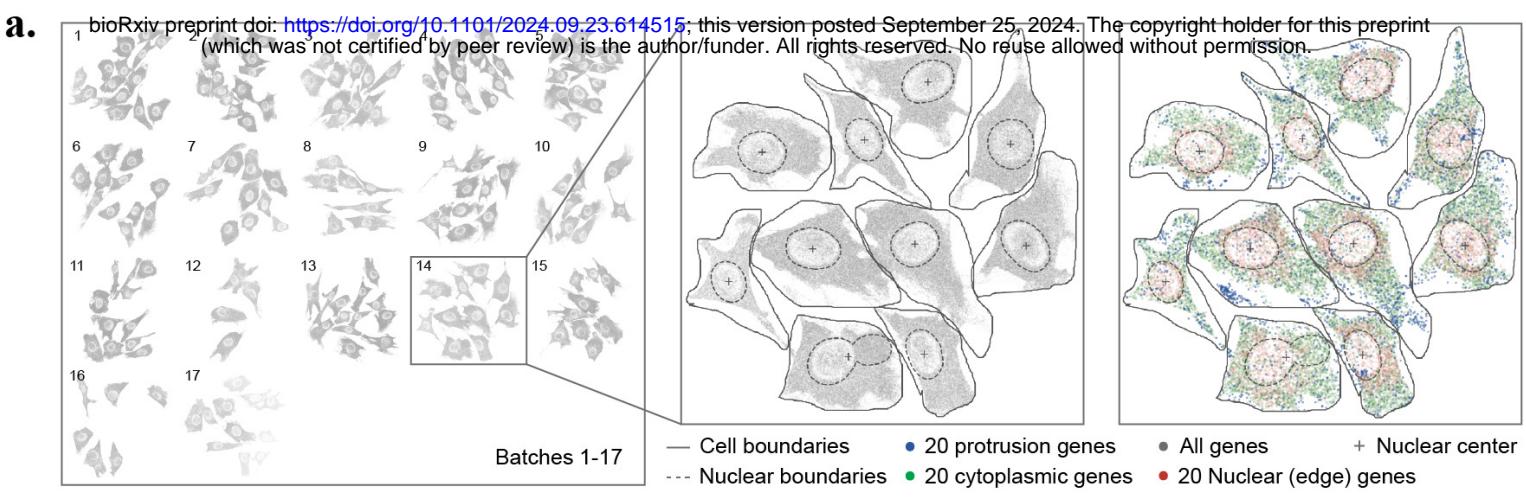
- 1264 21. Xia, C., et al., *Spatial transcriptome profiling by MERFISH reveals subcellular RNA*
1265 *compartmentalization and cell cycle-dependent gene expression*. Proceedings of the
1266 National Academy of Sciences, 2019. **116**(39): p. 19490-19499.
- 1267 22. Zhang, M., et al., *Molecularly defined and spatially resolved cell atlas of the whole*
1268 *mouse brain*. Nature, 2023. **624**(7991): p. 343-354.
- 1269 23. Cho, C.-S., et al., *Concurrent activation of growth factor and nutrient arms of mTORC1*
1270 *induces oxidative liver injury*. Cell Discovery, 2019. **5**(1): p. 60.
- 1271 24. mrichter23, *snRNA-seq2_young*. 2021.
- 1272 25. Nagai, K., et al., *Structure, function and evolution of the signal recognition particle*. The
1273 EMBO journal, 2003. **22**(14): p. 3479-3485.
- 1274 26. Abugessaisa, I., et al., *FANTOM5 transcriptome catalog of cellular states based on*
1275 *Semantic MediaWiki*. Database, 2016. **2016**: p. baw105.
- 1276 27. Solnestam, B.W., et al., *Comparison of total and cytoplasmic mRNA reveals global*
1277 *regulation by nuclear retention and miRNAs*. BMC genomics, 2012. **13**(1): p. 1-9.
- 1278 28. Singh, D.K. and K.V. Prasanth, *Functional insights into the role of nuclear-retained long*
1279 *noncoding RNAs in gene expression control in mammalian cells*. Chromosome research,
1280 2013. **21**: p. 695-711.
- 1281 29. Halpern, K.B., et al., *Nuclear retention of mRNA in mammalian tissues*. Cell reports,
1282 2015. **13**(12): p. 2653-2662.
- 1283 30. Kleshchevnikov, V. *Single-nucleus RNA-seq from adult mouse brain sections paired to*
1284 *10X Visium spatial RNA-seq*. 2021 [cited 2024; Available from:
1285 <https://www.ebi.ac.uk/biostudies/arrayexpress/studies/E-MTAB-1115>.
- 1286 31. Kingma, D.P. and J. Ba, *Adam: A method for stochastic optimization*. arXiv preprint
1287 arXiv:1412.6980, 2014.
- 1288 32. Miller, J.J., *Asymptotic properties of maximum likelihood estimates in the mixed model of*
1289 *the analysis of variance*. The Annals of Statistics, 1977: p. 746-762.
- 1290 33. Liu, Y., et al., *ACAT: a fast and powerful p value combination method for rare-variant*
1291 *analysis in sequencing studies*. The American Journal of Human Genetics, 2019. **104**(3):
1292 p. 410-421.
- 1293 34. Pillai, N.S. and X.-L. Meng, *An unexpected encounter with Cauchy and Lévy*. 2016.
- 1294 35. Raftery, A.E., D. Madigan, and J.A. Hoeting, *Bayesian model averaging for linear*
1295 *regression models*. Journal of the American Statistical Association, 1997. **92**(437): p.
1296 179-191.
- 1297 36. Kingman, J.F.C., *Poisson processes*. Vol. 3. 1992: Clarendon Press.
- 1298 37. Lab, L. *Seq-Scope Resources*. 2019 [cited 2024; Available from:
1299 <https://lee.lab.medicine.umich.edu/seq-scope>.
- 1300 38. Stringer, C., et al., *Cellpose: a generalist algorithm for cellular segmentation*. Nature
1301 methods, 2021. **18**(1): p. 100-106.
- 1302 39. Cheng, G. *Spatiotemporal transcriptomic atlas of mouse organogenesis using DNA*
1303 *nanoball patterned arrays*. 2021 [cited 2024; Available from:
1304 <https://db.cngb.org/search/project/CNP0001543/>.
- 1305 40. *MOSTA: Mouse Organogenesis Spatiotemporal Transcriptomic Atlas*. 2024 [cited 2024;
1306 Available from: <https://db.cngb.org/stomics/mosta/>.
- 1307 41. Mah, C. *Bento spatial AnnData formatted datasets*. 2021 [cited 2024; Available from:
1308 https://figshare.com/articles/dataset/Bento_spatial_AnnData_formatted_datasets/1510923_6/2.

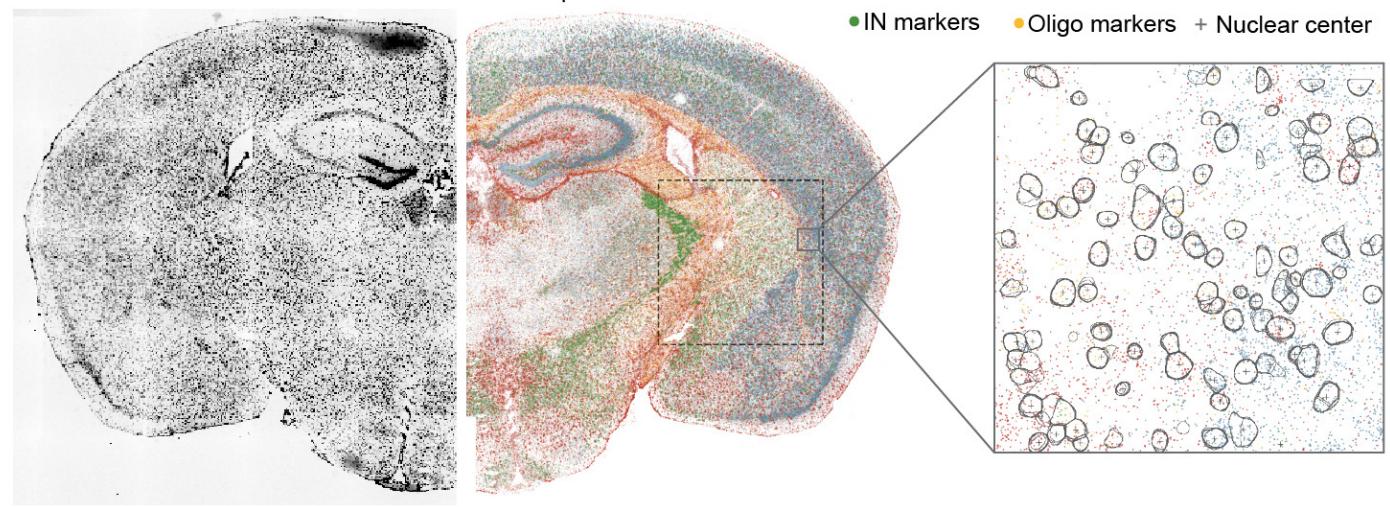
- 1310 42. Laboratory, X.Z. *Spatially resolved single-cell transcriptomics datasets acquired using*
1311 *MERFISH on the adult whole mouse brain*. 2023; Available from:
1312 <https://download.brainimagelibrary.org/29/3c/293cc39ceea87f6d/>.
- 1313 43. Cui, M., *Introduction to the k-means clustering algorithm based on the elbow method*.
1314 Accounting, Auditing and Finance, 2020. **1**(1): p. 5-8.
- 1315 44. Li, H.-D., C.-X. Lin, and J. Zheng, *GTFtools: a software package for analyzing various*
1316 *features of gene models*. Bioinformatics, 2022. **38**(20): p. 4806-4808.
- 1317 45. Savojardo, C., et al., *DeepSig: deep learning improves signal peptide detection in*
1318 *proteins*. Bioinformatics, 2018. **34**(10): p. 1690-1696.
- 1319 46. Gencode. *Gencode*. 2024 [cited 2024 13/05]; Available from:
1320 <https://www.gencodegenes.org/>.
- 1321
1322



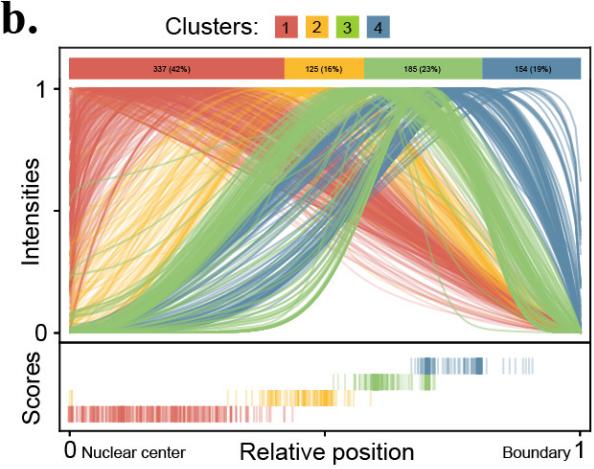




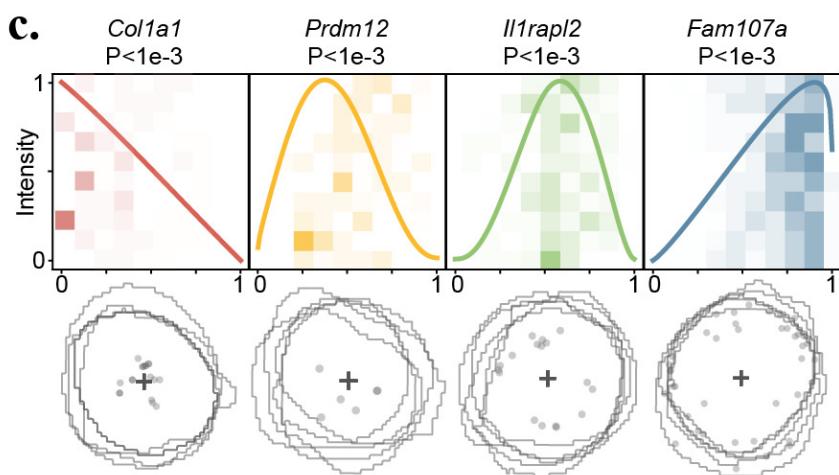




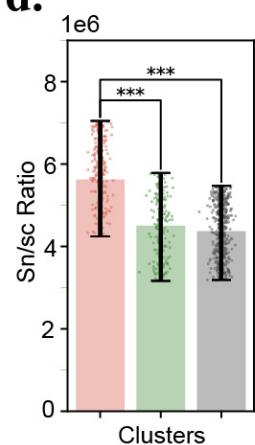
b.



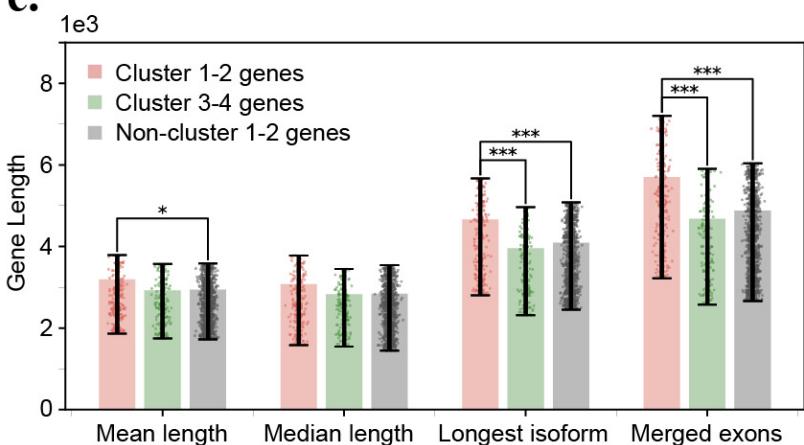
c.



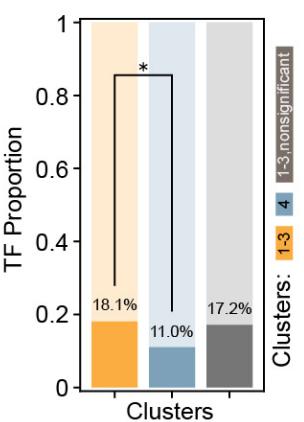
d.



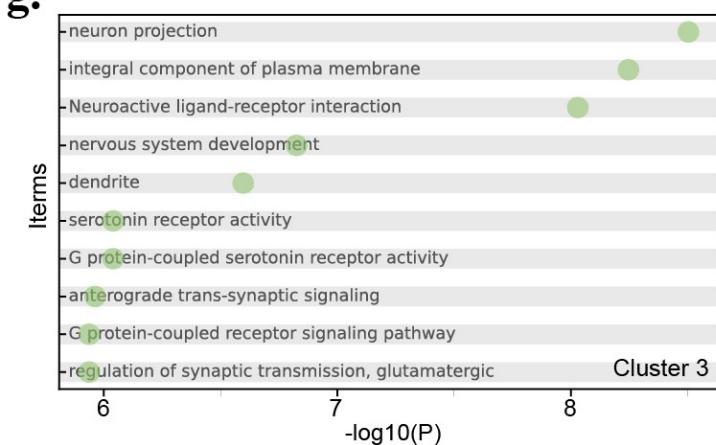
e.



f.



g.



h.

