

Intracellular Spatial Transcriptomic Analysis Toolkit (InSTAnT)

Anurendra Kumar (✉ akumar455@gatech.edu)

Georgia Institute of Technology

Alex Schrader (✉ alexws2@illinois.edu)

University of Illinois Urbana-Champaign

Ali Boroojeny (✉ ae20@illinois.edu)

University of Illinois Urbana-Champaign

Marisa Asadian (✉ asadian2@illinois.edu)

University of Illinois Urbana-Champaign <https://orcid.org/0000-0002-3554-2278>

Juyeon Lee (✉ julee@millikin.edu)

University of Illinois Urbana-Champaign

You Song (✉ yjsong2@illinois.edu)

University of Illinois Urbana-Champaign

Sihai Zhao (✉ sdzhao@illinois.edu)

University of Illinois Urbana-Champaign

Hee-Sun Han (✉ hshan@illinois.edu)

University of Illinois Urbana-Champaign

Saurabh Sinha (✉ saurabh.sinha@bme.gatech.edu)

Georgia Institute of Technology

Article

Keywords:

DOI: <https://doi.org/>

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

[Read Full License](#)

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

1 **Intracellular Spatial Transcriptomic Analysis Toolkit (InSTAnT)**

2

3 Anurendra Kumar¹, Alex W. Schrader⁴, Ali Ebrahimpour Boroojeny⁶, Marisa Asadian⁴, Juyeon Lee⁴, You
4 Jin Song⁷, Sihai Dave Zhao^{5,8*}, Hee-Sun Han^{4,8*}, Saurabh Sinha^{2,3*}

5

6 1 College of Computing, Georgia Institute of Technology, Atlanta, GA, 30332, USA

7 2 The Wallace H. Coulter Department of Biomedical Engineering, Georgia Institute of Technology,

8 Atlanta, GA, 30332, USA

9 3 H. Milton Stewart School of Industrial & Systems Engineering, Georgia Institute of Technology, Atlanta,
10 GA, 30318, USA

11 4 Department of Chemistry, University of Illinois Urbana-Champaign, Urbana, IL, 61801, USA

12 5 Department of Statistics, University of Illinois Urbana-Champaign, Urbana, IL, 61820, USA

13 6 Department of Computer Science, University of Illinois Urbana-Champaign, Urbana, IL, 61801, USA

14 7 Department of Cell and Developmental Biology, University of Illinois Urbana-Champaign, Urbana, IL,
15 61801, USA

16 8 Carl R. Woese Institute for Genomic Biology, University of Illinois Urbana-Champaign, Urbana, IL,
17 61801, USA

18 * Corresponding Authors

19

20

Abstract

21 Imaging-based spatial transcriptomics technologies such as MERFISH offer snapshots of cellular
22 processes in unprecedented detail, but new analytic tools are needed to realize their full potential.
23 We present InSTAnT, a computational toolkit for extracting molecular relationships from spatial
24 transcriptomics data at the intra-cellular resolution. InSTAnT detects gene pairs and modules with
25 interesting patterns of mutual co-localization within and across cells, using specialized statistical
26 tests and graph mining. We showcase the toolkit on datasets profiling a human cancer cell line
27 and hypothalamic preoptic region of mouse brain. We performed rigorous statistical assessment
28 of discovered co-localization patterns, found supporting evidence from databases and RNA
29 interactions, and identified subcellular domains associated with RNA-colocalization. We identified
30 several novel cell type-specific gene co-localizations in the brain. Intra-cellular spatial patterns
31 discovered by InSTAnT mirror diverse molecular relationships, including RNA interactions and
32 shared sub-cellular localization or function, providing a rich compendium of testable hypotheses
33 regarding molecular functions.

34

Introduction

36 A grand challenge in biology is to understand how molecules and cells cooperatively perform
37 higher-level processes and how these processes are coordinated to perform life functions. An
38 emerging approach to this question involves using single-cell sequencing technologies which
39 allows profiling of cellular composition and states at unprecedented resolution^{1,2}. Spatial omics
40 technologies further bolster this approach by characterizing the spatial organization of molecules
41 and cells, providing insights into their functional organization. Various analytic tools have been
42 developed to extract biological insights from spatial data, such as detecting spatially variable
43 genes^{3,4}, identifying spatial domains and their cellular compositions⁵⁻⁷, reconstructing spatial
44 gradients in developing organs⁸, or inferring cell-cell interactions^{9,10}. Most of these efforts,
45 however, have focused on cell-level or coarser resolution analyses. For grid-based spatial
46 encoding technologies, such as Visium¹¹ or DBiT-seq¹², the resolution is limited by grid size, which
47 is often larger than cells. Even with single-molecule resolution technologies¹³⁻¹⁷, tissue-scale
48 analyses mostly set the unit of analysis to be a cell^{6,7}. The focus on cell-level analyses is likely
49 due to the straightforward interpretations they provide, such as cellular arrangements around
50 diseased phenotypes¹⁸, cellular interactions⁹, and spatial context-dependent cell functions¹⁹.

51

52 Analyzing subcellular patterns of transcriptome expression can add new dimensions to our
53 understanding of cell functions. RNA localization underlies important cellular processes such as
54 transcriptional regulation²⁰⁻²², translational regulation²³, and protein localization^{24,25}. The few
55 studies that perform subcellular analyses on spatial transcriptomics data show exciting potential.
56 For instance, Xia et al²⁶ estimated RNA velocity based on the relative distribution of genes in
57 nuclei versus cytoplasm^{27,28} while Bento²⁹, a recently proposed analytical toolkit, identifies

58 subcellular domains where a gene tends to appear and was used to explore molecular
59 interactions involving RNA Binding Proteins (RBP)^{30,31}.

60
61 Despite this initial progress, the subcellular spatial landscape of RNA molecules remains largely
62 unexplored, especially for single-molecule resolution maps, which tap into a new dimension of
63 spatial architecture: spatial organization of molecules in a cell. A facet of the subcellular spatial
64 landscape that naturally merits attention is RNA-RNA proximity. Molecular interactions are
65 mediated by physical contacts; thus, the distance profile of molecular pairs can be used to infer
66 potentially interacting pairs. More broadly, RNA-RNA proximity may arise due to various reasons:
67 direct interactions between molecules, interactions with common mediator molecules, and
68 interactions with a subcellular structure, etc. Each of these sources of proximity in turn implicates
69 biological functions and molecular mechanisms. Even though single-molecule resolution spatial
70 transcriptomics data offer an unprecedented window into this world of sub-cellular organization,
71 there are no analysis tools to probe these phenomena in a large-scale and unbiased fashion.

72
73 Here, we introduce Intra-cellular Spatial Transcriptomics Analysis Toolkit (InSTAnT), a set of
74 methods for extracting subcellular localization patterns of RNA. It identifies gene pairs whose
75 transcripts tend to appear within distance d significantly more than by chance (“ d -colocalized
76 pairs”) and reports the cellular domains where they appear. Additional modules characterize the
77 d -colocalized pairs by their cell-type specificity and tissue-scale spatial modulation, and also
78 identify colocalizing gene modules. InSTAnT employs formal statistical procedures to account for
79 various sources of confounding such as overall transcript abundance, which is critical for
80 highlighting gene pairs whose transcript-proximity has biological implications. Demonstrative
81 applications of the InSTAnT toolkit to MERFISH data on a human osteosarcoma cell line and on
82 mouse hypothalamic preoptic region identified hundreds of d -colocalized gene pairs with low
83 estimated false positive rates and high reproducibility between replicates and data sources. The
84 identified gene pairs exhibit biologically relevant higher order characteristics such as specificity to
85 cell types or non-random spatial distribution in the tissue sample. We also found evidence of their
86 possible relationship to RNA-RNA or RNA-protein interactions, pathway-level co-functionality,
87 and localization to domains such as nuclear speckles. Our results suggest that InSTAnT can
88 recover known biology and generate new hypotheses about the functional role of RNA spatial
89 localization. We believe that the statistical concept of d -colocalization introduced in this work will
90 serve as a fundamental unit of subcellular spatial transcriptomics analyses, similar to how co-
91 expression analysis has served as a core concept of transcriptomics analysis.

92

93 RESULTS

94

95 Overview of InSTAnT

96 InSTAnT is a suite of statistical tools for spatial transcriptomics analysis at sub-cellular resolution.
97 It can discover intracellular spatial patterns involving transcripts of multiple genes, leading to
98 hypotheses regarding their functional relationships. At its heart is a statistical test to detect
99 "proximal pairs" of genes by analyzing the spatial coordinates of transcripts of a set of genes
100 within that cell, available from single-molecule resolution spatial transcriptomics technologies¹³⁻¹⁷.
101 Specifically, the "Proximal Pairs" (PP) test determines if transcripts of a gene pair, in a given cell,
102 are located within a distance threshold d significantly more often than expected by chance (**Figure**
103 **1a**). The null expectation may vary from cell to cell, depending on cell size and RNA density, so
104 it is calculated empirically based on the distances between all detected pairs of transcripts in a
105 cell regardless of gene identities. The test provides a p-value for each gene pair, representing its
106 departure from this expectation (Methods). The scale parameter d is user-configurable, allowing
107 the user to probe the spatial texture at different scales. The PP test can be implemented in either
108 two- or three-dimensions (PP-3D), depending on whether or not data are available from multiple
109 z-planes (Methods).

110

111 We define a " d -colocalized" gene pair to be a pair that is detected as proximal pair by the PP test
112 in significantly many cells. This gives us increased confidence in a spatial relationship between
113 the two genes. Like other statistical phenomena such as differential expression of a gene or co-
114 expression of a gene pair, d -colocalization may serve as a starting point for discovery of
115 underlying biological relationships. To detect d -colocalization, InSTAnT provides a test called
116 "Conditional Poisson Binomial" (CPB) test that assigns a p-value to a gene pair based on the
117 number of cells in which it is found to be a proximal pair. This test is based on a Poisson Binomial
118 distribution and allows for the fact that different cells have varying numbers of proximal pairs due
119 to varying transcript counts and spatial distributions (**Figure 1b**, Methods). Initially, we noticed
120 certain genes to feature among the reported d -colocalized pairs far more frequently, due to their
121 high expression (**Supplementary Figure 1**). The CPB test de-emphasizes pairs involving such
122 genes by adjusting the null distribution of each pair to account for the global d -colocalization
123 frequency of the involved genes (Methods).

124 Through the PP and CPB tests, InSTAnT unbiasedly identifies gene pairs with a tendency for
125 spatial proximity, at the level of individual cells (proximal pairs) and at the level of all cells (*d*-
126 colocalized pairs), respectively. The InSTAnT suite is available as a python package with routines
127 that return PP test results for every cell and CPB test results across all cells, for each gene pair.
128 To assist with biological interpretation of the detected spatial relationships, it can annotate each
129 *d*-colocalized gene pair with the cellular regions where its proximal transcripts tend to be found:
130 nuclear, peri-nuclear, cytosolic and peri-membrane. InSTAnT reports the primary and secondary
131 regions that has most PP counts for each gene pair across all cells (**Figure 1c**, Methods).
132 InSTAnT also implements additional analyses to study *d*-colocalization in intact tissue, where a
133 number of complex biological factors such as heterogeneity of cell types and interactions among
134 neighboring cells are at play. These factors may influence, or be influenced by, RNA-RNA
135 proximity patterns. InSTAnT can assess the cell-type specificity of *d*-colocalized gene pairs,
136 characterize tissue-level spatial modulation of *d*-colocalization patterns, and identify modules of
137 genes that are all frequently colocalized across multiple cells (**Figure 1c**).

138

139 **InSTAnT finds gene-gene relationships with high accuracy**

140 We first applied InSTAnT to the published MERFISH data on human osteosarcoma cells (U2-
141 OS), which profiles 130 genes in 3237 cells with an average of 1243 transcripts per cell³²
142 (Methods). Through the analysis, we identified 'proximal pairs' within each cell and '*d*-colocalized
143 pairs' across all cells with high accuracy. We calculated false positive rates (FPRs) by applying
144 InSTAnT to a random baseline dataset established by permuting the gene labels of all transcripts
145 within each cell, which recapitulates the spatial patterns of the original data but not the gene-gene
146 relationships. As shown in **Figure 2b** (blue), the PP test identifies hundreds of significant proximal
147 pairs with an estimated FPR below 10%. Smaller values of the scale parameter *d* yielded larger
148 FPR values (red and lemon, Figure 2B), suggesting lower sensitivity of the test and/or lesser
149 frequency of proximal pairs in this regime. We found similar operating characteristics for the CPB
150 test (**Figure 2c**). Throughout our paper, we use FPR to select p-value threshold for PP
151 (FPR<10%) and CPB Test (FPR<1%). We arrived at similar estimates of accuracy through an
152 entirely different approach that exploits presence of "blank" gene probes in the data (Methods and
153 **Supplementary Figure 2**). Overall, our tests suggested that hundreds of gene pairs exhibit the
154 *d*-colocalization phenomenon, out of all ~8,500 pairs possible with 130 genes.

155

156 The CPB test had sufficient power to identify 404 *d*-colocalized gene pairs at an FPR of < 1% (p
157 < 0.001), with $d = 4 \mu\text{m}$ (~5% of the diameter of an average cell) (**Supplementary Table 1**). An
158 example of a highly significant pair thus found is *THBS1-COL5A1*, with a p-value below ~1E-300,
159 the smallest number reportable by the program. This pair appeared as a proximal pair (PP test p-
160 value < 0.01) in ~74% of the 3,147 cells where both genes were detected. **Figure 2a** shows the
161 distribution of PP test p-values for this gene pair in all cells, compared to the distribution of the
162 strongest p-value in each cell after shuffling gene labels. The comparison illustrates how the CPB
163 test detects the persistent appearance of a proximal pair across many cells.

164
165 Our next assessment focused on the replicability of *d*-colocalization findings across four biological
166 replicates of the U2OS data set available from Moffitt et al.³². We identified the most significant
167 gene pairs (CPB test, $d = 4 \mu\text{m}$) in each replicate and observed that ~80% of the top 50 – 400
168 gene pairs are common between replicates (**Figure 2d**), supporting the reproducibility of the
169 reported pairs. The same assessment performed after randomizing each of the four replicate data
170 sets yielded a baseline level of ~5% or less for the replicability expected by chance.

171
172 We also tested the extent to which *d*-colocalization phenomena persist across independent
173 MERFISH experiments. For this, we generated the spatial transcriptome map of U2OS cells using
174 our home built MERFISH platform (Methods). We used InSTAnT to identify *d*-colocalized gene
175 pairs from our dataset and compared the top K (for varying values of K) gene pairs between the
176 Moffitt et al. and our data. As shown in **Figure 2e**, about 30-40% of the identified gene pairs are
177 shared between these two studies, across the range of K examined. The same analysis with
178 randomized versions of the two datasets reveals < 5% of the gene pairs to be shared between
179 studies. As another reference point, a similar comparison of the top co-expressed gene pairs
180 (detected using correlation of cellular transcript counts) shows similar or lesser extent of
181 commonality between the two studies (**Supplementary Figure 3**). Taken together, these
182 reproducibility analyses suggest that the *d*-colocalized gene pairs reported by InSTAnT capture
183 real biological phenomena or relationships.

184

185

186 **InSTAnT constructs global *d*-colocalization maps**

187

188 The 404 *d*-colocalized gene pairs found using CPB test at $d = 4 \mu\text{m}$ (**Supplementary Table 1**)
189 constitute the global *d*-colocalization map. InSTAnT provides annotations of the cellular regions
190 where each gene pair tends to colocalize, revealing perinuclear and nuclear colocalization as

191 most frequent (**Figure 3a-c, Supplementary Figure 4**). We also noted many gene pairs to
192 colocalize in the cytosolic (23) or cell periphery (16) regions (see **Figures 3d,e** for examples),
193 though far less often than the other two categories.

194
195 A d -colocalization map is expected to capture different biology at different values of d . The maps
196 created from the published U2OS data at $d = 1 \mu\text{m}$ (**Supplementary Table 2**) and $4 \mu\text{m}$ revealed
197 substantial complementarity (**Figure 3f**): while 152 pairs were common to the top 404 significant
198 pairs of either map, 197 of the pairs in the $d = 4$ map had CPB test p-value > 0.1 in the $d = 1$ map,
199 and 167 gene pairs were similarly exclusive to the $d = 1$ map. Two examples of such scale-specific
200 pairs are *FASN-DYNC1H1* (only with $d = 4$) and *CENPF-PRKCA* (only with $d = 1$). (See
201 **Supplementary Figure 5** for a more detailed report of their scale-dependence.) These results
202 illustrate scale-dependence of the colocalization phenomenon and suggest that multiple types of
203 biological relationships may underlie its detection.

204
205 The d -colocalization map probes a new type of information and may represent yet-to-be-explored
206 phenomena. Reconstructing gene-gene co-expression networks is a common analysis performed
207 with non-spatial single cell RNA-seq data³³. To test if the global d -colocalization map reflects such
208 co-expression networks or if it reveals a different type of relationship, we derived a co-expression
209 network from cell-level transcript counts in the same MERFISH data and found it to share $\sim 30\%$
210 of gene pairs with the colocalization map (Hypergeometric test p-value $6.3\text{e-}70$) (**Figure 3g**,
211 **Supplementary Table 3**). Over 70% of the pairs in either “co-expressed” or “colocalized” set
212 were exclusive to that set, suggesting that d -colocalization relationships are not revealed through
213 conventional co-expression analysis.

214
215 In addition to constructing a basic global map, InSTAnT can run the PP test in a “intra-nucleus”
216 mode where the analysis, including null distribution estimation, is limited to subnuclear transcripts.
217 This mode is critical for detecting subnuclear phenomenon. The default (whole-cell) mode
218 assumes the null distribution as uniform throughout a cell, disregarding the selective enrichment
219 of certain genes in subcellular regions. Thus, nucleus-enriched genes, such as long noncoding
220 RNAs (lncRNAs), often dominate detected co-localized pairs. For example, 89 of the 404 pairs in
221 the U2-OS global co-localization map involved the lncRNA *MALAT1*, which is the most nucleus-
222 enriched gene (89% in nucleus). The intra-nucleus mode effectively removes such bias.
223 As expected, many gene pairs detected by the whole-cell mode have far stronger p-values than
224 the intra-nucleus mode due to the greater number of transcripts examined. (**Supplementary**

225 **Figure 6**). However, we also observed a significant number of gene pairs that were assigned
226 greater statistical significance in the intra-nucleus analysis. Such pairs promise to reveal
227 biologically meaningful spatial patterns within nuclei, as might arise for instance from
228 colocalization of a gene pair to subnuclear structures, organelles and domains.

229

230 ***d*-colocalization maps suggest functional relationships in U2OS cells**

231

232 One plausible mechanism for *d*-colocalization is direct or indirect interaction between two RNAs.
233 To test this, we computed an RNA interaction score (“RRI score”) for all gene pairs using
234 RNAplex³⁴. To capture the greater proximity expected of interacting RNAs, we set *d* to 200 nm
235 (MERFISH resolution between pixels is 167nm). For each gene, we tested if its transcript tends
236 to have a higher RRI score for the RNAs of its *d*-colocalization partners (Methods) and found this
237 to be the case for eight genes out of 130 (FDR ≤ 0.2) (**Supplementary Table 4**). An example is
238 shown in **Figure 3i**, focusing on *USP9X*. In summary, this analysis suggests that RNA-RNA
239 interactions may underlie some of the relationships in a global *d*-colocalization map at a suitably
240 small value of the scale parameter.

241

242 Furthermore, we found that the *d*-colocalized gene pairs were enriched with functionally related
243 gene pairs, where we define a gene pair to be functionally related if both genes are present in the
244 same KEGG pathway or are annotated with same biological process, molecular function, or
245 cellular component GO terms (Methods) (**Figure 3h**). The highest enrichment happened with
246 molecular function GO terms, where 461 functionally related pairs and 403 *d*-colocalized pairs
247 had an overlap of 67 pairs. Interestingly, all 67 pairs in this intersection were annotated with the
248 term “protein binding”. Overall, these results suggest that *d*-colocalization of a gene pair may have
249 biological consequences such as colocalization of their protein products or protein binding to form
250 a ribonucleoprotein (RNP) complex.

251

252 The intra-nuclear analysis shows that a *d*-colocalization map can detect RNA-protein interactions
253 as well as identify subnuclear domains. The most prominent pair in the intra-nucleus analysis at
254 $d=2\ \mu\text{m}$ is *MALAT1-SRRM2*, with a CPB test p-value of $2.50\text{e-}16$ (see **Figure 3j**), while the
255 corresponding p-value in the whole-cell analysis is 0.51 (see marked point in Supplementary
256 Figure 1c). It is detected as a proximal pair in 11% of the nuclei, the most for any pair involving
257 either *SRRM2* or *MALAT1*. Notably, the *SRRM2* protein is a key marker of nuclear speckles (NS),
258 organizing NS formation via liquid condensation³⁵, and the lncRNA *MALAT1* is well known to be

259 localized to NS³⁶, suggesting that the detected intra-nuclear *d*-colocalization of these two RNAs
260 may be related to their colocalization in NS. This is an intriguing possibility though, since NS
261 localization of SRRM2 protein does not imply or necessitate a similar localization of its mRNA. To
262 see whether lncRNA *MALAT1* and mRNA *SRRM2* colocalize near NS, we co-stained *MALAT1*,
263 *SRRM2* mRNA, and SON in U2-OS cells using single molecule FISH and immunostaining (**Figure**
264 **3k-n**). For *SRRM2*, the probes were designed separately for intron and exon to distinguish pre-
265 mRNA and mRNA. As expected, all the *SRRM2* intron signals directly overlap with the *SRRM2*
266 exon signals. Consistent with the InSTAnT result, most *SRRM2* RNAs are *d*-colocalized with
267 *MALAT1* in *SRRM2* positive cells (99±1%, N=13 cells). The overlaid SON signals show that most
268 *d*-colocalized *MALAT1-SRRM2* pairs are within 1 µm distance from NS (92±4% for *SRRM2* exon
269 and 88±8% for *SRRM2* intron). It is well known that SRRM2 protein signals overlaps with SON
270 signals³⁵; thus, our result shows the *d*-colocalization of *SRRM2* mRNA and pre-mRNA with
271 SRRM2 proteins in nucleus. Further, these results suggest that *d*-colocalization maps can be
272 used to infer subcellular domains, such as NS.

273

274 **InSTAnT analysis of brain MERFISH data reveals cell type-specific spatial patterns**

275

276 We next used InSTAnT to analyze MERFISH data³⁷ on 5149 cells from the hypothalamic preoptic
277 region in mouse. This brain dataset includes nine different cell types (**Figure 4a**), so InSTAnT's
278 *d*-colocalization maps can be used to give additional insight into cell type differences. The data
279 feature seven z-planes and were thus analyzed with the PP-3D test of proximal pairs. We set the
280 scale parameter *d* to 2 µm, corresponding to ~5% of average cell diameter. The analysis identified
281 474 gene pairs with CPB test p-value < 1e-5 (**Supplementary Table5**) (estimated FPR < 1%).
282 This map was further processed with downstream InSTAnT modules for cell type specificity and
283 spatial modulation.

284

285 InSTAnT uses a sequence of statistical tests (**Methods, Figure 4b**) to place *d*-colocalized gene
286 pairs into one of three categories of cell type-specificity. Category 3 comprises pairs that were not
287 associated with any cell type (Bonferroni corrected hypergeometric p-value ≥ 0.05 ,
288 **Supplementary Table 5**). Pairs that did appear as proximal pairs more frequently in some cell
289 types than expected were further divided into two classes – those where cell type specificity may
290 arise simply because one of the genes in the pair is expressed specifically in that cell type
291 (Category 1) and those whose association goes beyond what would be expected from the cell-
292 type specificity of either gene's expression (Category 2) (Methods). We identified 5 gene pairs in

293 Category 2, specific to inhibitory neurons, excitatory neurons, and endothelial cells
294 (**Supplementary Table 6**), while 203 pairs fell in Category 1.

295
296 Gene pairs with strong *d*-colocalization signal in each category captured interesting biological
297 processes involving their counterpart protein-protein interactions. In Category 1, the genes *Aqp4*
298 (Aquaporin 4), *Cxcl14* (CXC motif chemokine ligand 14) and *Mlc1* (Modulator of VRAC current 1)
299 show strong pairwise *d*-colocalization associated with astrocytes (CPB test p-value < 1.9E-149,
300 Hypergeometric p-value of cell type association < 2.23E-12). As illustrated for the pair *Cxcl14*-
301 *Mlc1* in **Figure 4c**, these pairs are frequently colocalized in cells of most types, but with a higher
302 frequency in astrocytes, leading to the statistically detected specificity. *Cxcl14* transcripts are
303 known to be enriched in and possibly locally translated in peripheral astrocyte processes
304 (PAPs)³⁸. We speculate that *Mlc1* transcripts are also subject to local translation in PAPs, leading
305 to the *d*-colocalization of *Cxcl14* and *Mlc1*. Additionally, MLC1 protein forms a complex with AQP4
306 in cultured astrocytes³⁹ and localizes to the cell membrane^{38,40} providing the functional implication
307 of *Mlc1-Aqp4* RNA *d*-colocalization.

308
309 In Category 2 (**Figure 4d**), transmembrane proteins *Gpr165* (G protein-coupled receptor 165) and
310 *uc011zyl.1* (adhesion molecule with Ig like domain 2) form the *d*-colocalized pair most significantly
311 associated with inhibitory neurons, while *Gpr165* and *Omp* (Olfactory marker protein, known to
312 be involved in olfactory signaling processes⁴¹ form a *d*-colocalized pair specific to excitatory
313 neurons (**Supplementary Figure 7**). This example illustrates that different *d*-colocalized pairs
314 involving a common gene (*Gpr165*) can statistically mark different cell types. We observed 61 *d*-
315 colocalized pairs in Category 3. *Aldh111-Mlc1* is the strongest pair (CPB test p-value: 1.6E-169),
316 detected as a proximal pair in 7% of all cells, but these cells are not enriched for any one cell type
317 (**Figure 4e**). This example suggests that *d*-colocalization can capture biological relationships that
318 transcend any cell type-specific function of the constituent genes.

319
320
321 We also identified *d*-colocalized gene pairs that marked cellular function, in particular inhibitory
322 versus excitatory neurons (**Supplementary Table 6**). For instance, *Esr1* (estrogen receptor 1)
323 and *Npy2r* (Neuropeptide Y receptor Y2) are *d*-colocalized specifically in inhibitory neurons
324 compared to excitatory neurons (p-value 5.9E-8, see **Figure 4f**). Prior work shows that the
325 expression of these two genes underlies a social behavioral switch in virgin mice via activation of

326 a specific subtype of neurons⁴², suggesting the functional implication of *Esr1-Npy2r* *d*-
327 colocalization.

328

329

330 **InSTAnT reveals tissue-level spatial modulation of *d*-colocalization patterns**

331

332 Brain tissue is well-known to be spatially heterogeneous, so we applied InSTAnT's spatial
333 modulation analyses to study how *d*-colocalization varies across the mouse hypothalamic preoptic
334 region. Such tissue-level spatial modulation has been reported for individual gene expression^{3,4}.
335 In contrast, here we used InSTAnT to identify spatial patterns of transcript colocalization.

336

337 The analysis is based on a probabilistic model for calculating data likelihood under the hypothesis
338 of spatial modulated *d*-colocalization, for a specific gene pair. The probabilistic model (**Figure 5a**)
339 examines whether the PP test detects significant colocalization in a cell and assumes that the
340 probability of this happening depends on observed colocalization in neighboring cells, rewarding
341 spatially clustered distributions of cells that support colocalization. Such a model is then
342 contrasted with a null model lacking spatial dependence, resulting in a log likelihood ratio (LLR)
343 score being assigned to each gene pair in the *d*-colocalization map. Pairs above a threshold
344 (obtained using randomization of data) are then designated as spatially modulated. This yielded
345 99 spatially modulated pairs out of the 474 pairs in the global map (**Supplementary Table 7**). A
346 similar analysis for U2OS data yielded 11 gene pairs out of 404 *d*-colocalized pairs in the
347 corresponding global map. The stark difference in extent of spatial modulation detected is
348 expected, since intercellular communication plays a greater role in the biology underlying the brain
349 data compared to cell line data.

350

351 Forty nine of the 99 spatially modulated pairs in the brain data exhibited *d*-colocalization in a cell
352 type-specific manner (p-value 5E-6, Bonferroni corrected p-value < 0.05). For instance, the gene
353 pair *Sgk1-Ttyh2* – the strongest spatially modulated pair (LLR 305, **Figure 5f**) – colocalizes far
354 more frequently in mature oligodendrocytes than others (Hypergeometric test p-value 1.5e-248,
355 **Figure 5b**). *Sgk1* is a serine/threonine-protein kinase that mediates oligodendrocyte plasticity in
356 mouse in response to stress^{43,44} and regulates several ion channels⁴⁵, while *Ttyh2* is a chloride
357 channel noted for its transcriptional response to chronic stress in mouse oligodendrocytes⁴⁶. It is
358 plausible that the oligodendrocyte-specific *d*-colocalization results from a co-functional
359 relationship between these two genes. The pair *Slc17a6-Syt4* is the second strongest spatially
360 modulated *d*-colocalized pair (LLR 191), detected in six different cell types but highly specific to

361 excitatory neurons (**Supplementary Figure 8**). In contrast to these two examples where *d*-
362 colocalization is significant in multiple cell types but more frequent in one cell type, the pair *Cd24a*-
363 *Mlc1* exhibits spatially modulated *d*-colocalization (LLR 79, **Figure 5g**) that is significant only in
364 ependymal cells (**Figure 5c**).

365
366 We also found 15 spatially modulated gene pairs whose *d*-colocalization is not specific to any cell
367 type (Hypergeometric test p -value > 0.05 for every cell type), the strongest being *Col25a1-Gad1*
368 (LLR 97, **Figures 5d,h**). *Col25a1* is generated by different types of neurons, i.e., inhibitory as well
369 as excitatory, and interneurons in retino-recipient regions of the mouse brain, in a *Gad1*-
370 dependent pattern⁴⁷. In summary, the above examples of spatially modulated *d*-colocalization
371 provide a rich pool of potential functional relationships for future exploration.

372

373 **InSTAnT reveals modules of genes colocalizing with each other**

374
375 We asked if the significant gene pairs found by InSTAnT point to the existence of *d*-colocalization
376 “modules”, i.e., sets of genes whose transcripts tend to occur in subcellular proximity, across
377 many cells, drawing inspiration from co-expression module discovery⁴⁸. Colocalized gene
378 modules, if found, may reflect ribonucleoprotein complex formation^{29,49} or other shared functional
379 relationships³⁰.

380
381 InSTAnT provides two complementary routines for gene module discovery. The first routine,
382 called Global Colocalization Clustering (GCC), identifies modules by representing the CPB test
383 results as a matrix of gene-gene *d*-colocalization strengths and clustering rows and columns of
384 this matrix (Methods). **Figure 6a** shows the results of such clustering for U2OS data, revealing
385 two modules (top left) whose compositions are shown in **Figure 6b**. Module M1 (spatially
386 illustrated in **Figures 6d,e**) consists of 14 genes, with 85 of 91 pairs being significantly *d*-
387 colocalized and all but one of these significant gene pairs being assigned a perinuclear region
388 annotation. Gene Ontology (GO) enrichment analysis of the module revealed shared annotations
389 (p -value < 0.05 , **Figure 6c**) related to cytoskeleton and ribonucleoprotein complexes. mRNA-
390 cytoskeletal associations have been long known to play a key role in mRNA transport and
391 targeting to specific subcellular locations, partly mediated by RBPs and ribonucleoprotein
392 complexes^{50,51}. Module M1 includes gene pairs whose protein products are known to interact,
393 e.g., FASN-SPTBN1⁵² and PRPF8-SRRM2^{53,54}. Module M1 also shares four genes with the nine-
394 gene module called “Group II” found to colocalize (at a coarser resolution) in fibroblast MERFISH
395 data¹³. The second module (M2) comprises eight genes, with 23 of 28 pairs being significantly *d*-

396 colocalized, mostly with perinuclear annotation. The module is significantly enriched for several
397 GO terms, e.g., positive regulation of cell death and receptor complex (**Supplementary Figure**
398 **9**), and its sub-cellular colocalization may thus mirror a co-functioning of its protein products.

399
400 A module reported by GCC comprises gene pairs whose *d*-colocalization is supported by many
401 cells, but these supporting cells differ for different gene pairs and very few cells may have the
402 entire module colocalized. Motivated by this, InSTAnT includes a second module discovery
403 routine, called “Frequent Subgraph Mining” (FSM)⁵⁵, that seeks a network of genes “colocalized”
404 in many cells. (Colocalization of a network in a cell means that every edge in that network is a
405 proximal gene pair in that cell (**Figure 6f**.) FSM can be used to find networks with a pre-specified
406 minimum size (numbers of nodes and edges) that are supported by a large number of cells
407 (Methods). For illustration, we used FSM to search for fully connected networks (“cliques”) with
408 at least four genes and found a single module – *Sgk1*, *Ttyh2*, *Ndr1* and *Ermn* (**Figure 6g**) – that
409 is colocalized in 72 cells, far greater than the support of the next most frequent four-gene clique
410 (12 cells) (**Figures 6i-k**). The six gene pairs comprising this module are *d*-colocalized individually,
411 is specific to mature oligodendrocytes and the module is significantly associated (p-value 8.3e-3)
412 with myelin sheath⁵⁶⁻⁵⁹ (**Figure 6h**). We speculate that their co-localization in specific partitions
413 inside cell reflects coordinated transport and translation in mature oligodendrocytes.

414 415 **Discussions**

416 In this work, we present the InSTAnT toolkit to screen for subcellular colocalization patterns of RNA
417 pairs and modules in an unbiased manner, through rigorous statistical analysis of single-molecule
418 resolution spatial transcriptomics data. We define *d*-colocalization as a new statistical phenomenon
419 that may point to biological relationships such as RNA-RNA interactions, formation of condensates
420 and shared subcellular localization. InSTAnT is a suite of statistical tests, at the heart of which lie the
421 Proximal Pair (PP) test that finds colocalized gene pairs in a single cell and the Conditional Poisson
422 Binomial (CPB) test that aggregates results of PP test across cells and reports *d-colocalized* gene
423 pairs. InSTAnT provides spatial region annotations for the reported gene pairs to aid biological
424 interpretation. It also includes procedures to characterize a *d*-colocalized gene pair based on its cell
425 type specificity or spatial modulation and to identify colocalized gene modules.

426
427 We employed InSTAnT to detect hundreds of gene pairs with low false positive rate and high
428 reproducibility on human U2OS cell line and mouse brain data. The InSTAnT analysis results suggest
429 that *d*-colocalization map can provide insights into various types of molecular interactions: RNA-RNA
430 interactions (**Figure 3i**), protein-protein interaction or shared pathway membership (**Figure 3h**) and

431 RNA-protein interactions (**Figure 3j-n**). The RNA d-colocalized pairs can be used to infer detailed
432 subcellular structures or characterize membrane-less organelles such as NS. These results indicate
433 that the spatial distribution of RNAs has “texture” rather than being relatively random as previously
434 perceived. Our brain data analysis shows that some RNA d-colocalized pairs have cell-type specificity,
435 are spatially modulated, and share functional annotation with other colocalizing pairs. All these results
436 suggest that RNA colocalization likely has biological consequences.

437
438 InSTAnT allows us to represent a cell as a graph where nodes represent genes and edges represent
439 proximal gene pairs. Such a graph, along with the transcript count vector commonly used to represent
440 an individual cell, may prove powerful in single cell analytics, allowing us to discover novel cell types
441 through a more nuanced clustering of cells than possible using count vectors alone. It will be exciting
442 to apply InSTAnT functionalities on future data sets that profile orders of magnitude more genes²⁶
443 (~10K). There are straight-forward ways to adapt the toolkit to efficiently handle this scenario, such as
444 by sampling of transcript pairs to estimate background probabilities in the PP test and by using a
445 greedy approach to testing only a subset of gene pairs. We expect such applications to help us better
446 characterize intracellular compartmentalization and provide complementary axes of information for
447 discovering regulatory and signaling interactions with and between cells.

448
449

450 **Online Methods**

451

452 **Code Availability**

453 The code is available at <https://github.com/anurendra/InSTAnT>.

454

455 **InSTAnT user guide**

456 InSTAnT tools have tunable parameters that can be selected based on the user’s requirement. We
457 selected the scale parameter d based on the average cell’s diameter and threshold for CPB test based
458 on False Positive Rate (1%) estimates. The user can also obtain region annotations of a gene pair’s
459 colocalization if the data include masks for cell and nucleus boundaries. Similarly, they may run cell
460 type specificity analysis if the data include cell type information. We advise caution when using
461 InSTAnT with small distance thresholds, such as 1 μm or less, as the false positive rates in this regime
462 can be high. This is due to the fact that colocalization with small distance is relatively rare in MERFISH
463 data and the estimate of null probability of a pair of transcripts being proximal, a key aspect of the PP
464 test, is error-prone in such cases. We believe that higher number of transcripts and improved optical
465 resolution¹⁷ may alleviate this problem.

466

467 **U2OS Dataset**

468 We obtained MERFISH data³² on a human osteosarcoma cell line (U2-OS) from
469 http://zhuang.harvard.edu/MERFISHData/data_for_release.zip . We used the authors' Matlab code to
470 extract and output the data in table format. We filtered the data to retain transcripts having minimum area
471 of 3 and intensity of $10^{0.75}$. The dataset had 7 replicates. We were able to extract data for four replicates –
472 *rep2*, *rep3*, *rep4*, *rep5*; the other replicates presented severe memory management challenges and were
473 not analyzed. Most of the reported results are from analysis of *rep3*, which profiles 130 genes in 3237 cells
474 with an average of 1243 transcripts per cell. Global *d*-colocalization maps were constructed for all four
475 replicates and compared to assess reproducibility.

476

477 **Brain Dataset**

478 Data reported in Moffit et al.³⁷ were obtained through personal communication with Dr. Jeffrey Moffitt. The
479 dataset contained 6325 cells with 553 average number of transcripts across 7 z-planes. We obtained cell
480 type assignment from Supplementary Table1 from Moffit et al.³⁷. We removed ambiguous cells leading to
481 5149 cells with 9 cell types. Proximal pairs were detected in cells that have at least one z-plane with 20 or
482 more transcripts.

483

484 **MERFISH imaging and Analysis**

485 **General cell culture conditions:** U2 OS cells were cultured in minimal essential medium (MEM) from
486 ATCC with 1 mM sodium pyruvate, 10% fetal bovine serum (FBS), and 1% penicillin-streptomycin (Pen-
487 Strep). The cells were obtained from ATCC and maintained using the recommended protocol.

488 **MERFISH sample preparation:** U2 OS MERFISH samples were prepared using a previously published
489 method⁶⁰. In brief, U2 OS cells were plated on a salinized 40mm #1.5 coverslip (Fisher Scientific). Plated
490 cells were transferred to a 37 °C and 5% CO₂ incubator overnight to grow. Cells were then fixed with 4%
491 paraformaldehyde (Electron Microscopy Sciences) and permeabilized with 0.5% (vol/vol) Triton X-100
492 (Sigma Aldrich). Samples were stained with encoding probes (10nM/probe) and anchor probes (1μM) for
493 36 hours in a humidified incubator at 37 °C. To stabilize the cells during clearing, the stained cells were
494 embedded in a thin, 4% polyacrylamide (PA) gel. Fiducial beads (Spherotech, FP-0245-2) were also
495 included in the gel to align rounds of MERFISH images.

496 **Commonly used imaging solutions:** The following solutions were used during imaging experiments
497 described in this work. Readout wash buffer was adapted from Moffit et al.⁶⁰ and contained 10% (v/v)
498 ethylene carbonate (Sigma Aldrich), 0.1% Triton X-100 in 2x SSC. Imaging buffer adapted from Moffit et
499 al.⁶⁰ and contained 5mM 3,4-dihydroxybenzoic acid (PCA; Sigma Aldrich), 2 mM trolox (Sigma Aldrich), 50
500 μM trolox quinone, 1:500 of recombinant protococatechuate 3,4-dioxygenase (rPCO; OYC Americas),
501 adjusted to a pH of 7-7.2 using 1 N NaOH (VWR International) in 2x SSC. Cleavage buffer was adapted

502 from⁶⁰ and contained 0.05 M TCEP HCl, adjusted to a pH of 7-7.2 using 1 N NaOH, in 2x SSC. Stripping
503 buffer was adapted from Eng. et al.¹⁴ and contained 55% formamide, and 0.1% Triton X-100 in 2x SSC.

504 **MERFISH imaging:** All images were acquired using a Zeiss Axiovert-200m widefield microscope (Carl
505 Zeiss AG) located in the IGB core imaging facility. The sample was placed into a flow cell (Bioprotech,
506 FCS2), filled with RNase free 2x SSC, and connected to a lab built automated flow system. Briefly,
507 computer-controlled valves (Hamilton, MVP/4, 8-5 valve) are used to select which solution was pulled
508 across the sample by a computer controlled pump (Gilson, Minipuls 3). All systems are controlled by a
509 custom designed Python script that can communicate with the microscope to start imaging or start flowing
510 after an imaging round is done. In brief, a single round of imaging involves staining with fluorescently labeled
511 readout probes (0.4 mL/min for 6 minutes, and 0.34 mL/min for 6 minutes), washing with readout wash
512 buffer (0.23 mL/minute for 9 minutes) to remove unbound probes, and imaging buffer was flowed into the
513 flow cell prior to imaging (0.34 mL/minutes for 6 minutes) to reduce photobleaching. A single quad band
514 excitation filter (Chroma, ZET402/468/555/638x) and dichroic (Chroma, ZT405/470/555/640rpc-UF1) were
515 used to image all samples. Excitation was provided by a 7 laser system (LDI WF, 89 North). Alexa Fluor
516 647 (Fisher scientific) labeled probes were excited using a 647 nm laser (0.5 W) with a ET700/75m
517 (Chroma) emission filter, and 1.5 second exposure time. Atto 565 (Atto tec) labeled probes were excited
518 using a 555 nm laser (1 W) with a ET610/75m (Chroma) emission filter, and a 0.75 second exposure time.
519 Fiducial beads were imaged with a 405 nm laser (0.3 W) with a ET440/40m emission filter, and a 1 second
520 exposure time. Samples were imaged with a 63x oil immersion objective (Carl Zeiss AG, 420782-9900-
521 000), and focus was maintained between imaging rounds using Definite Focus (Carl Zeiss AG). 9 z planes
522 with 0.7 μ m steps were taken for each FOV, and a total of 100 FOVs were acquired. After imaging is
523 complete, a cleavage buffer (0.2 mL/minute for 15 minutes) was flowed across the sample to remove the
524 fluorophores from the probes. The cleavage buffer was washed away using RNase free 2x SSC (0.5
525 mL/minute for 10 minutes). This process was repeated for a total of 8 rounds of imaging. PolyA probes
526 were stained after the final imaging round using the same method as described above.

527 **MERFISH data processing:** Individual FOVs were exported from czi format into 16 bit tiff format using Zen
528 (Carl Zeiss AG) using the image export method. Images then were reformatted into image stacks by FOV
529 and round. A modified copy of MERLIN⁶¹ was used to decode MERFISH spots. In brief, for each FOV,
530 images from different rounds are aligned using fiducial beads that were imaged in each round. Aligned
531 images are then normalized, decoded, and identified spots filtered using previously published methods²⁶.
532 Cell segmentation was done separately from MERLIN using Cellpose⁶² on PolyA and DAPI images for each
533 FOV. To improve FOV alignment to neighboring FOVs, the DAPI channel was used with the restitching
534 function found in Zen (Edge detection: on, minimal overlap: 5%, maximal shift: 15%, comparer: best, Global
535 optimizer: best). Using the aligned images, segmented cells that cross FOV boundaries were merged into
536 single cells, and global positions were generated for each spot. Spots are then assigned to cells based on
537 their spatial coordinates. Spots were then filtered to remove any spot smaller than 3 pixels in size.

538 **smFISH probe design:** All smFISH probes were designed using the Stellaris probe designer (Biosearch
539 technologies). Probes were designed using the following settings: Masking level: 5, max number of probes:
540 48, oligo length: 20, minimum spacing length: 2. SRRM2 exon probes were designed against SRRM2
541 isoform ENST00000301740 (GRCh38.p13). SRRM2 intron probes were randomly selected from probes
542 designed for three different introns defined by ensemble (SRRM2-230 intron 1, SRRM2-230 intron 2, and
543 SRRM2-230 intron 10) (GRCh38.p13). MALAT1 probes were designed against MALAT1 isoform
544 ENST00000534336 (GRCh38.p13). All probes were purchased from Biosearch modified with mdC (TEG-
545 Amino) at the 3' terminus. The probes were dissolved in TE buffer and labeled using AF488/Cy3/Cy5 NHS
546 esters for MALAT1, SRRM2 intron, and SRRM2 exon, respectively. The labeled probes were purified using
547 the Bio-Rad Bio-Spin P-6 purification columns (Cat # 732-6221).

548 **smFISH sample preparation:** Approximately 1.5-1.8 million U2OS cells were plated on a #1.5, 40 mm
549 coverslip (Fisher Scientific) that has been UV treated before plating. The cells were then transferred to an
550 incubator at 37 °C and 5% CO₂, overnight for 12-16 hours.

551 Modified from Fei et al.⁶³, the sample was rinsed with 1x PBS (Corning), followed by fixation using 4%
552 paraformaldehyde (PFA; Electron Microscopy Sciences) in 1x PBS for 10 minutes at room temperature
553 (RT). The sample was then washed three times with 1x PBS and permeabilized with 0.5% Triton X-100
554 (Sigma Aldrich), 2 mM vanadyl ribonucleoside complexes (VRC; Sigma Aldrich) in 1x PBS for 10 minutes
555 on ice, followed by three quick washes with 1x PBS. At this point, the sample can be stored in 70% Ethanol
556 at 4 °C if the experiment needs to be paused temporarily.

557 To prepare for smFISH hybridization, sample was rinsed with 10% formamide (Sigma Aldrich) in 2x saline
558 sodium citrate (SSC; Fisher Scientific). smFISH probe hybridization buffer was prepared with 0.2 mg/mL of
559 bovine serum albumin (BSA; Fisher Scientific), 2 mM VRC, 10% dextran sulfate (Sigma Aldrich), 1 mg/mL
560 yeast tRNA (Fisher Scientific), 10% formamide, 1% murine RNase inhibitor (New England BioLabs) in 2x
561 SSC. Avoid light exposure from this point forward. smFISH probes were then added to the FISH
562 hybridization buffer at a final concentration of 14 nM for each targeted RNA (MALAT1, SRRM2 intron, and
563 SRRM2 exon).

564 A humidified chamber was made using an empty pipette box filled halfway with nuclease-free water
565 (Corning) at the base and a UV-treated glass slide covered with a parafilm layer on top. A 100 µl drop of
566 the FISH probe hybridization buffer was then added on top of the parafilm layer and the sample was casted
567 over the drop with the cell side facing down. The chamber was then placed in an incubator in dark and
568 wrapped entirely with aluminum foil overnight at 37 °C for at least 16 hours. The sample was quickly rinsed
569 two times with 10% formamide in 2x SSC then stained with 4',6-diamidino-2-phenylindole (DAPI; Invitrogen
570 by Fisher Scientific) 1:1000 of 1 mg/mL stock solution and 1:5000 of Fluoro-Max Blue Aqueous Fluorescent
571 Particles (fluorescent beads; Fisher Scientific) in 2x SSC. The sample was incubated with the DAPI and
572 fluorescent beads solution for 5 minutes while rocking at RT, followed by a quick wash with 2x SSC, then
573 stored in 2x SSC at 4 °C until ready for imaging.

574 **Protein staining:** After smFISH imaging, the sample can be stored in 1x PBS at 4 °C for up to a week
575 before protein staining. Samples were fixed a second time with 4% PFA in 1x PBS for 5 minutes at RT,
576 then rinsed three times with 1x PBS. This was followed by incubation with a blocking solution of 1% BSA
577 in 1x PBS for three consecutive times with 10 minutes each time at RT.

578 The SON primary antibody (Anti-SON, Sigma Aldrich, HPA023535) was kept at -20 °C until ready for use.
579 The primary antibody stock solution of 1:1000 was prepared with 1x PBS and kept on ice. A 1:5000 primary
580 antibody dilution was prepared in blocking solution and the sample was incubated with 200 µl of the primary
581 antibody solution for approximately 1 hour at RT in the dark.

582 The sample was washed with blocking solution three consecutive times with a 10-minute incubation each
583 time at RT, followed by three washes with 1x PBS, for 10 minutes each time at RT.

584 Secondary antibody was conjugated to Alexa Fluor 647 (Goat anti-rabbit, Invitrogen, A21245). The
585 concentrated secondary antibody was kept at 4 °C until ready for use. Sample staining was accomplished
586 by 1:1000 dilution of the secondary antibody in blocking solution and casting of the sample on a 200 µl drop
587 of the secondary antibody solution, with the cell side facing down. The sample was then incubated for 1
588 hour in the dark at RT. The sample was re-stained with DAPI in 1x PBS with the same concentration and
589 incubation time described in smFISH staining section. This was followed by a quick rinse with 1x PBS and
590 the sample was stored in 1x PBS at 4 °C until ready for imaging.

591 **smFISH image acquisition:** smFISH and protein imaging were done on the same MERFISH imaging and
592 fluidic system described above (MERFISH imaging). After placing the sample into the flow cell, imaging
593 buffer was flowed through the system (0.34 mL/minute for 5 minutes). Excitation and dichroic filters were
594 the same as used above. The following dyes, lasers, and emission filters were used for smFISH imaging.

Channel	Target	Laser line (power)	Exposure time	Emission filter
DAPI	Fiducial beads, nuclei	405 nm (0.3 W)	0.075 seconds	ET440/40m
Alexa Fluor 488	MALAT1 lncRNA	470 nm (1 W)	2 seconds	ET525/50m
Cy3	SRRM2 intron RNA	555 nm (1 W)	2 seconds	ET610/75m
Cy5	SRRM2 exon mRNA	640 nm (0.5 W)	3 seconds	ET700/75m

595
596 Samples were imaged with the same 63x oil immersion objective as above, and focus was maintained
597 between imaging rounds using Definite Focus. 9 z planes were imaged with a step size of 0.7 µm. After
598 imaging, smFISH probes were removed using a stripping buffer that was flowed through the system (0.34
599 mL/minutes for 5 minutes) without removing the sample from the microscope. After stripping the sample
600 was washed with 2x SSC (0.5 mL/minutes for 5 minutes). The sample was imaged a second time using the
601 same settings as above. After imaging the sample was removed from the flow cell and placed into 1x PBS
602 prior to protein staining (Protein staining).

603 After protein staining was complete, sample was placed into the flow cell and filled with imaging buffer. The
604 same region imaged during the smFISH experiment was found and reimaged using the same objective and
605 z stack settings as above. The following imaging settings were used.

Channel	Target	Laser line (power)	Exposure time	Emission filter
DAPI	Fiducial beads, nuclei	405 nm (0.3 W)	0.05 seconds	ET440/40m

Alexa Fluor 647	SON protein	640 nm (0.5 W)	1.5 seconds	ET700/75m
-----------------	-------------	----------------	-------------	-----------

606

607 **SRRM2 image registration and alignment:** Individual FOVs were exported from czi format into 16 bit tiff
608 format using Zen's (Carl Zeiss AG) image export method. To align images from the same FOV across
609 multiple rounds of imaging or experiment, blue fluorescent beads imaged in the DAPI channel were used
610 as fiducial markers. We found that aligning images from the same experiment required a simple translation.
611 To align protein images with mRNA images, an iterative rotation and translation process was developed.
612 For each iterative round of alignment, the protein DAPI channel was rotated, then translated to best align
613 with the mRNA image, this warped image was then used as the starting protein DAPI image for the next
614 round of alignment. We found that it took between 2 and 5 rounds of alignment to align protein images to
615 mRNA images. Chromatic aberration was corrected by aligning all channels to the Cy5 channel. Multicolor
616 beads (Multi-speck bead slide, Carl Zeiss AG, 1783-455) that included dyes in the Alexa Fluor 488, Cy3,
617 and Cy5 channels were used to correct Alexa Fluor 488 and Cy3 channels. The DAPI channel was
618 corrected to the Cy5 channel using the fiducial bead cross talk between the DAPI and Alexa Fluor 488
619 channels. This was done by calculating the shift between non-nuclear regions of the DAPI and Alexa Fluor
620 488 channels, then adding the Alexa Fluor 488 to Cy5 shift to the DAPI to Alexa Fluor 488 shift.

621 **SRRM2 image preprocessing:** To remove cross talk in DAPI and Alexa Fluor 488 channels caused by
622 the fiducial beads, stripped Alexa Fluor 488 mRNA channel was subtracted from the stained Alexa Fluor
623 488 channel. As fiducial beads are not affected by the mRNA stripping conditions, any spots that remain in
624 the stripped Alexa Fluor 488 channel would be from the beads, not from MALAT1 mRNA. In order to reduce
625 background in other images, round subtraction was also done on the other channels of the mRNA FOV.

626 **SRRM2 co-localization analysis:** Co-localization analysis was done on a single z plane from each
627 experiment stack. Images were then filtered using a high pass filter (5 pixel sigma) and Lucy-Richardson
628 deconvolution (10 iterations, 9 pixel filter size, 1.4 pixel sigma). Filtered images are then converted to binary
629 masks with manually defined thresholds. To remove false positives in the MALAT1 channel, the MALAT1
630 mask was multiplied with the inverse of the stripped MALAT1 mask. Cell nuclei were identified using the
631 DAPI channel, and segmented using a manually defined threshold.

632 The co-localization rate was calculated for each nucleus defined from the DAPI channel. To calculate the
633 co-localization rate between two channels, each channel is multiplied against the nuclei mask. For each
634 spot in the first mask, the spot was dilated by 2 μm and then compared against the second mask. If the
635 dilated spot overlaps any spot in the second mask, it is considered to be colocalized. The colocalization
636 rate was then calculated to be the following:

$$637 \quad \text{colocalization percent} = \frac{\text{Co-localized spots ct}}{\text{Total spots ct}} * 100\%$$

638 The colocalization percent was averaged across 13 cells.

639 **SRRM2 figure and generation (Figures 3k-n):** SRRM2 exon and intron images were filtered using a high
640 pass filter with 2 pixel sigma, while MALAT1 was filtered using high pass filter with 5 pixel sigma. Raw SON
641 images were used in panels **Figures 3 m,n**.

642

643 **False Positive Rate (FPR)**

644 We generate random baseline dataset established by permuting the gene labels of all transcripts within
645 each cell, which recapitulates the spatial patterns of the original data but not the gene-gene relationships.

646 FPR is obtained by comparing the number of detected pairs obtained on randomized data with number of
647 detected pairs on real data.

648 Ten of the 140 genes probed in the U2OS MERFISH data set were “blanks”, meaning that they
649 do not represent any particular RNA or other molecule. Any gene pair involving such blank “genes”, if found
650 to d-colocalize, is clearly a false positive. This provided us another opportunity to assess
651 the false positive errors in our global co-localization map. We recorded the fraction of such false positives
652 among predicted pairs at varying levels of significance (Supplement Figure 2, blue).

653

654 **Hyperparameter selection**

655 Scale parameter d was chosen to be 4 microns in U2OS dataset and 2 microns for Brain dataset, as it
656 corresponded to ~5% of average diameter of a cell in the respective datasets. The p-value threshold for PP
657 test was chosen to be 0.01 for both the datasets which resulted in FPR~5%. p-value threshold for CPB test
658 was chosen to be $1e-3$ for U2OS and $1e-5$ for Brain dataset as it resulted in FPR<1%. p-value threshold
659 for frequent subgraph mining on Brain dataset was chosen to be 0.05 as threshold of 0.01 didn't yield any
660 subgraph.

661

662 **Proximal Pair (PP) test**

663 PP test reports proximal pairs of genes in a particular cell. A gene pair g_i, g_j is a proximal pair in a cell if
664 their transcripts are proximally located (separated by distance d or less) significantly more often than
665 expected by chance. The null probability p is estimated from the distances between all pairs of transcripts
666 (regardless of gene identities) in the cell, by calculating the fraction of transcript pairs that are proximally
667 located. Let t_i and t_j denote the transcript counts of genes g_i, g_j respectively in the cell, let $T = t_i t_j$ and let
668 K be the number of proximally located transcript pairs of these genes. The PP test performs a Binomial
669 test providing a p-value for g_i, g_j as

670

$$671 \text{p-value}(g_i, g_j) = \text{Binomial}(T, p, K)$$

672

673 **PP-3D test**

674 PP-3D is an extension of PP test to handle three-dimensional data in the form of 2D (x-y) locations of
675 transcripts in each of multiple z-planes. We assume that data from different planes are independent and
676 identically distributed. The new distribution is the sum of independent Binomial distributions (with the same
677 parameter), which is also a Binomial distribution. The null probability of two transcripts being proximal is
678 estimated as a weighted combination of estimated null probability for each of the z-planes,

679

680

$$p \equiv \frac{\sum_z l_z p_z}{\sum_z l_z}$$

681 where, p_z denotes the null probability for z -th plane, l_z denotes the total number of transcripts in z -th slice.

682 T and K are also aggregated across z -planes:

683

$$T = \sum_z T_z$$

684

$$K = \sum_z K_z$$

685 where K_z is total number of proximal transcript pairs and T_z is total number of transcript pairs (of g_i, g_j) in

686 z -th plane. PP-3D calculates a p-value for each gene pair as $p\text{-value}(g_i, g_j) = \text{Binomial}(T, p, K)$.

687

688 **Conditional Poisson Binomial (CPB) test**

689 CPB test detects a d -colocalized gene pair, i.e., a gene pair that is a proximal pair in significantly many

690 cells. It assigns a p-value to the number of cells in which a gene pair is found to be proximal pair detected

691 using PP test. We first describe a simpler version of the test (“unconditional Poisson Binomial” or UPB) test

692 that assumes that all gene pairs are equally likely to be proximal pair in a cell but allows for the fact that

693 different cells may have different number of proximal pairs. Let X_{ij}^c be a binary variable denoting if g_i, g_j are

694 a proximal pair in c -th cell. X_{ij}^c is assumed to follow a Bernoulli distribution with parameter p_0^c , which is

695 estimated as the fraction of proximal gene pairs in the cell:

696

$$p_0^c \equiv \frac{\sum_{k \leq l} X_{k,l}^c}{\sum_{k \leq l} 1} = \frac{\sum_{k \leq l} X_{k,l}^c}{\binom{n}{2}}$$

697 where n denotes total number of genes. This estimate of p_0^c assumes that all gene pairs can be a proximal

698 pair. To incorporate the fact that a gene pair cannot be a proximal pair if either of the genes is not expressed

699 in the cell, the above estimate is modified as,

700

$$p_0^c \equiv \frac{\sum_{k \leq l} X_{k,l}^c}{\sum I_{k \leq l}(g_k, g_l)}$$

701 where $I(g_k, g_l)$ is an indicator function that equals to 1 iff both g_k and g_l are expressed.

702

703 CPB test is a modified version of the UPB test that accounts for the possibility that all gene pairs are not

704 equally likely to be colocalized in a cell and sets the Bernoulli parameter (p_0^c above) to be gene pair-

705 dependent. Let z_i denote total number of proximal pairs having gene i as one of the genes, aggregated

706 across all cells, i.e.,

707

$$z_i = \sum_{j \leq c} X_{ij}^c$$

708

709 We use these global summary statistics to model the prior probability Π_{ij} that a proximal pair detected in a

710 cell is the gene pair g_i, g_j , as follows:

711

$$\Pi_{ij} \equiv \frac{z_i z_j}{\sum_{i \leq j} z_i z_j}$$

712 This model de-emphasizes gene pairs comprising genes that are frequently found to be in proximal pairs
713 across cells. Now, the Bernoulli parameter for variable X_{ij}^c is estimated as

714

$$p_{ij}^c \equiv 1 - (1 - \Pi_{ij})^{\sum_{i \leq j} X_{ij}^c}$$

715

716 The total number of cells where g_i, g_j is a proximal pair follows a Poisson Binomial distribution

717

$$\sum_{c=1}^m X_{ij}^c \sim \text{Poisson Binomial} (p_{ij}^1, \dots, p_{ij}^m)$$

718 **Spatial Annotation**

719 A d-colocalized pair is annotated by cellular region where the gene pair's proximal pairs tend to be found.
720 We define four categories – Nucleus (Nuc), Peri-Nucleus (PN), Cytosol (Cyto) and Cell Periphery (CP).
721 Proximal pairs in each cell are annotated by cellular region and is aggregated across cells to yield primary
722 and secondary category. Perinuclear (PN) region is defined as including x microns on either side of the
723 nuclear membrane, while Cell Periphery (CP) is defined as regions within y microns of the cell membrane.
724 Remaining regions are designated as Cytosol (Cyt) or Nucleus (Nuc). We chose x = 2.5 micron which
725 corresponded to ~43% of nucleus transcripts being annotated as perinuclear, and y = 4 micron which
726 corresponds to ~35% cytosolic transcripts being annotated as cell periphery.

727

728 **RNA-RNA Interaction (RRI)**

729 For RRI, we set distance d to be equal to the resolution of MERFISH data (200 nm). The small distance
730 was chosen to capture gene pairs whose d -colocalization may be explained due to the binding of their
731 transcripts. We used RNAplex³⁴ to compute the RRI scores. For this, we retrieved the nucleotide
732 sequences from the Ensembl database⁶⁴ and got the specific transcript id to get the correct spliced form.
733 RNAplex has been shown to be among the most accurate tools while being fast enough to compute the
734 scores for gene pairs with their full transcripts. Finally, we perform a gene-centric analysis for each of the
735 130 genes. For each gene, we ask if top 10 d -colocalized pairs (out of 130) has significantly higher number
736 of pairs with RRI score greater than a fixed threshold (RRI>35). We perform a Binomial test whose success
737 probability is obtained as follows. We model background distribution by fitting a Gaussian distribution to the
738 RRI scores of the pairs with d -colocalization score greater than 0.01. The survival probability of RRI scores
739 higher than the fixed threshold (RRI>35) serves as the success probability of Binomial test. Finally, we
740 perform an FDR correction using the Benjamini-Hochberg procedure⁶⁵. 8 of the genes pass this FDR
741 correction showing that RRI may be a plausible mechanism for their d-colocalized pairs.

742

743 **Enrichment Analysis**

744 To understand the biological mechanism or consequences of d-colocalization, we tested if the compendium
745 of d-colocalized gene pairs has significant overlap with functionally related gene pairs. We define a gene

746 pair to be functionally related if both genes are present in same KEGG
 747 pathway or are annotated with same GO terms more than K times. K was chosen such that
 748 number of gene pairs is similar across d-colocalized and functionally related set. In our analysis, K (MF) =
 749 2 , K (BP) =1 , K (CC) = 3, K (pathway) = 1. We performed a hypergeometric test between d-colocalized
 750 pairs and functionally related set.

751

752 **Cell Type Specificity of a d-colocalized Gene Pair**

753 InSTAnT employs a series of statistical tests to categorize a d-colocalized pair based on its cell type
 754 specificity. First, it tests the association between cells where a gene pair was deemed a significant proximal
 755 pair and cells of a particular type (e.g., inhibitory neurons), using a Hypergeometric test. (This process is
 756 repeated for every cell type.) If such an association is found to be statistically significant, it is subjected to
 757 further tests to determine if the cell type specificity arises simply because one of the genes in the pair is
 758 expressed specifically in that cell type. For this, InSTAnT utilizes a version of the generalized
 759 Hypergeometric test that tests for an association between two sets conditional on a third set⁶⁶, as described
 760 below. In this case, the third set comprises the cells with high expression of one of the genes in the pair.

761

762 Let U be the set of all cells, M be the set of cells of a particular cell type, O be the set of cells where a gene
 763 pair is deemed a proximal pair and E be the set of cells with high expression of one of the genes in the pair.
 764 M , O and E are subsets of U . The threshold for high gene expression used in defining E is chosen such
 765 that $\text{size}(E) = \text{size}(M)$. Let $|M \cap E| = \gamma, |M \cap O| = \lambda, |E \cap O| = \alpha$. The Hypergeometric test p-value of
 766 association between M and O is given by the probability that a random set of size $|O|$ has an overlap
 767 (intersection) of size greater than or equal to λ with M . However, we wish to test if the overlap between M
 768 and O is significant beyond what is expected not from a random set of size $|O|$ but a random set of this size
 769 that respects the known overlap between M and E and between E and O . For this, we calculate probability
 770 of the overlap between M and a random set of $|O|$ being greater than or equal to λ conditional on the
 771 observed overlap between M and E and that between E and O , as follows:

772

$$773 \frac{\sum_{k=\lambda}^{\min(|M|,|O|)} \sum_{\beta=0}^k \binom{\gamma}{\beta} \binom{m-\gamma}{k-\beta} \binom{n_1-\gamma}{\alpha-\beta} \binom{|U|-|M|-|E|+\gamma}{|O|-\alpha-k+\beta}}{\binom{|E|}{\alpha} \binom{|U|-|E|}{|O|-\alpha}}$$

774

775 This is an example of multivariate hypergeometric distribution. We use *scipy.stats.multivariate_hypergeom*
 776 package for multivariate hypergeometric distribution.

777

778 For each gene pair that is associated with a cell type, InSTAnT performs the above test twice, each time
 779 conditioning on a set E defined by the high expression cells for one of the genes of the pair. Significant p-
 780 values in both tests thus performed indicate that the cell type-specificity of the d-colocalized gene pair is

781 significant beyond what is expected from the specificity of either gene's expression. Furthermore, InSTAnT
782 tests if either gene of the pair is a marker of the cell type, defined as any gene among the top 10 by
783 association between their expression and the cell type. A marker gene is found by conducting
784 Hypergeometric test of overlap between O and E .

785
786 Using the above tests, InSTAnT categorizes a d -colocalized gene pair vis-à-vis its cell type specificity as
787 follows: If the gene pair is significantly associated with a cell type (first test above), then it belongs to
788 Category 1 if the association is significant by the Hypergeometric test conditional on high expression cells
789 of both genes and neither gene is a marker of the cell type, otherwise it belongs to Category 2. Category 3
790 comprises d -colocalized gene pairs that are not associated with any cell type (Bonferroni corrected
791 hypergeometric p-value ≥ 0.05 , Supplement Table 5).

792

793 **Probabilistic graphical model for Spatial Modulation**

794

795 InSTAnT uses a likelihood ratio test to determine if sub-cellular colocalization of a d -colocalized gene pair
796 is spatially modulated at the tissue level. Informally, this means that the cells in which the gene pair is
797 deemed to be a proximal pair are non-randomly distributed in the physical space.

798 The probabilistic model is formulated around a graph with a node for each cell and edges between
799 neighboring cells. Two cells are neighboring cells if they are located within a configurable distance (set to
800 100 micron in our tests). Each node is associated with a binary variable s_c that indicates whether the
801 specific gene pair (say g_i, g_j) is a proximal pair in the corresponding cell c , as detected by the PP test. The
802 variable s_c is assumed to be a Bernoulli-distributed variable. The null hypothesis is that the Bernoulli
803 parameter is a global constant p^{global} shared across all cells, i.e., it does not depend on the cell c and thus
804 on its spatial location:

$$805 \quad H_0: s_c \sim Ber(p^{global})$$

806 p^{global} is estimated as the fraction of cells where the gene pair g_i, g_j is a proximal pair, which is its
807 maximum likelihood estimate. In the alternative hypothesis, the model assumes that the distribution of
808 variable s_c depends on the fraction of cells c' in the neighborhood of c for which $s_{c'} = 1$. Let p^{local} be the
809 fraction of cells c' in the neighborhood of c for which $s_{c'} = 1$.

$$810 \quad H_1: s_c \sim Ber(w p^{local} + (1 - w)p^{global})$$

$$811 \quad 0 < w < 1$$

812 The parameters p^{global}, p^{local}, w are learnt by maximizing likelihood. Weight w controls the contribution of
813 local neighborhood. InSTAnT calculates the log likelihood ratio (LLR) for each gene pair in the d -
814 colocalization map and pairs with LLR above a threshold are designated as spatially modulated. The
815 threshold is obtained by random permutation of the of s_c values of cells, repeating the above test and
816 selecting the highest LLR score (over all gene pairs) seen on the randomized data. This allows us to detect
817 spatially clustered distributions of cells supporting g_i, g_j colocalization.

818

819

820 **Module Discovery: Global Colocalization Clustering (GCC)**

821 GCC is a procedure to analyze a *d*-colocalization map to identify subsets of genes that exhibit a high
822 frequency of pairwise *d*-colocalization relationships. To this end, it represents the *d*-colocalization map as
823 an $n \times n$ matrix (n = number of genes) whose entries are the negative logarithm of p-values of gene pairs
824 from the CPB test and performs a hierarchical clustering of rows and columns using Euclidean distance
825 with Ward criterion. (The constant 1e-64 is added to all the p-values to handle zero p-values prior to taking
826 logarithms.)

827

828

829 **Module Discovery: Frequent subgraph mining (FSM)**

830 FSM seeks a network of genes that is “colocalized” in many cells, where colocalization of a network in an
831 individual cell means that every gene pair connected by an edge in that network is a proximal pair in that
832 cell. It constructs a *colocalization graph* for each cell with genes as nodes and edges representing proximal
833 gene pairs from PP test. It then uses an efficient graph mining tool called gSPAN⁵⁵ to detect subgraphs
834 with a pre-specified minimum size (numbers of nodes and edges) that are supported by a pre-specified
835 minimum number of cells.

836

837

838

839

840 **Acknowledgements**

841 We thank Dr. Jeffrey Moffitt for sharing the data from Moffit et al.³⁷, Zijun Wu for assistance in formatting of
842 figures, Alton S. Barbehenn for helpful discussion for the statistical analysis, Prof. Prasanth V.
843 Kannanganattu for advice on sample preparation. **Funding:** This work was supported by the National
844 Institutes of Health (R35GM131819 to S.S., 1R35GM147420 to H.-S.H, and T32- 842 GM136629 to M.A.),
845 Johnson & Johnson (WiSTEM2D Award for Science to H.-S.H.), Cancer Center at Illinois (Seed grant to
846 H.-S.H), and Georgia Institute of Technology (Wallace H. Coulter Distinguished Faculty Chair:
847 S.S.) **Facilities:** We acknowledge Core Facilities at the Carl R. Woese Institute for Genomic Biology for
848 their microscope and staff support.

849

850 **References**

851

- 852 1. Rao, A., Barkley, D., França, G.S. & Yanai, I. Exploring tissue architecture using spatial
853 transcriptomics. *Nature* **596**, 211-220 (2021).
- 854 2. Marx, V. Method of the Year: spatially resolved transcriptomics. *Nature methods* **18**, 9-
855 14 (2021).

- 856 3. Svensson, V., Teichmann, S.A. & Stegle, O. SpatialDE: identification of spatially variable
857 genes. *Nature methods* **15**, 343-346 (2018).
- 858 4. Zhu, J., Sun, S. & Zhou, X. SPARK-X: non-parametric modeling enables scalable and
859 robust detection of spatial expression patterns for large spatial transcriptomic studies.
860 *Genome biology* **22**, 1-25 (2021).
- 861 5. Pham, D. *et al.* stLearn: integrating spatial location, tissue morphology and gene
862 expression to find cell types, cell-cell interactions and spatial trajectories within
863 undissociated tissues. *BioRxiv* (2020).
- 864 6. Dries, R. *et al.* Giotto: a toolbox for integrative analysis and visualization of spatial
865 expression data. *Genome biology* **22**, 1-31 (2021).
- 866 7. Hu, J. *et al.* SpaGCN: Integrating gene expression, spatial location and histology to
867 identify spatial domains and spatially variable genes by graph convolutional network.
868 *Nature methods* **18**, 1342-1351 (2021).
- 869 8. Hildebrandt, F. *et al.* Spatial Transcriptomics to define transcriptional patterns of
870 zonation and structural components in the mouse liver. *Nature communications* **12**, 1-14
871 (2021).
- 872 9. Liu, Z., Sun, D. & Wang, C. Evaluation of cell-cell interaction methods by integrating
873 single-cell RNA sequencing data with spatial information. *Genome Biology* **23**, 1-38
874 (2022).
- 875 10. Li, D., Ding, J. & Bar-Joseph, Z. Identifying signaling genes in spatial single-cell
876 expression data. *Bioinformatics* **37**, 968-975 (2021).
- 877 11. Rao, N., Clark, S. & Habern, O. Bridging genomics and tissue pathology: 10x genomics
878 explores new frontiers with the visium spatial gene expression solution. *Genetic
879 Engineering & Biotechnology News* **40**, 50-51 (2020).
- 880 12. Liu, Y. *et al.* High-spatial-resolution multi-omics sequencing via deterministic barcoding
881 in tissue. *Cell* **183**, 1665-1681. e18 (2020).
- 882 13. Chen, K.H., Boettiger, A.N., Moffitt, J.R., Wang, S. & Zhuang, X. Spatially resolved,
883 highly multiplexed RNA profiling in single cells. *Science* **348**, aaa6090 (2015).
- 884 14. Eng, C.-H.L. *et al.* Transcriptome-scale super-resolved imaging in tissues by RNA
885 seqFISH+. *Nature* **568**, 235-239 (2019).
- 886 15. Lee, J.H. *et al.* Highly multiplexed subcellular RNA sequencing in situ. *Science* **343**,
887 1360-1363 (2014).
- 888 16. Wang, X. *et al.* Three-dimensional intact-tissue sequencing of single-cell transcriptional
889 states. *Science* **361**, eaat5691 (2018).
- 890 17. Alon, S. *et al.* Expansion sequencing: Spatially precise in situ transcriptomics in intact
891 biological systems. *Science* **371**, eaax2656 (2021).
- 892 18. Chen, W.-T. *et al.* Spatial transcriptomics and in situ sequencing to study Alzheimer's
893 disease. *Cell* **182**, 976-991. e19 (2020).
- 894 19. Dong, K. & Zhang, S. Deciphering spatial domains from spatially resolved
895 transcriptomics with an adaptive graph attention auto-encoder. *Nature communications*
896 **13**, 1-12 (2022).
- 897 20. Doyle, F. *et al.* Bioinformatic Tools for Studying Post-Transcriptional Gene Regulation.
898 *Post-Transcriptional Gene Regulation*, 39.
- 899 21. Parton, R.M., Davidson, A., Davis, I. & Weil, T.T. Subcellular mRNA localisation at a
900 glance. *Journal of cell science* **127**, 2127-2133 (2014).
- 901 22. Kloc, M., Zearfoss, N.R. & Etkin, L.D. Mechanisms of subcellular mRNA localization. *Cell*
902 **108**, 533-544 (2002).
- 903 23. Besse, F. & Ephrussi, A. Translational control of localized mRNAs: restricting protein
904 synthesis in space and time. *Nature reviews Molecular cell biology* **9**, 971-980 (2008).
- 905 24. Martin, K.C. & Ephrussi, A. mRNA localization: gene expression in the spatial dimension.
906 *Cell* **136**, 719-730 (2009).

- 907 25. Blower, M.D. Molecular insights into intracellular RNA localization. *International review of*
908 *cell and molecular biology* **302**, 1-39 (2013).
- 909 26. Xia, C., Fan, J., Emanuel, G., Hao, J. & Zhuang, X. Spatial transcriptome profiling by
910 MERFISH reveals subcellular RNA compartmentalization and cell cycle-dependent gene
911 expression. *Proceedings of the National Academy of Sciences* **116**, 19490-19499
912 (2019).
- 913 27. Bergen, V., Lange, M., Peidli, S., Wolf, F.A. & Theis, F.J. Generalizing RNA velocity to
914 transient cell states through dynamical modeling. *Nature biotechnology* **38**, 1408-1414
915 (2020).
- 916 28. Bergen, V., Soldatov, R.A., Kharchenko, P.V. & Theis, F.J. RNA velocity—current
917 challenges and future perspectives. *Molecular systems biology* **17**, e10282 (2021).
- 918 29. Mah, C.K. *et al.* Bento: A toolkit for subcellular analysis of spatial transcriptomics data.
919 *BioRxiv* (2022).
- 920 30. Engel, K.L., Arora, A., Goering, R., Lo, H.Y.G. & Taliaferro, J.M. Mechanisms and
921 consequences of subcellular RNA localization across diverse cell types. *Traffic* **21**, 404-
922 418 (2020).
- 923 31. Fazal, F.M. & Chang, H.Y. Subcellular spatial transcriptomes: Emerging frontier for
924 understanding gene regulation. in *Cold Spring Harbor symposia on quantitative biology*
925 Vol. 84 31-45 (Cold Spring Harbor Laboratory Press, 2019).
- 926 32. Moffitt, J.R. *et al.* High-throughput single-cell gene-expression profiling with multiplexed
927 error-robust fluorescence in situ hybridization. *Proceedings of the National Academy of*
928 *Sciences* **113**, 11046-11051 (2016).
- 929 33. Van Dam, S., Vosa, U., van der Graaf, A., Franke, L. & de Magalhaes, J.P. Gene co-
930 expression analysis for functional classification and gene–disease predictions. *Briefings*
931 *in bioinformatics* **19**, 575-592 (2018).
- 932 34. Tafer, H. & Hofacker, I.L. RNAplex: a fast tool for RNA–RNA interaction search.
933 *Bioinformatics* **24**, 2657-2663 (2008).
- 934 35. Ilik, İ.A. *et al.* SON and SRRM2 are essential for nuclear speckle formation. *Elife* **9**,
935 e60579 (2020).
- 936 36. Miyagawa, R. *et al.* Identification of cis-and trans-acting factors involved in the
937 localization of MALAT-1 noncoding RNA to nuclear speckles. *Rna* **18**, 738-751 (2012).
- 938 37. Moffitt, J.R. *et al.* Molecular, spatial, and functional single-cell profiling of the
939 hypothalamic preoptic region. *Science* **362**, eaau5324 (2018).
- 940 38. Sakers, K. *et al.* Astrocytes locally translate transcripts in their peripheral processes.
941 *Proceedings of the National Academy of Sciences* **114**, E3830-E3838 (2017).
- 942 39. Lanciotti, A. *et al.* Megalencephalic leukoencephalopathy with subcortical cysts protein 1
943 functionally cooperates with the TRPV4 cation channel to activate the response of
944 astrocytes to osmotic stress: dysregulation by pathological mutations. *Human molecular*
945 *genetics* **21**, 2166-2180 (2012).
- 946 40. Hwang, J., Vu, H.M., Kim, M.-S. & Lim, H.-H. Plasma membrane localization of MLC1
947 regulates cellular morphology and motility. *Molecular brain* **12**, 1-14 (2019).
- 948 41. Fleischer, J., Schwarzenbacher, K., Besser, S., Hass, N. & Breer, H. Olfactory receptors
949 and signalling elements in the Grueneberg ganglion. *Journal of neurochemistry* **98**, 543-
950 554 (2006).
- 951 42. Liu, M., Kim, D.-W., Zeng, H. & Anderson, D.J. Make war not love: The neural substrate
952 underlying a state-dependent switch in female social behavior. *Neuron* **110**, 841-856. e6
953 (2022).
- 954 43. Miyata, S. *et al.* Plasma corticosterone activates SGK1 and induces morphological
955 changes in oligodendrocytes in corpus callosum. *PloS one* **6**, e19859 (2011).

- 956 44. Miyata, S. *et al.* Sgk1 regulates desmoglein 1 expression levels in oligodendrocytes in
957 the mouse corpus callosum after chronic stress exposure. *Biochemical and biophysical*
958 *research communications* **464**, 76-82 (2015).
- 959 45. Dattilo, V., Amato, R., Perrotti, N. & Gennarelli, M. The emerging role of SGK1 (Serum-
960 and Glucocorticoid-Regulated Kinase 1) in major depressive disorder: Hypothesis and
961 mechanisms. *Frontiers in Genetics* **11**, 826 (2020).
- 962 46. Cathomas, F. *et al.* Oligodendrocyte gene expression is reduced by and influences
963 effects of chronic social stress in mice. *Genes, Brain and Behavior* **18**, e12475 (2019).
- 964 47. Monavarfeshani, A., Knill, C.N., Sabbagh, U., Su, J. & Fox, M.A. Region- and cell-specific
965 expression of transmembrane collagens in mouse brain. *Frontiers in Integrative*
966 *Neuroscience* **11**, 20 (2017).
- 967 48. Lemoine, G.G., Scott-Boyer, M.-P., Ambroise, B., Périn, O. & Droit, A. GWENA: gene
968 co-expression networks analysis and extended modules characterization in a single
969 Bioconductor package. *BMC bioinformatics* **22**, 1-20 (2021).
- 970 49. Cassella, L. & Ephrussi, A. Subcellular spatial transcriptomics identifies three
971 mechanistically different classes of localizing RNAs. *Nature Communications* **13**, 1-16
972 (2022).
- 973 50. Jansen, R.-P. RNA–cytoskeletal associations. *The FASEB Journal* **13**, 455-466 (1999).
- 974 51. Singer, R.H. The cytoskeleton and mRNA localization. *Current opinion in cell biology* **4**,
975 15-19 (1992).
- 976 52. Huang, J. *et al.* Identification of the fatty acid synthase interaction network via iTRAQ-
977 based proteomics indicates the potential molecular mechanisms of liver cancer
978 metastasis. *Cancer Cell International* **20**, 1-14 (2020).
- 979 53. Mering, C.v. *et al.* STRING: a database of predicted functional associations between
980 proteins. *Nucleic acids research* **31**, 258-261 (2003).
- 981 54. Szklarczyk, D. *et al.* The STRING database in 2021: customizable protein–protein
982 networks, and functional characterization of user-uploaded gene/measurement sets.
983 *Nucleic acids research* **49**, D605-D612 (2021).
- 984 55. Yan, X. & Han, J. gspan: Graph-based substructure pattern mining. in *2002 IEEE*
985 *International Conference on Data Mining, 2002. Proceedings.* 721-724 (IEEE, 2002).
- 986 56. Okura, A. *et al.* SGK1 in Schwann cells is a potential molecular switch involved in axonal
987 and glial regeneration during peripheral nerve injury. *Biochemical and Biophysical*
988 *Research Communications* **607**, 158-165 (2022).
- 989 57. King, R.H. *et al.* Ndr1 in development and maintenance of the myelin sheath.
990 *Neurobiology of disease* **42**, 368-380 (2011).
- 991 58. Dugas, J.C., Tai, Y.C., Speed, T.P., Ngai, J. & Barres, B.A. Functional genomic analysis
992 of oligodendrocyte differentiation. *Journal of Neuroscience* **26**, 10967-10983 (2006).
- 993 59. Ziaei, A. *et al.* Ermin deficiency leads to compromised myelin, inflammatory milieu, and
994 susceptibility to demyelinating insult. *Brain Pathology*, e13064 (2022).
- 995 60. Moffitt, J.R. *et al.* High-performance multiplexed fluorescence in situ hybridization in
996 culture and tissue with matrix imprinting and clearing. *Proceedings of the National*
997 *Academy of Sciences* **113**, 14456-14461 (2016).
- 998 61. Emanuel, G., Eichhorn, S. & Zhuang, X. MERlin-Scalable and extensible MERFISH
999 analysis software, v0. 1.6. *Zenodo doi* **10**(2020).
- 1000 62. Stringer, C., Wang, T., Michaelos, M. & Pachitariu, M. Cellpose: a generalist algorithm
1001 for cellular segmentation. *Nature methods* **18**, 100-106 (2021).
- 1002 63. Fei, J. *et al.* Quantitative analysis of multilayer organization of proteins and RNA in
1003 nuclear speckles at super resolution. *Journal of cell science* **130**, 4180-4192 (2017).
- 1004 64. Hubbard, T. *et al.* The Ensembl genome database project. *Nucleic acids research* **30**,
1005 38-41 (2002).

- 1006 65. Ferreira, J. & Zwinderman, A. On the benjamini–hochberg method. *The Annals of*
1007 *Statistics* **34**, 1827-1849 (2006).
- 1008 66. Kazemian, M., Zhu, Q., Halfon, M.S. & Sinha, S. Improved accuracy of supervised CRM
1009 discovery with interpolated Markov models and cross-species comparison. *Nucleic acids*
1010 *research* **39**, 9463-9472 (2011).
1011

Figures



Figure 1

Figure 2



Figure 2

Figure 6



Figure 3

Figure 4



Figure 4

Figure 5



Figure 5

Figure 1



Figure 6

Figure 3

Supplementary Files

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

- [SupplementaryFigures.pdf](#)
- [SupplementaryFigures.pdf](#)