

1 **Benchmarking computational methods to identify spatially variable genes and** 2 **peaks**

3 Zhijian Li^{1,2,3#}, Zain M. Patel^{1,2,3#}, Dongyuan Song⁴, Guanao Yan⁵, Jingyi Jessica Li⁵, and Luca
4 Pinello^{1,2,3*}

5 ¹Broad Institute of Harvard and MIT, Cambridge, MA, USA.

6 ²Molecular Pathology Unit, Center for Cancer Research, Massachusetts General Hospital,
7 Boston, MA, USA.

8 ³Department of Pathology, Harvard Medical School, Boston, MA, USA.

9 ⁴Interdepartmental Program of Bioinformatics, University of California, Los Angeles, CA, USA

10 ⁵Department of Statistics and Data Science, University of California, Los Angeles, CA, USA

11 #Authors contributed equally

12 *Corresponding author: Luca Pinello (lpinello@mgh.harvard.edu)

16 **Abstract**

17 Spatially resolved transcriptomics offers unprecedented insight by enabling the profiling of gene
18 expression within the intact spatial context of cells, effectively adding a new and essential
19 dimension to data interpretation. To efficiently detect spatial structure of interest, an essential
20 step in analyzing such data involves identifying spatially variable genes. Despite researchers
21 having developed several computational methods to accomplish this task, the lack of a
22 comprehensive benchmark evaluating their performance remains a considerable gap in the field.
23 Here, we present a systematic evaluation of 14 methods using 60 simulated datasets generated
24 by four different simulation strategies, 12 real-world transcriptomics, and three spatial ATAC-seq
25 datasets. We find that spatialDE2 consistently outperforms the other benchmarked methods,
26 and Moran's I achieves competitive performance in different experimental settings. Moreover,
27 our results reveal that more specialized algorithms are needed to identify spatially variable peaks.

28

29

30

31

32 Introduction

33 Recent years have witnessed significant progress in spatially-resolved transcriptome profiling
34 techniques that simultaneously characterize cellular gene expression and their physical position,
35 generating spatial transcriptomic (ST) data. The application of these techniques has dramatically
36 advanced our understanding of disease and developmental biology, for example, tumor-
37 microenvironment interactions¹, tissue remodeling following myocardial infarction², and mouse
38 organogenesis³, among others.

39

40 Spatial transcriptome profiling methods are broadly categorized into two groups, i.e., next-
41 generation sequencing (NGS)-based (including 10x Visium⁴; Slide-seq^{5,6}; HDST⁷; STARmap⁸) and
42 imaging-based (including seqFISH⁹ and MERFISH¹⁰) (**Fig. 1a**). They vary in terms of the number of
43 genes and spatial resolution. Specifically, NGS-based assays usually provide genome-wide gene
44 expression through spots profiling multiple cells, thus precluding the possibility of delineating
45 expression at the single-cell level. At the same time, the imaging-based methods can generate
46 sub-cellular resolution data but can only detect a subset of genes (30-300). Due to these
47 differences in the number of genes and spatial resolution, distinct computational methods and
48 algorithms are required for the downstream analysis of each data type. In the case of NGS-based
49 profiles, an important task involves associating cell types with spatial locations through cell-type
50 deconvolution, often leveraging paired single-cell RNA-seq data to compensate for the low
51 spatial resolution¹¹⁻¹³. On the other hand, for imaging-based profiles, the initial step involves
52 performing cell segmentation to accurately delineate the boundaries of individual cells¹⁴.

53

54 One common task for all ST profiles, regardless of the employed protocols, is to identify genes
55 that exhibit spatial patterns¹⁵ (**Fig. 1a**). These genes, defined as spatially variable genes (SVGs),
56 contain additional information about the spatial structure of the tissues of interest, compared to
57 highly variable genes (HVGs). Examples of SVGs include genes involved in developmental
58 gradients¹⁶, cell signaling pathways¹⁷, and tumor micro-environment interface¹. Additionally,
59 SVGs may be critical to downstream tasks such as detecting spatial domains¹⁸ and inferring
60 spatially aware gene regulatory networks (GRNs)¹⁹. To detect SVGs, researchers have developed
61 various computational methods by incorporating the spatial context into the analysis. As the
62 number of methods keeps increasing, it becomes difficult for users to choose the best
63 approaches effectively. Previous benchmarking studies have typically compared no more than
64 seven computational methods²⁰⁻²², significantly fewer than the currently available methods ($n >$
65 14). Furthermore, since obtaining ground truth from real-world ST profiles is not feasible, these
66 studies have relied on simulation data to evaluate the accuracy of each method in detecting SVGs.
67 However, the simulation data were generated either only using the predefined spatial

68 clusters^{20,22} or with a very limited number of spatial patterns (e.g., *spots* where the expression
69 forms round contours and *linear* where the expression forms rectangular shapes)²¹.
70 Consequently, the limitations of the simulation strategies may introduce inflating performance
71 metrics compared to realistic settings. Therefore, there is a clear need for a comprehensive
72 benchmarking study incorporating more methods and employing enhanced simulation strategies
73 to capture biologically plausible patterns of interest. Such a study would provide a more robust
74 and unbiased evaluation of the available methods for detecting SVGs in ST profiles, enabling
75 researchers to make informed decisions when selecting the most appropriate computational
76 methods for their analyses.

77

78 In this work, we comprehensively evaluated 14 methods (see **Table 1**) for identifying SVGs (the
79 selection of the 14 methods is discussed in the Discussion). We created multiple benchmarking
80 datasets ($n = 60$) with verified ground truths and compared the methods in terms of prediction
81 accuracy, sensitivity, specificity, statistical calibration, and scalability. We also investigated the
82 impact of identified SVGs on spatial domain detection. Finally, we explored the applicability of
83 the methods to other spatial modalities, specifically examining their effectiveness on spatial
84 ATAC-seq data. Our benchmark results indicate that *SpatialDE2*²³ generally outperformed the
85 other tested methods. Furthermore, *Moran's I*²⁴, despite its simplicity, consistently exhibited
86 performance either comparable to or superior to most methods in our benchmark evaluations.
87 Our results provide a detailed comparison of SVG detection methods and serve as a reference
88 for both users and method developers.

89

90 **RESULTS**

91

92 **Overview of computational methods for detection of spatially variable genes**

93 In contrast to the identification of HVGs solely from genes expression levels (i.e., mRNA molecular
94 abundance) in single-cell RNA sequencing (scRNA-seq) data, detecting SVGs requires the
95 additional consideration of spatial information at the cellular or subcellular level. A common and
96 straightforward approach is to build a k-nearest-neighbor (KNN) graph where each node
97 represents a spatial spot, and the edges between nodes represent the spatial proximity of spots.
98 SVGs are identified by combining this spatial neighbor graph with gene expression profiles. For
99 instance, *Moran's I* estimates the correlation coefficient of the expression between a spot and
100 its neighbors^{24,25}. Similarly, *Spanve* quantifies the divergence in gene expression distributions
101 between randomly and spatially sampled locations using Kullback-Leibler (KL) divergence²⁶. A
102 higher correlation or distribution divergence indicates that the gene is more likely to have a non-
103 random spatial pattern. Moreover, *scGCO* utilizes a hidden Markov random field (HMRF) to
104 capture the spatial dependence of each gene's expression levels and uses a graph cuts algorithm

105 to identify the SVGs²⁷. *SpaGCN* first builds a graph by integrating gene expression, spatial
106 location, and histology information (when available) and then clusters the spots using a graph
107 convolutional network (GCN)²⁸; then SVGs are identified by differential expression (DE) analysis
108 on the obtained clusters²⁹. *SpaGFT* constructs a KNN graph of spots based on their spatial
109 proximity and then transforms each gene's expression to the frequency domain; genes with low-
110 frequency signals tend to have less random spatial patterns³⁰. *Sepal* uses a diffusion model to
111 assess the degree of randomness of each gene's spatial expression pattern and ranks the genes
112 accordingly³¹.

113

114 Another strategy to incorporate spatial information involves utilizing a kernel function that takes
115 spatial distance as input and computes a covariance matrix to capture the spatial dependency of
116 gene expression across locations. This covariance matrix represents a prior of the underlying
117 spatial pattern. One of the pioneer methods is *SpatialDE*³², which models the normalized
118 expression data using non-parametric Gaussian Process (GP) regression and tests the significance
119 of the spatial covariance matrix for each gene by comparing the fitted models with and without
120 the spatial covariance matrix. *SpatialDE2*²³ further extends this framework by providing technical
121 innovations and computational speedups. *SPARK*³³ proposes another extension by modeling the
122 raw counts with a generalized linear model based on the over-dispersed Poisson distribution. It
123 provides a more robust statistical approach (Cauchy combination rule³⁴) to assess the significance
124 of the identified SVGs. In contrast, *BOOST-GP* uses a zero-inflated negative binomial (ZINB)
125 distribution to model the read counts and infers the model parameters via a Markov Chain Monte
126 Carlo (MCMC) algorithm³⁵. Similarly, *GPcounts*³⁶ models the counts with a negative binomial (NB)
127 distribution and estimates the model parameters using variational Bayesian inference to improve
128 computational efficiency. Notably, *SPARK-X* stands as an exception by directly comparing the
129 expression covariance matrix and the spatial distance covariance matrix, yielding substantial
130 computational efficiency gains³⁷.

131

132 In addition, two hybrid methods, namely *nnSVG* and *SOMDE*, have been developed to integrate
133 graph and kernel approaches to capture the spatial dependence between spatial spots. The
134 *nnSVG* method utilizes a hierarchical nearest-neighbor GP to model the large-scale spatial data³⁸,
135 providing computational efficiency gains over the standard Gaussian process used in *SpatialDE*.
136 On the other hand, *SOMDE* employs a self-organizing map (SOM) to cluster neighboring cells into
137 nodes and subsequently fits node-level spatial gene expression using a Gaussian process to
138 identify SVGs³⁹. Both methods reduce the computational complexity of kernel approaches by
139 leveraging a spatial graph, which significantly improves their scalability. We summarized the key
140 features of the 14 methods in **Table 1**.

141

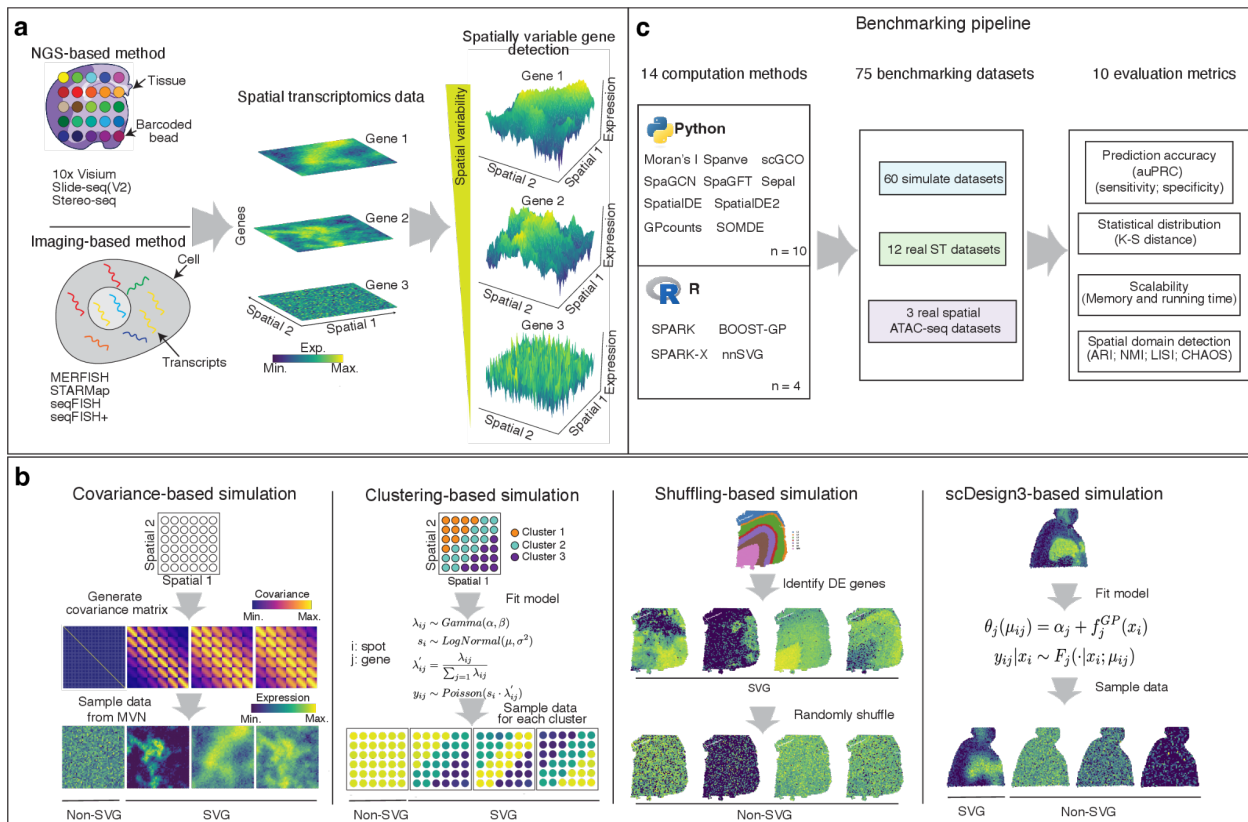
142 **Table 1| Overview of computational methods for identification of spatially variable genes.**

Method	Spatial model	Core methodology	Significance test	Input	Gene ranking	Language	Refs.
Moran's I	Graph	Correlation	Permutation	Norm.	Moran's I	Python	25, 24
Spanve	Graph	Sampling	G-test	Norm.	KL divergence	Python	26
scGCO	Graph	Graph cuts	CSR model	Norm.	FDR	Python	27
SpaGCN	Graph	Clustering	Wilcoxon test	Norm.	FDR	Python	29
SpaGFT	Graph	Fourier transform	Wilcoxon test	Norm.	GFT score	Python	30
Sepal	Graph	Diffusion model	NA	Norm.	Sepal score	Python	31
SpatialDE	Kernel	GP	Chi-square	Norm.	FSV	Python	32
SpatialDE2	Kernel	GP	NA	Norm.	FSV	Python	23
SPARK	Kernel	GP	Chi-square	Counts	Adj. p-value	R	33
SPARK-X	Kernel	Covariance test	Chi-square	Counts	Adj. p-value	R	37
BOOST-GP	Kernel	GP	BFDR	Counts	PPI	R	35
GPcounts	Kernel	GP	Chi-square	Counts	LLR	Python	36
nnSVG	Graph & Kernel	GP	LR test	Norm.	FSV	R	38
SOMDE	Graph & Kernel	GP	Chi-square	Counts	Adj. p-value	Python	39

143 *We grouped the methods based on the underlying spatial model. KL, Kullback-Leibler; GP,*
 144 *Gaussian Process; FDR, false discovery rate; HFRM, hidden Markov random field; CSR, complete*
 145 *spatial randomness, FSV, fraction of spatial variance; BFDR, Bayesian false discovery rate; PPI,*
 146 *posterior probabilities of inclusion; LR, likelihood ratio.*

147

148



149
 150 **Fig. 1 Overview of spatial transcriptome profiling protocols, benchmarking datasets with**
 151 **simulation designs, and benchmarking workflow.** **a**, Left: a schematic showing the NGS-based
 152 and imaging-based technologies for profiling spatially resolved transcriptomes. Middle:
 153 Visualization of gene expression with various patterns in spatial space. Colors refer to the
 154 expression levels of genes. Right: 3D plots showing the expression of the genes with different
 155 spatial patterns. A gene with a highly spatially correlated expression pattern is defined as a
 156 spatially variable gene (SVG; shown on the top), otherwise as a non-SVG (shown on the bottom).
 157 The x-axis and y-axis represent spatial coordinates, and the z-axis represents the expression of
 158 that gene. **b**, Schematics showing four approaches to simulate spatial transcriptomics datasets
 159 with ground truths. In the covariance-based simulation, we sampled data from a multivariate
 160 normal (MVN) with different covariance matrices for SVGs and non-SVGs. In the clustering-based
 161 simulation, we generated SVGs as differentially expressed genes for pre-defined spatial clusters.
 162 In the shuffling-based simulation, we first identified cluster-specific DE genes as SVGs and then
 163 generated the non-SVGs through data shuffling. In scDesign3-based simulation, we modeled a
 164 gene's expression as a function of spatial locations via Gaussian Process regression. **c**,
 165 Benchmarking workflow. We compared 14 computational methods on 60 simulated, 12 spatial
 166 transcriptomics, and three spatial ATAC-seq datasets. The evaluation metrics include prediction
 167 accuracy (measured by auPRC, sensitivity, and specificity), statistical distribution similarity
 168 (measured by K-S distance), scalability (measured by memory and running time), and spatial
 169 domain detection accuracy (measured by ARI, NMI, LISI, and CHAOS). K-S, Kolmogorov-Smirnov.

170
 171
 172 **Benchmarking datasets and pipeline**

173 In this study, the primary challenge we faced while benchmarking the 14 methods for detecting
174 SVGs was the lack of established datasets with verified true labels (i.e., true SVGs and non-SVGs)
175 in real-world scenarios. Hence, we focused on simulated data, an approach grounded in
176 precedent studies^{26,27,33,37,38}. Addressing the limitation of previous simulations that
177 predominantly utilized pre-defined spatial clusters — a strategy failing to mirror the rich diversity
178 of spatial patterns —, we formulated three innovative strategies and employed an recent
179 simulation framework to foster a more representative simulation dataset: covariance-based,
180 clustering-based, shuffling-based, and scDesign3-based simulation, illustrated in **Fig. 1b**.

181
182 In the covariance-based simulations, we sampled gene expression data from a multivariate
183 normal (MVN) distribution where the covariance matrix was pre-defined based on the spatial
184 coordinates (see Methods). To generate SVGs with various spatial patterns, we employed
185 multiple Gaussian kernels with diverse length scales to define the covariance matrix. We also
186 controlled the noise levels and covariance amplitudes to introduce varying degrees of complexity
187 (**Supplementary Fig. 1**). For non-SVGs, we simply used the identity matrix as the covariance
188 matrix. In the clustering-based simulations, we first fitted a Gamma-Poisson mixture model on
189 real-world spatial transcriptomics profiles from breast tumors with annotated spatial clusters⁴⁰.
190 We then generated synthetic data for each cluster by manipulating the log-fold change for each
191 gene to simulate different gene expression levels and to assess the sensitivity of each SVG
192 detection method (**Supplementary Fig. 2**). In the shuffling-based simulation, we downloaded a
193 spatial transcriptomics dataset generated from the human dorsolateral prefrontal cortex
194 (DLPFC)⁴¹ with distinct and well-annotated spatial clusters, and we obtained “true labels” based
195 on differential expression analysis and data shuffling. Briefly, we considered the cluster-specific
196 markers as true SVGs and randomly shuffled the spots to remove spatial correlation, creating
197 non-spatially variable expressions (**Supplementary Fig. 3a-c**). Finally, we used scDesign3⁴², a
198 recent simulation framework for generating realistic spatial transcriptomics datasets with pre-
199 specified true SVGs (**Supplementary Fig. 4**). Details of the simulated datasets were provided in
200 **Supplementary Table 1**.

201
202 Using the simulated datasets as described above, we benchmarked 14 SVG detection methods
203 to identify spatially variable genes, covering most of the currently available methods for this task
204 as detailed in **Table 1** (for the method selection, see Discussion). We evaluated their prediction
205 performance based on the area under the Precision-Recall curve (AUPRC), a widely accepted
206 metric for assessing classification accuracy. Additionally, we compared the sensitivity, specificity,
207 and statistical calibration of the methods. Using simulated data, we also investigated the memory
208 requirements and time scalability of the methods in relation to the number of spatial spots.
209 Importantly, considering potential downstream applications, we evaluated the impact of the

210 detected SVGs on spatial domain detection and measured the performance of this task against
211 the true labels using commonly used metrics such as the Adjusted Rand Index (ARI) and
212 Normalized Mutual Information (NMI). Finally, we explored the possibilities of applying these
213 methods, which were developed for spatial transcriptomics data, to spatial ATAC-seq data for
214 detecting spatially variable peaks (SVPs). The results were evaluated based on clustering analysis
215 using the local inverse Simpson's index (LISI) and the spatial chaos score (CHAOS) metrics¹⁸. The
216 overall benchmarking workflow is presented in **Fig. 1c**.

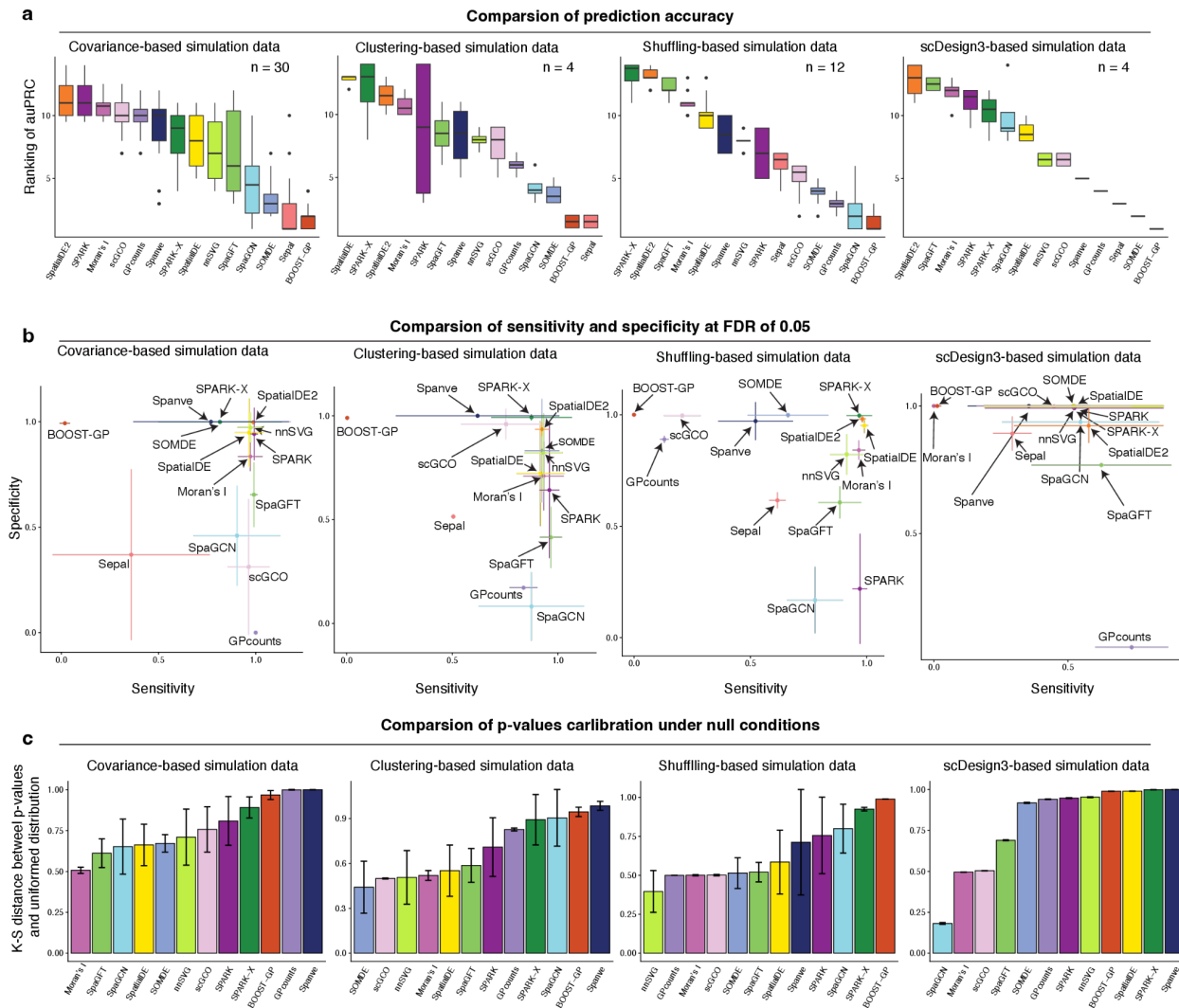
217

218 **Benchmarking prediction performance of the methods**

219 We reasoned that identifying SVGs can be considered as a binary classification problem where
220 the task is distinguishing SVGs from non-SVGs based on the statistical significance of a calculated
221 score that should capture the degree of spatially variable pattern. Currently available methods
222 typically provide different scores to rank the genes. For example, *SpatialDE* and *SpatialDE2* use
223 the fraction of spatial variance (FSV) estimated by the GP regression model, while *SpaGFT* defined
224 a GFT score as the sum of the low-frequency Fourier coefficients. We first assessed if these scores
225 could correctly separate SVGs from non-SVGs. To this end, we applied the 14 algorithms to 50
226 simulated datasets generated using different strategies and evaluated the results using auPRC.
227 We observed that the methods exhibited various accuracies across the benchmarking datasets
228 (**Supplementary Fig. 5a-b**). Specifically, for the covariance-, clustering-, and scDesign3-based
229 simulation datasets, most algorithms achieved a high auPRC at a modest noise level, and their
230 performance declined as the noise level increased. On the other hand, for shuffling-based
231 simulation data, we found instead that *SPARK-X*, *SpatialDE2*, *SpaGFT*, and *Moran's I* showed
232 consistently higher auPRC than alternative methods.

233

234 To compare the performance, we ranked the methods based on auPRC for each experimental
235 setting and visualized the overall results within each dataset and across all the datasets (**Fig. 2a**;
236 **Supplementary Fig. 5c**). Remarkably, we found that *SpatialDE2* outperformed all other methods
237 on two datasets (i.e., covariance- and scDesign3-based simulations), performed the second-best
238 on the shuffling-based simulation data, and performed the third-best on clustering-based
239 simulation data, demonstrating its robust performance. Interestingly, our evaluation revealed
240 that *Moran's I* statistic, which solely relies on auto-correlation between spots and their
241 neighbors, showed the second-best performance despite its relative simplicity compared to
242 other methods (**Supplementary Fig. 5c**). Moreover, *SPARK* and *SPARK-X* displayed similar
243 performance, likely because they used the same kernel functions to capture spatial dependency.



244
245

246 **Fig. 2 Comparison of the methods using simulated datasets.** **a**, Box plot comparing the
 247 prediction performance of the methods for covariance-based (n=30), clustering-based (n=4),
 248 shuffling-based (n=12), and scDesign3-based (n=4) simulation datasets. The y-axis represents the
 249 rank of the method based on auPRC. A higher rank denotes a higher auPRC. **b**, Evaluation of
 250 sensitivity and specificity of each method for a false discovery rate (FDR) of 0.05. Each dot
 251 represents an average true positive rate and a true negative rate. The error bar represents the
 252 standard deviation of the corresponding values. **c**, Bar plot comparing the statistical calibration
 253 evaluated by K-S distance between the distribution of empirical p-values and uniformed
 254 distribution for the null hypothesis. A lower K-S distance represents a more calibrated model. K-
 255 S: Kolmogorov–Smirnov.

256

257

258

259

260

261

Sensitivity, specificity, and statistical distribution under null conditions

262 Many of the benchmarked methods also calculate statistical significance, enabling users to
263 identify the most relevant SVGs *ad hoc*. However, because they utilize distinct statistical tests
264 based on different null hypotheses, it is unclear how sensitive and specific the results are. To
265 address this, we subsequently analyzed the methods' sensitivity (true positive rate) and
266 specificity (true negative rate) at a false discovery rate (FDR) of 0.05. Of note, *SpatialDE2* and
267 *Sepal* were excluded from this analysis because they do not provide statistical significance results.

268

269 In terms of sensitivity, we found that most methods achieved high values on the datasets with
270 low noise. However, the performance decreased when the noise level increased, a trend that
271 mirrored our findings in the accuracy evaluation (**Supplementary Fig. 5a**). Intriguingly, our
272 analysis showed that no single method consistently outperformed the rest in both sensitivity and
273 specificity across all benchmark datasets (**Fig. 2b**). For example, *SPARK* and *SpaGFT* displayed high
274 sensitivity but low specificity. In contrast, *Spanve* and *SOMDE* showed high specificity but low
275 sensitivity (**Supplementary Fig. 5-6**). These findings suggest that more sophisticated statistical
276 approaches are needed to control both false positives and false negatives. Nevertheless, we
277 found that *SpatialDE* exhibited the best balance between sensitivity and specificity.

278

279 Additionally, we evaluated the p-value distribution of the methods under null conditions for each
280 dataset. To this end, we measured the Kolmogorov–Smirnov (K-S) distance between the
281 distribution of the computed p-values and the uniform distribution (ranging from 0 to 1). The
282 intuition is that a well-calibrated model should produce uniformly distributed p-values between
283 0 and 1 under the null condition. Therefore, a smaller distance represents a better-calibrated
284 approach. Our analysis revealed that the methods demonstrated various degrees of calibrations
285 in different datasets (**Fig. 2c; Supplementary Fig. 8a-b**). For instance, *SOMDE* showed the best
286 calibration on the clustering-based simulation dataset but was not well-calibrated on the other
287 three datasets. Next, we aggregated the results across all the benchmarking datasets to compare
288 the methods comprehensively. We found that *Moran's I* exhibited the best calibration among
289 the selected methods. This can potentially be attributed to that this method used permutation
290 to estimate the background distribution, thereby accurately recapitulating the true negatives
291 (i.e., non-SVGs).

292

293

294 **Benchmarking the scalability of the methods**

295 Subsequently, we evaluated the space and time scalability of the analyzed methods. Given that
296 all methods independently estimate the spatial variability for each gene, the scalability, in theory,
297 is primarily influenced by the number of spatial locations. To benchmark this aspect, we
298 generated ten simulation datasets, each consisting of the same number of genes ($n = 100$) but

299 varying the number of spots, ranging from 100 to 40000. We applied every method to each of
300 the ten simulation datasets and recorded the memory consumption and running time as
301 performance metrics (**see Methods**).

302

303 Our initial examination of memory usage revealed that most methods displayed moderate
304 memory requirements, typically staying below 32 GB, even when confronted with datasets
305 containing 40000 spots (**Fig. 3a**). For example, we observed that *Moran's I* consumed less than
306 4GB for all datasets. This favorable outcome suggests that these methods can be executed on
307 modern laptops without encountering memory constraints. Among them, *SOMDE* exhibited the
308 most efficient memory usage across all benchmarking datasets, followed by *Spanve* and *SPARK-*
309 *X* (**Supplementary Fig. 9a**). In contrast, both *SPARK* and *SpatialDE* exhibited significant increases
310 in memory demand as the number of spots in the dataset increased. For instance, when applied
311 to a dataset with 20000 spots, *SPARK* necessitated approximately 250 GB of memory, while
312 *SpatialDE* consumed roughly 150 GB when dealing with a dataset containing 40000 spots. These
313 observations can be attributed to the fact that both *SPARK* and *SpatialDE* are based on Gaussian
314 Process regression, requiring the estimation of a covariance matrix across all spots.
315 Consequently, this leads to a cubic scaling relationship with the number of spots, resulting in the
316 pronounced memory consumption we observed for these two methods.

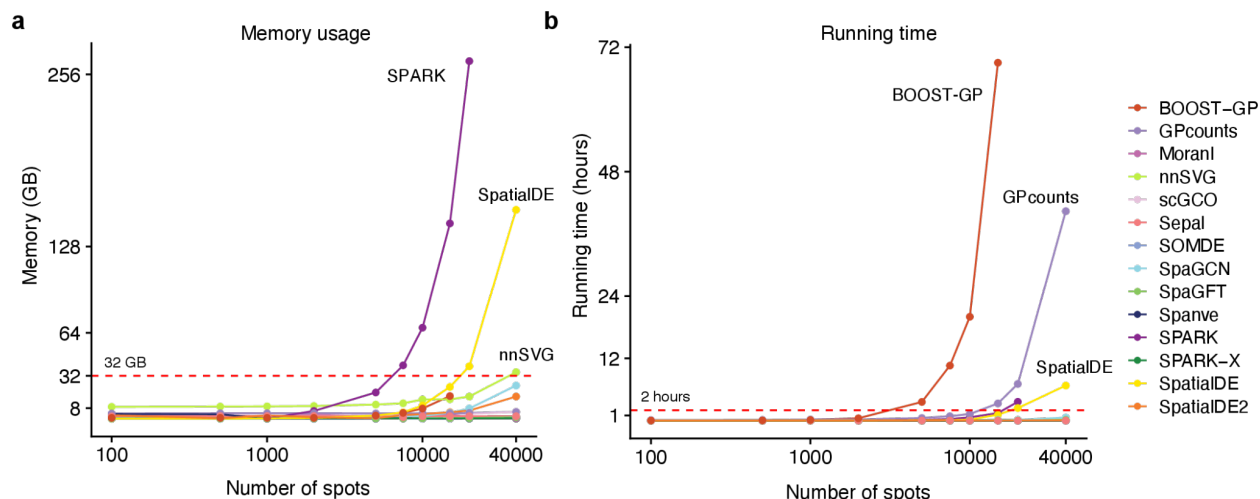
317

318 Regarding running time, we observed that *SOMDE* achieved the best scalability again, closely
319 followed by *SPARK-X* and *scGCO*. Notably, most methods completed their computations within a
320 reasonable timeframe of about 2 hours (**Fig. 3b; Supplementary Fig. 9b**), making them suitable
321 for practical usage. Both *BOOST-GP* and *GPcounts* exhibited poor scalability with increasing
322 numbers of spots. For instance, *BOOST-GP*'s computational time escalated significantly, requiring
323 three days to process a dataset containing 20000 spots and failing to produce results within five
324 days for a dataset with 40000 spots. Similarly, despite running on a GPU, *GPcounts* still require
325 approximately 45 hours to process the largest datasets. In summary, our analysis revealed that
326 *SOMDE* and *SPARK-X* exhibited the most favorable scalability when handling datasets with an
327 increasing number of spots. In addition, *SpatialDE2* and *Moran's I* statistics, the top two
328 performers in the evaluation of prediction accuracy, also demonstrated competitive scalability.

329

330

331



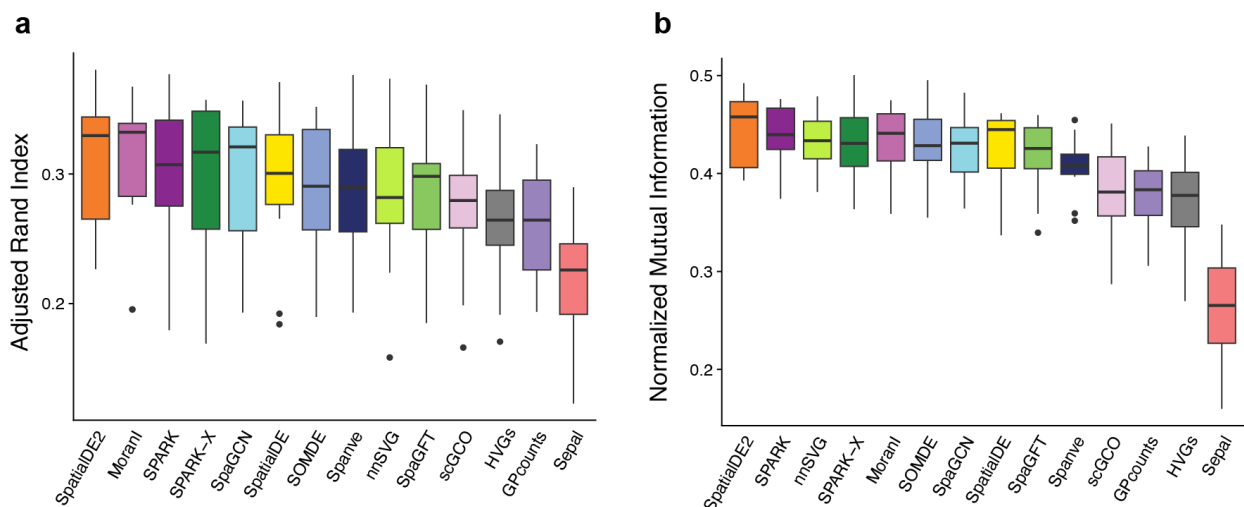
332
333

334 **Fig. 3 Scalability of the methods.** **a**, Line plot showing the memory scalability of the methods.
335 The x-axis represents the number of spots (log10) of the input datasets with 100 genes. The y-
336 axis represents consumed memory (measured as GB) by each method. The red dash line denotes
337 32 GB. We labeled the top four methods. **b**, Same as **a** for time scalability. The y-axis represents
338 the consumed time (measured as hours) of each method.

339
340

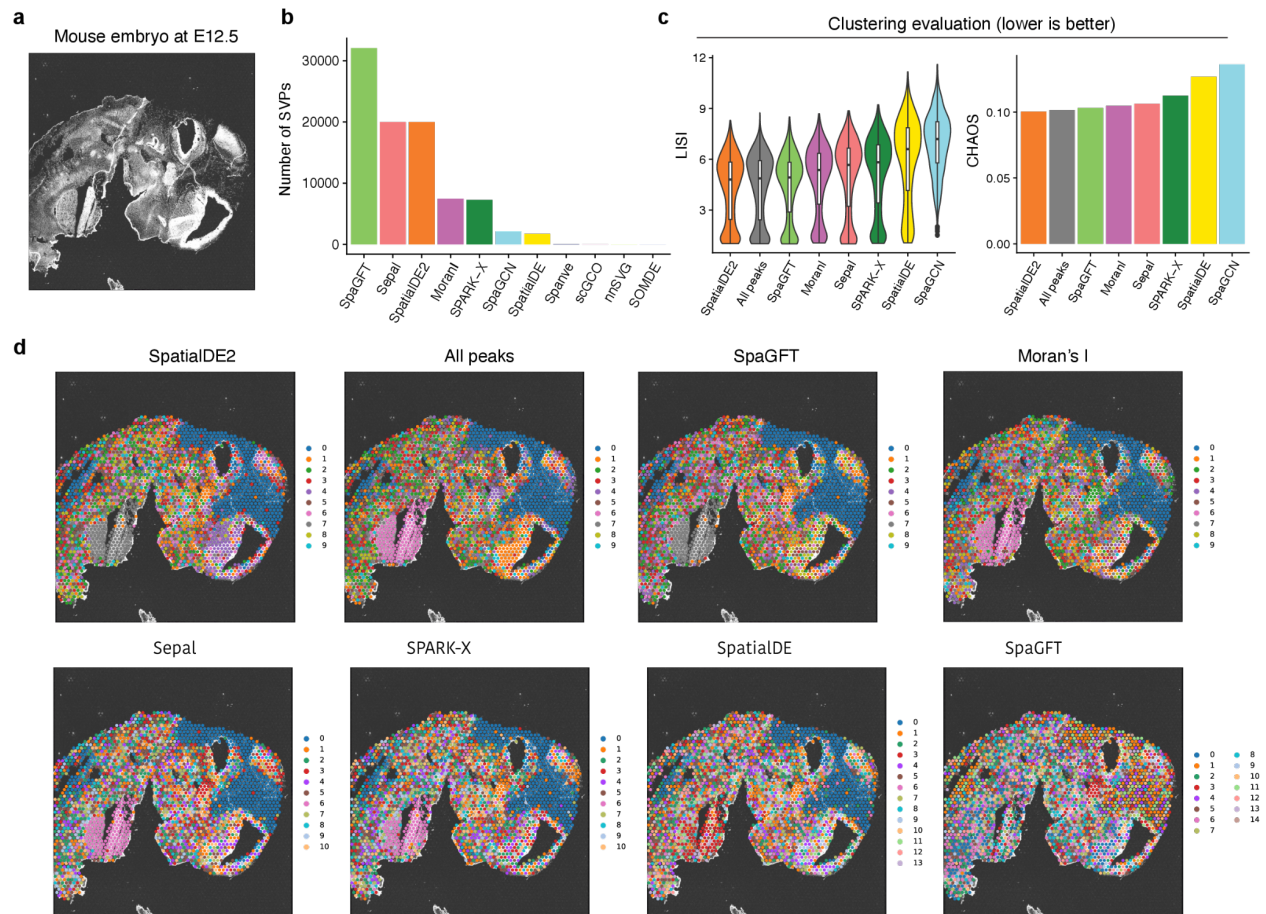
341 Benchmarking the impact of identified SVGs on spatial domain detection

342 One of the important applications of spatially resolved transcriptomics is the identification of
343 tissue or region substructures through domain detection analysis. In non-spatially resolved
344 scRNA-seq data, it is a standard practice to utilize HVGs as features for cell clustering⁴³. Therefore,
345 we hypothesized that employing SVGs could similarly be beneficial for spatial domain detection.
346 Using the human DLPFC datasets, we first evaluated which method might capture the most
347 informative features for this task. To this end, we ran the methods to identify SVGs for each
348 dataset and observed significant variations in the number of detected SVGs, highlighting
349 discrepancies between the methods (**Supplementary Fig. 10a**). Specifically, *scGCO*, *SOMDE*,
350 *GPcounts*, and *Spanve* tended to yield a low number of SVGs (<1000) across all datasets. In
351 contrast, *SpaGFT*, *Moran's I*, *SpaGCN*, and *SpatialDE* generated more SVGs. Because *SpatialDE2*
352 and *Sepal* do not perform significance test, we here used the top 2000 genes based on the FSV
353 and *Sepal* scores, respectively (see Methods).



354
355 **Fig. 4 Impact of detected SVGs on spatial domain detection analysis. a,** Box plots showing the
356 clustering performance as evaluated using ARI. Methods are ranked by the average value. **b,**
357 Same as **a** for NMI.
358

359 Subsequently, we used the graph-based Leiden algorithm⁴⁴ (resolution = 1) to cluster the spatial
360 spots based on the detected SVGs as input features for each method and dataset. To establish a
361 baseline for comparison, we also selected the top 2000 HVGs, therefore discarding spatial
362 information in this feature selection procedure. The detection results were evaluated against the
363 annotated spatial domains using two metrics: Adjusted Rand Index (ARI) and Normalized Mutual
364 Information (NMI) (see Methods for details). Remarkably, we observed that all methods, except
365 for *GPcounts* and *Sepal*, exhibited improved accuracy when utilizing SVGs compared to using only
366 HVGs (**Fig. 4a-b; Supplementary Fig. 10b**). This finding underscores the importance and power
367 of incorporating spatial information into this analysis, which can better capture the spatial
368 organization and tissue structures. Among the evaluated methods, *SpatialDE2* demonstrated the
369 highest average ARI (0.31) and NMI (0.44), further confirming its superior performance.
370 Additionally, *Moran's I* achieved the second-highest average ARI (0.303), closely followed by
371 *SPARK* (0.301) and *SPARK-X* (0.296). Concerning the NMI metric, *SPARK* ranked second-best with
372 an average value of 0.438, followed by *nnSVG* (0.434) and *SPARK-X* (0.43). In conclusion, our
373 analysis highlighted that incorporating SVGs can notably enhance clustering accuracy in spatial
374 transcriptomic analysis. Furthermore, it revealed that *SpatialDE2*, *SPARK*, and *SPARK-X* generally
375 outperformed other methods in this context, showcasing their effectiveness in capturing
376 meaningful spatial patterns and facilitating the discovery of tissue structures in spatial
377 transcriptomics data.



378
 379 **Fig. 5 Benchmarking the methods on spatial ATAC-seq data.** **a**, Image of a mouse embryo at days
 380 of E12.5. **b**, Number of detected spatially variable peaks by each method. **c**, Left: violin plot
 381 showing the LISI scores. Methods are sorted by the median values. Right: The bar plot shows the
 382 CHAOS score. For both metrics, a lower value represents a better performance. **d**, Visualization
 383 of obtained clusters by using spatially variable peaks identified by different methods.

384
 385

386 Benchmarking the methods with spatial ATAC-seq profiles

387 Recent technological advances have allowed for profiling spatially-resolved chromatin
 388 accessibility^{45,46}. However, specific methods for detecting spatially variable open chromatin
 389 regions (i.e., spatially variable peaks, abbreviated as SVPs) are currently lacking. In this section,
 390 we aimed to investigate the feasibility of applying methods developed for SVG detection to
 391 analyze spatial chromatin accessibility profiles. For this, we downloaded spatial ATAC-seq data
 392 from mouse gestational development at embryonic days of E12.5⁴⁶ (**Fig. 5a**). Following data
 393 processing, we obtained a dataset consisting of 2246 spatial spots and 34460 peaks representing
 394 open chromatin regions (see **Methods**). Subsequently, we tried to employ each of the 14
 395 methods to detect SVPs. However, given that these methods were not specifically designed for
 396 this task, we encountered several challenges. *BOOST-GP* and *GPcounts* failed to produce results

397 even after 120 hours of running, due to the fact the number of peaks exceeded substantially the
398 number of genes, highlighting the limitation of these two methods in terms of scalability.
399 Additionally, *SPARK* encountered memory issues and did not yield any results. As in the previous
400 section, we wanted to investigate if SVPs recovered from these procedures could boost spatial
401 clustering. Since *SpatialDE2* and *Sepal* do not provide statistical results, we here used the top
402 20000 peaks. For other methods, we determined the peaks at the FDR of 0.05. Importantly, we
403 observed considerable variation in the number of SVPs detected by different methods (**Fig. 5b**).
404 For example, *nnSVG* and *SOMDE* did not identify any significant peaks, indicating their limitations
405 in capturing spatial variability in this context. In contrast, *SpaGFT* identified almost all the peaks
406 as significant as SVPs (32079 out of 34460)

407

408 In the subsequent step, we used Leiden-based clustering analysis—utilizing the SVPs to group the
409 spots—to evaluate the quality of SVPs discovered by each method. We excluded *Spanve* and
410 *scGCO* for this analysis as they only detected 39 and 26 SVPs. Because the ground truth is
411 unavailable in this dataset, we measured the spatial locality and continuity of the clusters using
412 two metrics: the local inverse Simpson's index (LISI) and the spatial chaos score (CHAOS)¹⁸. The
413 underlying assumption is that a more accurate identification of SVPs would yield more
414 continuous and cohesive clusters¹⁸. We also included the results generated using all the peaks as
415 a baseline. Interestingly, we observed that *SpatialDE2* outperformed other methods (median LISI
416 = 4.8; CHAOS = 0.1), indicating that it has good potential to identify SVPs (**Fig. 5c-d**). Surprisingly,
417 our analysis revealed that using all peaks yielded the second-best performance (median LISI =
418 4.87; CHAOS = 0.102). This finding suggests that more specialized methods are required to
419 analyze spatial chromatin accessibility data. Similar results were also observed from the spatial
420 ATAC-seq from embryos at E13.5 and E15.5 (**Supplementary Fig. 11**).

421

422 **Discussion**

423 Recently, over a dozen computational methods have been developed to identify spatially variable
424 genes for spatial transcriptomics data. These methods diverge substantially in several aspects,
425 including the assumptions in modeling spatial relationships between cells (graph vs. kernel), the
426 algorithms to estimate spatial variation (e.g., auto-correlation vs. Gaussian Process regression vs.
427 graph cut), the statistical tests to determine significances (e.g., permutation test vs. Wilcoxon
428 test vs. Chi-square test), the choice of input data (raw counts vs. normalized data), and the
429 programming languages (Python vs. R) (**Table 1**). These factors complicate the selection of
430 methods for users, a situation exacerbated by the current absence of systematic benchmarking
431 of the methods' performance.

432

433 In this study, we systematically evaluated the performance of 14 methods for detecting SVGs
434 using simulated and real-world data. Compared to previous works^{20,22}, we used four different
435 approaches to generate simulation data. The rationale behind this quadruple simulation strategy
436 is to minimize potential biases and prevent the undue advantages of certain methods on specific
437 types of simulated data. For instance, methods like *SpatialDE*, which models spatial covariance
438 in their algorithmic framework, might overestimate performance when evaluated on covariance-
439 based simulation data. On the other hand, methods such as *SpaGCN* can benefit from clustering-
440 based simulation as this method utilizes clusters to identify SVGs. Moreover, the shuffling-based
441 simulation enables the testing of methods against real-world spatial transcriptomics data. In
442 addition, scDesign3-based simulation can generate *in silico* spatial transcriptomics data,
443 enhancing our simulation with a method capable of explicitly accounting for dependencies
444 between genes. Overall, our benchmark datasets covered a variety of scenarios and represented
445 a useful resource for developing and testing methods in the future.

446

447 Our evaluation results revealed that *SpatialDE2* generally outperformed other methods by
448 providing a high average auPRC across various experimental settings, however, this method does
449 not provide any statistical significance for the recovered genes. Interestingly, we found that
450 *Moran's I* achieved the second-best prediction performance, although it is simply based on auto-
451 correlation between spots and their spatial neighbors, which has been neglected in previous
452 benchmarks^{23,26,27,30,32,33,37}. Going forward, it would be prudent to include *Moran's I* as a baseline
453 in future SVG benchmarking. Additionally, we observed comparable effectiveness between
454 kernel-based and graph-based methods on simulation data and real-world datasets, suggesting
455 their capability to effectively capture similar spatial dependencies. Regarding the sensitivity and
456 specificity, we observed that no single method consistently outperformed the others for both
457 metrics on all benchmarking datasets, indicating that robustly estimating statistical significance
458 remains a difficult problem. Our analysis highlighted the superior p-value calibration of *Moran's*
459 *I*, which is attributable to their use of permutation tests that produce well-calibrated statistical
460 significance.

461

462 Scalability is another crucial aspect to consider, especially with the emergence of large-scale
463 spatially-resolved profiling methods capable of capturing sub-cellular resolution and
464 accommodating an increasing number of spots, exemplified by Stereo-seq³ ($> 10^4$ spots). Our
465 investigations revealed notable distinctions in scalability between graph-based and kernel-based
466 methods, with the former generally outperforming the latter. Among the methods we examined,
467 *SOMDE* stood out as the most efficient in both memory utilization and running time
468 (**Supplementary Fig. 9**). It is important to note that *SOMDE* initially clusters adjacent data points
469 into graph nodes and then employs GP regression to identify SVGs in a node-centric manner. This

470 strategy significantly mitigates the complexities associated with both time and memory. SPARK-
471 X, as a kernel-based method, demonstrated comparable performance to *SOMDE* by directly
472 comparing the expression and spatial distance covariance matrix rather than using GP regression
473 to estimate spatial variation, unlike its predecessor *SPARK*. Moreover, we found that *SpatialDE2*
474 demonstrated reasonable scalability. Given its superior prediction performance, we envision that
475 this method could be the default one to use in practice. Of note, this method provides no
476 statistical significance for its SVG identification results. Therefore, we recommend selecting the
477 top genes based on the fraction of spatial variation (FSV) for downstream analysis. In summary,
478 our findings not only underscore the significance of scalability in the context of SVG detection
479 but also shed light on the relative advantages of different analytical methods when processing
480 large-scale datasets. Future methods should consider scalability alongside prediction
481 performance as advanced spatial profiling techniques produce better quality and larger quantity
482 of data.

483

484 In addition, we also demonstrated that the incorporation of SVGs identified by 12 out of the 14
485 methods led to a notable enhancement in spatial domain detection when applied to real data
486 with annotated clusters, as opposed to relying solely on HVGs. These results imply the
487 significance of capturing spatial information in improving clustering analysis by incorporating
488 more comprehensive information on the architecture of complex tissues and tumors. As novel
489 technologies like Slide-Tag⁴⁷ emerge, enabling the simultaneous acquisition of single-cell
490 measurements and spatial data, we anticipate a surge in the adoption and popularity of SVG
491 identification tools in various downstream analysis tasks.

492

493 Finally, we also showed that some methods can be applied to other modalities like Spatial-ATAC
494 seq, facilitating the identification of potential SVPs. We use the term “potential” due to the
495 absence of a ground truth; instead, we leveraged the SVPs for clustering and evaluated the spatial
496 locality and continuity of the obtained clusters. It is essential to note that not all methods were
497 capable of detecting SVPs due to limitations in memory or algorithmic complexity. Therefore,
498 there is a pressing need to develop novel methodologies or modify existing ones to make them
499 applicable to spatial-ATAC seq data. Tools focused on discerning SVPs have the potential to reveal
500 the regulatory elements that govern gene expression profiles within specific spatial sub-regions.
501 This, in turn, can enhance our understanding of the regulatory mechanisms governing SVGs and,
502 consequently, the spatial organization of tissues and tumors. In the future, integrating SVGs and
503 SVPs through novel algorithms holds tremendous potential to facilitate the construction of
504 accurate spatially aware gene regulatory networks.

505

506 Although we have covered a large number of available methods (n = 14) in the present study,
507 there are still some methods that are not included. This is because either the repository has not
508 been maintained for a long time, resulting in outdated dependencies that make it difficult to
509 install and execute, for instance, *trendsceek*⁴⁸, or the method was unavailable during the
510 preparation of our manuscript, for example, *BSP*⁴⁹. Another limitation of our work is its exclusive
511 focus on spatial transcriptomics and spatial ATAC-seq, despite the advent of other spatially-
512 resolved omics data, including spatial proteomics. Future directions may also include testing and
513 or adapting SVG detection methods on these modalities. Nonetheless, our benchmarking study
514 provides a detailed evaluation of various SVG detection methods across simulated and real-world
515 datasets of spatial transcriptomics and spatial-ATAC-seq.

516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547

Methods

Simulation datasets

We used for different approaches to generate simulated spatial transcriptomics data with ground

548 truth. The details are described below.

549

550 **Covariance-based simulation.** This simulation is based on a pre-defined covariance matrix.
551 Specifically, given m spatial locations where each location is represented by its coordinates x ,
552 for each gene, we first calculated a covariance matrix $K \in R^{m \times m}$ based on multiple Gaussian
553 kernels:

$$554 \quad K(a, b) = \sum_{n=1}^N \beta_n \cdot \exp\left(\frac{\|x(a) - x(b)\|^2}{2 \cdot l_n^2}\right)$$

$$555 \quad (\beta_1, \beta_2, \dots, \beta_N) \sim \text{Dirichlet}(1/N, \dots, 1/N)$$

556 where $x(a)$ and $x(b)$ denote two spatial locations, N is the number of kernels, β_n is the weight
557 of the n th kernel and is sampled from a Dirichlet distribution, l_n denotes the length scale. By
558 sampling β for each gene, we obtained different spatial covariance matrices.

559

560 We next sampled expression $\lambda_j \in R^m$ for gene j across all locations from a multivariate normal
561 distribution (MVN) as follows:

$$562 \quad \log(\lambda_j) \sim (1 - \alpha) \cdot \text{MVN}(\mu, \sigma^2 \cdot K) + \alpha \cdot \text{MVN}(\mu, \sigma^2 \cdot I)$$

563 where σ^2 represents the amplitude of the spatial covariance; I is an identity matrix (i.e. with zeros
564 everywhere except on the diagonal); $\alpha \in [0, 1]$ denotes the noise level in simulated gene
565 expression. When $\alpha = 1$, signals are sampled from an MVN without spatial correlation, thus they
566 are considered non-spatially variable genes. Because some methods can only work on raw
567 counts, we next converted the data to counts as follows:

$$568 \quad \lambda_{ij}' = \frac{\lambda_{ij}}{\sum_{j=1} \lambda_{ij}}$$

$$569 \quad y_{ij} \sim \text{Poisson}(s \cdot \lambda_{ij}')$$

570 where s denotes the library size and is set to 10,000 for all locations.

571

572 To evaluate the prediction accuracy, we generated simulation data on a 50 by 50 grid layout (in
573 a total of 2500 spots) by setting $N = 5$ and using $l \in [1, 3.25, 5.5, 7.75, 10]$ to generate different
574 covariance matrices. Moreover, we used five noise levels $\alpha \in [0, 0.2, 0.4, 0.6, 0.8]$ and six different
575 amplitudes of the covariance matrix $\sigma \in [0.5, 1, 1.5, 2, 2.5, 3]$ to generate 30 simulation datasets
576 (**Supplementary Fig. 1**). We generated 100 SVGs and 100 non-SVGs for each dataset using the
577 above process.

578

579 To benchmark the scalability of the methods with the number of spatial spots, we generated ten
580 simulation datasets as described above. Each dataset had the same number of genes ($n = 100$)
581 and a different number of spots ($n = 100, 500, 1000, 2000, 5000, 7500, 10000, 15000, 20000,$
582 40000).

583

584 **Clustering-based simulation.** This simulation is based on spatial data with annotated clusters. To
585 generate a simulation that can recapitulate a real ST dataset, we followed the two-step strategy
586 proposed in SRTsim⁵⁰, i.e., first estimating the parameters required for the simulation from a real
587 ST dataset and then generating a synthetic dataset based on the estimated parameters.
588 Specifically, given a real-world dataset with a count matrix $Y \in R^{m \times n}$ where m is the number of
589 spatial locations, n is the number of genes, and y_{ij} is the expression of gene j in spatial location
590 i , we modeled the count y_{ij} using the following process:

591
$$y_{ij} \sim \text{Poisson}(s_i \cdot \lambda_{ij}')$$

592
$$\lambda_{ij}' = \frac{\lambda_{ij}}{\sum_{j=1} \lambda_{ij}}$$

593
$$s_i \sim \text{LogNormal}(\mu, \sigma^2)$$

594
$$\lambda_{ij} \sim \text{Gamma}(\alpha, \beta)$$

595 We denoted s_i the total number of reads in location i and assumed that it follows a Log-Normal
596 distribution parameterized by μ and σ . We denoted λ_{ij} the log normalized mean expression of
597 gene j sampled from a Gamma distribution parameterized by α and β . λ_{ij}' represents the
598 proportion of gene j at the location i . We estimated the parameters μ , σ , α , and β using the
599 maximum likelihood algorithm based on the count matrix Y and log normalized matrix Y_{norm}
600 which was generated using functions `pp.normalize_total` and `pp.log1p` from the `scanpy`
601 package⁵¹.

602

603 Once we inferred the parameters, we used them to generate synthetic data by sampling data
604 using the above process. To obtain SVGs, we randomly selected a number of genes for each input
605 spatial cluster and multiplied the sampled mean expression by a differential factor. Since the
606 clusters are spatially associated, these marker genes are considered as SVGs. For non-SVGs, the
607 differential factors were set to one. Using breast tumors as input, we generated simulation data
608 with 100 SVGs and 100 non-SVGs. We varied the differential expression levels from 0.5 to 2,
609 generating four simulation datasets (**Supplementary Fig. 2**).

610

611 **Shuffling-based simulation.** To test the methods against real-world data, we here created true
612 labels through data shuffling. For this, we downloaded the LIBD human dorsolateral prefrontal
613 cortex (DLPFC) spatial transcriptomics data from <http://research.libd.org/spatialLIBD>. The data
614 was generated with the 10x Genomics Visium platform and included 12 samples. Each sample
615 was manually annotated as one of the six prefrontal cortex layers (L1-6) and white matter (WM).
616 We filtered the genes by the number of detected spots (>500). Next, we identified marker genes
617 for each cluster with differential expression analysis (t-test, p-value < 0.01, and logFC > 1). These
618 marker genes were considered true positives (i.e., spatially variable genes). Next, we randomly
619 permuted the spots to remove spatial correlation to generate uniformly distributed gene
620 expression profiles. We considered these genes as true negative (i.e., non-spatially variable
621 genes). This resulted in an average number of 549 true labels across all the samples
622 (**Supplementary Fig. 3**).

623

624 **scDesign3-based simulation.** scDesign3 aims to generate realistic in silico data by first learning
625 interpretable parameters from real data and then generating synthetic data. We installed
626 scDesign3 (v0.99.6) from <https://github.com/SONGDONGYUAN1994/scDesign3> and followed
627 the tutorial to generate four datasets with different numbers of true positives ranging from 50
628 to 200 (**Supplementary Fig. 4**).

629

630

631 **Identify SVGs with computational methods.**

632 We described below the details of running the methods to identify SVGs.

633

634 **Moran's I.** *Moran's I* measures the correlation of gene expression between a spatial location and
635 its neighbors²⁵. We computed *Moran's I* score using Squidpy (v1.2.3)²⁴ by following the tutorial:
636 https://squidpy.readthedocs.io/en/stable/auto_examples/graph/compute_moran.html. Spatial
637 neighbors were found using the function `spatial_neighbors`, and scores were estimated using the
638 function `spatial_autocorr`. We set parameter `n_perms` to 100 to obtain the statistical significance
639 and used 0.05 as the threshold for the adjusted p-value to identify significant SVGs. To compute
640 auPRC, we used the *Moran's I* score to rank genes.

641

642 **Spanve.** *Spanve* (Spatial Neighborhood Variably Expressed Genes) is a non-parametric statistical
643 approach for detecting SVGs²⁶. Similar to *Moran's I*, this method uses the difference between a
644 location and its spatial neighbors to estimate the spatial variation. Specifically, for each gene, it
645 computes Kullback-Leibler divergence between space-based and randomly sampled expressions.
646 The significance is calculated by the G-test. We installed *Spanve* (v0.1.0) and ran the method by

647 following the tutorial: <https://github.com/zjuggx/Spanve/blob/main/tutorial.ipynb>. Genes were
648 ranked by FDR to compute auPRC. We used 0.05 as the threshold for FDR to select significant
649 SVGs.

650

651 ***SpaGFT***. *SpaGFT* is a hypothesis-free Fourier transform model to identify SVGs³⁰. It decomposed
652 the signal from the spatial domain to the frequency domain based on a spatial KNN graph. and
653 estimated a GFTscore per gene on the Fourier coefficient for low-frequency signals. We installed
654 *SpaGFT* (v0.1.1.4) and ran it by following the tutorial: <https://spagft.readthedocs.io/en/latest>.
655 We computed auPRC for this method using GFTscore and selected significant SVGs using q-value
656 < 0.05.

657

658 ***SpaGCN***. *SpaGCN* is a graph convolutional network (GCN)-based approach that integrates gene
659 expression, spatial location, and histology to identify SVGs²⁹. It first identifies spatial domains
660 through clustering and then detects SVGs that are enriched in each domain. We installed *SpaGCN*
661 (v1.2.5) and ran the method by following the tutorial:
662 <https://github.com/jianhuupenn/SpaGCN/blob/master/tutorial/tutorial.ipynb>. We used the
663 adjusted p-values to rank genes for computing auPRC and select significant SVGs (<0.05).

664

665 ***scGCO***. *scGCO* (single-cell graph cuts optimization) utilizes a hidden Markov random field (HMRF)
666 with graph cuts to identify SVGs²⁷. We installed *scGCO* (v1.1.0) and executed the method by
667 following the tutorial: [https://github.com/WangPeng-](https://github.com/WangPeng-Lab/scGCO/blob/master/code/Tutorial/scGCO_tutorial.ipynb)
668 [Lab/scGCO/blob/master/code/Tutorial/scGCO_tutorial.ipynb](https://github.com/WangPeng-Lab/scGCO/blob/master/code/Tutorial/scGCO_tutorial.ipynb). To compute auPRC, we used FDR
669 to rank the genes. To select significant SVGs, we used 0.05 of FDR as threshold.

670

671 ***Sepal***. *Sepal* assesses the degree of randomness exhibited by the expression profile of each gene
672 through a diffusion process and ranks the genes accordingly³¹. We computed the *Sepal* score
673 using Squidpy (v1.2.3) by following the tutorial:
674 https://squidpy.readthedocs.io/en/stable/auto_examples/graph/compute_sepal.html. We used
675 the *sepal* score to rank the genes to calculate auPRC.

676

677 ***SpatialDE***. *SpatialDE* is one of the pioneer methods for identifying SVGs³². It models the
678 normalized spatial gene expression using the Gaussian process regression and estimates the
679 significance by comparing the models with and without spatial covariance. We installed
680 *SpatialDE*~(v1.1.3) with pip and processed the data with the functions `NaiveDE.stabilize` and
681 `NaiveDE.regress_out`. We ran the function `SpatialDE.run` to obtain results and used the fraction

682 of spatial variance (FSV) to compute auPRC. To select significant SVGs, we used the adjusted p-
683 values (< 0.05).

684

685 ***SpatialDE2***. *SpatialDE2* is a flexible framework for modeling spatial transcriptomics data that
686 refines *SpatialDE* by providing technical innovations and computational speedups²³. We obtained
687 the source code from <https://github.com/PMBio/SpatialDE> and estimated spatial variance using
688 the function `SpatialDE.fit`. Similar to *SpatialDE*, we ranked the genes by FSV to compute auPRC.

689

690 ***SPARK***. *SPARK* extended the computation framework proposed in *SpatialDE* by directly modeling
691 the raw count data using a generalized linear spatial model (GLSM) based on Poisson
692 distribution³³. We obtained *SPARK* (v1.1.1) from <https://github.com/xzhoulab/SPARK> and ran the
693 method by following the tutorial [https://xzhoulab.github.io/SPARK/02 SPARK Example](https://xzhoulab.github.io/SPARK/02_SPARK_Example). We
694 used the adjusted p-values to compute auPRC and select significant SVGs (< 0.05).

695

696 ***SPARK-X***. *SPARK-X* is a non-parametric covariance test method based on multiple spatial kernels
697 for modeling sparse count data from spatial transcriptomic experiments³⁷. We ran *SPARK-X*
698 (v1.1.1) by following the tutorial: [https://xzhoulab.github.io/SPARK/02 SPARK Example](https://xzhoulab.github.io/SPARK/02_SPARK_Example). We
699 used the adjusted p-values to compute auPRC and select significant SVGs (< 0.05).

700

701 ***BOOST-GP***. *BOOST-GP* is a Bayesian hierarchical model to analyze spatial transcriptomics data
702 based on zero-inflated negative binomial distribution and Gaussian process³⁵. We downloaded
703 the source codes of *BOOST-SP* from <https://github.com/Minzhe/BOOST-GP> and ran the function
704 `boost.gp` by setting the parameters `iter` to 100 and `burn` to 50. We used p-values to compute
705 auPRC and select significant SVGs using 0.05 as the threshold.

706

707 ***GPcounts***. *GPcounts* implemented Gaussian process regression for modeling counts data using a
708 negative binomial likelihood function³⁶. We obtained the source codes of *GPcounts* from
709 <https://github.com/ManchesterBioinference/GPcounts> and followed the tutorial
710 [https://github.com/ManchesterBioinference/GPcounts/blob/master/demo notebooks/GPcount](https://github.com/ManchesterBioinference/GPcounts/blob/master/demo_notebooks/GPcounts_spatial.ipynb)
711 [ts_spatial.ipynb](https://github.com/ManchesterBioinference/GPcounts/blob/master/demo_notebooks/GPcounts_spatial.ipynb). To compute auPRC, we ranked the genes by the log-likelihood ratio (LLR),
712 representing the ratio between the dynamic and constant (null) models. Significant SVGs were
713 selected based on the q-values with 0.05 as the threshold. We noted that *GPcounts* sometimes
714 failed to generate results for certain genes, especially when applied to real-world datasets. In this
715 case, we set the LLR as 0 and the q-value as 1.

716

717 **nnSVG.** *nnSVG* is a method built on nearest-neighbor Gaussian processes to identify SVGs³⁸. We
718 installed the package (v1.2.0) from Bioconductor and ran the method by following the tutorial
719 <https://bioconductor.org/packages/release/bioc/vignettes/nnSVG/inst/doc/nnSVG.html>. We
720 used the fraction of spatial variance estimated by the method to compute auPRC. Significant
721 SVGs were selected based on adjusted p-values using 0.05 as the threshold.

722

723 **SOMDE.** *SOMDE* uses a self-organizing map (SOM) to cluster neighboring locations into nodes
724 and then uses a Gaussian process to fit the node-level spatial gene expression to identify SVGs³⁹.
725 We installed *SOMDE* (v0.1.7) with pip and followed the tutorial
726 https://github.com/WhirlFirst/somde/blob/master/slide_seq0819_11_SOM.ipynb to run the
727 method. We ranked the genes by FSV to compute auPRC and selected significant SVGs based on
728 the q-values using 0.05 as the threshold.

729

730 **Benchmarking prediction performance, sensitivity, and specificity**

731 We applied each method on the simulated datasets to identify SVGs. For comparison, we
732 computed the auPRC using the function `pr.curve` from the R package `PRROC`⁵² by ranking the
733 prediction for each method accordingly (see **Table 1**). We calculated the true positive rate
734 (sensitivity) and true negative rate (specificity) at the false discovery rate of 0.05 as follows:

$$735 \quad TPR = \frac{TP}{TP + FN}$$

$$736 \quad TNR = \frac{TN}{TN + FP}$$

737 where *TP* denotes the number of true positives, *FN* denotes the number of false negatives, *TN*
738 denotes the number of true negatives and *FP* denotes the number of false positives. For
739 *SpatialDE2* and *Sepal*, we selected the top *n* genes (*n* = the number of true positives) as significant
740 SVGs to sensitivity and specificity.

741

742 **Benchmarking scalability with the number of spatial spots**

743 We used the Snakemake⁵³ workflow (v7.25.2) management system to evaluate the scalability of
744 each method with the number of spatial spots. For this, we generated simulation datasets with
745 100 genes and various numbers of spots from 100 to 40000 (see above). Next, we ran each
746 method on a dedicated HPC node with AMD EPYC 7H12 64-Core Processor using the same
747 computational resource (1TB memory, 120 hours, and 10 CPUs) defined by the Snakemake
748 pipeline. For methods (i.e., *GPcounts* and *SpatialDE2*) that require a graphics processing unit
749 (GPU) for running, we used an A100 with 40GB of memory. We measured the memory usage and
750 running time using the benchmark directive provided by the Snakemake tool (`--benchmark`).

751 Notably, we could not run *SPARK* for datasets with 40000 spots because of memory issues.
752 Moreover, *BOOST-GP* did not generate results for datasets with 20000 and 40000 spots within
753 120 hours.

754

755 **Benchmarking impact of identified SVGs on spatial domain detection analysis**

756 We utilized the human DLPFC datasets to evaluate the impact of **identified** SVGs on spatial
757 domain detection. We ran the methods and determined the SVGs based on an FDR 0.05. Since
758 *SpatialDE2* and *Sepal* did not provide statistical significance, we selected the top 2000 genes
759 based on the FSV and Sepal scores, respectively. Because *BOOST-GP* failed to produce any results
760 after 120 hours of running, we excluded it from this evaluation. In addition, we also identified
761 the highly variable genes using the function `scanpy.pp.highly_variable_genes` and used the top
762 2000 as our baseline for comparison. We next used these genes to perform dimension reduction
763 using the function `scanpy.tl.pca` and generated a k-nearest-neighbor graph with
764 `scanpy.pp.neighbors`. The clustering was conducted using the function `scanpy.tl.leiden`
765 (resolution = 1). We next compared the obtained clusters (denoted by X) with the annotated
766 layers (denoted by Y). We assessed the clustering quality with Adjusted Rand Index (ARI):

$$767 \quad ARI(Y, X) = \frac{\sum_{ij} \frac{c_{ij}}{2} - (\sum_i \frac{a_i}{2} \sum_j \frac{b_j}{2}) / \frac{n}{2}}{\frac{1}{2} (\sum_i \frac{a_i}{2} + \sum_j \frac{b_j}{2}) - (\sum_i \frac{a_i}{2} \sum_j \frac{b_j}{2}) / \frac{n}{2}}$$

768 where c_{ij} denotes the number of common spots for each obtained cluster i and ground truth j ,
769 $a_i = \sum_j^s c_{ij}$, and $b_i = \sum_i^r c_{ij}$. We also calculated the Normalized Mutual Information
770 (NMI) for comparison as follows:

$$771 \quad NMI(Y, X) = \frac{2 \cdot I(Y, X)}{H(Y) + H(X)}$$

772 where H represents entropy of the partition and I represents the mutual information between
773 clusters and the true labels. Both ARI and NMI have values from 0 to 1, with 1 indicating that the
774 two partitions are the same and 0 indicating that the two are independent.

775

776 **Benchmarking the methods for spatial ATAC-seq data**

777 We downloaded spatial ATAC-seq data of mouse embryos at stages E12.5, E13.5, and E15.5 from
778 GEO with accession number GSE214991. We first identified open chromatin regions by peak
779 calling on all the spots using `MACS2`⁵⁴ (`--nomodel --nolambda --shift -75 --extsize 150`) and
780 obtained 34460 (E12.5), 31099 (E13.5), and 69896 (E15.5) peaks for each sample, respectively.
781 We next built a cell-by-peak count matrix using the fragments and peaks as input based on the
782 function `FeatureMatrix` from the `Signac`⁵⁵ package. We only retained the spots that were located
783 on the tissue.

784

785 We ran each method on the cell-by-peak matrix of spatial ATAC-seq data from mouse embryos
786 to detect spatially variable peaks. For those methods that require normalized data as input, we
787 used TF-IDF (Term Frequency - Inverse Document Frequency) for normalization. Of note, *BOOST-*
788 *GP* and *GPcounts* failed to produce results after 120 hours, and we could not obtain results from
789 *SPARK* due to memory issues. Therefore, these three methods were excluded from this
790 evaluation. We selected the significant variable peaks using an FDR of 0.05 for the rest of the
791 methods. For *SpatialDE2* and *Sepal*, we opted for the top 20000 peaks since they did not provide
792 statistical significance. We next used the spatially variable peaks to cluster the spots. Because the
793 true clusters were unavailable, we evaluated the clustering performance by following ref.¹⁸ based
794 on the local inverse Simpson's index (LISI) and the spatial chaos score (CHAOS). The LISI score was
795 calculated as follows:

$$796 \quad S = \frac{1}{\sum_{k=1}^K p(k)}$$

797 where $p(k)$ is the probability that the cluster label k is in the local neighborhood, and K is the
798 total number of clusters. A lower LISI score indicates more homogeneous neighborhood clusters
799 of the spots. The CHAOS score was calculated as follows:

$$800 \quad CHAOS = \frac{\sum_{k=1}^K \sum_{i,j}^{n_k} d_{kij}}{N}$$

801 where d_{kij} is the Euclidean distance between the spots i and j in the cluster k and N is the total
802 number of spots. A lower CHAOS indicates better spatial continuity.

803
804

805 References

- 806 1. Hunter, M. V., Moncada, R., Weiss, J. M., Yanai, I. & White, R. M. Spatially resolved
807 transcriptomics reveals the architecture of the tumor-microenvironment interface. *Nat.*
808 *Commun.* **12**, 6278 (2021).
- 809 2. Kuppe, C. *et al.* Spatial multi-omic map of human myocardial infarction. *Nature* **608**, 766–
810 777 (2022).
- 811 3. Chen, A. *et al.* Spatiotemporal transcriptomic atlas of mouse organogenesis using DNA
812 nanoball-patterned arrays. *Cell* **185**, 1777–1792.e21 (2022).
- 813 4. Ståhl, P. L. *et al.* Visualization and analysis of gene expression in tissue sections by spatial

- 814 transcriptomics. *Science* **353**, 78–82 (2016).
- 815 5. Rodriques, S. G. *et al.* Slide-seq: A scalable technology for measuring genome-wide
816 expression at high spatial resolution. *Science* **363**, 1463–1467 (2019).
- 817 6. Stickels, R. R. *et al.* Highly sensitive spatial transcriptomics at near-cellular resolution with
818 Slide-seqV2. *Nat. Biotechnol.* **39**, 313–319 (2021).
- 819 7. Vickovic, S. *et al.* High-definition spatial transcriptomics for in situ tissue profiling. *Nat.*
820 *Methods* **16**, 987–990 (2019).
- 821 8. Wang, X. *et al.* Three-dimensional intact-tissue sequencing of single-cell transcriptional
822 states. *Science* **361**, (2018).
- 823 9. Lubeck, E., Coskun, A. F., Zhiyentayev, T., Ahmad, M. & Cai, L. Single-cell in situ RNA
824 profiling by sequential hybridization. *Nature methods* vol. 11 360–361 (2014).
- 825 10. Chen, K. H., Boettiger, A. N., Moffitt, J. R., Wang, S. & Zhuang, X. RNA imaging. Spatially
826 resolved, highly multiplexed RNA profiling in single cells. *Science* **348**, aaa6090 (2015).
- 827 11. Kleshchevnikov, V. *et al.* Cell2location maps fine-grained cell types in spatial
828 transcriptomics. *Nat. Biotechnol.* (2022) doi:10.1038/s41587-021-01139-4.
- 829 12. Cable, D. M. *et al.* Robust decomposition of cell type mixtures in spatial transcriptomics.
830 *Nat. Biotechnol.* **40**, 517–526 (2022).
- 831 13. Biancalani, T. *et al.* Deep learning and alignment of spatially resolved single-cell
832 transcriptomes with Tangram. *Nat. Methods* **18**, 1352–1362 (2021).
- 833 14. Petukhov, V. *et al.* Cell segmentation in imaging-based spatial transcriptomics. *Nat.*
834 *Biotechnol.* **40**, 345–354 (2022).
- 835 15. Palla, G., Fischer, D. S., Regev, A. & Theis, F. J. Spatial components of molecular tissue

- 836 biology. *Nat. Biotechnol.* **40**, 308–318 (2022).
- 837 16. Lohoff, T. *et al.* Integration of spatial and single-cell transcriptomic data elucidates mouse
838 organogenesis. *Nat. Biotechnol.* **40**, 74–85 (2022).
- 839 17. Cang, Z. *et al.* Screening cell-cell communication in spatial transcriptomics via collective
840 optimal transport. *Nat. Methods* **20**, 218–228 (2023).
- 841 18. Shang, L. & Zhou, X. Spatially aware dimension reduction for spatial transcriptomics. *Nat.*
842 *Commun.* **13**, 7203 (2022).
- 843 19. de Luis Balaguer, M. A. *et al.* Predicting gene regulatory networks by combining spatial and
844 temporal gene expression data in *Arabidopsis* root stem cells. *Proc. Natl. Acad. Sci. U. S. A.*
845 **114**, E7632–E7640 (2017).
- 846 20. Chen, C., Kim, H. J. & Yang, P. Evaluating spatially variable gene detection methods for
847 spatial transcriptomics data. *bioRxiv* 2022.11.23.517747 (2022)
848 doi:10.1101/2022.11.23.517747.
- 849 21. Jiang, X. *et al.* Spatial Transcriptomics Arena (STAr): an Integrated Platform for Spatial
850 Transcriptomics Methodology Research. *bioRxiv* (2023) doi:10.1101/2023.03.10.532127.
- 851 22. Charitakis, N. *et al.* Disparities in spatially variable gene calling highlight the need for
852 benchmarking spatial transcriptomics methods. *bioRxiv* 2022.10.31.514623 (2023)
853 doi:10.1101/2022.10.31.514623.
- 854 23. Kats, I., Vento-Tormo, R. & Stegle, O. SpatialDE2: Fast and localized variance component
855 analysis of spatial transcriptomics. *bioRxiv* 2021.10.27.466045 (2021)
856 doi:10.1101/2021.10.27.466045.
- 857 24. Palla, G. *et al.* Squidpy: a scalable framework for spatial omics analysis. *Nat. Methods* **19**,

- 858 171–178 (2022).
- 859 25. Moran, P. A. P. Notes on continuous stochastic phenomena. *Biometrika* **37**, 17–23 (1950).
- 860 26. Cai, G., Chen, Y., Gu, X. & Zhou, Z. Spanve: an Effective Statistical Method to Detect
861 Spatially Variable Genes in Large-scale Spatial Transcriptomics Data.
862 <https://europepmc.org/article/ppr/ppr613993>.
- 863 27. Zhang, K., Feng, W. & Wang, P. Identification of spatially variable genes with graph cuts.
864 *Nat. Commun.* **13**, 5488 (2022).
- 865 28. Kipf, T. N. & Welling, M. Semi-Supervised Classification with Graph Convolutional
866 Networks. *arXiv [cs.LG]* (2016).
- 867 29. Hu, J. *et al.* SpaGCN: Integrating gene expression, spatial location and histology to identify
868 spatial domains and spatially variable genes by graph convolutional network. *Nat. Methods*
869 **18**, 1342–1351 (2021).
- 870 30. Chang, Y. *et al.* Spatial omics representation and functional tissue module inference using
871 graph Fourier transform. *bioRxiv* 2022.12.10.519929 (2023)
872 doi:10.1101/2022.12.10.519929.
- 873 31. Andersson, A. & Lundeberg, J. sepal: identifying transcript profiles with spatial patterns by
874 diffusion-based modeling. *Bioinformatics* **37**, 2644–2650 (2021).
- 875 32. Svensson, V., Teichmann, S. A. & Stegle, O. SpatialDE: identification of spatially variable
876 genes. *Nat. Methods* **15**, 343–346 (2018).
- 877 33. Sun, S., Zhu, J. & Zhou, X. Statistical analysis of spatial expression patterns for spatially
878 resolved transcriptomic studies. *Nat. Methods* **17**, 193–200 (2020).
- 879 34. Liu, Y. & Xie, J. Cauchy combination test: a powerful test with analytic p-value calculation

- 880 under arbitrary dependency structures. *arXiv [stat.ME]* (2018).
- 881 35. Li, Q., Zhang, M., Xie, Y. & Xiao, G. Bayesian modeling of spatial molecular profiling data via
882 Gaussian process. *Bioinformatics* **37**, 4129–4136 (2021).
- 883 36. BinTayyash, N. *et al.* Non-parametric modelling of temporal and spatial counts data from
884 RNA-seq experiments. *Bioinformatics* **37**, 3788–3795 (2021).
- 885 37. Zhu, J., Sun, S. & Zhou, X. SPARK-X: non-parametric modeling enables scalable and robust
886 detection of spatial expression patterns for large spatial transcriptomic studies. *Genome*
887 *Biol.* **22**, 184 (2021).
- 888 38. Weber, L. M., Saha, A., Datta, A., Hansen, K. D. & Hicks, S. C. nnSVG for the scalable
889 identification of spatially variable genes using nearest-neighbor Gaussian processes. *Nat.*
890 *Commun.* **14**, 4059 (2023).
- 891 39. Hao, M., Hua, K. & Zhang, X. SOMDE: a scalable method for identifying spatially variable
892 genes with self-organizing map. *Bioinformatics* **37**, 4392–4398 (2021).
- 893 40. Andersson, A. *et al.* Spatial deconvolution of HER2-positive breast cancer delineates
894 tumor-associated cell type interactions. *Nat. Commun.* **12**, 6012 (2021).
- 895 41. Maynard, K. R. *et al.* Transcriptome-scale spatial gene expression in the human
896 dorsolateral prefrontal cortex. *Nat. Neurosci.* **24**, 425–436 (2021).
- 897 42. Song, D. *et al.* scDesign3 generates realistic in silico data for multimodal single-cell and
898 spatial omics. *Nat. Biotechnol.* (2023) doi:10.1038/s41587-023-01772-1.
- 899 43. Heumos, L. *et al.* Best practices for single-cell analysis across modalities. *Nat. Rev. Genet.*
900 1–23 (2023).
- 901 44. Traag, V. A., Waltman, L. & van Eck, N. J. From Louvain to Leiden: guaranteeing well-

- 902 connected communities. *Sci. Rep.* **9**, 5233 (2019).
- 903 45. Deng, Y. *et al.* Spatial profiling of chromatin accessibility in mouse and human tissues.
904 *Nature* **609**, 375–383 (2022).
- 905 46. Llorens-Bobadilla, E. *et al.* Solid-phase capture and profiling of open chromatin by spatial
906 ATAC. *Nat. Biotechnol.* (2023) doi:10.1038/s41587-022-01603-9.
- 907 47. Russell, A. J. C. *et al.* Slide-tags: scalable, single-nucleus barcoding for multi-modal spatial
908 genomics. *bioRxiv* (2023) doi:10.1101/2023.04.01.535228.
- 909 48. Edsgård, D., Johnsson, P. & Sandberg, R. Identification of spatial expression trends in
910 single-cell gene expression data. *Nat. Methods* **15**, 339–342 (2018).
- 911 49. Wang, J. *et al.* Dimension-agnostic and granularity-based spatially variable gene
912 identification. *bioRxiv* (2023) doi:10.1101/2023.03.21.533713.
- 913 50. Zhu, J., Shang, L. & Zhou, X. SRTsim: spatial pattern preserving simulations for spatially
914 resolved transcriptomics. *Genome Biol.* **24**, 39 (2023).
- 915 51. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data
916 analysis. *Genome Biol.* **19**, 15 (2018).
- 917 52. Grau, J., Grosse, I. & Keilwagen, J. PRROC: computing and visualizing precision-recall and
918 receiver operating characteristic curves in R. *Bioinformatics* **31**, 2595–2597 (2015).
- 919 53. Mölder, F. *et al.* Sustainable data analysis with Snakemake. *F1000Res.* **10**, 33 (2021).
- 920 54. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, R137 (2008).
- 921 55. Stuart, T., Srivastava, A., Madad, S., Lareau, C. A. & Satija, R. Single-cell chromatin state
922 analysis with Signac. *Nat. Methods* **18**, 1333–1341 (2021).

923

924 **Data Availability**

925 All the datasets used in this study are publicly available. The human DLPFC data were downloaded
926 from <http://research.libd.org/spatialLIBD>. The breast tumor data were downloaded from
927 <https://github.com/almaan/her2st>. Spatial-ATAC-seq data were obtained from the Gene
928 Expression Omnibus (GEO) with the accession number GSE214991.

929

930 **Code Availability**

931 The code for running the benchmarked methods is available on GitHub:
932 https://github.com/pinelloab/SVG_Benchmarking. The code for generating the simulation
933 datasets is available on Github: <https://github.com/pinelloab/simstpy>.

934

935 **Acknowledgment**

936 The authors would like to thank the members of Pinello Lab for the discussion. L.P. is partially
937 supported by the National Human Genome Research Institute (NHGRI) Genomic Innovator Award
938 (R35HG010717).

939

940 **Author contributions statement**

941 Z.L. and L.P. conceived the study. Z.L. and Z.P. conducted the experiments and analyzed the
942 results. D.S. and G.Y. supported the simulation results using scDesign3. Z.L. and Z.P. wrote the
943 manuscript, revised by L.P. and J.J.L. All authors reviewed the manuscript.