

1 Identifying spatially variable genes by projecting to morphologically 2 relevant curves

3 Phillip B. Nicol^{1,2}, Rong Ma^{1,2}, Rosalind J. Xu^{3,4},
4 Jeffrey R. Moffitt^{3,4,5}, and Rafael A. Irizarry^{1,2,*}

5 ¹*Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA 02115, USA*

6 ²*Department of Data Science, Dana Farber Cancer Institute, Boston, MA 02215, USA*

7 ³*Program in Cellular and Molecular Medicine, Boston Children's Hospital, Boston MA 02115, USA*

8 ⁴*Department of Microbiology, Blavatnik Institute, Harvard Medical School, Boston, MA 02115, USA*

9 ⁵*Broad Institute of Harvard and MIT, Cambridge, MA 02142, USA*

10 * Correspondence: rafael.irizarry@dfci.harvard.edu

11 November 21, 2024

12 **Abstract**

13 Spatial transcriptomics enables high-resolution gene expression measurements while preserving the
14 two-dimensional spatial organization of the sample. A common objective in spatial transcriptomics
15 data analysis is to identify spatially variable genes within predefined cell types or regions within the
16 tissue. However, these regions are often implicitly one-dimensional, making standard two-dimensional
17 coordinate-based methods less effective as they overlook the underlying tissue organization. Here we
18 introduce a methodology grounded in spectral graph theory to elucidate a one-dimensional curve that
19 effectively approximates the spatial coordinates of the examined sample. This curve is then used to
20 establish a new coordinate system that better reflects tissue morphology. We then develop a generalized
21 additive model (GAM) to detect genes with variable expression in the new *morphologically relevant* coord-
22 inate system. Our approach directly models gene counts, thereby eliminating the need for normalization
23 or transformations to satisfy normality assumptions. We demonstrate improved performance relative to
24 current methods based on hypothesis testing, while also accurately estimating gene expression patterns
25 and precisely identifying spatial loci where deviations from constant expression are observed. We validate
our approach through extensive simulations and by analyzing experimental data from multiple platforms

26 such as Slide-seq and MERFISH.

27 1 Introduction

28 Spatial transcriptomics (ST) technologies permit high-resolution measurement of gene expression while main-
29 taining the spatial coordinates of the samples (Rao et al., 2021; Moses and Pachter, 2022). These technologies
30 have the potential to improve our understanding of the influence of cellular spatial organization on important
31 biological processes and disease. One of the starting points for ST analysis is the identification of spatially
32 variable genes (SVGs) (Adhikari et al., 2024). Since spatial variability can often be explained by differences
33 in cell type (Cable et al., 2022b), it is common to test for SVGs within a predefined cell type or spatial
34 domain (Yu and Luo, 2022).

35 Current statistical approaches for SVG detection perform a hypothesis test for each gene, quantifying
36 the evidence of spatial variability using a p -value (Svensson et al., 2018; Sun et al., 2020; Hao et al., 2021;
37 Zhu et al., 2021; Weber et al., 2023). However, these methodologies are incapable of distinguishing genes
38 whose spatial expression patterns manifest in fundamentally different ways, such as along distinct anatomical
39 features within the tissue. For example, MERFISH measurements of a healthy mouse colon revealed two
40 dominant patterns of spatial variability, denoted here as *localized* (Fig 1a, left) and *radial gradient* (Fig 1a,
41 right), respectively. Genes exhibiting localized variation, such as *Ddx58*, are characterized by a distinct patch
42 of expression in one region of the colon, whereas genes exhibiting radial gradient variation, such as *Apob*, are
43 characterized by a gradual change in expression between the outside and inside of the mucosa. Importantly,
44 while both of these examples are illustrations of spatially variable genes, their distinct spatial distributions
45 have important implications for their biological interpretation. *Ddx58* is an interferon-stimulated gene,
46 and the localized expression observed in the mucosa is representative of local patches of interferon activity
47 and interferon-stimulated gene expression described previously (Van Winkle et al., 2022). By contrast, the
48 radial distribution of *Apob*—a marker of the final stages of mature enterocytes—shows the known variation
49 of epithelial cells from the base to the tip of colonic crypts (Moor et al., 2018). Although current leading
50 approaches successfully identified these genes as spatially variable (Fig S1, S2), they lack the capability of
51 distinguishing between these two modes of spatial variation. Additionally, these existing methodologies do
52 not allow for precise pinpointing of the locations where spatially relevant gene expression occurs.

53 Although numerous statistical techniques exist for the estimation of two-dimensional surfaces (Wood,
54 2003; Schulz et al., 2018), and some of these used for ST (Cable et al., 2022a), it is noteworthy that in many
55 applications the primary interest is in genes that exhibit variation along one-dimensional paths. For example,
56 the exercise of visually detecting the localized pattern shown in *Ddx58* could be described as searching for

57 deviation from a baseline expression level along a curve tracing through the mucosa. The radial gradient
58 pattern exhibited by *Apob* could be described as change in expression in the direction perpendicular to the
59 curve. This implies that a curve-based coordinate transform could help separate genes with a localized burst
60 from those with radial gradient, which, in turn, could facilitate new biological findings. In addition to the
61 colon (**Fig 1a**), cell types in the brain also commonly exhibit distinct one-dimensional spatial structure. In
62 this paper we also consider two Slide-seq datasets: granule cells from the mouse cerebellum ([Cable et al.,](#)
63 [2022b](#)) (**Fig 1b**) and CA3 cells from the mouse hippocampus ([Stickels et al., 2021](#)) (**Fig 1c**), although
64 numerous additional examples exist.

65 In this paper, we introduce a statistical framework that estimates a one-dimensional curve passing through
66 the ST coordinates and then uses that estimated curve to define a *morphologically relevant* coordinate system.
67 Although similar curve-estimation methods have been used for pseudotime analysis in single-cell RNA-seq
68 ([Street et al., 2018](#)), we find that our methodology grounded in spectral graph theory yields better results
69 on two-dimensional ST data. Moreover, pseudotime methods do not measure variation *orthogonal* to the
70 curve which is critical to distinguish between localized and radial gradient patterns.

71 Upon estimating the curve, we employ a generalized additive model (GAM) to model expression as
72 a (possibly non-linear) function of the morphologically relevant coordinates. We refer to our approach
73 as *MorphoGAM*. Unlike previously published hypothesis tests for SVGs, MorphoGAM identifies the exact
74 location and mode of the pertinent expression pattern, thereby resulting in more interpretable findings. An
75 additional advantage of summarizing the two-dimensional coordinates using one-dimensional projections is
76 increased statistical power to detect relevant SVGs due to reduced complexity of model fit to estimate spatial
77 effects.

78 2 Results

79 **MorphoGAM identifies morphologically relevant coordinates in spatial transcriptomics data.**

80 We begin by modeling the spatial location of cell j , $1 \leq j \leq n$, using *morphologically relevant* coordinates t_j
81 and r_j . Specifically, we assume that the standard two-dimensional spatial coordinates $x_j \in \mathbb{R}^2$ lie close to a
82 latent curve:

$$\mathbb{E}(x_j) = f(t_j) \tag{2.1}$$

83 where $f : [a, b] \rightarrow \mathbb{R}^2$ is a smooth one-dimensional parametric curve. We write $f(t) = (f_1(t), f_2(t))$ to denote
84 the two component functions of the curve. The first morphologically relevant coordinate t_j describes the
85 position of cell j along the curve. In the [Methods](#) we describe in detail our approach based on spectral

86 graph theory to estimate t_j and f . Briefly, our approach relates the distance between coordinates $|t_i - t_j|$
87 to the shortest path in a k -nearest neighbor graph G_k and then shows that t_j can be estimated through an
88 eigendecomposition of a centered shortest path matrix. Upon obtaining an estimate \hat{t}_j , the curve f can be
89 estimated by smoothing each dimension separately. We plug in \hat{t}_j to (2.1) to obtain

$$\mathbb{E}(x_{j1}) = f_1(\hat{t}_j) \quad (2.2)$$

$$\mathbb{E}(x_{j2}) = f_2(\hat{t}_j) \quad (2.3)$$

90 We thus obtain \hat{f}_1 and \hat{f}_2 by using regression splines as implemented in *mgcv* (Wood, 2017).

91 Our methodology to estimate t_j in the case of a linear curve (i.e., when $f(a) \neq f(b)$) is motivated
92 by the ISOMAP (Tenenbaum et al., 2000) technique for non-linear dimensionality reduction. We extend
93 this approach to allow our method to address scenarios wherein f constitutes a closed curve (i.e., when
94 $f(a) = f(b)$).

95 A detailed visual assessment indicates that, when applied to the slice of healthy mouse colon, Mor-
96 phoGAM excels in estimating $f(t)$ (Fig 2a) and the morphologically relevant coordinate t_j (Fig 2b).

97 The second *morphologically relevant* coordinate, denoted here by r_j , is defined by using the distance from
98 the cell's coordinates to the position on the curve $f(t)$. Explicitly, the magnitude of r_j is given by

$$|r_j| = \|x_j - \hat{f}(\hat{t}_j)\|_2. \quad (2.4)$$

99 To determine the sign of r_j , we set

$$\text{sign}(\hat{r}_j) = \text{sign} \left[\langle x_j - \hat{f}(\hat{t}_j), R\hat{f}'(\hat{t}_j) \rangle \right] \quad (2.5)$$

100 where $R : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ is a counter-clockwise 90 degrees rotation: $R(v_1, v_2) = (-v_2, v_1)$. The conceptual
101 framework behind equation (2.5) can be understood by envisioning a traversal along the curve, where the
102 velocity vector at time t is $f'(t)$. Residuals exhibiting a positive sign would be placed on the left-hand side
103 of the curve as one progresses, whereas those with a negative sign would be on the right-hand side. The
104 left-hand side can be ascertained through a counter-clockwise rotation R of the velocity vector $f'(t)$. This
105 coordinate for each cell is also morphologically pertinent, as illustrated in Figure (Fig 2c).

106 After transforming to the morphologically relevant coordinate system, the difference between localized
107 and radial gradient patterns becomes immediately clear. SVGs with localized patterns show variation in
108 the first morphologically relevant coordinate t_j (Fig 2d) whereas SVGs exhibiting a radial gradient pattern

109 show variation in the second morphologically relevant coordinate r_j (**Fig 2e**).

110 **MorphoGAM outperforms existing curve estimation approaches.** Hastie and Stuetzle (1989) in-
111 troduced model (2.1) as a general approach to estimate a curve passing through a set of points (in arbitrary
112 dimension). This method, known as *Principal Curves*, is used by the popular pseudotime method *Sling-*
113 *shot* (Street et al., 2018). Hastie and Stuetzle (1989) employ an iterative algorithm that alternates between
114 updating f and updating t_j . However, we find that this iterative approach is unsuitable for the highly non-
115 linear structures observed in spatial transcriptomics data. To demonstrate this, we applied *Principal Curves*
116 to granule cells from the mouse cerebellum, measured using Slide-seqV2 (Cable et al., 2022b) (see Figure
117 1b). We found that this approach did not accurately estimate the curve for a variety of tuning parameter
118 choices (**Fig 3a-c**). In contrast, if we set the number of nearest neighbors k to 5, MorphoGAM accurately
119 identified the path and the first morphologically relevant coordinate (**Fig 3d, Fig S3**). The performance is
120 similar for $k = 10$ (**Fig 3e**) and only begins to degrade once $k = 30$ (**Fig 3f**).

121 To quantitatively evaluate the robustness of these methods with respect to their tuning parameters, we
122 defined *ground truth* t_j by carefully hand-drawing a path (**Fig S4**) and compared estimates \hat{t}_j obtained
123 with different values of the tuning parameter to the ground truth t_j . We observed that the spearman
124 correlation between the estimates $\hat{t}_j(\text{df})$ and t_j was below 0.6 for Principal Curves (**Fig 3g**). In contrast,
125 with MorphoGAM, the correlation exceeded 0.95 for $k < 10$ and remained above 0.9 for values up to $k = 30$
126 (**Fig 3h**). We also applied both methods to simulated swiss roll data (**Fig S5**) and the mouse colon
127 dataset (**Fig S6**) and found similar performance gains from using MorphoGAM. An additional advantage
128 of MorphoGAM is that the tuning parameter k (number of nearest neighbors) is more interpretable than
129 degrees of freedom (df) for a smoothing spline, which makes it easier to find a reasonable value in practice.

130 **MorphoGAM allows for interpretable detection of spatially variable genes.** Following the estima-
131 tion of morphologically relevant coordinates \hat{t}_j and \hat{r}_j , MorphoGAM identifies spatially variable genes using
132 a generalized additive model (GAM) (Hastie and Tibshirani, 1986). Specifically, we denote the
133 count for gene g in cell j with Y_{gj} and model it with

$$Y_{gj} \sim \text{NegBinom}(n_j \mu_{gj}, \theta_g) \tag{2.6}$$
$$\log \mu_{gj} = \beta_{g0} + h_g(\hat{t}_j) + s_g(\hat{r}_j)$$

134 where h_g and s_g are unknown smooth functions, β_{g0} is an unknown intercept, and n_j is the total counts
135 for cell j . θ_g is the inverse dispersion parameter because $\text{Var}(Y_{gj}) = n_j \mu_{gj} + (n_j \mu_{gj})^2 / \theta_g$. In this model,
136 the gene g is spatially variable if $h_g \neq 0$ or $s_g \neq 0$. Estimating the parameters of model (2.6) is achieved

137 by writing h_g and s_g as the sum of basis functions and then adding a penalty to encourage smoothness
138 (Methods). Although p -values can be computed by testing the null hypothesis $h_g = 0$ or $s_g = 0$, we
139 recommend inspecting the estimated functions \hat{h}_g and \hat{s}_g along with estimates of their covariance (to measure
140 uncertainty). We specifically use *adaptive shrinkage* (Stephens, 2017) to further regularize the functions with
141 higher uncertainty (see Methods).

142 Although we can examine the entire function estimate, we also introduce two summaries useful for ranking
143 genes automatically. Specifically, we consider the *peak*

$$\hat{P}_g := \sup_t \hat{h}_g(t) \quad (2.7)$$

144 which estimates the maximum log-fold change from the baseline log expression $\hat{\beta}_{g0}$. Because this measure-
145 ment could prioritize large multiplicative changes in small genes we also define the *range*

$$\hat{R}_g := \sup_t \left[n_{\text{med}} \exp(\hat{\beta}_{g0} + \hat{h}_g(t)) \right] - \inf_t \left[n_{\text{med}} \exp(\hat{\beta}_{g0} + \hat{h}_g(t)) \right] \quad (2.8)$$

146 to account for genes that have large differences on the original scale of the counts. Here n_{med} is defined
147 as the median of the n_j , so that \hat{R}_g can be directly interpreted as a count difference. We note $\hat{s}_g(t)$ could
148 replace $\hat{h}_g(t)$ in both equations (2.7) and (2.8).

149 We emphasize that the model (2.6) can easily be modified depending on the particular scientific question.
150 For example, if only variation along the curve is of interest then we only need to examine $\hat{h}_g(t)$. Moreover,
151 the model is flexible enough to account for other potential confounders in the linear predictor.

152 **MorphoGAM improves power to detect relevant spatially variable genes.** We applied Mor-
153 phoGAM to the CA3 mouse hippocampus cells (Fig 1c) to estimate a one-dimensional curve \hat{f} and mor-
154 phologically relevant coordinates \hat{t}_j (Fig 4a). In the original analysis of this dataset, Cable et al. (2022b)
155 used 2D locally weighted regression to identify genes with a high coefficient of variation (CV). This analysis
156 identified two genes *Rgs14* and *Cpne9* that exhibited variable expression at different ends. We applied the
157 GAM model (2.6) with s_g removed to identify genes varying along the curve \hat{f} . Our approach corroborated
158 the finding of *Rgs14* ($\hat{P}_g = 2.77$, $p < 10^{-16}$) and *Cpne9* ($\hat{P}_g = 1.41$, $p < 10^{-16}$) (Fig 4b).

159 We hypothesized that MorphoGAM increases statistical power to detect SVGs by projecting the two-
160 dimensional ST coordinates to a one-dimensional morphologically relevant coordinate. To demonstrate this,
161 we simulated a gene such that $\mu_{gj} = 1 + \kappa \exp(-\sigma(\hat{t}_j - 0.5)^2)$, where \hat{t}_j is as above. In order to compare to
162 the approaches based on hypothesis testing, we labeled a gene as spatially variable if the p -value was below
163 the transcriptome wide significance level of 0.05/20000. MorphoGAM had consistently higher power than

164 two state-of-the-art methods for detecting SVGs, nnSVG (Weber et al., 2023) and SPARK-X (Zhu et al.,
165 2021) (Fig 4c). To show that the increase in power did not come at the price of an inflated type I error rate,
166 we randomly permuted all spatial locations (to generate a null dataset with no spatially variable genes) and
167 found that our method was conservative at a variety of significance levels (Fig S7). When ranking by genes
168 with a large peak or range, our approach identified genes that were not reported in the original analysis of
169 Cable et al. (2022b) such as *Fxyd6* ($\hat{P}_g = 3.28$, $p < 10^{-16}$) and *Hpca* ($\hat{R}_g = 12.67$, $p < 10^{-16}$) (Fig 4d).

170 **MorphoGAM identifies spatially variable genes in the mouse colon data.** We applied MorphoGAM
171 to the MERFISH measurements of a slice of healthy mouse colon (Fig 1) with the goal of separating genes
172 with localized and radial gradient patterns of expression. Because localized genes are characterized by a
173 burst in expression along the curve, we used the peak of estimated functions $\hat{h}_g(\hat{t}_j)$ to rank genes (Fig
174 5a). Radial gradient genes, on the other hand, are characterized by a smooth transition along the second
175 morphologically relevant coordinate, so for this we found genes with a large range in \hat{s}_g (Fig 5b). Figure
176 5 also lists the ranking of each gene of both nnSVG and SPARK-X, showing that the targeted analysis
177 of MorphoGAM prioritizes genes that could have been missed if only hypothesis-based tests were used for
178 SVGs. In particular, *Ddx58* was found to have a large peak in the first morphologically relevant coordinate
179 ($\hat{P}_g = 2.01$, $p < 10^{-16}$) and *Apob* was found to have a large peak in the direction of the second morphologically
180 relevant coordinate ($\hat{P}_g = 1.24$, $p < 10^{-16}$). We also plot the genes with the largest range in the direction
181 of the first morphologically relevant coordinate and the genes with the largest peak in the direction of the
182 second morphologically relevant coordinate in Figure S8.

183 3 Discussion

184 We introduced an approach to estimate the curve passing through spatial transcriptomics coordinates and
185 leveraged this curve to define *morphologically relevant* coordinates. A GAM is used to model spatial variation
186 along these morphologically relevant coordinates, which we have shown to be an interpretable and powerful
187 approach to find relevant spatially variable genes. Importantly, we have advocated to directly use summaries
188 of the estimated functions rather than relying on a null hypothesis test, as p -values do not provide information
189 about the mode of spatial variation and are in general sensitive to misspecification in the assumed model
190 (Greenland et al., 2016).

191 The proposed methodology presents certain limitations. First, the final results depend on the accurate
192 annotation of cell types or spatial domains, and inaccuracies at this stage may propagate to MorphoGAM.
193 Furthermore, the approach is not inherently applicable in scenarios where the tissue structure cannot be
194 adequately represented by a one-dimensional framework. As a result, as part of the software package sup-

195 porting the implementation of MorphoGAM, we have developed a tool allowing manual curve drawing $f(t)$,
196 as shown in **Figure S4**. This facilitates the application of MorphoGAM, although manually, in instances
197 where variation along a predetermined trajectory is to be identified. More broadly, our GAM methodology,
198 accompanied by summaries of estimated functions (such as range and peak), can be easily extended to the
199 two-dimensional domain by employing thin-plate splines ([Wood, 2003](#)).

200 Morphologically relevant coordinates may offer considerable utility beyond the scope of spatially vari-
201 able genes. For instance, the alignment of multiple spatial transcriptomics (ST) slices may be enhanced
202 by leveraging these morphologically relevant coordinates instead of conventional two-dimensional coordi-
203 nates. Future research could profitably explore the application of morphologically pertinent coordinates in
204 conducting multi-sample ST analyses.

205 References

- 206 S. D. Adhikari, J. Yang, J. Wang, and Y. Cui. Recent advances in spatially variable gene detection in spatial
207 transcriptomics. *Computational and Structural Biotechnology Journal*, 2024.
- 208 D. M. Cable, E. Murray, V. Shanmugam, S. Zhang, L. S. Zou, M. Diao, H. Chen, E. Z. Macosko, R. A.
209 Irizarry, and F. Chen. Cell type-specific inference of differential expression in spatial transcriptomics.
210 *Nature methods*, 19(9):1076–1087, 2022a.
- 211 D. M. Cable, E. Murray, L. S. Zou, A. Goeva, E. Z. Macosko, F. Chen, and R. A. Irizarry. Robust decom-
212 position of cell type mixtures in spatial transcriptomics. *Nature biotechnology*, 40(4):517–526, 2022b.
- 213 P. Cadinu, K. N. Sivanathan, A. Misra, R. J. Xu, D. Mangani, E. Yang, J. M. Rone, K. Tooley, Y.-C. Kye,
214 L. Bod, et al. Charting the cellular biogeography in colitis reveals fibroblast trajectories and coordinated
215 spatial remodeling. *Cell*, 187(8):2010–2028, 2024.
- 216 R. Cannoodt. princurve 2.0: Fit a principal curve in arbitrary dimension, Jun 2018. URL [https://doi.](https://doi.org/10.5281/zenodo.3351282)
217 [org/10.5281/zenodo.3351282](https://doi.org/10.5281/zenodo.3351282).
- 218 C. Eckart and G. Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):
219 211–218, 1936.
- 220 B. Ghojogh, M. Crowley, F. Karray, and A. Ghodsi. Multidimensional scaling, sammon mapping, and
221 isomap. In *Elements of Dimensionality Reduction and Manifold Learning*, pages 185–205. Springer, 2023.

- 222 S. Greenland, S. J. Senn, K. J. Rothman, J. B. Carlin, C. Poole, S. N. Goodman, and D. G. Altman.
223 Statistical tests, p values, confidence intervals, and power: a guide to misinterpretations. *European journal*
224 *of epidemiology*, 31(4):337–350, 2016.
- 225 M. Hao, K. Hua, and X. Zhang. Somde: a scalable method for identifying spatially variable genes with
226 self-organizing map. *Bioinformatics*, 37(23):4392–4398, 2021.
- 227 T. Hastie and W. Stuetzle. Principal curves. *Journal of the American statistical association*, 84(406):502–516,
228 1989.
- 229 T. Hastie and R. Tibshirani. Generalized additive models. *Statistical Science*, 1(3):297–318, 1986.
- 230 B. Kégl, A. Krzyzak, T. Linder, and K. Zeger. Learning and design of principal curves. *IEEE transactions*
231 *on pattern analysis and machine intelligence*, 22(3):281–297, 2000.
- 232 A. Lastra. *Parametric Geometry of Curves and Surfaces*. Springer, 2021.
- 233 A. E. Moor, Y. Harnik, S. Ben-Moshe, E. E. Massasa, M. Rozenberg, R. Eilam, K. B. Halpern, and
234 S. Itzkovitz. Spatial reconstruction of single enterocytes uncovers broad zonation along the intestinal
235 villus axis. *Cell*, 175(4):1156–1167, 2018.
- 236 L. Moses and L. Pachter. Museum of spatial transcriptomics. *Nature methods*, 19(5):534–546, 2022.
- 237 A. Rao, D. Barkley, G. S. França, and I. Yanai. Exploring tissue architecture using spatial transcriptomics.
238 *Nature*, 596(7871):211–220, 2021.
- 239 D. Righelli, L. M. Weber, H. L. Crowell, B. Pardo, L. Collado-Torres, S. Ghazanfar, A. T. Lun, S. C. Hicks,
240 and D. Risso. Spatialexperiment: infrastructure for spatially-resolved transcriptomics data in r using
241 bioconductor. *Bioinformatics*, 38(11):3128–3131, 2022.
- 242 E. Schulz, M. Speekenbrink, and A. Krause. A tutorial on gaussian process regression: Modelling, exploring,
243 and exploiting functions. *Journal of mathematical psychology*, 85:1–16, 2018.
- 244 M. Stephens. False discovery rates: a new deal. *Biostatistics*, 18(2):275–294, 2017.
- 245 R. R. Stickels, E. Murray, P. Kumar, J. Li, J. L. Marshall, D. J. Di Bella, P. Arlotta, E. Z. Macosko,
246 and F. Chen. Highly sensitive spatial transcriptomics at near-cellular resolution with slide-seq2. *Nature*
247 *biotechnology*, 39(3):313–319, 2021.
- 248 K. Street, D. Risso, R. B. Fletcher, D. Das, J. Ngai, N. Yosef, E. Purdom, and S. Dudoit. Slingshot: cell
249 lineage and pseudotime inference for single-cell transcriptomics. *BMC genomics*, 19:1–16, 2018.

- 250 S. Sun, J. Zhu, and X. Zhou. Statistical analysis of spatial expression patterns for spatially resolved tran-
251 scriptomic studies. *Nature methods*, 17(2):193–200, 2020.
- 252 V. Svensson, S. A. Teichmann, and O. Stegle. Spatialde: identification of spatially variable genes. *Nature*
253 *methods*, 15(5):343–346, 2018.
- 254 J. B. Tenenbaum, V. d. Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality
255 reduction. *science*, 290(5500):2319–2323, 2000.
- 256 W. S. Torgerson. Multidimensional scaling: I. theory and method. *Psychometrika*, 17(4):401–419, 1952.
- 257 J. A. Van Winkle, S. T. Peterson, E. A. Kennedy, M. J. Wheadon, H. Ingle, C. Desai, R. Rodgers, D. A.
258 Constant, A. P. Wright, L. Li, et al. Homeostatic interferon-lambda response to bacterial microbiota
259 stimulates preemptive antiviral defense within discrete pockets of intestinal epithelium. *elife*, 11:e74072,
260 2022.
- 261 L. M. Weber, A. Saha, A. Datta, K. D. Hansen, and S. C. Hicks. nnsvg for the scalable identification of
262 spatially variable genes using nearest-neighbor gaussian processes. *Nature communications*, 14(1):4059,
263 2023.
- 264 S. N. Wood. mgcv: Gams and generalized ridge regression for r. *R news*, 1(2):20–25, 2001.
- 265 S. N. Wood. Thin plate regression splines. *Journal of the Royal Statistical Society Series B: Statistical*
266 *Methodology*, 65(1):95–114, 2003.
- 267 S. N. Wood. *Generalized additive models: an introduction with R*. chapman and hall/CRC, 2017.
- 268 J. Yu and X. Luo. Identification of cell-type-specific spatially variable genes accounting for excess zeros.
269 *Bioinformatics*, 38(17):4135–4144, 2022.
- 270 J. Zhu, S. Sun, and X. Zhou. Spark-x: non-parametric modeling enables scalable and robust detection of
271 spatial expression patterns for large spatial transcriptomic studies. *Genome biology*, 22(1):184, 2021.

272 4 Methods

273 **Statistical model for latent curve.** Let $x_j \in \mathbb{R}^2$ denote the spatial coordinates of cell $1 \leq j \leq n$. We
274 assume that

$$\mathbb{E}(x_j) = f(t_j) \tag{4.1}$$

275 where $f : [a, b] \rightarrow \mathbb{R}^2$ is a smooth parametric curve. We will write $f(t) = (f_1(t), f_2(t))$ to denote the two
 276 component functions of the curve. For the moment, we assume that f does not intersect itself, so that $t_i \neq t_j$
 277 implies $f(t_i) \neq f(t_j)$. Note that both f and t_j are unknown and must be estimated. See Section S1.2 for
 278 detailed discussions on the identifiability conditions for f and t_j .

279 The arclength of f can be expressed in terms of the first derivative $f'(t) := (f'_1(t), f'_2(t))$:

$$\int_a^b \|f'(t)\|_2 dt \quad (4.2)$$

280 We will assume that f has a *unit-speed parametrization*, which means that $\|f'(t)\|_2 = 1$ for all t . Any
 281 parametric curve such that $f'(t) \neq 0$ for all t can be reparametrized to satisfy this requirement (Section
 282 S1.2). In particular, this means the arc-length between two points on the curve is equal to the difference in
 283 their coordinates:

$$t_i - t_j = \int_{t_i}^{t_j} \|f'(t)\|_2 dt \quad (t_j < t_i) \quad (4.3)$$

284 **First morphologically relevant coordinate.** Our approach to estimate t_j in the case of a non-intersecting
 285 curve (i.e., $x_1 \neq x_2 \Rightarrow f(x_1) \neq f(x_2)$) leverages the relationship in (4.3). Interestingly, it turns out that the
 286 estimated \hat{t}_j is the same as the first component produced by the ISOMAP algorithm for manifold learning
 287 (Tenenbaum et al., 2000). Based on equation (4.3), we can estimate the arclength between t_i and t_j using
 288 shortest paths in a k -nearest neighbor (KNN) graph G_k . If k is chosen sufficiently small, the shortest path
 289 between two vertices (cells) will have a similar shape as f . To make this precise, we define $d_{G_k}(i, j)$ to be
 290 the shortest path between x_i and x_j in G_k , so that

$$d_{G_k}(i, j) \approx |t_i - t_j| \quad (4.4)$$

291 Our estimate \hat{t}_j will be chosen to satisfy the approximation in (4.4). To construct this, we follow the
 292 steps of classical multidimensional scaling (cMDS) (Torgerson, 1952). Squaring both sides of (4.4) yields

$$d_{G_k}^2(i, j) \approx t_i^2 + t_j^2 - 2t_i t_j \quad (4.5)$$

293 Viewing $d_{G_k}^2$ as an $n \times n$ matrix, the operation of double centering (Lemma S1.1) yields $b_{G_k} \in \mathbb{R}^{n \times n}$ such
 294 that

$$-\frac{1}{2} b_{G_k}(i, j) \approx t_i t_j \quad (4.6)$$

295 Given this, we set

$$\hat{t} = \operatorname{argmin}_{t \in \mathbb{R}^n} \left\| -\frac{1}{2}b_{G_k} - t \otimes t \right\|_F^2 \quad (4.7)$$

296 where \otimes is the outer product defined in Section S1.1. The optimization problem in (4.7) has (under mild
297 conditions) a closed form solution given by the leading eigenvector of $-\frac{1}{2}b_{G_k}$ (scaled by the square root of
298 the leading eigenvalue), see Lemma S1.2. In practice, we standardize the resulting \hat{t} so that it takes values
299 between 0 and 1.

300 Once \hat{t}_j is obtained, the curve f can be estimated by smoothing each component function separately. We
301 plug in \hat{t}_j to (2.1) to obtain

$$\mathbb{E}(x_{j1}) = f_1(\hat{t}_j) \quad (4.8)$$

$$\mathbb{E}(x_{j2}) = f_2(\hat{t}_j) \quad (4.9)$$

302 We obtain \hat{f}_1 and \hat{f}_2 by using regression splines as implemented in *mgcv* (Wood, 2017).

303 **Extending the method to closed curves.** For closed curves, we have $f(a) = f(b)$, violating the non-
304 intersecting condition required above. In this case, the approximation in (4.3) no longer holds because there
305 could be a shorter path passing over the endpoint. However, we can still obtain an explicit form for the
306 shortest path between t_i and t_j by applying the law of cosines:

$$\begin{aligned} \theta_i &:= 2\pi \left(\frac{t_i - b}{b - a} \right) \\ d_{G_k}^2(i, j) &\approx (b - a) \arccos \left(1 - \frac{(\cos(\theta_i) - \cos(\theta_j))^2 - (\sin(\theta_i) - \sin(\theta_j))^2}{2} \right). \end{aligned} \quad (4.10)$$

307 As $d_{G_k}^2$ appears to have no simple form, we make a second-order Taylor approximation:

$$d_{G_k}^2(i, j) \approx c [(\cos(\theta_i) - \cos(\theta_j))^2 - (\sin(\theta_i) - \sin(\theta_j))^2] \quad (4.11)$$

308 where c is some constant (the entire derivation is in Section S1.4). Applying the double centering operation
309 to (4.11) yields

$$-\frac{1}{2}b_{G_k}(i, j) \approx c \cos(\theta_i) \cos(\theta_j) + c \sin(\theta_i) \sin(\theta_j). \quad (4.12)$$

310 Because $\cos(\theta), \sin(\theta) \in \mathbb{R}^n$ are expected to be approximately orthogonal, a reasonable approximation θ_j
311 is given by

$$\hat{\theta}_j = \operatorname{atan2}(v_2(-b_{G_k}/2)_j, v_1(-b_{G_k}/2)_j) \quad (4.13)$$

312 where $v_k(\cdot)$ denotes the k -th leading eigenvector of a matrix and atan2 is the 2-argument arctangent function.
 313 Thus $\hat{\theta}_j$ can be converted back to \hat{t}_j via equation (4.13), although our downstream analysis of SVG detection
 314 will be invariant to this scaling.

315 **Second morphologically relevant coordinate.** The second morphologically relevant coordinate $\hat{r}_j \in \mathbb{R}$
 316 describes how far from the estimated curve a cell's coordinates are. The magnitude of the coordinate is
 317 defined as

$$\|\hat{r}_j\|_2 := \|x_j - \hat{f}(\hat{t}_j)\|_2 \quad (4.14)$$

318 The sign of the second coordinate is determined by

$$\text{sign}(\hat{r}_j) = \text{sign}\langle x_j - \hat{f}(\hat{t}_j), R\hat{f}'(\hat{t}_j) \rangle \quad (4.15)$$

319 where $R : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ is a counter-clockwise rotation by 90 degrees: $R(v_1, v_2) = (-v_2, v_1)$.

320 The intuition behind this equation is that cells/spots with a positive sign would be on the left-hand side
 321 if one was driving along the curve. The left-hand side is identified by a counter-clockwise rotation R of the
 322 velocity vector $f'(t)$. Again, in practice we standardize \hat{r}_j to be in the interval $[0, 1]$ although this could be
 323 modified depending on the specific scientific question.

324 **Disconnected graphs.** The estimation procedure described above requires G_k to be connected. However,
 325 if k is chosen large enough to ensure the graph is fully connected, then d_{G_k} may not capture more subtle
 326 morphological features. For this reason, we permit the procedure to be applied separately to disconnected
 327 components of G_k and then *stitched* together to create the final curve. Given G_k has C connected com-
 328 ponents, let $x_1^{(c)}$ and $x_2^{(c)}$ denote endpoints of the curve describing the c -th component. We then identify
 329 the connections between $x_i^{(c)}$ and $x_{i'}^{(c')}$, $i, i' \in \{1, 2\}$, $c \neq c'$ of minimum Euclidean distance that produce a
 330 single (connected) curve. Note that this may require reversing the direction of a curve fit to one particular
 331 component. We identify the optimal connections through a brute force search of the $C! \cdot 2^C$ possibilities.
 332 Because this is computationally infeasible for large C , we require that k is at least large enough to ensure
 333 that $C \leq 5$.

334 **Generalized additive model to identify spatially variable genes.** Let Y_{gj} denote the count for gene
 335 g ($1 \leq g \leq G$) in cell/spot j ($1 \leq j \leq n$). As noted before, we consider the model

$$\begin{aligned} Y_{gj} &\sim \text{NegBinom}(n_j \mu_{gj}, \theta_g) \\ \log \mu_{gj} &= \beta_{g0} + h_g(\hat{t}_j) + s_g(\hat{r}_j) \end{aligned} \quad (4.16)$$

336 where h_g and s_g are unknown smooth functions, β_{g0} is an unknown intercept, and $n_j := \sum_{g=1}^G Y_{gj}$ is a known
 337 offset. Note that we are using the following standard parameterization of the negative binomial distribution:

$$\mathbb{P}(Y_{gj} = y) = \frac{\Gamma(y + \theta)}{\Gamma(\theta)y!} \left(\frac{\theta}{\theta + n_j \mu_{gj}} \right)^\theta \left(\frac{n_j \mu_{gj}}{\theta + n_j \mu_{gj}} \right)^y \quad y = 0, 1, 2, \dots \quad (4.17)$$

338 For identifiability, we also assume that $\sum_{j=1}^n h_g(\hat{t}_j) = \sum_{j=1}^n s_g(\hat{r}_j) = 0$. We use *mgcv* (Wood, 2017) to
 339 estimate the functions in model (2.6); this method writes h_g and s_g as a linear combination of known basis
 340 functions

$$\log \mu_{gj} = \beta_{g0} + \sum_{\ell=1}^{L_t} \beta_{g\ell}^{(t)} \phi(\hat{t}_j) + \sum_{\ell=1}^{L_r} \beta_{g\ell}^{(r)} \psi(\hat{r}_j) \quad (4.18)$$

341 Although there is flexibility in the choice of ϕ and ψ , we use cubic regression splines (cyclic for ϕ when f is a
 342 closed curve). Estimation proceeds by maximizing the log-likelihood of the model parameters $\ell(\beta_{g0}, \beta_g^{(t)}, \beta_g^{(r)})$
 343 (here $\beta_g^{(t)} \in \mathbb{R}^{L_t}$ and $\beta_g^{(r)} \in \mathbb{R}^{L_r}$ are vectors of coefficients) subject to a smoothness penalty:

$$\operatorname{argmin}_{\beta} \left\{ -\ell(\beta_{g0}, \beta_g^{(t)}, \beta_g^{(r)}) + \lambda_t (\beta_g^{(t)})^\top S_t \beta_g^{(t)} + \lambda_r (\beta_g^{(r)})^\top S_r \beta_g^{(r)} + \lambda \left((\beta_g^{(t)})^\top \beta_g^{(t)} + (\beta_g^{(r)})^\top \beta_g^{(r)} \right) \right\} \quad (4.19)$$

344 where S_t and S_r are (known) matrices depending on the second derivative of the chosen basis functions
 345 (Wood, 2001). *mgcv* performs a procedure to select the best choice of λ_t and λ_r and we set $\lambda = 1$ by default.
 346 Upon estimating the coefficients, *mgcv* returns a Bayesian covariance matrix for uncertainty quantification.
 347 We use this to obtain the posterior standard deviation $\operatorname{sd}(\hat{\beta}_{g\ell}^{(c)})$ for each coefficient. For further shrinkage
 348 towards 0, we apply adaptive shrinkage (ash) (Stephens, 2017) to the estimated coefficients and their standard
 349 deviations to obtain the final estimate of \hat{h} and \hat{s} .

350 **Data and code availability.** The following datasets were used:

- 351 • The granule cells in **Figure 3** were obtained from the data provided by Cable et al. (2022b). Cells
 352 such that the 5-th nearest neighbor was 2 times greater than the median 5-th nearest neighbor were
 353 excluded. This procedure removed outlier cells.
- 354 • The CA3 cells were obtained from *STexampleData* (Righelli et al., 2022). Cells such that the 20-th
 355 nearest neighbor was 3 times greater than the median 20-th nearest neighbor were excluded. These
 356 values were used so that the retained set of cells visually matched Figure 5 of Cable et al. (2022b).
- 357 • MERFISH measurements of the adult healthy colon are available upon request from the authors.
 358 Briefly, these measurements were performed using standard MERFISH protocols (Cadimu et al., 2024)
 359 targeting a custom set of 1,920 genes.

360 MorphoGAM is available as an R package at <https://github.com/phillipnicol/MorphoGAM>. The reposi-
361 tory also includes scripts to reproduce all results in the paper.

362 **Acknowledgements**

363 PBN is supported by the National Institutes of Health grant T32CA009337.

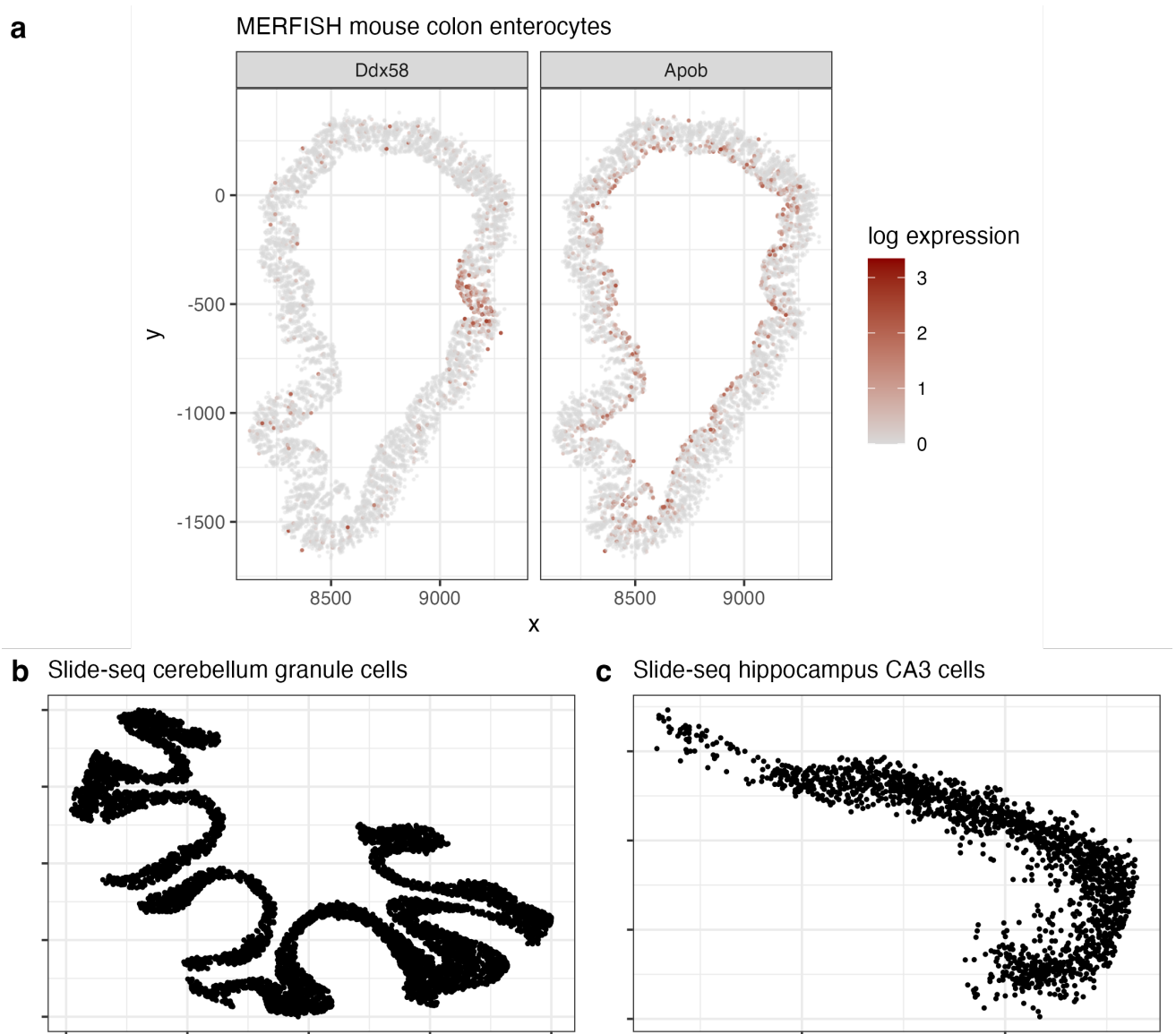


Figure 1: **a**. The spatial location of enterocytes identified in a cross section of the healthy mouse colon as measured with MERFISH. Cells are colored by the log-transformed expression of the two listed genes. The expression of *Ddx58* is called localized whereas the expression pattern of *Apob* is called radial gradient. The spatial locations of the plotted enterocytes lie close to a one-dimensional circular manifold. Additional examples of cell types with coordinates that lie close to a one-dimensional curve can be found in **b**. granule cells from the mouse cerebellum (Cable et al., 2022b) and **c**. CA3 cells in the mouse hippocampus (Stickels et al., 2021)

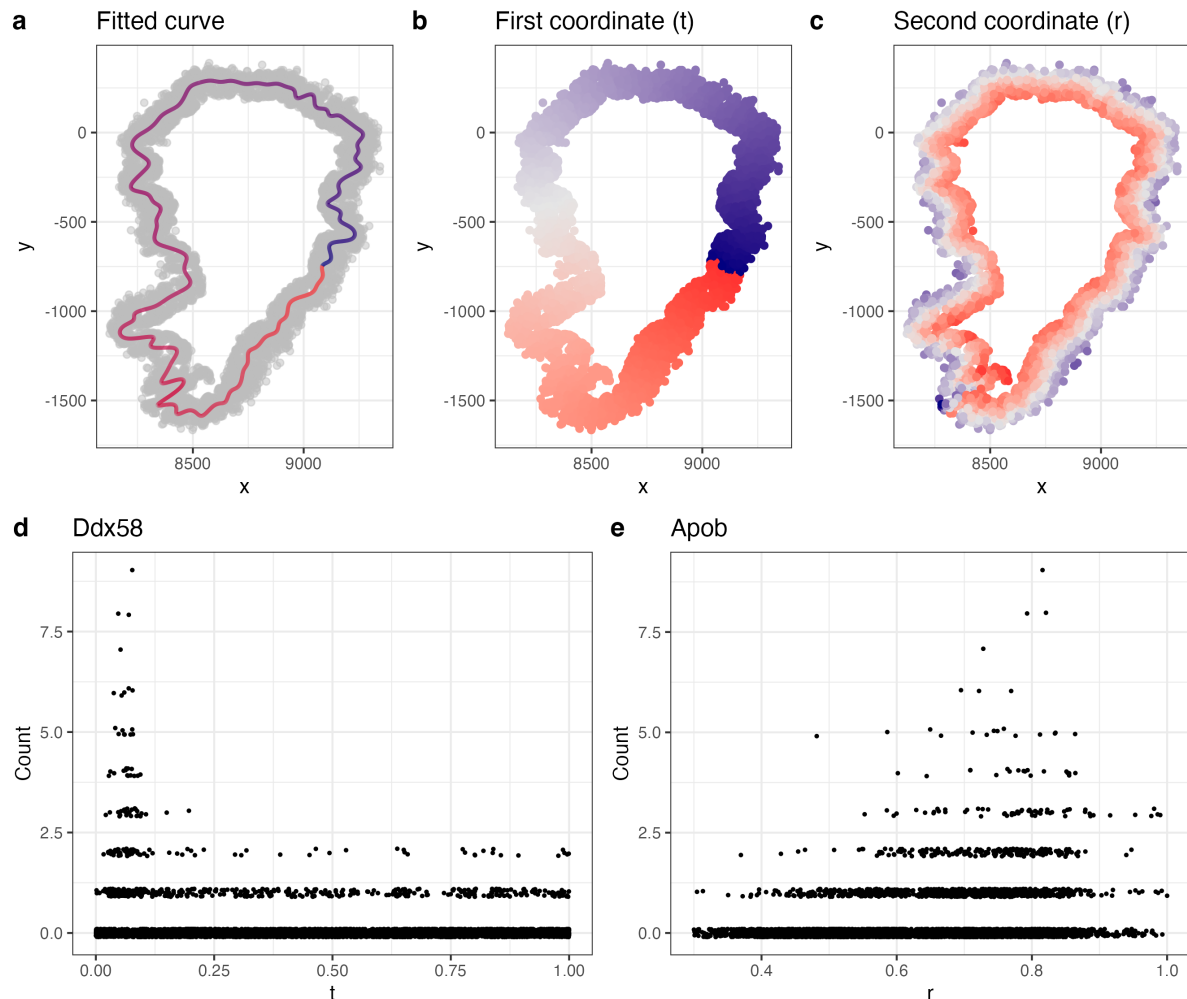


Figure 2: **Overview of MorphoGAM.** **a.** MorphoGAM begins by estimating a smooth parametric curve passing through the spatial transcriptomic sample coordinates. **b.** The first morphologically relevant coordinate is defined as the position of each cell (or more generally, sample) along the estimated curve from the previous step. **c.** The second morphologically relevant coordinate is defined as the position of the cell in the direction orthogonal to the curve at a given point. **d.** Genes with localized variation such as *Ddx58* show strong expression variation as a function of the first morphologically relevant coordinate. **e.** Genes with a radial gradient pattern such as *Apob* show expression variation as a function of the second morphologically relevant coordinate.

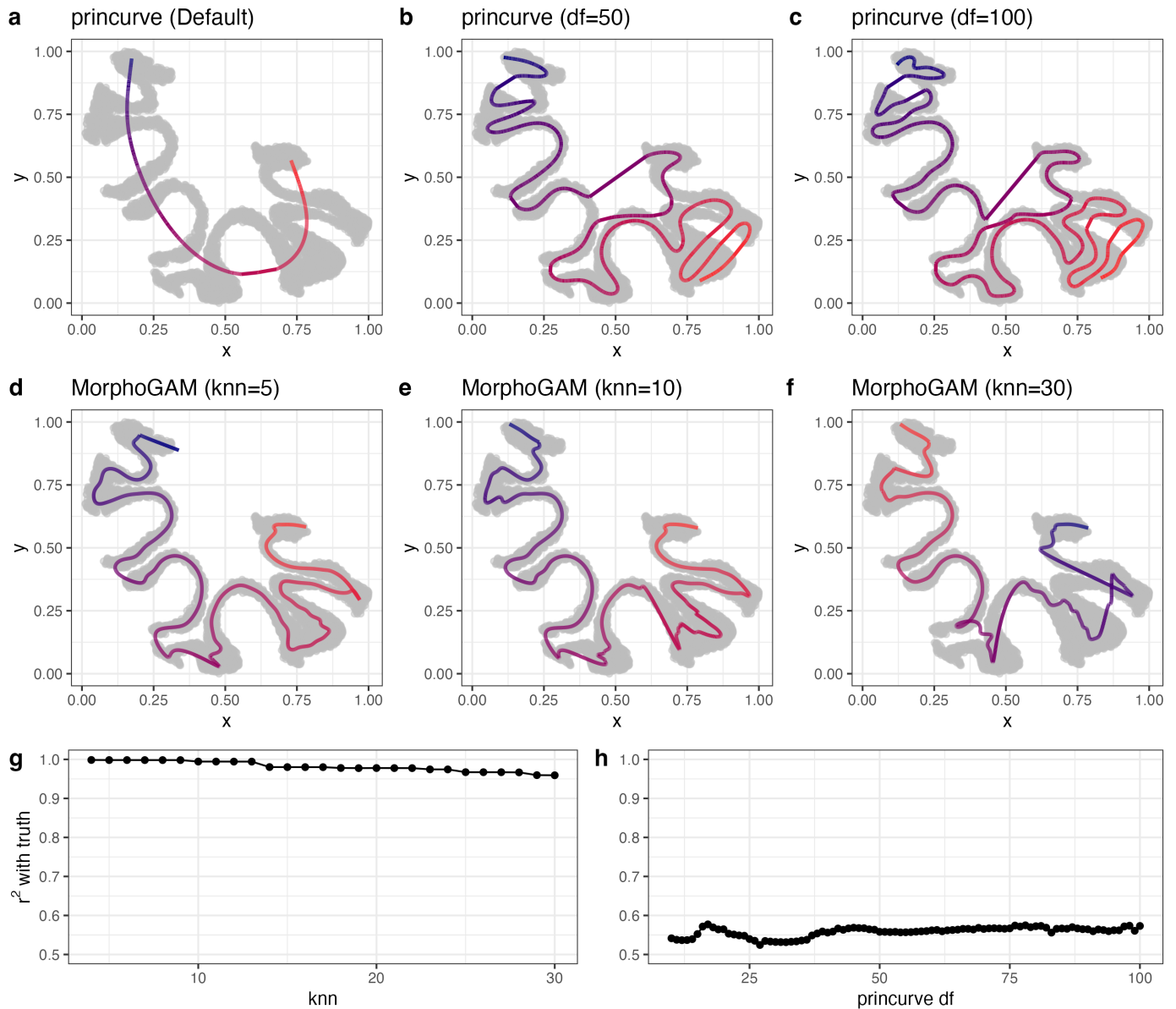


Figure 3: **MorphoGAM outperforms existing curve estimation approaches.** **a.** Applying the principal curves method (Hastie and Stuetzle, 1989) as implemented by the *princurve* (Cannoodt, 2018) R package. **a.**, **b.**, and **c.** show the estimated curve for three different values of the tuning parameter, which is degrees of freedom (df) of the smoothing spline. **d.**, **e.** and **f.** show the estimated curve from MorphoGAM with three different values of its tuning parameter k-nearest neighbor (kNN). **g.**, **h.** Using the estimated coordinate t_j from a hand-drawn ground truth (Fig S3) we compute the squared spearman correlation between this and the estimated coordinate from both methods as the tuning parameters vary.

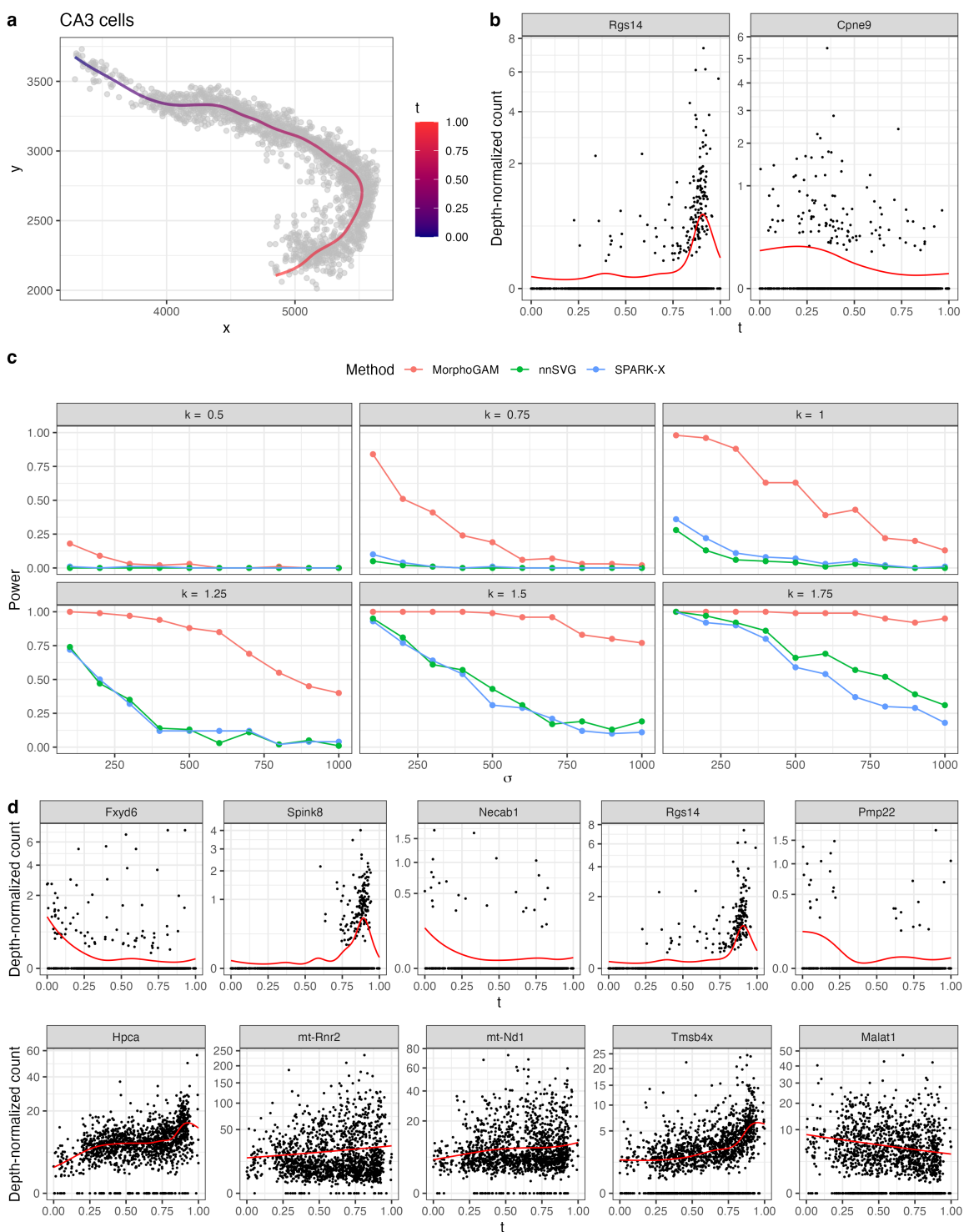


Figure 4: **MorphoGAM increases power to detect relevant spatially variable genes.** **a.** The estimated curve \hat{f} on the CA3 cells from mouse hippocampus (see Fig 1d). **b.** The estimated functions $\hat{h}(\hat{t}_j)$ for two previously reported spatially variable genes *Rgs14* and *Cpne9*. **c.** Comparing hypothesis-based frameworks to detect SVGs; a gene with $\mu_{gj} = 1 + \kappa \exp(-\sigma(\hat{t} - 0.5)^2)$ and $\theta = 5$ was simulated and labeled as SVG is the corresponding p -value was smaller than $0.05/20000$. The power reflects the proportion of 100 trials where the null hypothesis was correctly rejected. **d.** Plotting the genes with the largest peak and range summaries.

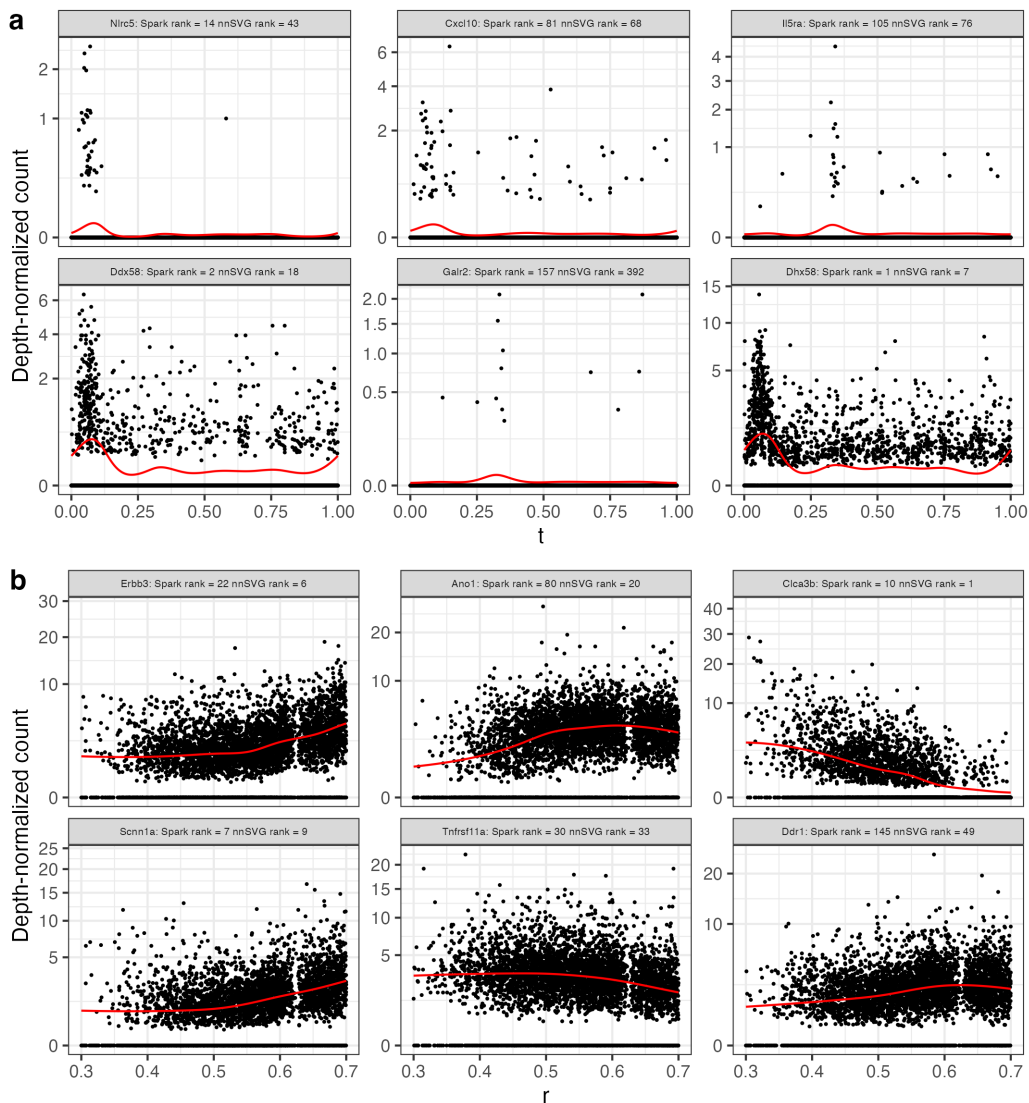


Figure 5: MorphoGAM identifies additional genes with localized and radial gradient pattern. **a.** The top six genes identified when ranking by the peak of h_g . That is, genes with a high log fold-change relative to baseline in the direction of the first morphologically relevant coordinate. **b.** The top six genes identified when ranking by the range of s_g . That is, genes with a large changes on the scale of the counts in the direction of the second morphologically relevant coordinate. Each label shows the ranking of the gene from SPARK-X (Zhu et al., 2021) and nnSVG (Weber et al., 2023).

364

365

366

Supplementary Material for “Identifying spatially variable genes by projecting to morphologically relevant curves”

367

S1 Mathematical details of curve and coordinate estimation

368

S1.1 Notation

369

370

371

372

We use \otimes to denote outer-product: if $u, v \in \mathbb{R}^n$ then $u \otimes v := uv^\top \in \mathbb{R}^{n \times n}$. If $f : [a, b] \rightarrow \mathbb{R}^k$ denotes a parametric curve, then f can be written in terms of k component functions $f(t) = (f_1(t), \dots, f_k(t))$ and $f'(t) := (f'_1(t), \dots, f'_k(t))$. We say f is *smooth* if $f'_i(t)$ exists and is continuous for all t . For a matrix $A \in \mathbb{R}^{m \times n}$, $\|A\|_F^2 = \sum_{i=1}^n \sum_{j=1}^m A_{ij}^2$ denotes Frobenius norm.

373

S1.2 Assumptions

374

375

We assume the following conditions, which are necessary (but not sufficient) for the identifiability of model (2.1):

376

377

1. $\bar{t} := \frac{1}{n} \sum_{j=1}^n t_j = 0$ and $t_1 < 0$.

2. $\|f'(t)\|_2 = 1$ for all t .

378

379

380

381

382

For condition 1, note that for any constant c , $\tilde{f} : [a - c, b - c] \rightarrow \mathbb{R}^2$ defined by $\tilde{f}(t) = f(t - c)$ satisfies $f(t_j) = \tilde{f}(t_j + c)$ for every j . Similarly, $\tilde{f} : [-b, -a] \rightarrow \mathbb{R}^2$ defined by $\tilde{f}(t) = f(-t_j)$ satisfies $f(t_j) = \tilde{f}(-t_j)$. For condition 2, define $h(t) = \int_a^t \|f'(s)\|_2 ds$ and note that $h'(t) = \|f'(t)\|_2$. Then we can define the reparameterized curve $\tilde{f} := f(h^{-1}(t))$. h^{-1} exists and is differentiable by the inverse function theorem. In particular,

$$\|(f \circ h^{-1})'(t)\|_2 = \|f'(h^{-1}(t)) \cdot (h'(h^{-1}(t)))^{-1}\|_2 = \|f'(h^{-1}(t)) \cdot (f'(h^{-1}(t)))^{-1}\|_2 = 1. \quad (\text{S1.1})$$

383

384

385

A full introduction to parametric curves is given by Lastra (2021). We also note that these conditions are not sufficient to ensure the identifiability of model (2.1) as a fully identifiable model would likely need to specify a distribution or a procedure from which the t_j are obtained.

386 S1.3 Linear curve

387 We now describe how to estimate the first coordinate in the case of a linear curve (i.e., $x_1 \neq x_2 \Rightarrow f(x_1) \neq$
 388 $f(x_2)$). If we assume that the approximation in equation (4.3) is equality, i.e., $d_{G_k}(i, j) = |t_i - t_j|$, then
 389 $d_{G_k}^2 \in \mathbb{R}^{n \times n}$ can be written in matrix form as

$$d_{G_k}^2 = t^2 \otimes 1_n + 1_n \otimes t^2 - 2t \otimes t \quad (\text{S1.2})$$

390 where t^2 is applied entry-wise to $t := (t_1, \dots, t_n)$ and $1_n \in \mathbb{R}^n$ is a vector of 1's. Now define the centering
 391 matrix $H \in \mathbb{R}^{n \times n}$ as

$$H = I - \frac{1}{n} 1_n \otimes 1_n \quad (\text{S1.3})$$

392 Applying H on the right has the property of subtracting the row means while applying H on the left
 393 subtracts the column means.

394 **Lemma S1.1.** *The “double centered” matrix b_{G_k} satisfies*

$$-\frac{1}{2} b_{G_k} := -\frac{1}{2} H d_{G_k}^2 H = t \otimes t \quad (\text{S1.4})$$

395 *Proof.* The proof is derived from [Ghojogh et al. \(2023\)](#). Because $H(1_n \otimes t^2) = 0$ and $(t^2 \otimes 1_n)H = 0$, we
 396 have

$$\frac{-1}{2} H d_{G_k}^2 H = \left(-\frac{1}{2} H(t^2 \otimes 1_n) + H(t \otimes t) \right) H \quad (\text{S1.5})$$

$$= H(t \otimes t)H \quad (\text{S1.6})$$

$$= t \otimes t \quad (\text{S1.7})$$

397 where the last line follows because $\bar{t} = 0$ by assumption. □

398 The above result shows that $-\frac{1}{2} b_{G_k}$ is a rank 1 matrix with a positive eigenvalue $\langle t, t \rangle > 0$. In practice,
 399 however, $-\frac{1}{2} b_{G_k}$ will be expected to have higher rank due to noise. For this reason, we estimate t using
 400 the top eigenvector associated with the largest eigenvalue. The following theorem shows that, under some
 401 conditions, the top eigenvector of a symmetric matrix leads to the best rank-one approximation with the
 402 smallest reconstruction error.

403 **Lemma S1.2.** *Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix, and suppose that $\lambda_{\max}(A) > 0$ and $\lambda_{\max}(A) >$*

404 $|\lambda_{\min}(A)|$, where λ_{\max} and λ_{\min} denotes the largest and smallest eigenvalues, respectively. Then

$$\operatorname{argmin}_{t \in \mathbb{R}^n} \|A - t \otimes t\|_F^2 = \sqrt{\lambda_{\max}(A)} u_1 \quad (\text{S1.8})$$

405 where u_1 is the unit eigenvector corresponding to $\lambda_{\max}(A)$.

406 *Proof.* As A is symmetric, we may write

$$A = \sum_{i=1}^n \lambda_i (u_i \otimes u_i) \quad (\text{S1.9})$$

407 with $u_1, \dots, u_n \in \mathbb{R}^n$ orthonormal. Then

$$A = \sum_{i=1}^n |\lambda_i| (\operatorname{sign}(\lambda_i) u_i) \otimes u_i \quad (\text{S1.10})$$

408 is a singular value decomposition (SVD) of A . By [Eckart and Young \(1936\)](#), we have

$$\min_{u, v \in \mathbb{R}^n} \|A - u \otimes v\|_F^2 = \sum_{i=2}^n |\lambda_i| \quad (\text{S1.11})$$

409 Moreover,

$$\min_{u, v \in \mathbb{R}^n} \|A - u \otimes v\|_F^2 \leq \min_{t \in \mathbb{R}^n} \|A - t \otimes t\|_F^2 \quad (\text{S1.12})$$

410 so $\min_{t \in \mathbb{R}^n} \|A - t \otimes t\|_F^2 \geq \sum_{i=2}^n |\lambda_i|$ as well. This minimum is achieved by setting $t = \sqrt{\lambda_{\max}(A)} u_1$. \square

411 In practice, it seems to be the case that the condition $\lambda_{\max}(A) > 0$ and $\lambda_{\max}(A) > |\lambda_{\min}(A)|$ always
412 holds.

413 S1.4 Closed curve

414 In the case of a closed curve $f(a) = f(b)$ and the approximation in equation (4.10) must be used for d_{G_k} :

$$\begin{aligned} \theta_i &:= 2\pi \left(\frac{t_i - b}{b - a} \right) \\ d_{G_k}(i, j) &\approx (b - a) \arccos \left(1 - \frac{(\cos(\theta_i) - \cos(\theta_j))^2 - (\sin(\theta_i) - \sin(\theta_j))^2}{2} \right). \end{aligned} \quad (\text{S1.13})$$

415 Consider $(b - a)^2 \arccos^2(1 - x^2/2)$ as a function of x . We make a second order Taylor approximation around
416 $x = 0$, which yields

$$\arccos(1) + (b - a)^2 x^2 = (b - a)^2 x^2 \quad (\text{S1.14})$$

417 Taking $c = (b - a)^2$ yields the approximation in (4.12). Then by double centering,

$$-\frac{1}{2}b_{G_k}(i, j) \approx c \cos(\theta_i) \cos(\theta_j) + c \sin(\theta_i) \sin(\theta_j). \quad (\text{S1.15})$$

418 This implies that $-\frac{1}{2}b_{G_k}$ will be approximately rank 2. Moreover, if n is large and θ_i densely populated
419 within $[0, 2\pi]$ then we have

$$\frac{1}{n} \sum_{i=1}^n \cos(\theta_i) \sin(\theta_i) \approx \int_0^{2\pi} \cos(\theta) \sin(\theta) d\theta = 0 \quad (\text{S1.16})$$

420 which shows that (S1.15) is also an (approximate) eigendecomposition of $-\frac{1}{2}b_{G_k}$. In particular, if the two
421 eigenvectors recover $\cos(\theta_i)$ and $\sin(\theta_i)$, respectively, then taking the arctangent function of the ratio should
422 be a reasonable approximation to θ_i . Because both $\cos(\theta)$ and $\sin(\theta)$ are approximate eigenvectors with
423 eigenvalue c , the top two eigenvectors could be invariant to rotation. However, in any case, the top two
424 eigenvectors would still represent the location of each cell on a circle, and the two-argument arctangent
425 function would still recover the angle along that circle.

426 **S2** Supplementary figures

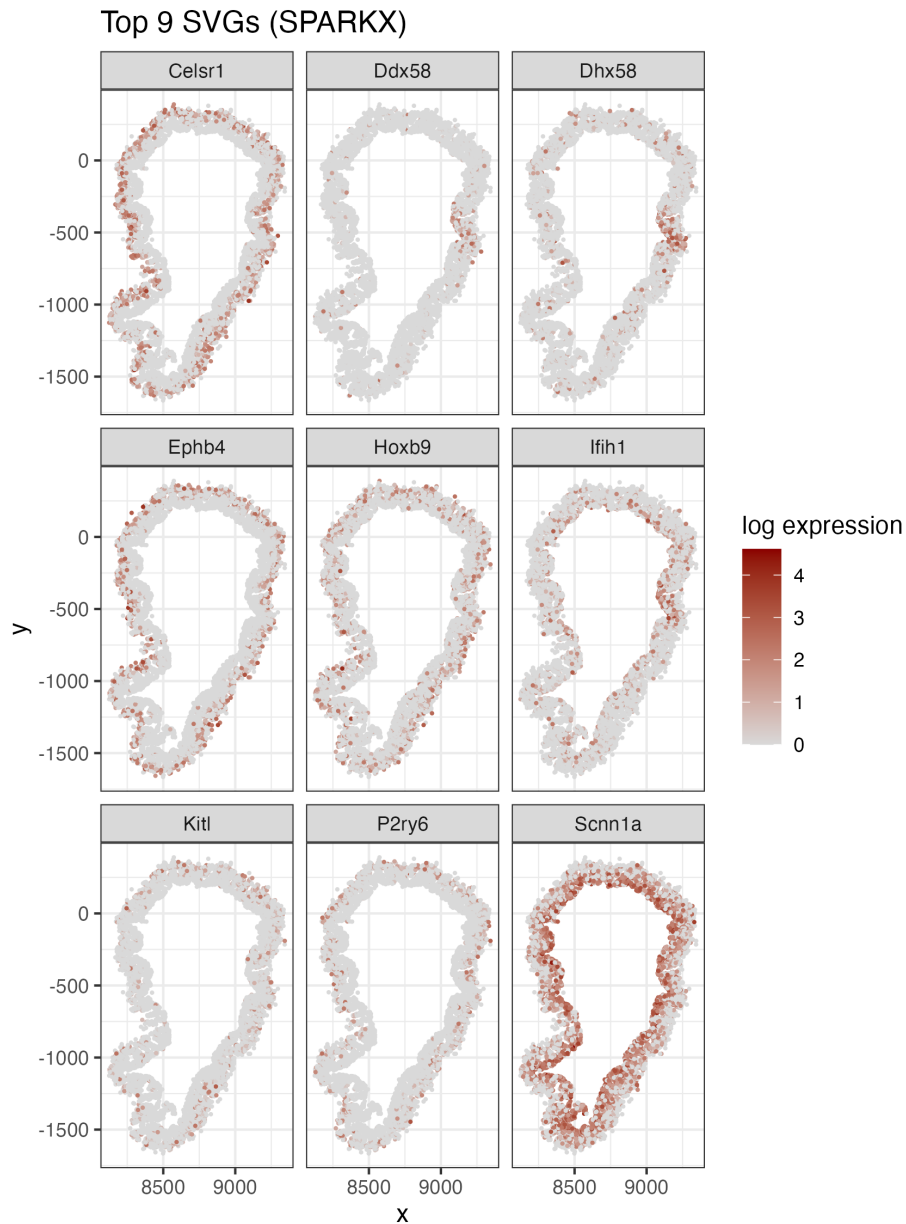


Figure S1: The top 9 SVGs identified by SPARK-X (Zhu et al., 2021) in a MERFISH measurement of a slice of the healthy mouse colon. Specifically, *Ddx58* had a reported (adjusted) p -value of 6.38×10^{-67} and *Apob* had a reported (adjusted) p -value of 8.04×10^{-9} .

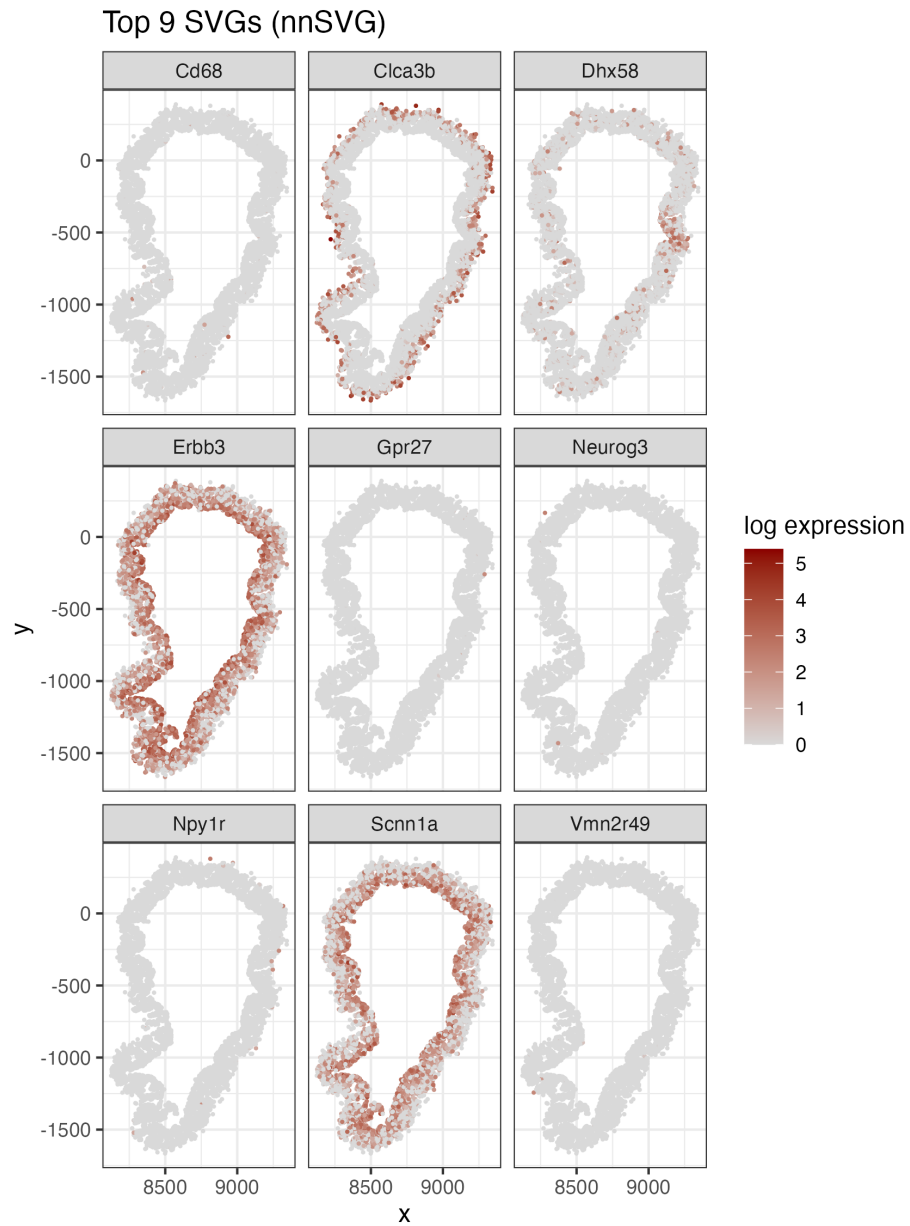


Figure S2: The top 9 SVGs identified by nSVG (Weber et al., 2023) in a MERFISH measurement of a slice of the healthy mouse colon. Specifically, *Ddx58* and *Apob* both had reported (adjusted) p -values of 0.

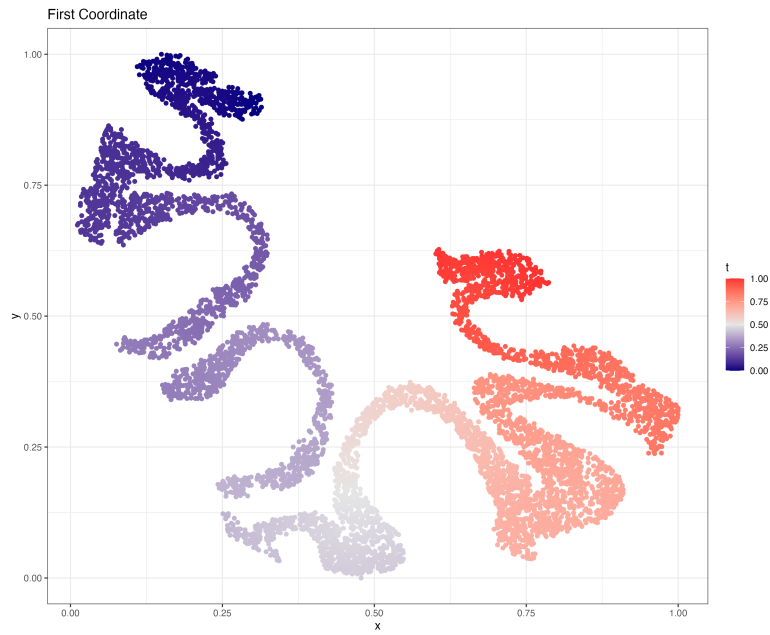


Figure S3: The estimated coordinate from the hand-drawn path on the granule cells.

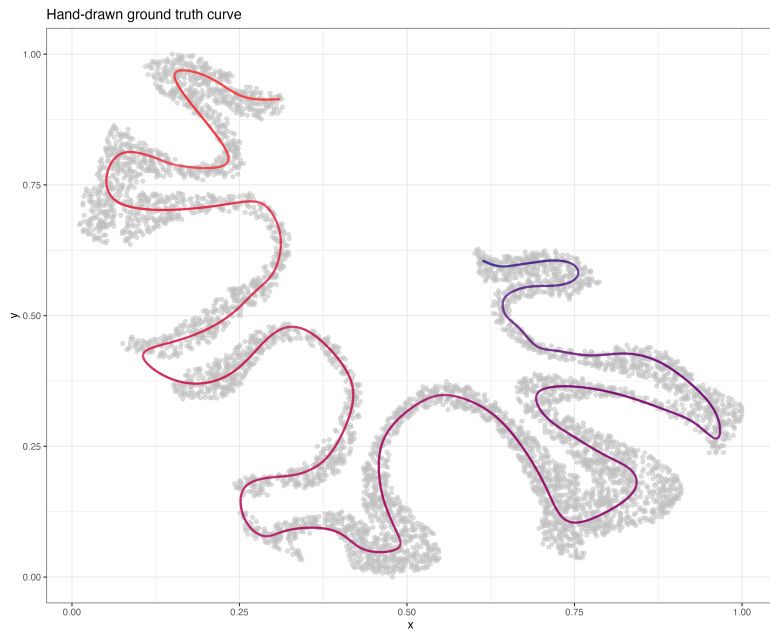


Figure S4: The hand-drawn curve on the granule cells.

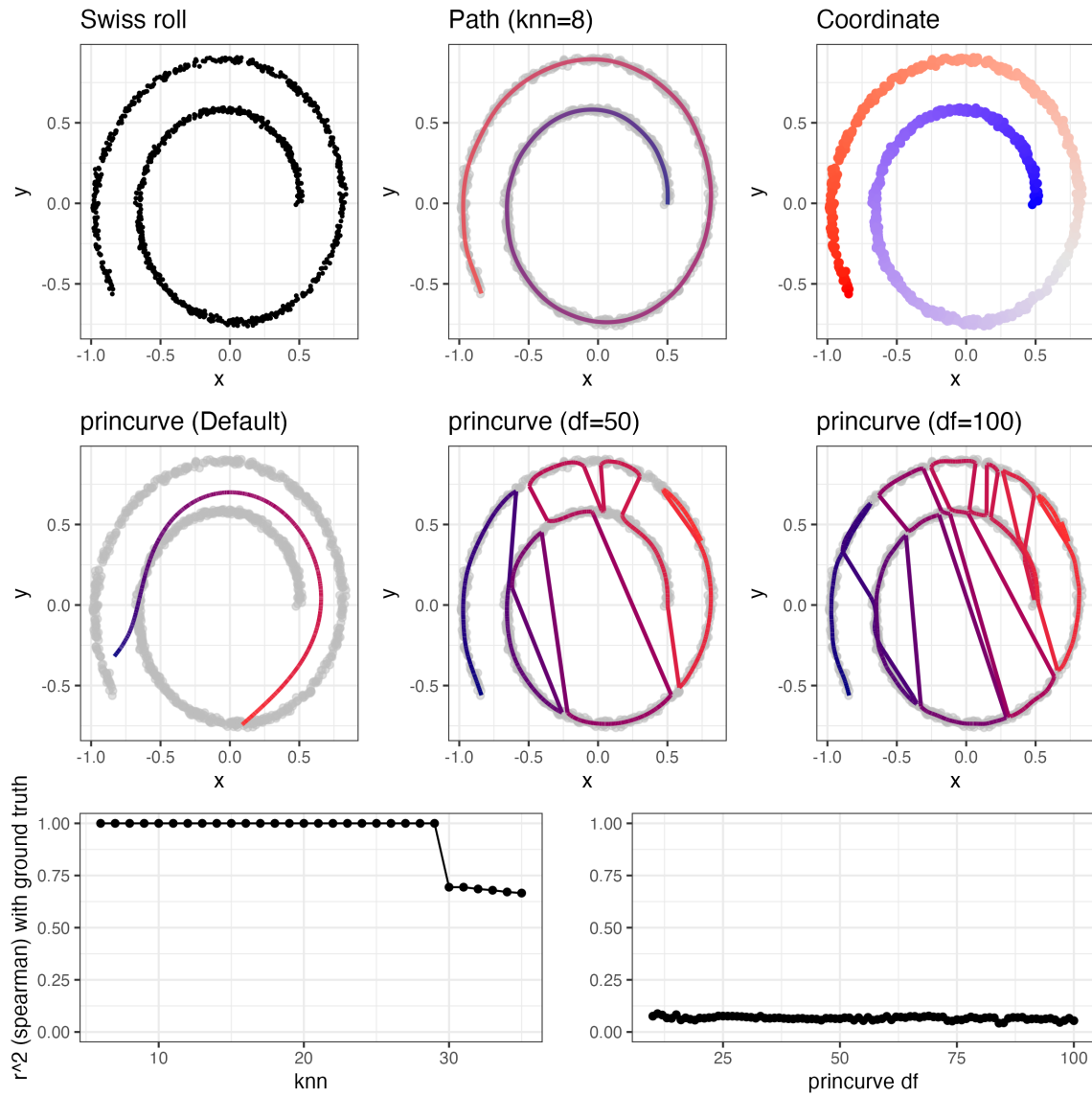


Figure S5: Repeating the analysis in Figure 3 instead using a simulated swiss roll. The inability of the standard principal curves algorithm to accurately reconstruct the swiss roll was discussed in (Kégl et al., 2000).

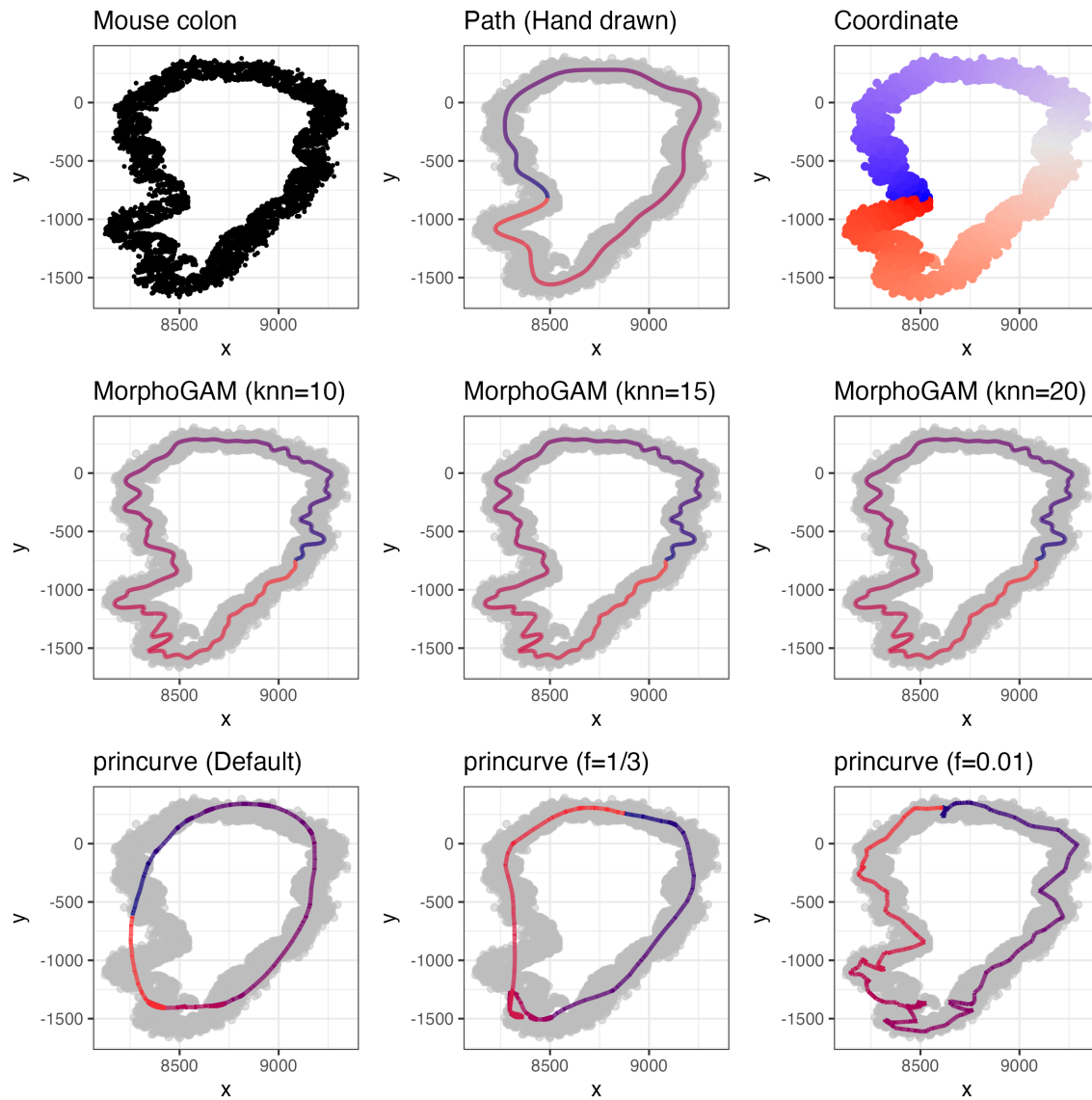


Figure S6: Repeating the analysis in Figure 3 instead using the mouse mucosa data. In this case a periodic smoother was used in princurve.

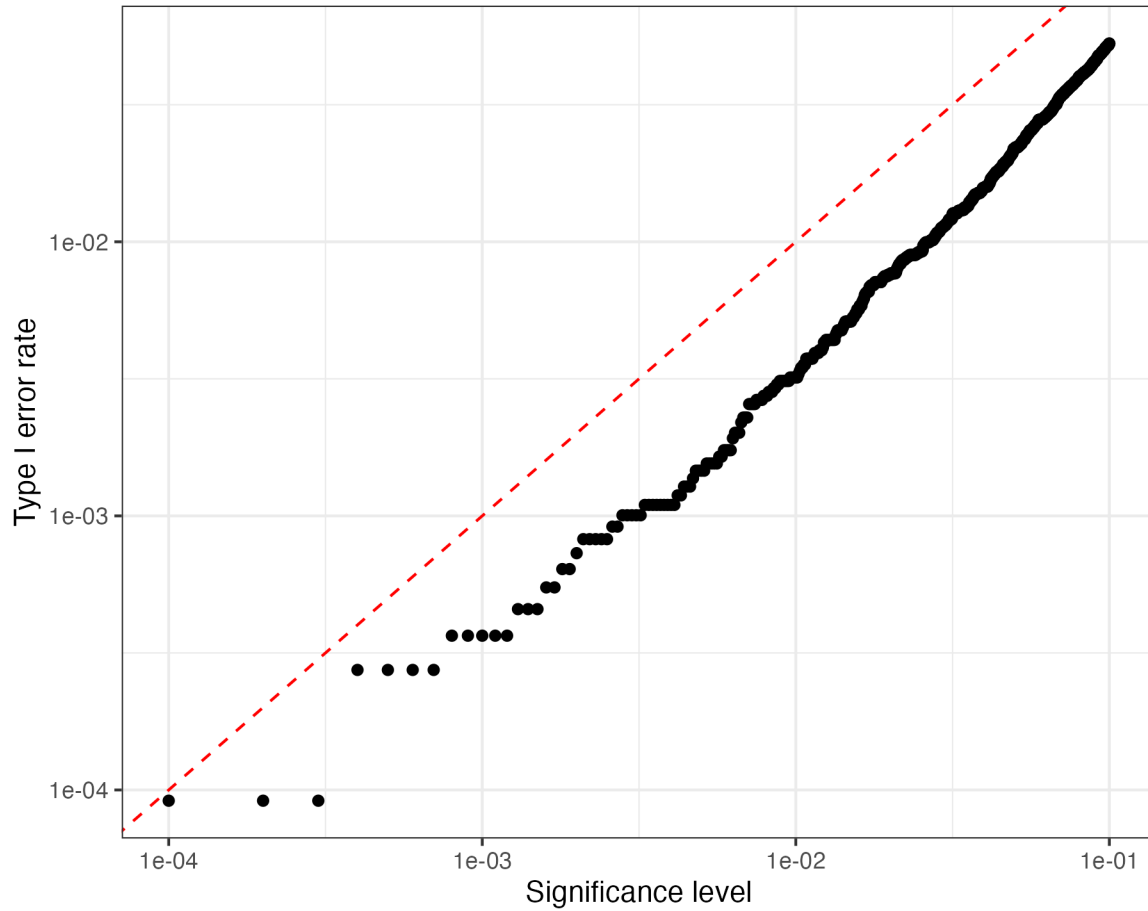


Figure S7: Spatial locations in the CA3 data were randomly permuted to produce a null dataset where there should be no SVGs. The proportion of genes with a p -value smaller each significance level was computed (the Type I error rate). The red-dashed line indicates the nominal type I error rate.

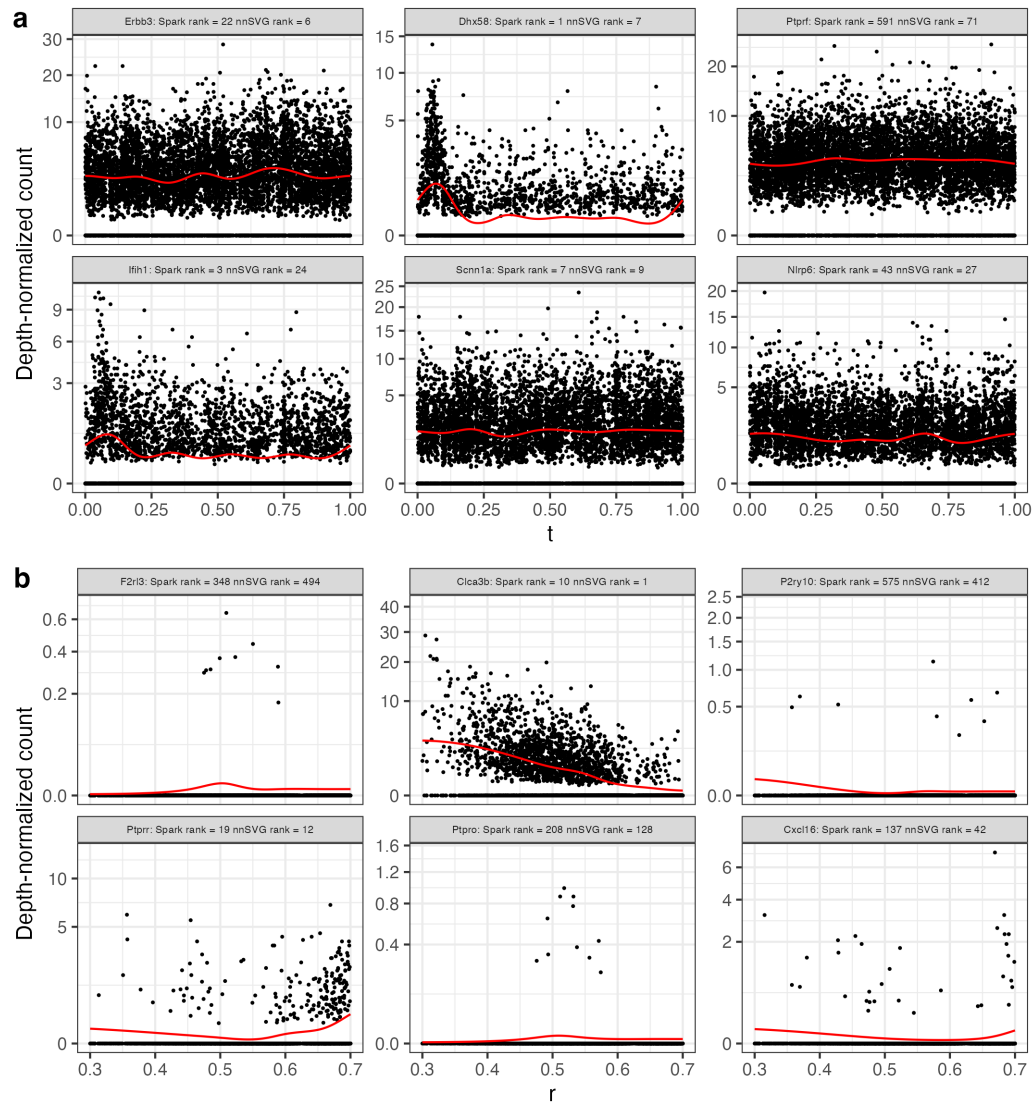


Figure S8: Repeating the analysis of Figure 5 plotting the genes with the largest range in the direction of the first morphologically relevant coordinate t_j and the genes with the largest peak in the direction of the second morphologically relevant coordinate r_j .