

Order under uncertainty: robust differential expression analysis using probabilistic models for pseudotime inference

Kieran Campbell^{1,2} and Christopher Yau^{2,3}

¹Department of Physiology, Anatomy and Genetics, University of Oxford, UK

²Wellcome Trust Centre for Human Genetics, University of Oxford, UK

³Department of Statistics, University of Oxford, UK

April 5, 2016

Abstract

Single cell gene expression profiling can be used to quantify transcriptional dynamics in temporal processes, such as cell differentiation, using computational methods to label each cell with a ‘pseudotime’ where true time series experimentation is too difficult to perform. However, owing to the high variability in gene expression between individual cells, there is an inherent uncertainty in the precise temporal ordering of the cells. Preexisting methods for pseudotime ordering have predominantly given point estimates precluding a rigorous analysis of the implications of uncertainty. We use probabilistic modelling techniques to quantify pseudotime uncertainty and propagate this into downstream differential expression analysis. We demonstrate that reliance on a point estimate of pseudotime can lead to inflated false discovery rates compared and that probabilistic approaches provide greater robustness and measures of the temporal resolution that can be obtained from pseudotime inference.

Background

The emergence of high-throughput single cell genomics as a tool for the precision study of biological systems Kalisky and Quake (2011); Shapiro et al. (2013); Macaulay and Voet (2014); Wills and Mead (2015) has given rise to a variety of novel computational and statistical modelling challenges Stegle et al. (2015); Trapnell (2015). One particular area of interest has been the

25 study of transcriptional dynamics in temporal processes, such as cell differentiation or prolifera-
26 tion Treutlein et al. (2014); Tsang et al. (2015), in order to understand the coordinated changes
27 in transcription programming that underlie these processes. In the study of such systems, prac-
28 tical experimental designs that can allow the collection of real time series data maybe difficult or
29 impossible to achieve. Instead, investigators have adopted computational methods to identify
30 temporal signatures and trends from unordered genomic profiles of single cells, a process known
31 as *pseudotemporal ordering* Qiu et al. (2011); Bendall et al. (2014); Marco et al. (2014); Trapnell
32 et al. (2014); Moignard et al. (2015); Reid and Wernisch (2015). Computational approaches for
33 this problem were first tackled in the context of gene expression microarray analysis of bulk cell
34 populations Magwene et al. (2003); Gupta and Bar-Joseph (2008); Qiu et al. (2011) but the
35 recent availability of single cell technology overcomes the limitations of measuring population
36 averaged signals in bulk analyses.

37 Pseudotemporal ordering of whole-transcriptome profiles of single cells with unsupervised
38 computational methods has an advantage over cytometry-based assays in that it does not rely on
39 *a priori* knowledge of marker genes. The principle underlying these methods is that each single
40 cell RNA sequencing experiment constitutes a time series in which each cell represents a distinct
41 time point along a continuum representing the underlying degree of temporal progress (Figure
42 1A). During the single cell capture process, the *true* temporal label that identifies the stage
43 of the cell is lost (Figure 1B) and these parameters become latent, unobserved quantities that
44 must be statistically inferred from the collection of single cell expression profiles (Figure 1C).
45 Importantly, absolute physical time information will in general be irretrievably lost and it is only
46 possible to assign a “pseudotime” for each cell that provides a relative quantitative measure of
47 progression. Consequently, whilst the correspondence between physical and pseudotime ordering
48 maybe conserved, the pseudotimes themselves are not necessarily calibrated to actual physical
49 times. Pseudotime ordering can potentially be used to recapitulate temporal resolution in an
50 experiment that does not explicitly capture labelled time series data. The pseudotimes could
51 then be used to identify genes that are differentially expressed across pseudotime (Figure 1D)
52 providing insight into the evolution of transcription programming.

53 Practically, current methods for pseudotime inference proceed via a multi-step process.
54 First, gene selection and dimensionality reduction techniques are applied to compress the infor-

55 mation held in the high-dimensional gene expression profiles to a small number of dimensions
56 (typically two or three for simplicity of visualisation). The identification of an appropriate
57 dimensionality reduction technique is a *subjective* choice and a number of methods have been
58 adopted such as Principal and Independent Components Analysis (P/ICA) and highly non-
59 linear techniques such as diffusion maps Haghverdi et al. (2015) or stochastic neighbourhood
60 embedding (SNE) Hinton and Roweis (2002); Van der Maaten and Hinton (2008); Amir et al.
61 (2013). This choice is guided by whether the dimensionality reduction procedure is able to
62 identify a suitable low-dimensional embedding of the data that contains a relatively smooth
63 trajectory that might plausibly correspond to the temporal process under investigation.

64 Next, the pseudotime trajectory of the cells in this low-dimensional embedding is charac-
65 terised. In Monocle Trapnell et al. (2014) this is achieved by the construction of a minimum
66 spanning tree (MST) joining all cells. The diameter of the MST provides the main trajectory
67 along which pseudotime is measured. Related graph-based techniques (Wanderlust) have also
68 been used to characterise temporal processes from single cell mass cytometry data Bendall et al.
69 (2014). In SCUBA Marco et al. (2014) the trajectory itself is directly modelled using principal
70 curves Hastie and Stuetzle (2012) and pseudotime is assigned to each cell by projecting its
71 location in the low-dimensional embedding on to the principal curve. The estimated pseudo-
72 times can then be used to order the cells and to assess differential expression of genes across
73 pseudotime.

74 A limitation of these approaches is that they provide only a single *point estimate* of pseu-
75 dotimes concealing the full impact of gene expression variability and technical noise. As a
76 consequence, the statistical uncertainty in the pseudotimes is not propagated to downstream
77 analyses precluding a thorough treatment of robustness and stability. To date, the impact of
78 this pseudotime uncertainty has not been explored and its implications are unknown as the
79 methods applied do not possess a probabilistic interpretation. However, we can examine the
80 stability of the pseudotime estimates by taking multiple random subsets of a dataset and re-
81 estimating the pseudotimes for each subset. For example, we have found that the pseudotime
82 assigned to the same cell can vary considerably across random subsets in Monocle (details given
83 in Supplementary Materials and Supplementary Figure S1).

84 In order to address pseudotime uncertainty in a formal and coherent framework, probabilistic

85 approaches using Gaussian Process Latent Variable Models (GPLVM) have been used recently
86 as non-parametric models of pseudotime trajectories Reid and Wernisch (2015); Campbell and
87 Yau (2015). These provide an explicit model of pseudotimes as latent embedded one-dimensional
88 variables and can be fitted within a Bayesian statistical framework allowing posterior uncertainty
89 in the pseudotimes to be derived using Markov Chain Monte Carlo (MCMC) simulations. In
90 this article we adopt this framework based to assess the impact of pseudotime uncertainty
91 on downstream differential analyses. We will show that pseudotime uncertainty can be non-
92 negligible and when propagated to downstream analysis may considerably inflate false discovery
93 rates. We demonstrate that there exists a limit to the degree of recoverable temporal resolution,
94 due to intrinsic variability in the data, with which we can make statements such as “this cell
95 precedes another”. Finally, we propose a simple means of accounting for the different possible
96 choices of reduced dimension data embeddings. We demonstrate that, given sensible choices
97 of low-dimensional representations, these can be combined to produce more robust pseudotime
98 estimates. Overall, we outline a modelling and analytical strategy to produce more stable
99 pseudotime based differential expression analysis.

100 **Results**

101 **Probabilistic pseudotime inference using Gaussian Process Latent Variable** 102 **Models**

103 We first provide a brief overview of the Gaussian Process Latent Variable Model Titsias and
104 Lawrence (2010). The GPLVM uses a Gaussian Process to define a stochastic mapping between
105 a low-dimensional latent space to a (typically) higher dimensional observation space. A Gaussian
106 Process is characterised by a mean function describing the expected mapping between the latent
107 and observation spaces and a covariance function that describes the covariance between the
108 mapping function evaluated at any two arbitrary latent positions. The covariance function
109 therefore acts to control the second-order statistics of the Gaussian Process and suitable choices
110 can be designed to encourage properties such as smoothness, periodicity or other second-order
111 features.

112 For this application, the latent space is one-dimensional, describing pseudotime progression
113 whilst the observations are the reduced dimensionality representations of the single cell expres-

114 sion data. We will use Bayesian inference to characterise the joint posterior distribution $p(\mathbf{t}|\mathbf{X})$
115 of the pseudotimes $\mathbf{t} = \{t_1, \dots, t_n\}$ given the expression data $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ for n single cells.
116 As the integrals involved are mathematically intractable, we will use Markov Chain Monte Carlo
117 simulations to obtain a numerical approximation to the posterior by drawing samples from the
118 posterior distribution. Each sample corresponds to one possible trajectory *and* ordering for the
119 cells with the set of samples providing an approximate distribution of pseudotimes. The pseu-
120 dotime values are between measured 0 and 1 where a value of 0 corresponds to one end state of
121 the temporal process and a value 1 to the other. In this work we focus only on non-bifurcating
122 processes. Figure 2 gives a diagrammatic representation of our proposed workflow and a more
123 detailed model descriptions is given in Methods.

124 Sources of uncertainty in pseudotime inference

125 We applied our probabilistic pseudotime inference to three published single-cell RNA-seq datasets
126 of differentiating cells: myoblasts in Trapnell et al. (2014) Trapnell et al. (2014), hippocampal
127 quiescent neural stem cells in Shin et al. (2015) Shin et al. (2015) and sensory epithelia from
128 the inner ear in Burns et al. (2015) Burns et al. (2015). For the Trapnell and Shin datasets
129 we used Laplacian Eigenmaps Belkin and Niyogi (2003) for dimensionality reduction prior to
130 pseudotime inference, while for the Burns dataset we used the PCA representation of the cells
131 from the original publication (for a detailed description of our analysis, see Supplementary
132 Methods: Data Analysis). These particular choices of reduced dimensionality representations
133 gave visually plausible trajectory paths in two dimensions.

134 An implicit assumption in pseudotime ordering is that proximity in pseudotime should reflect
135 proximity in the observation or data space. That is, two cells with similar pseudotime assign-
136 ments should have similar gene expression profiles but, in practice, cell-to-cell variability and
137 technical noise means that the location of the cells in the observation space will be variable even
138 if they truly do have the same pseudotime. We plotted posterior mean pseudotime trajectories
139 for the three datasets learned using the GPLVM in Figure 3A-C and the posterior predictive
140 data distribution $p(\mathbf{X}^*|\mathbf{X})$. The posterior predictive data distribution gives an indication of
141 where *future* data points might occur given the existing data. Notice that for all three data
142 sets, this distribution can be quite diffuse. This is due to a combination of actual cell-to-cell

143 expression variability, manifesting as a spread of data points around the mean trajectory, but
144 also model misspecification (the difference between what our “assumed” model and the “true”
145 but unknown data generating mechanism).

146 It is interesting to discuss the latter point as it is an issue that is often not adequately
147 addressed or fully acknowledged in the literature. The GPLVM applied assumes a homoscedastic
148 noise distribution which is uniform along the pseudotime trajectory. However, it is clear that
149 the variability of the data points can change along the trajectory and a heteroscedastic (non-
150 uniform) noise model may be more appropriate in certain scenarios. Unfortunately, whilst
151 models of heteroscedastic noise processes can be applied Le et al. (2005), these typically severely
152 complicate the statistical inference and require a model of how the variability changes over
153 pseudotime which is likely to be unknown. The important point here is that the posterior
154 probability calculations are always calibrated with respect to a given model. The better the
155 model represents the true data generating mechanism, the better calibrated the probabilities.
156 Model misspecification can also contribute to posterior uncertainty in inferred parameters.

157 Returning to the intrinsic cell-to-cell variability, we next considered the conditional posterior
158 predictive data distributions $p(\mathbf{X}^*|t^*, \mathbf{X})$ which are shown in Figure 3D-F. These distributions
159 show the possible distribution of future data points given the existing data *and* a theoretical
160 pseudotime t^* and, in this example, we condition on pseudotimes $t^* = 0.5$ and $t^* = 0.7$. Al-
161 though the two pseudotimes differ by a magnitude of 0.2, the conditional predictive distributions
162 are very close or overlapping. This means that cells with pseudotimes of 0.5 or 0.7 could have
163 given rise to data point occupying these overlapping regions. This variability is what ultimately
164 limits the temporal resolution that can be obtained.

165 It is important to note that the posterior mean trajectories correspond to certain *a priori*
166 or subjective smoothness assumptions (specified as hyperparameters in the model specification)
167 which dictate the curvature properties of the trajectory. Figure 4 shows three alternative
168 posterior mean pseudotime trajectories for the Trapnell data based on different hyperparameters
169 settings for the GPLVM. In a truly unsupervised scenario all three paths could be plausible as
170 we would have little information to inform us about the true shape of the trajectory. This would
171 become an additional source of uncertainty in the pseudotime estimates. However, we favoured
172 hyperparameter settings that gave rise to well-defined (unimodal) posterior distributions that

173 resulted in multiple independent Markov Chain Monte Carlo runs converging to the same mean
174 trajectory rather than settings that give rise to a “lumpy” posterior distribution with many
175 local modes corresponding to different interpretations of the data (see Supplementary Figure
176 S3). Later on, when we consider inference using multiple representations, the ability to specify
177 a wider choice of trajectories is useful as we will demonstrate how the correspondence between
178 pseudotime trajectories in different reduced dimension representations is not always obvious
179 from a visual analysis.

180 We next examined the posterior distributions in pseudotime assignment for four cells from
181 the Trapnell dataset in Figure 5A. Uncertainty in the estimate of pseudotime is assessed using
182 the highest probability density (HPD) credible interval (CI), the Bayesian equivalent of the
183 confidence interval. The 95% pseudotime CI typically covers around one quarter of the tra-
184 jectory, suggesting that pseudotemporal orderings of single-cells can potentially only resolve a
185 cell’s place within a trajectory to a coarse estimate (e.g. ‘beginning’, ‘middle’ or ‘end’) and
186 do not necessarily dramatically increase the temporal resolution of the data. One immediate
187 consequence of this is that it is unlikely that we can make definite statements such as whether
188 one cell comes exactly before or after another. This is illustrated in Figures 5B-D which dis-
189 plays the estimated pseudotime uncertainty for all three datasets. In all the datasets, the
190 general progression is apparent, but the precise ordering of the cells has a non-trivial degree of
191 ambiguity.

192 **Failure to account for pseudotime uncertainty leads to increased false discov-** 193 **ery rates**

194 The previous section addressed the sources of statistical uncertainty in the pseudotimes. We next
195 explored the impact of pseudotime uncertainty on downstream analysis. Specifically, we focused
196 on the identification of genes that are differentially expressed across pseudotime. Typically,
197 these analyses involve regression models that assume the input variables (the pseudotimes) are
198 both fixed and certain but, with our probabilistic model, we can use the posterior samples from
199 our Bayesian model to refit the regression model to each pseudotime estimate. In doing so we
200 can examine which genes are called as significant in each of the posterior samples and assess
201 the stability of the differential expression analysis to pseudotime uncertainty by recording how

202 frequently genes are designated as significant across the posterior samples. This allowed us to
203 re-estimate the false discovery rate (FDR) fully accounting for the variability in pseudotime. As
204 there are a multitude of sources of uncertainty on top of this (such as biological and technical
205 variability) this allows us to put a lower bound on the FDR of such analyses in general.

206 Precisely, we fitted the tobit regression model from Trapnell et al. (2014) for each gene
207 for each sample from the posterior pseudotime distribution, giving us a per-gene set of false-
208 discovery-rate-corrected Q -values. We then compared the proportion of times a gene is called
209 as differentially expressed (5% FDR) across all pseudotime samples to the Q -value using a point
210 pseudotime estimate based on the maximum *a posteriori* (or MAP) estimate. We reasoned that
211 if a gene is truly differentially expressed then such expression will be robust to the underlying
212 uncertainty in the ordering. Note for comparison, our MAP estimates with the GPLVM correlate
213 strongly with Monocle derived pseudotime point estimates (see Supplementary Figure S2).

214 Figure 6A shows two analyses for two illustrative genes (ITGAE and ID2) in the Trapnell
215 data set. Using the MAP pseudotime ordering, differential expression analysis of ITGAE over
216 pseudotime attained a q -value of 0.02. However, the gene was only called significant in only
217 9% of posterior pseudotime samples with a median q -value of 0.32. In contrast, ID1 - known to
218 be involved in muscle differentiation - had a q -value of 6.6×10^{-11} using the MAP pseudotime
219 ordering, but was also called significant in all the posterior pseudotime ordering samples having
220 a median q -value of 4.4×10^{-11} . This indicates that the significance of the temporal expression
221 variability of ID1 is robust with respect to posterior sampling of the pseudotime ordering whilst
222 the significance ITGAE is much more dependent on the ordering chosen.

223 As a conservative rule of thumb, we designated a putative temporal association as a false
224 positive if the gene has a Q -value less than 5% at the MAP estimate of pseudotime but is
225 significant in less than 95% of the posterior pseudotime samples. Looking across all genes in the
226 Trapnell data, Figure 6B shows that a significant number of genes that were found to have a Q -
227 value < 0.05 and deemed significant based on the MAP pseudotime ordering, did not replicate
228 consistently and were not robust to alternate orderings. In fact, across the three datasets we
229 analysed, we found that the false discovery rate, when adjusted for pseudotime uncertainty,
230 ranged from 4% to 20% (Figure 6C). This indicates the FDR can be up to much larger than
231 the expected 5% and crucially is variable across datasets meaning there is no simple rule of

232 thumb that can be applied *a priori* to account for pseudotime uncertainty. Such values remain
233 low enough that analyses examining the coexpression of gene sets across pseudotime (such as
234 in Trapnell et al. (2014)) will still be largely valid. However, if a set of robustly differentially
235 expressed genes is required or the FDR needs to be characterised then a full probabilistic
236 treatment of pseudotime is needed.

237 **A sigmoidal model of switch-like behaviour across pseudotime**

238 In the previous section, we examined differential expression across pseudotime by fitting gen-
239 eralized additive models to the gene expression profiles Trapnell et al. (2014). Their approach
240 used a Tobit regression model with a cubic smoothing spline link function. Hypothesis testing
241 using the likelihood ratio test is conducted against a null model of no pseudotime dependence.
242 This model provides a highly flexible but non-specific model of pseudotime dependence that
243 was not suited to the next question we wished to address.

244 Specifically, we were interested in whether we could identify if two genes switched behaviours
245 at the *same* (or similar) times during the temporal process and therefore an estimate of the time
246 resolution that can be gained from a pseudotime ordering approach. This requires estimation
247 of a parameter that can be directly linked to a switch on(/off) time that is not present in the
248 Tobit regression model. As a result, we propose a “sigmoidal” model of differential expression
249 across pseudotime that better captures switch-like gene (in)activation and has easy to interpret
250 parameters corresponding to activation time and strength. By combining such a parametric
251 model with the Bayesian inference of pseudotime we can then infer the resolution to which
252 we can say whether one gene switches on or off before another. Details of the sigmoidal gene
253 activation model are given in Methods and in Supplementary Methods.

254 We applied our sigmoidal model to learn patterns of switch-like behaviour of genes in the
255 Trapnell dataset. For each gene we estimated the *activation time* (t_0) as well as the *activation*
256 *strength* (k). We fitted these sigmoidal switching models to all posterior pseudotime samples
257 to approximate the posterior distribution for the time and strength parameters. We uncovered
258 a small set of genes whose median activation strength is distinctly larger than the rest and
259 had low variability across posterior pseudotime samples implying a population of genes that
260 exhibit highly switch-like behaviour (Figure 7A). Some genes showed high activation strength

261 for certain pseudotime orderings but low overall median levels across all the posterior samples.
262 We concluded that genes with large credible intervals on the estimates of activation strength
263 do not show robust switch-like behaviour and demonstrate the necessity of using probabilistic
264 methods to infer gene behaviour as opposed to point estimates that might give highly unstable
265 results.

266 Representative examples of genes with large and small activation strengths showed marked
267 differences in the gene expression patterns corresponding to strong and weak switch-like be-
268 haviour as expected (Figure 7B). In addition, we examined the posterior density activation
269 time t_0 for the five genes showing strong switching behaviour (Figure 7C). Under a point esti-
270 mate of pseudotime each gene would give a distinct activation time with which these genes can
271 be ordered. However, when pseudotime uncertainty is taken into account, a distribution over
272 possible activation times emerges. In this case, the five genes all have activation times between
273 0.3 and 0.5 precluding a precise ordering (if one exists) of activation. Visually, this seems sensi-
274 ble since there is considerable cell-to-cell variability in the expression of these genes and not all
275 cells express the genes during the “on” phase. We are therefore unable to determine whether
276 the “on” phase begins when the first cell with high expression is first observed in pseudotime or,
277 if it starts before, and the first few cells simply have null expression (for biological or technical
278 reasons).

279 We further explore this in Figure 8 which shows ten genes identified as having significant
280 switch-like pseudotime dependence but with a range of mean activation times t_0 . The switch-
281 like behaviour is stable to the different posterior pseudotime orderings that were sampled from
282 the GPLVM. It is clear that the two genes RARRES3 and C1S are activating at an earlier time
283 compared to the genes IL20RA and APOL4. However, we cannot be confident of the ordering
284 within the pairs RARRES3/C1S and IL20RA/APOL4 in pseudotime since the distributions
285 over the activation times are not well-separated and it is impossible to make any definitive
286 statements as to whether one of these genes (in)activates before another. If the probability of
287 a sequence of activation events is required, instead of examining each gene in isolation, we can
288 count the number of posterior samples in which one gene precedes another instead and evidence
289 may emerge of a possible ordering. These observations suggests a finite temporal resolution
290 limit that can be obtained using pseudotemporal ordering.

291 We note that we have deliberately avoid directly linking the sigmoidal gene activation and
292 GPLVM pseudotime models to derive a single, joint model. In a joint model, the inference
293 would attempt to order the cells in such a way as to maximise the fit of the sigmoidal and
294 GPLVM to the expression data. However, as the sigmoidal model is only intended to identify
295 genes with switch-like behaviour, it cannot explain other types of pseudotime dependence that
296 may and do exist. This model misspecification would potentially drive inference in ways that
297 cannot be foreseen.

298 **Learning trajectories from multiple reduced data representations**

299 Finally, we address the impact of the subjective choice of dimensionality reduction that is
300 normally applied to single cell gene expression data prior to pseudotime ordering and estimation.
301 Typically, the choice of dimensionality reduction approach is based on whether the method
302 gives rise to a putative pseudotime trajectory in the reduced dimensionality representation
303 from visual inspection followed by confirmational analysis by examining known marker genes
304 with established temporal association. This may lead to a number of possibilities since the same
305 trajectory may exist in a number of reduced dimensionality representations.

306 One characteristic of the GPLVM is that the likelihood is conditionally independent across
307 input dimensions. A consequence of this is that we can integrate heterogeneous data sources
308 to learn pseudotimes as there is no requirement that each input dimension should come from
309 the same representation or assay. We exploited this feature to examine the effect of the initial
310 dimensionality reduction stage and see if we can learn pseudotime trajectories from multiple
311 reduced dimension representations. Many dimensionality reduction algorithms have been ap-
312 plied to single-cell RNA-seq data, including PCA Shin et al. (2015), ICA Trapnell et al. (2014),
313 t-SNE Marco et al. (2014) and diffusion maps Haghverdi et al. (2015).

314 We applied our probabilistic pseudotime inference algorithm to Laplacian Eigenmaps, PCA
315 and t-SNE representations of the three datasets under consideration. We also applied the
316 algorithm using all three representations jointly the results of which can be seen in Figure 9.
317 While the pseudotime inference algorithm can fit trajectories to all three datasets individually,
318 combining multiple representations can lead to a clearer, better defined trajectory. This allows
319 us to track the same progression of cells through multiple reduced dimension representations

320 at once, providing an equivalence to the trajectories represented by different dimensionality
321 reduction algorithms. In fact there may be no correspondence between the trajectories in the
322 different reduced dimension representations at all if the analysis is not integrated.

323 It should be noted that this approach is not necessarily an *ideal* integration model and more
324 complex multi-view learning models Xu et al. (2013) should be investigated in future that will
325 resolve the potential dependencies between the input representations. However, as a number of
326 popular dimensionality reduction techniques (e.g. t-SNE) have no probabilistic interpretation
327 and possess no underlying generative model, it is challenging to incorporate these within a
328 coherent probabilistic framework (i.e. there is no likelihood function). The suggested technique,
329 though implying simplistic independence assumptions, has practical value for incorporating such
330 representations.

331 We caution though that this integration approach is not intended to contain any arbitrary
332 number of representations provided by the user. As each representation is ultimately drawn
333 from the same underlying raw data, a representations should only be included if it provides some
334 orthogonal (near-independent) information since the GPLVM assumes the representations are
335 independent. In practice, this means selecting a small number of representations drawn from
336 very different dimensionality reduction approaches. If the representations are not independent
337 and are related this can give rise to an artificial reduction in posterior variance since we would
338 be essentially doubling sample size by replicating the same data.

339 Discussion

340 Pseudotime ordering from gene expression profiling of single cells provides the ability to extract
341 temporal information about complex biological processes from which *true* time series experi-
342 mentation may be technically challenging or impossible. In our investigations we have sought
343 to characterise the utility of a probabilistic approach to the single cell pseudotime ordering
344 problem over approaches that only return a single point estimate of pseudotime. Our work is
345 significant since it has so far not been possible to assess the impact of this statistical uncertainty
346 in downstream analyses and to ascertain the level of temporal resolution that can be obtained.

347 In order to address this we adopted a Gaussian Process Latent Variable modelling framework
348 to perform probabilistic pseudotime ordering within a Bayesian inference setting. The GPLVM

349 allows us to probabilistically explore a range of different pseudotime trajectories within the
350 reduced dimensional space. We showed that in a truly unbiased and unsupervised analysis the
351 properties of the pseudotime trajectory will never purely be a product of the data alone and
352 can heavily depend on prior assumptions about the smoothness, length scales of the trajectory
353 and noise properties. Using samples drawn from the posterior distribution over pseudotime
354 orderings under the GPLVM we were able to assess if genes that showed a significant pseudotime
355 dependence under a point (MAP) pseudotime estimate would be robust to different possible
356 pseudotime orderings. In two of the three datasets we examined we discovered that, when
357 adjusted for pseudotime uncertainty, the false discovery rate may be significantly larger than
358 the target 5%. Our investigations show that reliance on a single estimate of pseudotime ordering
359 can lead to increased number of false discoveries but that it is possible to assess the impact of
360 such assumptions within a probabilistic framework.

361 It is important to note that the GPLVM used in our investigations is not intended to be
362 a single, all-encompassing solution for pseudotime modelling problems. For our purposes, it
363 provided a simple and relevant device for tackling the single trajectory pseudotime problem in a
364 probabilistic manner but clearly has limitations when the temporal process under investigation
365 contains bifurcations or heteroscedastic noise processes (as discussed earlier). Improved and/or
366 alternative probabilistic models are required to address more challenging modelling scenarios
367 but the general procedures we describe are generic and should be applicable to any problem
368 where statistical inference for a probabilistic model can give posterior simulation samples.

369 We also developed a novel sigmoidal gene expression temporal association model that en-
370 abled us to identify genes exhibiting a strong switch-like (in)activation behaviour. For these
371 genes we were then able to estimate the activation times and use these to assess the time reso-
372 lution that can be attained using pseudotime ordering of single cells. Our investigations show
373 that pseudotime uncertainty prevents precise characterisation of the gene activation time but a
374 probabilistic model can provide a distribution over the possibilities. In application, this uncer-
375 tainty means that it is challenging to make precise statements about when regulatory factors
376 will turn on or off and if they act in unison. This places an upper limit on the accuracy of
377 dynamic gene regulation models and causal relationships between genes that could be built
378 from the single cell expression data.

379 In conclusion, single cell genomics has provided a precision tool with which to interrogate
380 complex temporal biological processes. However, as widely reported in recent studies, the prop-
381 erties of single cell gene expression data are complex and highly variable. We have shown that
382 the many sources of variability can contribute to significant uncertainty in statistical inference
383 for pseudotemporal ordering problems. We argue therefore that strong statistical foundations
384 are vital and that probabilistic methods for provide a platform for quantifying uncertainty in
385 pseudotemporal ordering which can be used to more robustly identify genes that are differen-
386 tially expressed over time. Robust statistical procedures can also temper potentially unrealistic
387 expectations about the level of temporal resolution that can be obtained from computationally-
388 based pseudotime ordering. Ultimately, as the raw input data is not true time series data,
389 pseudotime ordering is only ever an attempt to solve a *missing data* statistical inference prob-
390 lem that we should remind ourselves involves quantities (pseudotimes) that are *unknown, never*
391 *can be known*.

392 Methods

393 In addition to the descriptions below, further methodological descriptions and links to code to
394 reproduce all our findings are given in Supplementary Methods.

395 Statistical model for probabilistic pseudotime

396 The hierarchical model specification for the Gaussian Process Latent Variable model is described
397 as follows:

$$\begin{aligned} \gamma &\sim \text{Gamma}(\gamma_\alpha, \gamma_\beta), \\ \lambda_j &\sim \text{Exp}(\gamma), \quad j = 1, \dots, P, \\ \sigma_j^2 &\sim \text{InvGamma}(\alpha, \beta), \quad j = 1, \dots, P, \\ t_i &\sim \text{TruncNormal}_{[0,1]}(\mu_t, \sigma_t^2), \quad i = 1, \dots, N, \\ \Sigma &= \text{diag}(\sigma_1^2, \dots, \sigma_P^2) \\ K^{(j)}(t, t') &= \exp(-\lambda_j(t - t')^2), \quad j = 1, \dots, P, \\ \mu_j &\sim \text{GP}(0, K^{(j)}), \quad j = 1, \dots, P, \\ \mathbf{x}_i &\sim \text{MultiNorm}(\boldsymbol{\mu}(t_i), \Sigma), \quad i = 1, \dots, N. \end{aligned} \tag{1}$$

398 where \mathbf{x}_i is the P -dimensional input of cell i (of N) found by performing dimensionality re-
399 duction on the entire gene set (for our experiments $P = 2$ following previous studies). The
400 observed data is distributed according to a multivariate normal distribution with mean func-
401 tion $\boldsymbol{\mu}$ and a diagonal noise covariance matrix $\boldsymbol{\Sigma}$. The prior over the mean function $\boldsymbol{\mu}$ in each
402 dimension is given by a Gaussian Process with zero mean and covariance function K given
403 by a standard double exponential kernel. The latent pseudotimes t_1, \dots, t_N are drawn from a
404 truncated Normal distribution on the range $[0, 1)$. Under this model $|\boldsymbol{\lambda}|$ can be thought of as
405 the arc-length of the pseudotime trajectories, so applying larger levels of shrinkage to it will
406 result in smoother trajectories passing through the point space. This shrinkage is ultimately
407 controlled by the gamma hyperprior on γ , whose mean and variance are given by $\frac{\gamma_\alpha}{\gamma_\beta}$ and $\frac{\gamma_\alpha}{\gamma_\beta^2}$
408 respectively. Therefore, adjusting these parameters allows curves to match prior smoothness
409 expectations provided by plotting marker genes.

410 The hyperparameters γ_α , γ_β , α , β , μ_t and σ_t^2 are fixed and values for specific experiments
411 for given in Supplementary Information. Inference was performed using the Stan probabilistic
412 programming language Gelman et al. (2015) and our implementation is available as an R package
413 at <http://www.github.com/kieranrcampbell/pseudogp>.

414 Integrating multiple representations

415 One feature of the GPLVM is that the likelihood is conditionally independent (given the pseu-
416 dotimes) across input dimensions. If we have a set of Q reduced dimension representations
417 of single-cell data $\{\mathbf{X}_i, i = 1, \dots, Q\}$ (be they multiple representations of the same assay, e.g.
418 RNA-seq, or multiple representations of multiple assays) then the likelihood factorises across
419 each representation. For example, if we have Laplacian Eigenmaps, PCA and t-SNE represen-
420 tations of the same data then the likelihood becomes

$$p(\{\mathbf{X}\}|\mathbf{t}) = p(\mathbf{X}_{\text{LE}}|\mathbf{t})p(\mathbf{X}_{\text{PCA}}|\mathbf{t})p(\mathbf{X}_{\text{tSNE}}|\mathbf{t}) \quad (2)$$

421 where \mathbf{t} is the pseudotime vector to be learned and inference proceeds straightforwardly using
422 this product likelihood.

423 Sigmoidal model for switch-like gene (in)activation behaviours across pseudo- 424 time

425 We detail the mathematical specification of the sigmoidal switch model below. Let y_{ij} debotes
426 the \log_2 gene expression of gene i in cell j at pseudotime t_j then

$$y_{ij}(t_j) \sim \text{Norm}(\mu_i(t_j), \sigma_i^2) \quad (3)$$

427 where

$$\mu_i(t_j) = \begin{cases} \mu_i^{(0)}, & \text{if gene } i \text{ not differentially expressed,} \\ \frac{2\mu_i^{(0)}}{1+\exp(-k_i(t_j-t_i^{(0)}))}, & \text{if gene } i \text{ differentially expressed.} \end{cases} \quad (4)$$

428 Under this model the parameter k_i can be thought of as an activation ‘strength’ relating to
429 how quickly a gene switches on or off, while $t_i^{(0)}$ relates to the pseudotime at which the gene
430 switches on or off.

431 The case of a gene not being differentially expressed is a nested model of the differential
432 expression case found by setting $k = 0$. Consequently we can use a likelihood ratio test with
433 no differential expression as the null hypothesis and differential expression as the alternative
434 and twice the difference in their log-likelihoods will form a χ^2 test statistic with 2 degrees of
435 freedom. The maximum likelihood estimates of the parameters under the differential expression
436 model have no analytical solution so L-BFGS-B optimisation was used (implemented in the R
437 package `switchde`, <http://github.com/kieranrcampbell/switchde>).

438 Competing interests

439 The authors declare that they have no competing interests.

440 Author’s contributions

441 K.C. and C.Y. conceived the study. K.C. developed software and performed computer simula-
442 tions. K.C. and C.Y. wrote the manuscript.

443 Acknowledgements

444 K.C. is supported by a UK Medical Research Council funded doctoral studentship. C.Y. is
445 supported by a UK Medical Research Council New Investigator Research Grant (Ref. No.
446 MR/L001411/1), the Wellcome Trust Core Award Grant Number 090532/Z/09/Z, the John
447 Fell Oxford University Press (OUP) Research Fund and the Li Ka Shing Foundation via a
448 Oxford-Stanford Big Data in Human Health Seed Grant.

449 References

- 450 Amir, E.-a. D., K. L. Davis, M. D. Tadmor, E. F. Simonds, J. H. Levine, S. C. Bendall,
451 D. K. Shenfeld, S. Krishnaswamy, G. P. Nolan, and D. Pe'er (2013, June). viSNE enables
452 visualization of high dimensional single-cell data and reveals phenotypic heterogeneity of
453 leukemia. *Nature biotechnology* 31(6), 545–52.
- 454 Belkin, M. and P. Niyogi (2003). Laplacian Eigenmaps for Dimensionality Reduction and Data.
455 1396, 1373–1396.
- 456 Bendall, S. C., K. L. Davis, E.-A. D. Amir, M. D. Tadmor, E. F. Simonds, T. J. Chen, D. K.
457 Shenfeld, G. P. Nolan, and D. Pe'er (2014, April). Single-cell trajectory detection uncovers
458 progression and regulatory coordination in human B cell development. *Cell* 157(3), 714–25.
- 459 Burns, J. C., M. C. Kelly, M. Hoa, R. J. Morell, and M. W. Kelley (2015). Single-cell RNA-Seq
460 resolves cellular complexity in sensory organs from the neonatal inner ear. *Nature Commu-*
461 *nications* 6, 8557.
- 462 Campbell, K. and C. Yau (2015). Bayesian gaussian process latent variable models for pseudo-
463 time inference in single-cell rna-seq data. *bioRxiv*, 026872.
- 464 Gelman, A., D. Lee, and J. Guo (2015). Stan a probabilistic programming language for
465 bayesian inference and optimization. *Journal of Educational and Behavioral Statistics*,
466 1076998615606113.
- 467 Gupta, A. and Z. Bar-Joseph (2008). Extracting dynamics from static cancer expression data.
468 *IEEE/ACM transactions on computational biology and bioinformatics / IEEE, ACM* 5(2),
469 172–82.

- 470 Haghverdi, L., F. Buettner, and F. J. Theis (2015). Diffusion maps for high-dimensional single-
471 cell analysis of differentiation data. *Bioinformatics* (May), 1–10.
- 472 Hastie, T. and W. Stuetzle (2012, March). Principal Curves.
- 473 Hinton, G. E. and S. T. Roweis (2002). Stochastic neighbor embedding. In *Advances in neural*
474 *information processing systems*, pp. 833–840.
- 475 Kalisky, T. and S. R. Quake (2011). Single-cell genomics. *Nature methods* 8(4), 311–314.
- 476 Le, Q. V., A. J. Smola, and S. Canu (2005). Heteroscedastic gaussian process regression. In
477 *Proceedings of the 22nd international conference on Machine learning*, pp. 489–496. ACM.
- 478 Macaulay, I. C. and T. Voet (2014, January). Single cell genomics: advances and future per-
479 spectives. *PLoS genetics* 10(1), e1004126.
- 480 Magwene, P. M., P. Lizardi, and J. Kim (2003). Reconstructing the temporal ordering of
481 biological samples using microarray data. *Bioinformatics* 19(7), 842–850.
- 482 Marco, E., R. L. Karp, G. Guo, P. Robson, A. H. Hart, L. Trippa, and G.-C. Yuan (2014, De-
483 cember). Bifurcation analysis of single-cell gene expression data reveals epigenetic landscape.
484 *Proceedings of the National Academy of Sciences of the United States of America* 111(52),
485 E5643–50.
- 486 Moignard, V., S. Woodhouse, L. Haghverdi, A. J. Lilly, Y. Tanaka, A. C. Wilkinson, F. Buett-
487 nner, I. C. Macaulay, W. Jawaid, E. Diamanti, S.-I. Nishikawa, N. Piterman, V. Kouskoff,
488 F. J. Theis, J. Fisher, and B. Göttgens (2015, February). Decoding the regulatory network
489 of early blood development from single-cell gene expression measurements. *Nature Biotech-*
490 *nology* 33(3).
- 491 Qiu, P., A. J. Gentles, and S. K. Plevritis (2011, April). Discovering biological progression
492 underlying microarray samples. *PLoS computational biology* 7(4), e1001123.
- 493 Qiu, P., E. F. Simonds, S. C. Bendall, K. D. Gibbs Jr, R. V. Bruggner, M. D. Linderman,
494 K. Sachs, G. P. Nolan, and S. K. Plevritis (2011). Extracting a cellular hierarchy from
495 high-dimensional cytometry data with spade. *Nature biotechnology* 29(10), 886–891.

- 496 Reid, J. E. and L. Wernisch (2015). Pseudotime estimation: deconfounding single cell time
497 series. *bioRxiv*, 019588.
- 498 Shapiro, E., T. Biezuner, and S. Linnarsson (2013, September). Single-cell sequencing-based
499 technologies will revolutionize whole-organism science. *Nature reviews. Genetics* *14*(9), 618–
500 30.
- 501 Shin, J., D. A. Berg, Y. Zhu, J. Y. Shin, J. Song, M. A. Bonaguidi, G. Enikolopov, D. W.
502 Nauen, K. M. Christian, G.-l. Ming, and H. Song (2015, August). Single-Cell RNA-Seq with
503 Waterfall Reveals Molecular Cascades underlying Adult Neurogenesis. *Cell Stem Cell* *17*(3),
504 360–372.
- 505 Stegle, O., S. a. Teichmann, and J. C. Marioni (2015, January). Computational and analytical
506 challenges in single-cell transcriptomics. *Nature Reviews Genetics* *16*(3), 133–145.
- 507 Titsias, M. and N. Lawrence (2010). Bayesian Gaussian Process Latent Variable Model. *Arti-*
508 *ficial Intelligence* *9*, 844–851.
- 509 Trapnell, C. (2015, Oct). Defining cell types and states with single-cell genomics. *Genome*
510 *Res* *25*(10), 1491–8.
- 511 Trapnell, C., D. Cacchiarelli, J. Grimsby, P. Pokharel, S. Li, M. Morse, N. J. Lennon, K. J.
512 Livak, T. S. Mikkelsen, and J. L. Rinn (2014, April). The dynamics and regulators of cell fate
513 decisions are revealed by pseudotemporal ordering of single cells. *Nature biotechnology* *32*(4),
514 381–6.
- 515 Treutlein, B., D. G. Brownfield, A. R. Wu, N. F. Neff, G. L. Mantalas, F. H. Espinoza, T. J.
516 Desai, M. A. Krasnow, and S. R. Quake (2014). Reconstructing lineage hierarchies of the
517 distal lung epithelium using single-cell rna-seq. *Nature* *509*(7500), 371–375.
- 518 Tsang, J. C., Y. Yu, S. Burke, F. Buettner, C. Wang, A. A. Kolodziejczyk, S. A. Teichmann,
519 L. Lu, and P. Liu (2015). Single-cell transcriptomic reconstruction reveals cell cycle and
520 multi-lineage differentiation defects in bcl11a-deficient hematopoietic stem cells. *Genome*
521 *biology* *16*(1), 1–16.
- 522 Van der Maaten, L. and G. Hinton (2008). Visualizing data using t-sne. *Journal of Machine*
523 *Learning Research* *9*(2579-2605), 85.

- 524 Wills, Q. F. and A. J. Mead (2015). Application of single cell genomics in cancer: Promise and
525 challenges. *Human molecular genetics*, ddv235.
- 526 Xu, C., D. Tao, and C. Xu (2013). A survey on multi-view learning. *arXiv preprint*
527 *arXiv:1304.5634*.

528 **Figures**

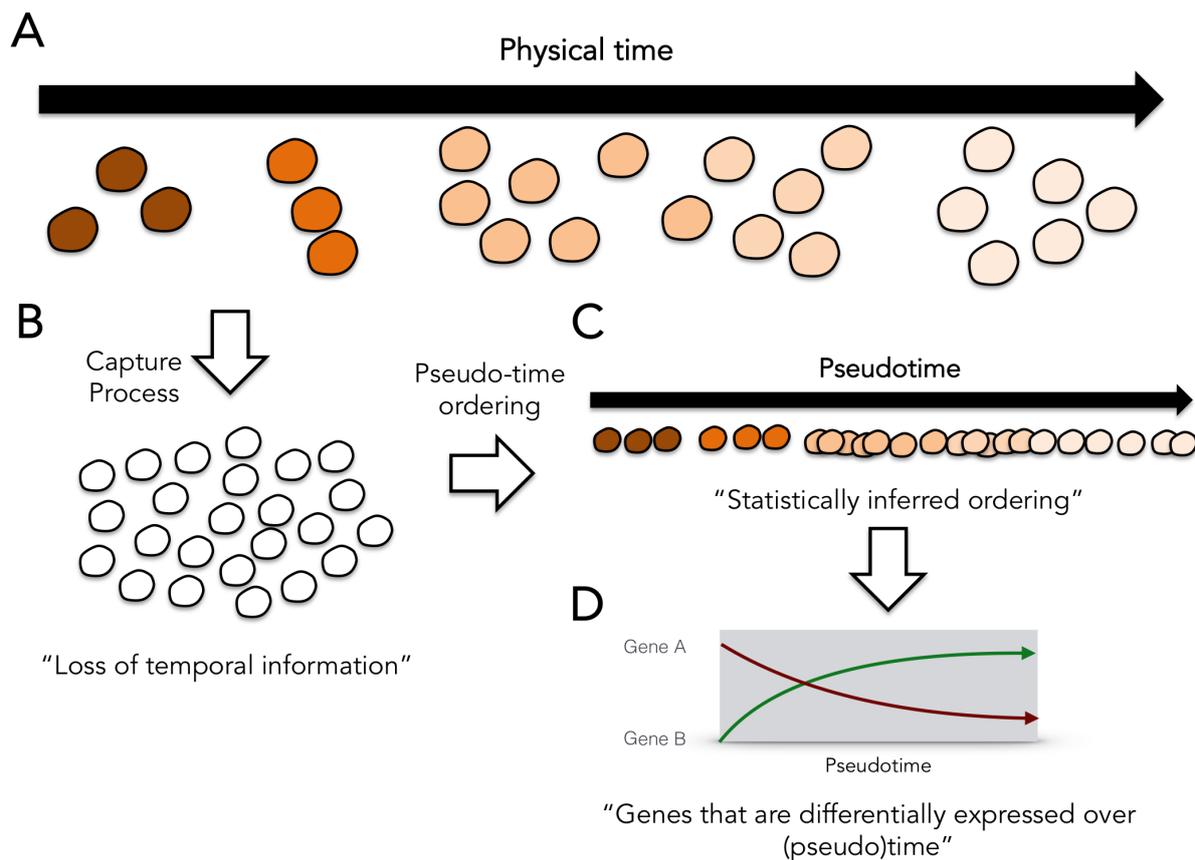


Figure 1: **The single cell pseudotime ordering problem.** (A) Single cells at different stages of a temporal process. (B) The temporal labelling information is lost during single cell capture. (C) Statistical pseudotime ordering algorithms attempt to reconstruct the relative temporal ordering of the cells but cannot fully reproduce physical time. (D) The pseudotime estimates can be used to identify genes that are differentially expressed over (pseudo)time.

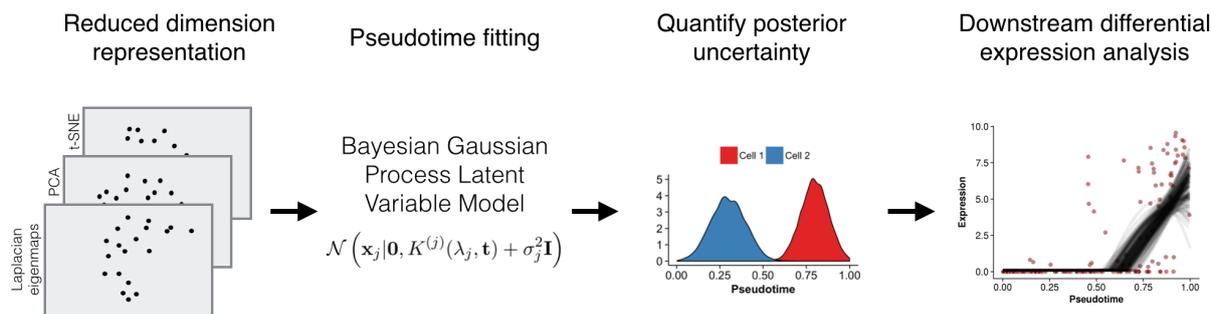


Figure 2: **Workflow for fitting Bayesian Gaussian Process Latent Variable Model pseudotime models.** Reduced-dimension representations of the gene expression data (from Laplacian eigenmaps, PCA and/or t-SNE) are created. The pseudotime can be fitted using one or more low dimensional representations of the data. Posterior samples of pseudotimes are drawn from a Bayesian GPLVM and these are used to obtain alternative pseudotime estimates. Downstream differential analyses can be performed on the posterior samples to characterise the robustness with respect to variation in pseudotime ordering.

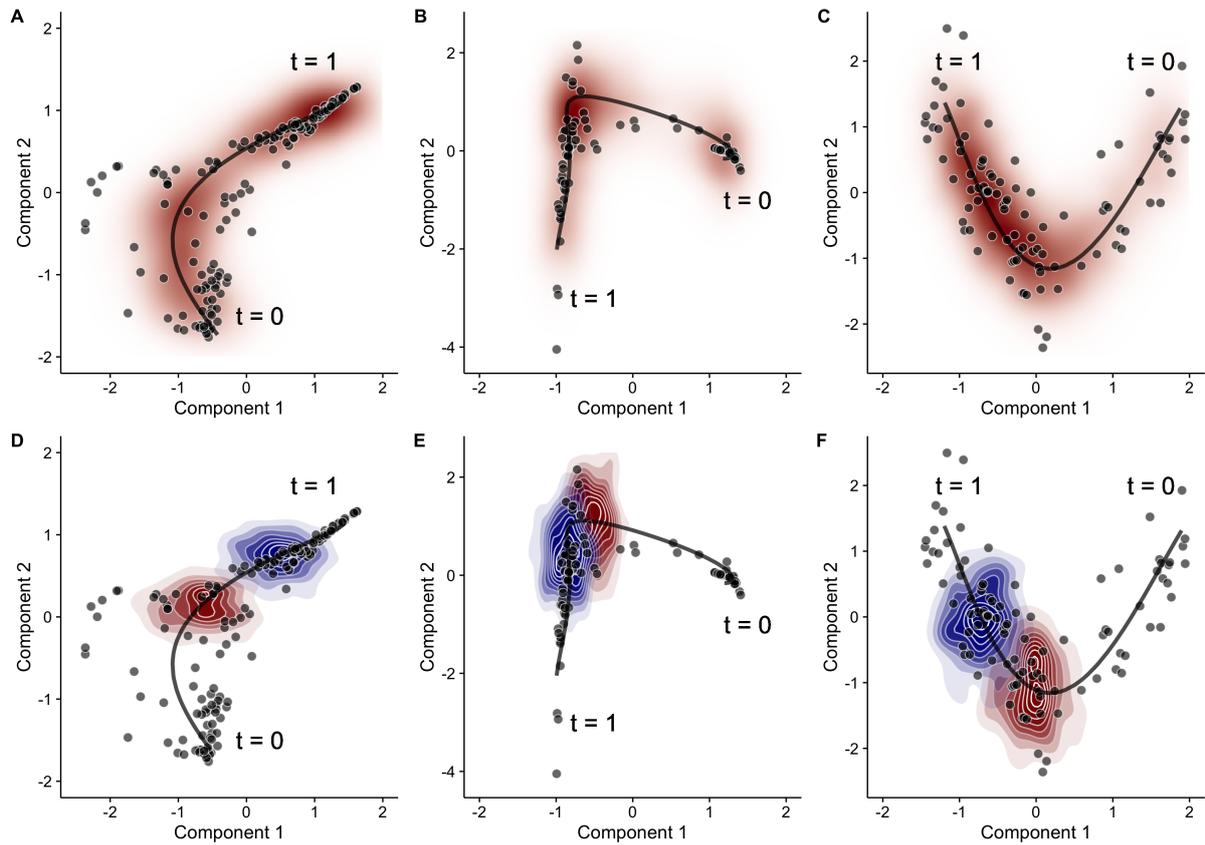


Figure 3: Posterior pseudotime trajectories for three single-cell RNA-seq datasets. Posterior pseudotime trajectories shown in a two-dimensional reduced representation space for (left) a Laplacian eigenmaps representation of Trapnell et al. (2014) Trapnell et al. (2014), (centre) Laplacian eigenmaps representation of Burns et al. (2015) Burns et al. (2015) and (right) PCA representation of Shin et al. (2015) Shin et al. (2015). Each point represents a cell and the black line represents the mean pseudotime trajectory. Plots (A-C) shows the overall posterior predictive data density (red) whilst (D-F) shows the conditional posterior predictive data density for $t = 0.5$ (red) and $t = 0.7$ (blue).

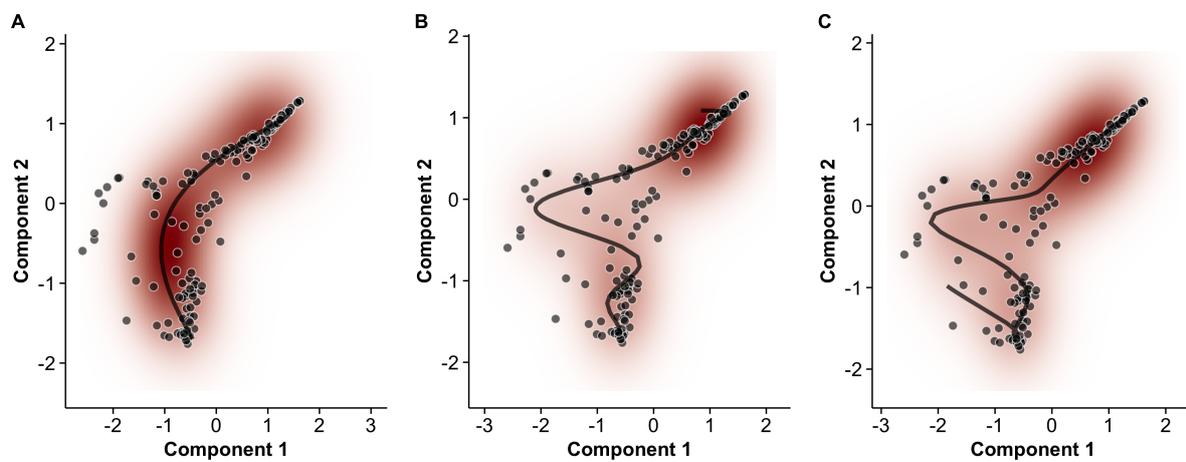


Figure 4: **Effect of prior expectations on pseudotime trajectories.** The prior probability distribution (defined in terms of hyperparameters $(\gamma_\alpha, \gamma_\beta)$ in our model) on the expected smoothness of pseudotime trajectories can fundamentally change the inferred progression path. Examples shown using the data of Trapnell et al. (2014) Trapnell et al. (2014). Red - shows the density of the posterior predictive data distribution. Black - shows the mean pseudotime trajectory. Shrinkage hyperparameters $(\gamma_\alpha, \gamma_\beta)$ of $(30, 5)$, $(5, 1)$ and $(3, 1)$ were used for **A**, **B** and **C** respectively.

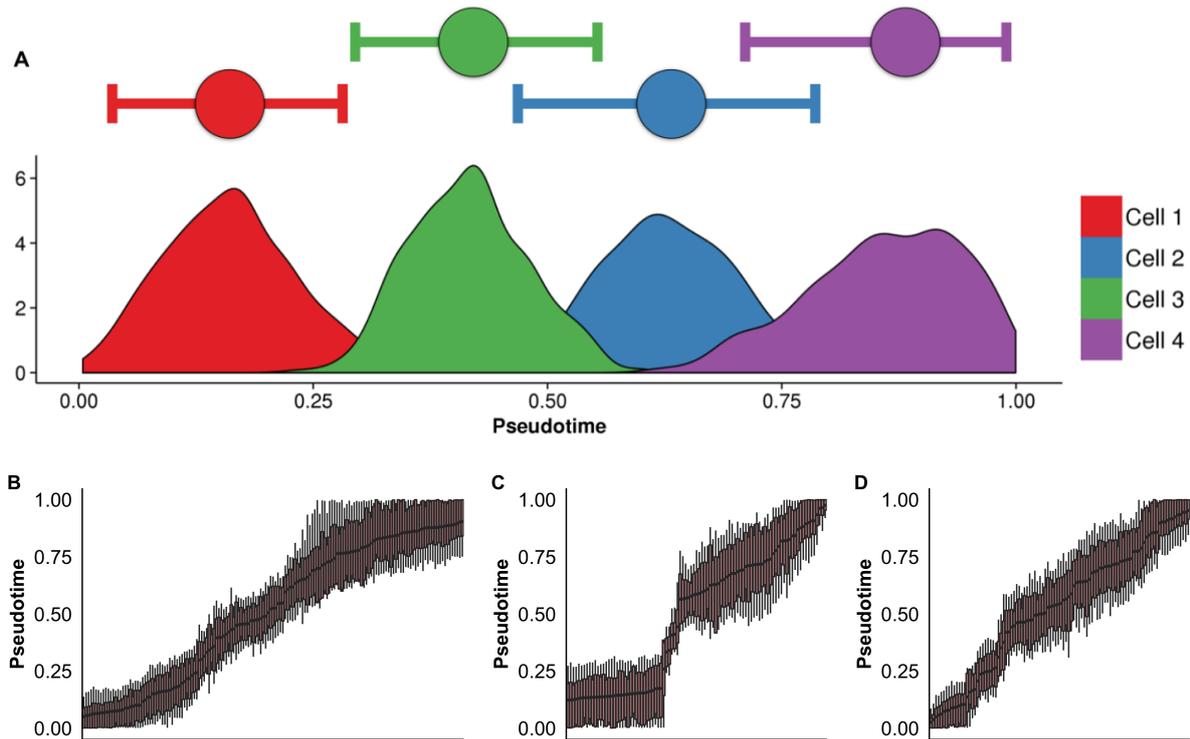


Figure 5: **Posterior uncertainty in pseudotime trajectories.** (A) Posterior uncertainty in pseudotimes for four randomly selected cells from the Trapnell et al. (2014) dataset. Horizontal bars represent the 95% highest probability density (HPD) credible interval (CI), which typically covers around a quarter of the pseudotime trajectory. (B-D) Boxplots showing the posterior uncertainty for each cell from the Trapnell et al. (2014) datasets. The edges of the boxes and tails correspond to the 75% and 95% HPD-CIs respectively.

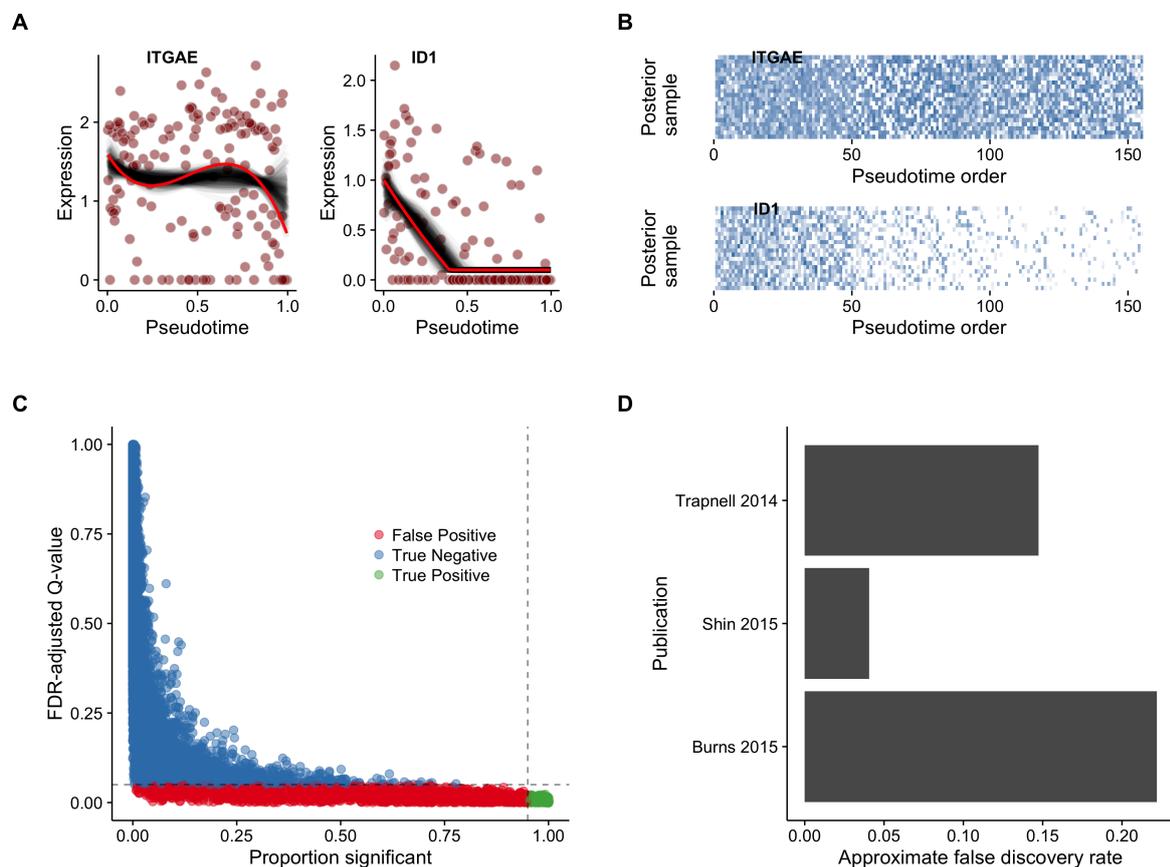


Figure 6: **Approximate FDR for differential expression across pseudotime.** (A) Gene expression plots across pseudotime, with black traces corresponding to models fitted to pseudotime samples while the red trace corresponds to the point (MAP) estimate for two exemplar genes and (B) corresponding posterior pseudotime orderings. (C) Scatter plot of point estimate q-values against proportion significant for all genes (Trapnell dataset). (D) Approximate false discovery rates (AFDR) for three datasets (Trapnell et al. 2014, Shin et al. 2015 and Burns et al. 2015).

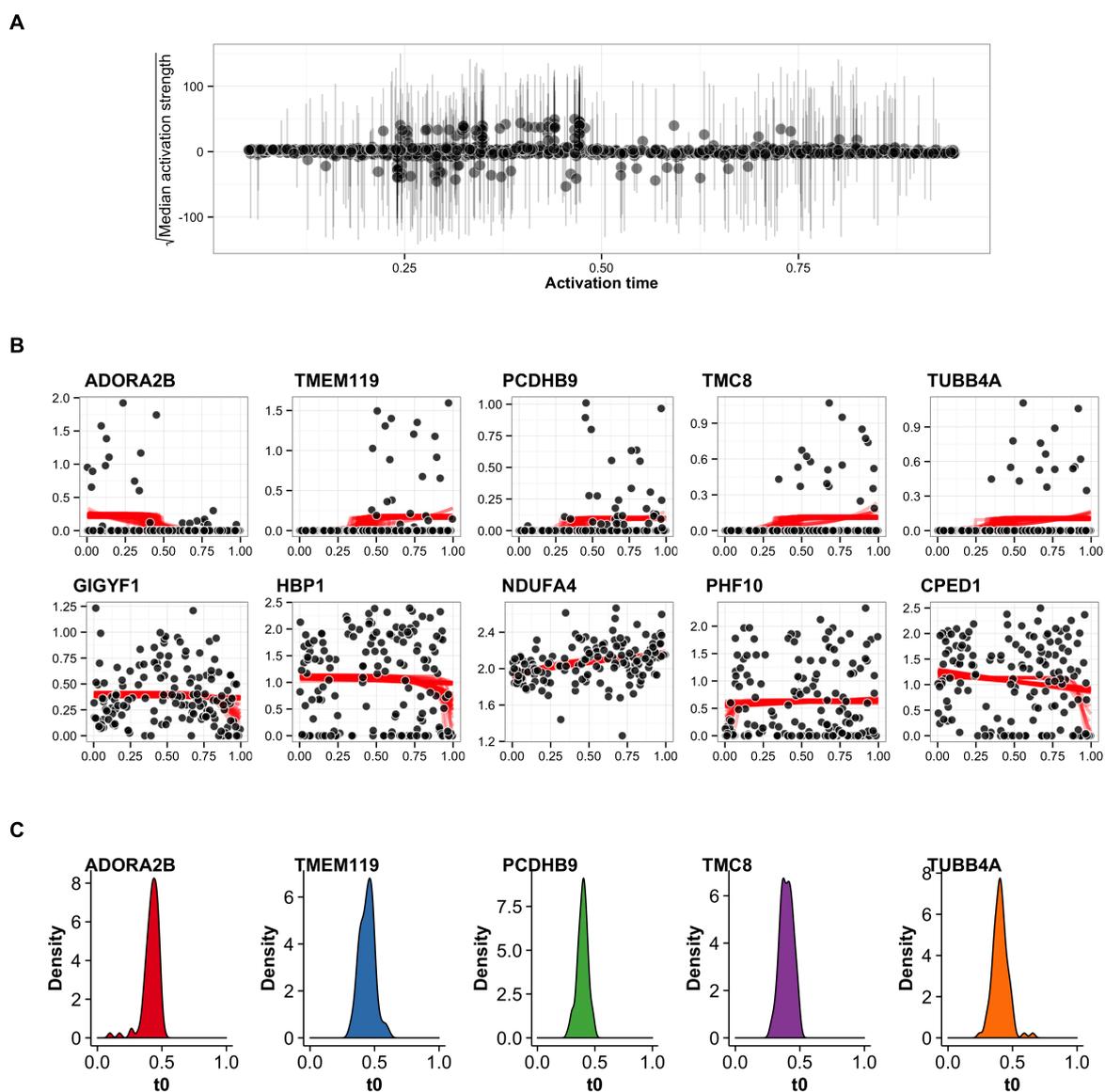


Figure 7: **Robust inference of switch-like behaviour in genes across pseudotime.** (A) The square-root of the median of the activation strength parameter k across all pseudotime samples as a function of activation time t_0 . The error bars show the 95% credible interval, demonstrating that point estimates can severely skew the apparent behaviour of genes and a requirement for a robust Bayesian treatment of gene expression. A distinct population of genes whose median activation strength sits separate from the majority close to the x-axis implies a subset of genes show true switch-like behaviour. (B) Representative examples of genes whose median activation strength is large (top row) compared to small (bottom row). Each black point represents the gene expression of the cell with red lines corresponding to posterior traces of the sigmoidal gene expression model. Genes with a large activation strength show a distinct gene expression pattern compared to those with a small activation strength. (C) A posterior density plot of the activation time for the five genes showing strong activation strength in (B).

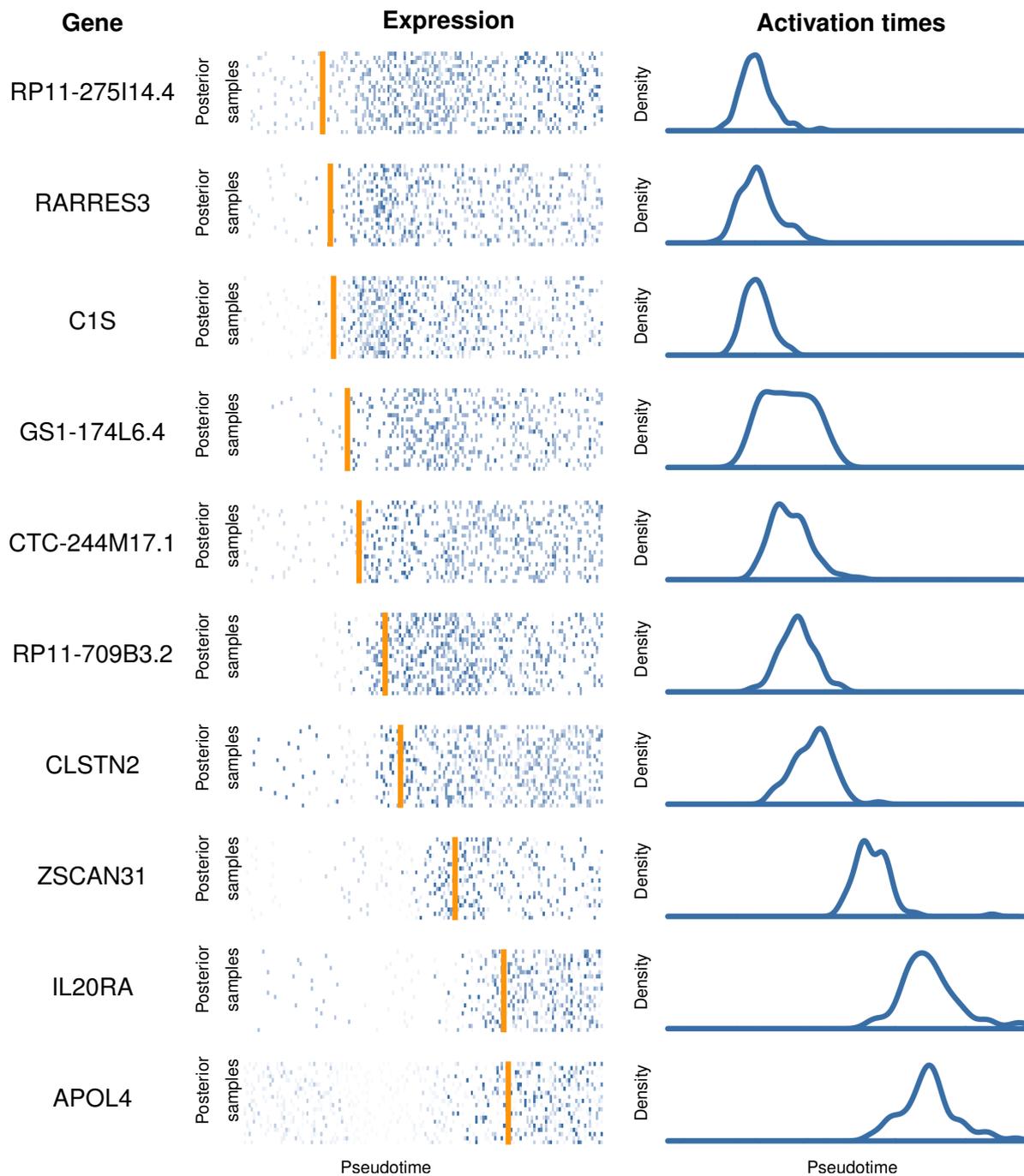


Figure 8: **Identifying pseudotime dependent gene activation behaviour.** Ten selected genes from Trapnell et al. (2014) found using our sigmoidal gene activation model exhibiting a range of activation times. For each gene, we show the expression levels of each cell (centre) where each row corresponds to an ordering according to a different posterior samples of pseudotime. The orange line corresponds to a point estimate of the activation time. The posterior density of the estimated activation time is also shown (right).

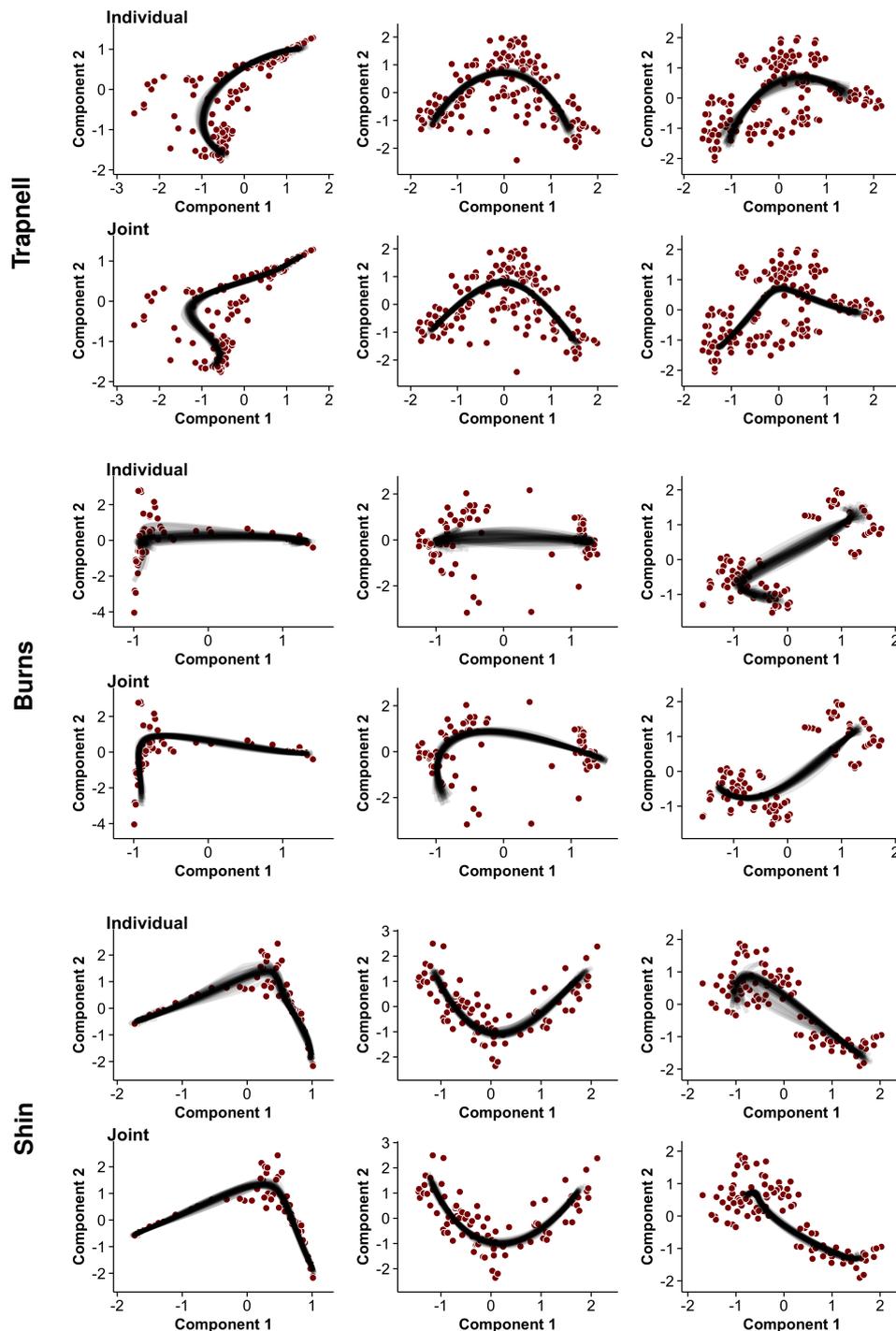


Figure 9: Learning pseudotime from three reduced dimension representations (Laplacian eigenmaps, PCA and t-SNE) of single-cell gene expression data from the three datasets studied (Trapnell, Burns and Shin). For each dataset the left column shows the Laplacian Eigenmaps representation, middle shows PCA and right shows t-SNE (Supplementary Methods). Pseudotime trajectories are fitted either on each representation individually (top row of each dataset) or jointly for all representations (bottom row). It can be seen that trajectory fits are more stable when the joint representations are used. Such analysis allows us to track cellular trajectories across multiple visualisations showing an equivalency of dimensionality reduction algorithms in the context of single-cell RNA-seq data.