

1 **DTFLOW: Inference and Visualization of Single-cell Pseudo-temporal
2 Trajectories Using Diffusion Propagation**

3

4 Jiangyong Wei^{1,a,#1}, Tianshou Zhou^{2,b}, Xinan Zhang^{3,c} and Tianhai Tian^{4,*d}

5

6 ¹College of Science, Huazhong Agricultural University, Wuhan, 430070, China,

7 ²School of Mathematics and Statistics, Sun Yat-sen University, Guangzhou, China,

8 ³School of Mathematics and Statistics, Central China Normal University, Wuhan,
9 430079, China,

10 ⁴School of Mathematical Sciences, Monash University, Melbourne VIC 3800, Australia

11 ^{#1}Previous address: School of Statistics and Mathematics, Zhongnan University of
12 Economics and Law, Wuhan, China

13 * Corresponding author.

14 E-mail: jiangyong.wei@mail.hzau.edu.cn (Wei J), mcszhtsh@mail.sysu.edu.cn (Zhou
15 T), xinanhang@mail.ccnu.edu.cn (Zhang X), tianhai.tian@monash.edu (Tian T)

16

17 **Running title: Wei J et al / Inference of Single-cell Pseudo-time Trajectories**

18

19 ^aORCID: 0000-0002-4608-8226.

20 ^bORCID: 0000-0002-0797-0531.

21 ^cORCID: 0000-0001-8614-7975.

22 ^dORCID: 0000-0001-6191-0209.

23

24 Total word counts (from “Introduction” to “Conclusions” or “Materials and methods”):

25 5662

26 Total figures: 8

27 Total tables: 0

28 Total supplementary figures: 10

29 Total supplementary tables: 0

30 Total supplementary files: 1

31 **ABSTRACT**

32 One of the major challenges in single-cell data analysis is the determination of cellular
33 developmental trajectories using single-cell data. Although substantial studies have
34 been conducted in recent years, more effective methods are still strongly needed to infer
35 the developmental processes accurately. In this work we devise a new method, named
36 DTFLOW, for determining the pseudo-temporal trajectories with multiple branches.
37 This method consists of two major steps: namely a new dimension reduction method
38 (i.e. Bhattacharyya kernel feature decomposition (BKFD)) and a novel approach,
39 named Reverse Searching on kNN Graph (RSKG), to identify the underlying
40 multi-branching processes of cellular differentiations. In BKFD we first establish a
41 stationary distribution for each cell to represent the transition of cellular developmental
42 states based on the random walk with restart algorithm and then propose a new distance
43 metric for calculating pseudo-times of single-cells by introducing the Bhattacharyya
44 kernel matrix. The effectiveness of DTFLOW is rigorously examined by using four
45 single-cell datasets. We compare the efficiency of the new method with two
46 state-of-the-art methods. Simulation results suggest that our proposed method has
47 superior accuracy and strong robustness properties for constructing pseudo-time
48 trajectories. Availability: DTFLOW is implemented in Python and available at
49 <https://github.com/statway/DTFLOW>.

50 **KEYWORDS:** Single-cell heterogeneity; Pseudo-time trajectory; Manifold learning;
51 Bhattacharyya kernel

52

53

54 **Introduction**

55 Recent advances in single-cell technology have provided powerful tools to measure
56 gene expression levels or protein activities of thousands single cells in one experiment.
57 Compared with the traditional experimental studies using bulk samples that average out
58 the responses from a large number of cells in a population, the analysis of cellular
59 aspects at the single-cell level offers promising advantages to investigate the
60 heterogeneity in cellular processes [1]. Since temporal data cannot be collected
61 straightforward, a major step in single-cell studies is to order individual cells according
62 to their progress along the differentiation pathways. The pseudo-temporal data based on
63 the ordered individual cells will ultimately lead to the reconstruction of regulatory
64 networks and cellular differentiation pathways [2].

65 Since the first algorithm Monocle for the pseudo-temporal ordering [3], a number
66 of data-driven computational methods have been developed to define the relative
67 position of each cell during the differentiation process. The methods for inferring
68 pseudo-time trajectories typically consist of two major steps: namely a dimensionality
69 reduction step and a trajectory modelling step. A class of methods based on the graph
70 theory uses the minimum-spanning tree (MST) or shortest path to construct the major
71 structure of trajectories, and then project all single cells onto the major structure to
72 obtain the pseudo-time trajectory. These methods include Wanderlust [4], Wishbone
73 [5], TSCAN [6], Monocle [3], Monocle2 [7], Waterfall [8], SCOUT [9], DensityPath
74 [10] and SoptSC [11]. Another type of algorithms employs the probabilistic model to
75 obtain the major structure of trajectories, such as Gpfates [12], DeLorean [13] and
76 PhenoPath [14]. In addition, other techniques have been used to develop effective
77 methods, include the methods based on differential equations (e.g., SCOUPE [15],
78 Pseudodynamics [16], and PBA [17]), and methods based on the principal curves such
79 as Embeddr [18] and Slingshot [19]. Usually the algorithms based on the graph theory
80 are more efficient, but the accuracy of the inference results is susceptible to the noise in
81 datasets. However, methods using the probabilistic model or differential equations
82 normally need high computational cost. Recently, a number of comparison studies have

83 been conducted to examine the performance of these algorithms [20], and more
84 effective methods can be found in the comprehensive review papers [21-24].

85 Network diffusion, also known as network propagation, has attracted much
86 attention in recent years for identifying disease genes, genetic modules and drug targets
87 [25]. It has also been used for manifold learning and pseudo-time calculation for
88 single-cell data. The nonlinear dimensionality reduction algorithms based on network
89 propagation include DCA, PHATE, etc. Among them, DCA obtains the low
90 dimensional representation of the high-dimensional dataset by minimizing the
91 Kullback-Leibler divergence between the observed diffusion states and
92 parameterized-multinomial logistic distributions [26]; and PHATE generates a Markov
93 transition matrix as the diffusion operator and then embeds the operator with the
94 non-metric multi-dimensional scaling (MDS) approach for the visualization of
95 single-cell datasets [27]. In addition, MAGIC alleviates the noises in single-cell
96 datasets and learns the intrinsic biological structure and gene interaction via data
97 diffusion [28]. Diffusion map, as a random walk approach, has also been used to
98 explore the developmental continuum of cell-fate transitions [29, 30]. The
99 pseudo-temporal trajectory algorithm DPT defines the diffusion pseudo-time distance
100 between two cells using the accumulated Markov transition matrix and determines the
101 ordering of cells based on the distances between a root cell and all other cells [31]. In
102 fact DPT can obtain the pseudo-temporal ordering results before the dimension
103 reduction step, and thus can detect subtle changes of gene expression.

104 Another important issue in single-cell studies is to identify branches in the
105 pseudo-time trajectories in order to explore the different developmental pathways. A
106 number of algorithms have been designed to determine the branches and optimal
107 bifurcation points. Among them, DPT determines the branching trajectories by the
108 correlation versus anti-correlation relationship of the dpt distances between cells [31],
109 and Wishbone identifies two post-bifurcation cell fates using the second eigenvector of
110 a mutual disagreement matrix [5]. In addition, SLICER uses the geodesic entropy
111 metric for branches assignment [32]; TSCAN finds the differentiation structure based
112 on the MST algorithm applied to the cluster centres [6]; and Monocle2 conducts the

113 branching assignment according to the branches of the DDRtree [7]. However, the
114 majority of these branching detection approaches can identify only one bifurcation
115 point. More sophisticated algorithms are strongly needed to determine the branching
116 processes with multiple bifurcation events.

117 In this paper, we propose a new method, named DTFLOW, for inferring the
118 pseudo-time trajectories using single-cell data. This method uses a new manifold
119 learning method, named Bhattacharyya kernel feature decomposition (BKFD), for the
120 visualization of underlying dataset structure. The innovation of this algorithm includes
121 the usage of the random walk with restart (RWR) method to transform each data point
122 into a discrete distribution and the Bhattacharyya kernel to calculate the similarities
123 between cells. Compared with DPT, RWR includes a free parameter that can be used to
124 tune for better inference results. More importantly, we propose a novel distance metric
125 based on the Bhattacharyya distance to preserve the distances along the manifold. In
126 addition, DTFLOW uses an innovative approach named Reverse Searching on kNN
127 Graph (RSKG) to identify the underlying multi-branching processes of cellular
128 differentiation. The effectiveness of our proposed algorithm is rigorously examined by
129 the application to analyse four single-cell datasets.

130 **Methods**

131 This section introduces the proposed DTFLOW for the inference of psuedo-time
132 ordering using single-cell data. **Figure 1** gives the framework of this algorithm and a
133 brief description of the major steps. The detailed steps can be found in Supplementary
134 Algorithm 1.

135 **Markov Adjacency Matrix Construction**

136 We denote N as the number of cells, D the number of genes, and $x_{ij} \in R^{N \times D}$ the
137 gene expression data. For each cell x_i ($i = 1, \dots, N$) with expression data $\{x_{i1}, \dots, x_{iD}\}$,
138 we first find its k most similar neighbors (include itself) through the kNN algorithm
139 based on the pairwise cell-cell Euclidean distance. To obtain a robust diffusion
140 operator, we construct the affinity matrix using the procedure in [33]. We first
141 transform the cell-cell Euclidean distances into the symmetric Gaussian kernel weights

142 to represent the affinities/similarities between cells. The transition probability between
143 any two neighbour cells is defined by the Gaussian kernel

144

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma_i\sigma_j}\right) \quad (1)$$

145 where σ_i and σ_j are the local kernel width of cell x_i and x_j , respectively. The value
146 of σ_i is based on the local density with its distance to the k -th nearest neighbour. If
147 cell x_i is a neighbour of x_j but x_j is not a neighbour of x_i , we let $K(x_j, x_i) =$
148 $K(x_i, x_j)$ to generate a symmetric kernel matrix. If x_i and x_j are not the neighbour of
149 each other, $K(x_i, x_j) = 0$. Then we normalize the kernel as

150

$$\tilde{K}(x_i, x_j) = \frac{K(x_i, x_j)}{\sqrt{Z(x_i)Z(x_j)}},$$

151

$$Z(x_i) = \sum_j K(x_i, x_j). \quad (2)$$

152 Finally, we define the Markov transition probability matrix using the normalization
153 over rows, defined by

154

$$M_{ij} = \frac{\tilde{K}(x_i, x_j)}{\sum_j \tilde{K}(x_i, x_j)}. \quad (3)$$

155 **Bhattacharyya kernel feature decomposition**

156 Random walk with restart (RWR), also known as Personalized PageRank (PPR) [34], is
157 a ranking algorithm that has a number of good properties such as the ability to capture
158 the global structure of a graph and multi-facet relationship between two nodes in the
159 graph [35]. It has a wide range of applications in link prediction, community detection,
160 and anomaly detection. In RWR, each cell is considered as a node, and RWR iteratively
161 calculates the relevance (proximity) score of each node with regard to a given seed node
162 i in the kNN graph. At each step t , RWR select a move from the current node either to
163 its neighbours with probability p , or return to itself with the restart probability $1 - p$.
164 Then the distribution vector satisfies the following equation:

165

$$s_i^t = ps_i^{t-1}M + (1 - p)e_i, \quad 0 < p < 1, \quad (4)$$

166 where s_i^t is an N -dimensional row distribution vector for the visiting probability of
167 each node from the seed node i after t steps, M is the transition matrix defined by
168 (3), and $s_i^0 = e_i$ is a unit direction row vector, which means that the propagation starts
169 from node i . Thus RWR can be regarded as a more general approach and DPT is a
170 special case of RWR (i.e. $p = 1$) (see Supplemental Information). We introduce the
171 stationary distribution for the visiting probabilities by iterating the updating step (4)
172 infinitely (i.e. $t \rightarrow \infty$) until the convergence of the distribution vector. Then the final
173 stationary distribution is

174
$$s_i = s_i^\infty = (1 - p)e_i(I - pM)^{-1}, \quad (5)$$

175 where I is the identity matrix. The diffusion matrix $S = [s_1, \dots, s_N]^T$ is written as

176
$$S = (1 - p)(I - pM)^{-1}, \quad (6)$$

177 The diffusion distribution of each node is a non-vanishing distribution, i.e. $s_{ij} > 0$ and
178 $\sum_{i=1}^N s_{ij} = 1$, where element s_{ij} of matrix S is the similarity score of node j towards
179 node i . Cells with a smaller similarity score should have ordering that is farther away
180 than those having larger similarity scores. Based on the properties of RWR, it is
181 assumed that the exact ordering of all nodes is determined by the diffusion process, and
182 each node is represented by a discrete distribution vector in the matrix S .

183 Suppose that p and q are two discrete probability distributions over the same
184 space $\Omega = \{x_1, \dots, x_N\}$, and let $p_i = p(x_i)$ and $q_i = q(x_i)$. Then Bhattacharyya
185 coefficient (BC) measures the similarity between p and q , given by

186
$$BC(p, q) = \sum_{i=1}^N \sqrt{p_i q_i}. \quad (7)$$

187 Based on the above definition, the Bhattacharyya kernel matrix [36] is defined by

188
$$G = \sqrt{S} \sqrt{S}^T = \left[\langle \sqrt{s_i}, \sqrt{s_j} \rangle \right]_{i,j=1,\dots,N} \quad (8)$$

189 where the square root operation is conducted for every element of the matrix, and $\langle \cdot, \cdot \rangle$
190 is the inner product of two vectors. Apparently the diagonal element of matrix G is the
191 inner product of vector $\sqrt{s_i}$ and has the value of unit one. Because G is a kernel
192 matrix, its eigenvalues are greater than or equal to 0.

193 According to Mercer's Theorem, there exists a kernel function k , satisfying that

194 $k(s_i, s_j) = G_{ij} = \langle \sqrt{s_i}, \sqrt{s_j} \rangle, \forall i, j \in [1, \dots, N].$ (9)

195 Based on the properties of kernel functions, we construct a new kernel k_1 with the
196 mapping operator ϕ , defined by

197 $G_{ij} = k(s_i, s_j) \triangleq e^{k_1(s_i, s_j)} = e^{\langle \phi(s_i), \phi(s_j) \rangle}$ (10)

198 Let $y_i = \phi(s_i)$, the above equation can be written as $\langle y_i, y_j \rangle = \log G_{ij}$. Then we
199 rewrite it in the matrix form

200 $\log G = Y^T Y,$ (11)

201 where $Y = [y_1, \dots, y_N]^T$, and the logarithm operation is applied to every element of
202 matrix G . Equation (11) is a linear transformation, and we perform SVD to obtain
203 decomposition

204 $\log G = V \Sigma V^T,$ (12)

205 where $V \in R^{N \times N}$ is an unitary matrix which satisfies $V^T V = I$, and Σ is a diagonal
206 matrix whose elements are the singular values of matrix $\log G$. For the purpose of
207 visualization, we use the largest d (positive) singular values to represent the major
208 feature of matrix $\log G$. The d low-dimensional embedding of Y , defined by

209 $Y_d = V_d \Sigma_d^{1/2}$ (13)

210 is used to represent the single-cell dataset. Here Σ_d is a matrix that includes only the
211 largest d singular values and V_d is the corresponding vectors. Normally we use
212 $d = 2$ or $d = 3$ for 2-dimensional or 3-dimensional visualization. Then we use the
213 low-dimensional dataset Y_d to visualize the underlying structure of the original
214 high-dimensional single-cell dataset. In this way we design a novel method for
215 dimensionality reduction by introducing the new Bhattacharyya kernel function.

216 Pseudo-time Ordering

217 Note that the Bhattacharyya distance [37], defined by

218 $D_B(i, j) = -\log \langle \sqrt{s_i}, \sqrt{s_j} \rangle = -\log G_{ij},$ (14)

219 is a measure of similarity between two probability distributions. Although it has been
220 widely used in engineering and statistical sciences, this metric is a measure of
221 divergence, and does not satisfy the triangle inequality in the inner product space. Thus,

222 it is not appropriate to use this metric to calculate the distances between cells. In this
223 work, we introduce a new distance metric to measure the distance between two cells.
224 From equation (11), we obtain the distance of two cells i and j as

225
$$\|y_i - y_j\|^2 = \|y_i\|^2 + \|y_j\|^2 - 2\langle y_i, y_j \rangle = -2 \log G_{ij}. \quad (15)$$

226 Since $\|y_i\|^2 = \log G_{ii} = 0$, we define the new distance metric as

227
$$D_{ij} = \|y_i - y_j\| = \sqrt{-2 \log G_{ij}}. \quad (16)$$

228 It can be shown that this new distance satisfies the triangle inequality in the inner
229 product space, which is essentially a kernel distance [38], and this distance metric is
230 defined on the manifold. Thus, we propose to use this metric to calculate the distance
231 between cells.

232 If the root cell x_r is known, the distance between the root cell and the i -th cell is
233 denoted as D_{ri} , and we use the vector $T_r = D_{r,:}$ to denote the pseudo-times of single
234 cells. However, if we do not know the root cell in advance, we can select a group of
235 cells as the root cells based on the sum of distances between a particular cell and all
236 other cells. Suppose we select R cells as the group of root cells, the pseudo-times of
237 single cells is given by $T_r = \sum_{r=1}^R D_{r,:}$. Finally we normalize the pseudo-times to
238 values between 0 and 1, given by

239
$$T = \frac{T_r - \min\{T_r\}}{\max\{T_r\} - \min\{T_r\}}. \quad (17)$$

240 **Reverse-searching in kNN graph for branch detection**

241 According to the Waddington's epigenetic landscape model [39], single cells will not
242 change their fates after differentiation. Based on the constructed kNN graph and
243 pseudo-time of each cell, we next propose a new method for branching detection using
244 the reverse-searching in kNN graph (RSKG). Supplementary Figure S1 shows a brief
245 description of RSKG for identifying multi-branching processes. In this algorithm, n is
246 the minimum number of cells required for forming one sub-branch, T the set of
247 pseudo-times of all cells, and A the set of indices array of the kNN graph of all cells.

248 For the id -th cell, $A[id]$ is the set of its k -nearest neighbours. In addition, we use

249 to store the reverse index ordering based on . We also use a nested list
250 *prop-groups* to store the candidate sub-branches/groups and a nested list *sub-branches*
251 to store the determined sub-branches. Initially these two nested lists are empty. The
252 major steps of this algorithm are described in Supplementary Algorithm 2 together with
253 Figure S1 (see supplemental materials).

254 This algorithm starts from the cell with the largest pseudo-time, whose index is
255 the first element in . We put the indexes of its neighbour in the nested list
256 *prop-groups* as the first candidate group. Then we consider the next element in
257 and its neighbour . If set has intersections with the list ,
258 then extend to the list in *prop-groups*. Otherwise, append the list
259 to *prop-groups* as a separate group. The similar procedure is applied to the
260 following elements with index .

261 For the following cells, if has intersections with two or more candidate
262 groups in *prop-groups*, and if the length of two or more intersected lists reaches ,
263 these lists will be moved from *prop-groups* to *sub-branches* and become a determined
264 branch; otherwise, if the length of the merged list does not reach *n*, merge these lists
265 together as one new list in *prop-groups*. If has intersections with lists in both
266 *prop-groups* and *sub-branches*, and if the length of and the intersected list in
267 *prop-groups* reaches , and the intersected list will be moved from *prop-groups*
268 to *sub-branches* to become a determined branch; otherwise, the elements of and
269 the intersected list in *prop-groups* will be assigned to the branches in *sub-branches* that
270 are closer to them.

271 Our method ensures that all the end points of sub-branches can be connected within
272 the kNN graph, and all pseudo-times of single cells in the previous sub-branch are less
273 than those in the subsequent sub-branch(es) along the development process. Thus this
274 new algorithm can provide more accurate inference results than the existing methods.

275 **Results**

276 In this section, we evaluate the robustness and accuracy of our proposed algorithm
277 DTFLOW for the inference of psuedo-time ordering using single-cell data. We apply
278 DTFLOW to analyse four real scRNA-seq datasets, namely the mouse embryo dataset
279 [40], myeloid progenitor MARS-seq dataset [41], female gonad scRNA-seq dataset
280 [42], and microwell-seq data [43]. This work does not include any work for the
281 pre-processing of experimental data. We use the datasets with the same input (namely
282 the same genes and same single-cells) from the published papers directly.

283 **Mouse embryo single-cell dataset**

284 The high-throughput reverse transcription polymerase chain reaction (RT-PCR) dataset
285 [40] describes the early-stages of the developmental process for mouse embryo. This
286 dataset includes the expression levels of 48 selected genes in 438 single cells at seven
287 different developmental stages, namely from the 1-cell zygote stage to the 64-cell
288 blastocyst stage. We first apply DTFLOW to project the 48-dimensional gene
289 expression data into the two-dimensional feature space by using the BKFD algorithm.

290 **Figure 2A** provides the visualization of single cells at different stages. It clearly reveals
291 the seven developmental stages/labels (i.e. 1-cell stage, 2-cell stage, ..., 64-cell stage),
292 which also validates the effectiveness of our proposed dimensional reduction technique.
293 Since not knowing the root cell in the dataset, we select a cell in the initial time stage,
294 which has the largest sum of distances to all other cells, as the root cell. The
295 differentiation process of single cells is characterized by the calculated pseudo-times in
296 **Figure 2B**, whose values range from zero of the first cell to the maximal value of unit
297 one of the last cell. These results suggest that the pseudo-times of individual cells are
298 recovered successfully.

299 We also test the influence of the minimal cell number n required for forming a
300 sub-branch. When we set a small value (i.e. $n \leq 11$), the individual cells in the lineage
301 process is divided into 5 sub-branches in **Figure 2C**. There are two bifurcation points
302 that separate cells into two distinct sub-branches along the differentiation process.
303 **Figure 2C** shows that the main lineage trajectory contains two major branches and one
304 of them further differentiates into two smaller branches. It also suggests that cell
305 differentiation does not occur in the early stages, but cells in the 32-cell stage

306 differentiate distinctly into trophectoderm (TE) and inner cell mass (ICM).
307 Subsequently, cells in the ICM stage further differentiate into epiblast (EPI) and
308 primitive endoderm (PE) in the 64-cell stage. After the second bifurcating event, the
309 embryo cells are divided into three distinct types: namely TE, PE, and EPI. However, if
310 we use a relatively large value of n ($= 12\sim 112$), the single cells will form only three
311 sub-branches with the first bifurcation event occurred in **Figure 2D**. The second
312 bifurcation event is not identified since the lengths of sub-branches are less than the
313 minimal cell number n . We use red triangles in **Figures 2C** and **2D** to indicate the
314 bifurcation points. Note that the distances between cells in our algorithm are calculated
315 based on the high-dimensional Bhattacharyya kernel matrix. However, the data
316 visualized in **Figure 2** are the low-dimensional data after the application of SVD.

317 In our proposed algorithm, there are two free parameters that should be determined
318 based on the datasets it applied. The first one is the number of closest neighbours k of
319 each data point, which is taken into account for the determination of affinity with
320 classes. To place greater emphasis on the local properties of the manifold structure, a
321 smaller value of k is preferred. Meanwhile, the value of k should also be large enough
322 for the connectivity of the kNN graph. The value of k in BKFD is usually smaller than
323 that in diffusion maps for dimension reduction, which implies that BKFD can capture
324 the local structure of manifold better than diffusion maps. We test different values of k
325 and find that the results are better if $k = 10$, which will be used in this work for
326 analyzing other datasets. The second parameter is the restart probability $1 - p$ that
327 controls the relative influence of both local and global topological structure. To smooth
328 the noise of data, a larger value of p (i.e. a smaller value of $1 - p$) may be preferred.
329 To test the influence of p , we calculate the pseudo-ordering of single cells using
330 different values of p . We use the Kendall rank correlation coefficient of the inference
331 results to compare the accuracy of the algorithms. Since knowing the stage number of
332 each cell in the experimental data, we determine the stage number of each cell in the
333 inferred trajectories and then calculate the Kendall rank correlation coefficient of these
334 two types of stage numbers. An algorithm has better accuracy if the value of this
335 correlation coefficient is larger. Numerical results in **Figure 3A** suggest that the

336 ordering accuracy is better when the value of p is around 0.9. Thus, we use $p = 0.9$ in
337 this work, including the results shown in **Figure 2**.

338 To demonstrate the effectiveness of our proposed algorithm, we compare the
339 performance of DTFLOW with two published state-of-the-art methods, namely DPT
340 and Monocle2 [7]. We use the Python toolkit Scanpy [44] for the implementation of
341 DPT. Supplementary Figures S2C and S3C show the branching detection results of
342 Scanpy and Monocle2, respectively. Figure S2C suggests that Scanpy detects only
343 three groups/sub-branches. It fails to identify the number of terminal states correctly,
344 and also obtains the wrong location of bifurcation point. Although Monocle2 identifies
345 three types of the terminal cells correctly in Figure S3C, it does not reveal the
346 intermediate state between state 1 and states 3/4 (i.e. the ICM stage) using the
347 dimensional reduction method DDRTree. For this dataset, we use Kendall rank
348 correlation coefficient to compare the accuracy of these three algorithms. The
349 calculated Kendall rank correlations are 0.862, 0.796, 0.761 for DTFLOW, Scanpy and
350 Monocle2, respectively, which suggests that our proposed method has better accuracy
351 than the two published methods.

352 Supplementary Figure S4 shows the expression levels of two genes, namely *Gata3*
353 and *Sox2*, based on the inferred pseudo-times using the three methods, which are
354 consistent with the results of visualization. It shows that only DTFLOW detects the
355 ICM stage correctly. The intermediate states of cell development in Monocle2 are not
356 revealed properly possibly because the differences between clusters are amplified by
357 the DDRTree method with the cluster centroids. In addition, DPT uses diffusion maps
358 for dimensional reduction, which may not be sensitive enough to the noise in dataset.

359 We further conduct the robustness analysis of each algorithm. We first use the whole
360 dataset to infer a trajectory and determine the position of each cell in this trajectory.
361 Then we sample part of the cells from the whole dataset and use the same algorithm to
362 determine the trajectory of cells in the sub-dataset. We calculate the Spearman rank
363 correlation coefficient between the positions of subset cells in the trajectory of the
364 whole dataset and those of the sub-dataset. An algorithm is more robust if the value of
365 the correlation coefficient is larger. We conduct 50 tests to measure the robustness

366 properties of these three algorithms. In each test we randomly sample 90% of cells (i.e.
367 394 cells) from the dataset and then calculate the Spearman rank correlation coefficient
368 of the pseudo-time ordering of the sub-dataset. Then we use the mean and standard
369 deviation of the correlation coefficient based on these 50 test results to measure the
370 robustness property of algorithms. **Figure 3B** shows that the robustness properties of
371 DTFOLW and Scanpy are better than that of Monocle2. In addition, the variance of
372 correlation coefficients obtained by DTFLOW is smaller than that of Scanpy, which
373 suggest that the performance of DTFLOW is more stable than the two published
374 methods. To examine the influence of the sample size, we conduct further robustness
375 test by randomly sampling 80% of cells (i.e. 350 cells) from the dataset. The Spearman
376 rank correlation coefficients in Supplementary Figure S5 are consistent with those
377 shown in **Figure 3B**.

378 **Mouse myeloid progenitors dataset**

379 We next apply our proposed algorithm to analyze the mouse myeloid progenitor
380 MARS-seq dataset that contains 2730 single cells and 3451 informative genes [41].
381 Note that 10 genes with corrupted names are removed from our analysis based on the
382 pre-processing of Scanpy. In this experimental study, 19 distinct, transcriptionally
383 homogeneous progenitor types/clusters have been identified through an EM-based
384 clustering approach. Among these clusters, clusters 1~6 represent Ery (erythroid
385 lineage progenitors) subpopulations, clusters 7~10 represent CMP (common myeloid
386 progenitors) subpopulations, cluster 11 is for the DC (dendritic cell) fate, clusters
387 12~18 correspond to GMP (granulocyte/macrophage progenitors) subpopulations, and
388 cluster 19 is the lymphoid lineage progenitors (outlier class) with only 31 cells. We
389 apply DTFLOW to project this dataset into the two-dimensional feature space. **Figure**
390 **4A** elucidates that CMP and its progenitors (namely Ery and GMP) are nearly separated
391 in three different regions, while DC and lymphoid cells deviate away from the main
392 differentiation progression process.

393 To reveal the cellular differentiation process, we select the same cell in [31] as the
394 root cell (i.e. the 840-th cell in cluster 8). **Figure 4B** demonstrates the pseudo-time
395 ordering results from the themyeloid progenitor stage. **Figure 4C** and **Figure 4D** show

396 different branching detection results that are determined by a smaller cell number of
397 $n = 6\sim 22$ and a relatively larger number of $n = 23\sim 124$ for forming sub-branches,
398 respectively. The first sub-branch in **Figure 4C** contains only a small number of cells.
399 DTFLOW ensures that the pseudo-time of each cell in the initial branch is less than that
400 of any other cells in the following sub-branches. Then cells differentiate into three
401 different terminal branches. Sub-branch two corresponds to the erythroid evolutionary
402 branch, sub-branch three is formed by cells within clusters 11 and 19, and sub-branch
403 four corresponds to the GMP branch. However, when a larger value of n is used,
404 sub-branches two and three merge together and form one large sub-branch in **Figure**
405 **4D**.

406 We next compare the branching detection results of DTFLOW, Scanpy and
407 Monocle2. Supplementary Figures S6C and S7C show the branching inference results
408 of Scanpy and Monocle2, respectively. Figure S6C suggests that Scanpy is also able to
409 identify 4 sub-branches. However, the root cell identified by Scanpy is not in the initial
410 group, which is unreasonable for the developmental process. Although Monocle2
411 successfully estimates 12 states in Figure S7C, which is consistent with the
412 experimental observation, it is difficult to analysis the changes of gene expression over
413 time based on this large branch number.

414 We then carry out robustness analysis of the three methods. For each method, we
415 randomly sample 2500 single cells from 2730 cells and then use the same methods to
416 infer the pseudo-times of the selected cells. Then we compare the pseudo-times of cells
417 in the sampled set with that of the corresponding cells in the whole dataset by using the
418 Spearman rank correlation coefficient. We conduct 50 repeated tests to measure the
419 robustness property of each method. **Figure 5A** shows that the robustness property of
420 DTFOLW and Scanpy is better than Monocle2. In addition, the variance of correlation
421 coefficients obtained by DTFLOW is smaller than that of Scanpy, which suggest that
422 the performance of DTFLOW is more stable than the two published methods.

423 Note that the gene expression levels in this dataset are not continuous and the three
424 terminal branches have different lengths. To illustrate this, Supplementary Figure S8
425 presents the expression visualization of two marker genes, namely *Elane* and *Klf1*.

426 These marker genes show similar significance in different branches for different
427 dimensionality reduction algorithms. Figure S8A shows that expression levels of gene
428 Elane are essential for the GMP process while gene Klf1 increase gradually on the
429 erythroid branch. The expression trends of these marker genes are different along the
430 constructed trajectories, which provides important information for developing gene
431 regulatory networks.

432 **Mouse female gonad scRNA-seq dataset**

433 The third dataset is the mouse female gonad scRNA-seq dataset that contains 563 single
434 cells and 822 genes at six developmental stages of gonadal differentiation (namely
435 E10.5, E11.5, E12.5, E13.5, E16.5, and post-natal day 6) [42]. We project this dataset
436 into the 3-dimensional space using our proposed algorithm BKFD. In **Figure 6A** all the
437 cells are presented by different colors for different embryonic stages. It shows that the
438 early progenitor cells subsequently lead to the differentiation to the granulosa cell
439 lineage and stromal progenitor cell lineage in around stages E11.5-E12.5. **Figure 6B**
440 gives the ordered pseudo-times of different single cells and **Figures 6C** presents the
441 inferred three sub-branches by our proposed DTFLOW.

442 We first compare the pseudo-time ordering accuracy of DTFLOW with Scanpy
443 and Monocle2. Figures S9 and S10 show the analysis results of Scanpy and Monocle2
444 for this dataset, respectively. The calculated Kendall rank correlations are 0.761, 0.702,
445 0.569 for DTFLOW, Scanpy and Monocle2, correspondingly. We also compare the
446 robustness properties of DTFLOW with those of Scanpy and Monocle2. We conduct 50
447 tests to measure the robustness properties of these three algorithms. In each test we
448 randomly sample 90% of cells (i.e. ~740 cells) from the dataset and find the
449 pseudo-times of these cells. Then we compare the pseudo-times of these cells in the
450 sampled set with that of the corresponding cells in the whole dataset by using the
451 Spearman rank correlation. **Figure 5B** gives the robustness property of these three
452 methods. Numerical results suggest that that the robustness properties of DTFLOW are
453 better than those of Scanpy and Monocle2.

454 **Comparison of dimensional reduction algorithms**

455 Now we compare our dimensional reduction algorithm BKFD with several popular and
456 widely used methods PCA, tSNE [45] and UMAP [46]. **Figure 7** shows the
457 visualization results of the mouse embryo single-cell dataset [40]. Based on the idea of
458 diffusion propagation, BKFD (**Figure 7A**) represents the cellular development
459 reasonably. **Figure 7B** suggests that PCA cannot distinguish differentiation stages very
460 well. Although tSNE (**Figure 7C**) and UMAP (**Figure 7D**) can separate
461 different cell types clearly for this dataset, the distance intervals of different cells types
462 are relatively large, which cannot be used to indicate cellular developmental processes
463 properly.

464 After the successful applications of our proposed algorithm to three relatively small
465 datasets, we next show the effectiveness of our method to a large dataset. We apply our
466 proposed method BKFD to a microwell-seq data that contains 51252 cells and 25912
467 genes [43]. After the data pre-processing, the dataset is reduced to 40210 cells with 100
468 approximate principal components [46]. Based on this dataset, we use four methods for
469 dimensional reduction. **Figure 8** shows the visualization results of this dataset with
470 eight major cell clusters. It suggests that BKFD can capture the developmental
471 trajectories in a better way in **Figure 8A**. In addition, tSNE and UMAP can also
472 distinguish different cell types clearly in **Figures 8C and 8D**, respectively. However,
473 PCA cannot show good visualization results with distinguishable cell clusters for this
474 dataset in **Figure 8B**.

475

476 **Discussion**

477 In this paper we propose a new method DTFLOW for conducting pseudo-time analysis
478 of single-cell data. This method has two major steps: namely a new dimension
479 reduction method BKFD and a novel approach RSKG to identify the underlying
480 multi-branching processes of cellular differentiation. In BKFD we first establish a
481 stationary distribution for each cell to represent the transition of cellular developmental
482 states based on the RWR algorithm, and then propose a new Bhattacharyya kernel
483 matrix to measure the distances between the distributions obtained by RWR. We use

484 this novel distance metric to calculate the pseudo-time distances between single cells
485 before dimension reduction. Thus, our method can reduce the information loss in data
486 processing and increase the inference accuracy. The combination of RWR and the
487 Bhattacharyya kernel matrix shows great power to explore the global structure of the
488 developmental processes using single-cell datasets. In addition, we design the RSKG
489 algorithm to identify the multi-branching of cellular processes. Four datasets are used
490 to compare the accuracy, robustness and branch detection of the proposed algorithm
491 with two popular published methods. Inference results suggest that our proposed
492 method is more accurate and robust than the published algorithms for developing the
493 pseudo-time trajectories of single cells.

494 The RWR algorithm is a popular method to estimate the global similarity between
495 a particular node point with other node points in the graph structure. We use this
496 method to transform the data of each node point to a stationary discrete distribution.
497 Thus, the input space becomes a set of distributions over the same space. The
498 performance of our proposed algorithm is affected by the choice of Gaussian kernel
499 function, the number of closest neighbours k and the restart probability $1 - p$ in the
500 RWR algorithm. Although we have examined the performance of the proposed
501 algorithm by using four datasets, the values of these parameters may vary from dataset
502 to dataset. In addition, BKFD uses the same restart probability for all the nodes, and this
503 may limit the effectiveness of random walk [47]. It is still a challenge to express each
504 cell by a distribution vector in a better way, which needs to be studied in the future.

505 The continuously topological structure of cellular developmental processes can be
506 analyzed by using the nearest-neighbour graph, which lays the basis of the DTFLOW
507 algorithm. The kNN graph describes the similarities between a cell and its neighbour
508 cells, and has been used twice in the proposed method, namely the definition of
509 transition probability matrix, which leads to the low-dimensional visualization via the
510 Bhattacharyya kernel matrix, and the determination of branching processes in the
511 RSKG algorithm. The new branch detection algorithm identifies the sub-branches
512 through reverse searching on the sequence of indices ordering and provides biological
513 insights into developmental bifurcations. It can ensure that the sub-branches can be

514 connected through the kNN graph, which in turn also verifies its consistency with the
515 pseudo-time inference and visualization results of BKFD.

516 Scalability is an important issue for the implementation of algorithms. Our
517 algorithm is connected to the dataset size (i.e. the number of cells) in two major steps:
518 the computation of matrix S by finding the inverse of matrix ($I-pM$) in (6), and the SVD
519 computation in (12). In this study we consider four datasets with cell numbers 438,
520 2730, 563 and 40210, respectively. The computational time of our algorithm is 0.224
521 second, 11.65s, 0.246s, and 3108.35s on a Lenovo ThinkPad P53 mobile workstation
522 with 2.6GHz CPU for these four datasets, which is close to the computing time of other
523 algorithms. In addition, the computing time is in the order of $O(N^2)$ in terms of the
524 dataset size N , which suggests our program is scalable to dataset size.

525 In summary, the proposed algorithm DTFLOW provides a new framework for
526 inferring the pseudo-time of single cells. Numerical results suggest that it is a power
527 tool for the inference and visualization of cellular developmental trajectories. Potential
528 future work may include the selection of parameters in the proposed method in order to
529 achieve optimal performance in single-cell data analysis.

530

531 **Data Accessibility**

532 The first dataset [40] is downloaded from [https://github.com/gcyuan/SCUBA/tree/
533 master/sample_data/guo2010](https://github.com/gcyuan/SCUBA/tree/master/sample_data/guo2010). The second dataset[41] is given by [https://github.com/
534 theislab/scanpy_usage/tree/master/170502_paul15](https://github.com/theislab/scanpy_usage/tree/master/170502_paul15). The third dataset[42] is download
535 from <https://github.com/IStevant/XX-XY-mouse-gonad-scRNA-seq>. The fourth
536 dataset [43] is download from https://github.com/ebecht/DR_benchmark.

537

538 **Authors' contributions**

539 JW designed the programs, analysed the data and drafted the manuscript. TZ and XZ
540 analysed and interpreted the data as well as revised manuscript critically. XZ helped
541 programming in revision. TT conceived of the study, participated in the programming
542 and drafting the manuscript. All authors read and approved the final manuscript.

543

544 **Competing interests**

545 The authors have declared no competing interests.

546

547 **Acknowledgments**

548 This work was supported by National Natural Science Foundation of China (Grant No.
549 11571368, 11931019, 11775314, 11871238) and the Fundamental Research Funds for
550 the Central Universities (Grant No. 2662019QD031).

551

552 **References**

- 553 [1] Farrell JA, Wang Y, Riesenfeld SJ, Shekhar K, Regev A, Schier AF. Single-cell
554 reconstruction of developmental trajectories during zebrafish embryogenesis. *Science*
555 2018; 360(6392): eaar3131.
- 556 [2] Laurenti E, Göttgens B. From haematopoietic stem cells to complex differentiation
557 landscapes. *Nature* 2018; 553(7689): 418.
- 558 [3] Trapnell C, Cacchiarelli D, Grimsby J, Pokharel P, Li S, Morse M, et al. The
559 dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering
560 of single cells. *Nature biotechnology* 2014; 32(4): 381–386.
- 561 [4] Bendall SC, Davis KL, Amir EaD, Tadmor MD, Simonds EF, Chen TJ, et al.
562 Single-cell trajectory detection uncovers progression and regulatory coordination in
563 human B cell development. *Cell* 2014; 157(3): 714–725.
- 564 [5] Setty M, Tadmor MD, Reich-Zeliger S, Angel O, Salame TM, Kathail P, et al.
565 Wishbone identifies bifurcating developmental trajectories from single-cell data.
566 *Nature biotechnology* 2016; 34(6): 637–645.
- 567 [6] Ji Z, Ji H. TSCAN: Pseudo-time reconstruction and evaluation in single-cell
568 RNA-seq analysis. *Nucleic acids research* 2016; 44(13): e117–e117.
- 569 [7] Qiu X, Mao Q, Tang Y, Wang L, Chawla R, Pliner HA, et al. Reversed graph
570 embedding resolves complex single-cell trajectories. *Nature methods* 2017; 14(10):
571 979.
- 572 [8] Shin J, Berg DA, Zhu Y, Shin JY, Song J, Bonaguidi MA, et al. Single-cell
573 RNA-seq with waterfall reveals molecular cascades underlying adult neurogenesis.
574 *Cell stem cell* 2015; 17(3): 360–372.
- 575 [9] Wei J, Zhou T, Zhang X, Tian T. SCOUT: A new algorithm for the inference of
576 pseudo-time trajectory using single-cell data. *Computational Biology and Chemistry*
577 2019; 80: 111–120.
- 578 [10] Chen Z, An S, Bai X, Gong F, Ma L, Wan L. DensityPath: an algorithm to visualize
579 and reconstruct cell state-transition path on density landscape for single-cell RNA
580 sequencing data. *Bioinformatics* 2019; 35(15): 2593–2601.

- 581 [11]Wang S, Karikomi M, MacLean AL, Nie Q. Cell lineage and communication
582 network inference via optimization for single-cell transcriptomics. Nucleic Acids
583 Research 2019 03; 47(11): e66–e66.
- 584 [12]Lönnberg T, Svensson V, James KR, Fernandez-Ruiz D, Sebina I, Montandon R,
585 et al. Single-cell RNA-seq and computational analysis using temporal mixture
586 modelling resolves Th1/Tfh fate bifurcation in malaria. Science immunology 2017;
587 2(9): 2192.
- 588 [13]Reid JE, Wernisch L. Pseudotime estimation: deconfounding single cell time
589 series. Bioinformatics 2016; 32(19): 2973–2980.
- 590 [14]Campbell K, Yau C. Uncovering pseudotemporal trajectories with covariates from
591 single cell and bulk expression data. Nature Communications 2018; 9(1): 2442.
- 592 [15]Matsumoto H, Kiryu H. SCOUPE: a probabilistic model based on the
593 Ornstein–Uhlenbeck process to analyze single-cell expression data during
594 differentiation. BMC bioinformatics 2016; 17(1): 232.
- 595 [16]Fischer DS, Fiedler AK, Kernfeld E, Genga RM, Hasenauer J, Maehr R, et al.
596 Beyond pseudotime: Following T-cell maturation in single-cell RNAseq time series.
597 bioRxiv 2017; p. 219188.
- 598 [17]Weinreb C, Wolock S, Tusi BK, Socolovsky M, Klein AM. Fundamental limits on
599 dynamic inference from single-cell snapshots. Proceedings of the National Academy of
600 Sciences 2018; 115(10): 2467–2476.
- 601 [18]Campbell K, Ponting CP, Webber C. Laplacian eigenmaps and principal curves for
602 high resolution pseudotemporal ordering of single-cell RNA-seq profiles. bioRxiv
603 2015; p. 027219.
- 604 [19]Street K, Risso D, Fletcher RB, Das D, Ngai J, Yosef N, et al. Slingshot: Cell
605 lineage and pseudotime inference for single-cell transcriptomics. BMC genomics 2018;
606 19(1): 477.
- 607 [20]Saelens W, Cannoodt R, Todorov H, Saeys Y. A comparison of single-cell
608 trajectory inference methods. Nature biotechnology 2019; 37(5): 547–554.
- 609 [21]Kester L, van Oudenaarden A. Single-Cell Transcriptomics Meets Lineage
610 Tracing. Cell Stem Cell 2018; 23(2): 166–179.

- 611 [22]Chen J, Renia L, Ginhoux F. Constructing cell lineages from single-cell
612 transcriptomes. *Molecular aspects of medicine* 2018; 59: 95–113.
- 613 [23]Cannoodt R, Saelens W, Saeys Y. Computational methods for trajectory inference
614 from single-cell transcriptomics. *European journal of immunology* 2016; 46(11):
615 2496–2506.
- 616 [24]Tritschler S, Büttner M, Fischer DS, Lange M, Bergen V, Lickert H, et al.
617 Concepts and limitations for learning developmental trajectories from single cell
618 genomics. *Development* 2019; 146(12): 170506.
- 619 [25]Cowen L, Ideker T, Raphael BJ, Sharan R. Network propagation: a universal
620 amplifier of genetic associations. *Nature Reviews Genetics* 2017; 18(9): 551.
- 621 [26]Wang S, Cho H, Zhai C, Berger B, Peng J. Exploiting ontology graph for
622 predicting sparsely annotated gene function. *Bioinformatics* 2015; 31(12): 357–364.
- 623 [27]Moon KR, van Dijk D, Wang Z, Gigante S, Burkhardt DB, Chen WS, et al.
624 Visualizing structure and transitions in high-dimensional biological data. *Nature
625 Biotechnology* 2019; 37(12): 1482–1492.
- 626 [28]Van Dijk D, Sharma R, Nainys J, Yim K, Kathail P, Carr AJ, et al. Recovering
627 gene interactions from single-cell data using data diffusion. *Cell* 2018; 174(3):
628 716–729.
- 629 [29]Angerer P, Haghverdi L, Buttner M, Theis FJ, Marr C, Buettner F. destiny:
630 diffusion maps for large-scale single-cell data in R. *Bioinformatics* 2016; 32(8):
631 1241–1243.
- 632 [30]Haghverdi L, Buettner F, Theis FJ. Diffusion maps for high-dimensional
633 single-cell analysis of differentiation data. *Bioinformatics*. 2015; 31(18): 2989–2998.
- 634 [31]Haghverdi L, Buettner M, Wolf FA, Buettner F, Theis FJ. Diffusion pseudotime
635 robustly reconstructs lineage branching. *Nature methods* 2016; 13(10): 845–848.
- 636 [32]Welch JD, Hartemink AJ, Prins JF. SLICER: inferring branched, nonlinear cellular
637 trajectories from single cell RNA-seq data. *Genome biology* 2016; 17(1): 106.
- 638 [33]Rohrdanz MA, Zheng W, Maggioni M, Clementi C. Determination of reaction
639 coordinates via locally scaled diffusion map. *The Journal of chemical physics* 2011;
640 134(12): 03B624.

- 641 [34] Tong H, Faloutsos C, Pan JY. Random walk with restart: fast solutions and
642 applications. *Knowledge and Information Systems* 2008; 14(3): 327–346.
- 643 [35] Yu AW, Mamoulis N, Su H. Reverse top-k search using random walk with restart.
644 Proceedings of the VLDB Endowment 2014; 7(5): 401–412.
- 645 [36] Jebara T, Kondor R. Bhattacharyya and expected likelihood kernels. In: *Learning*
646 theory and kernel machines. Springer 2003. p. 57–71.
- 647 [37] Kailath T. The divergence and Bhattacharyya distance measures in signal
648 selection. *IEEE transactions on communication technology*. 1967; 15(1): 52–60.
- 649 [38] Phillips JM, Venkatasubramanian S. A gentle introduction to the kernel distance.
650 arXiv preprint arXiv: 11031625. 2011.
- 651 [39] Ferrell Jr JE. Bistability, bifurcations, and Waddington’s epigenetic landscape.
652 *Current biology* 2012; 22(11): 458–466.
- 653 [40] Guo G, Huss M, Tong GQ, Wang C, Sun LL, Clarke ND, et al. Resolution of cell
654 fate decisions revealed by single-cell gene expression analysis from zygote to
655 blastocyst. *Developmental cell* 2010; 18(4): 675–685.
- 656 [41] Paul F, Arkin Y, Giladi A, Jaitin DA, Kenigsberg E, Keren-Shaul H, et al.
657 Transcriptional heterogeneity and lineage commitment in myeloid progenitors. *Cell*
658 2015; 163(7): 1663–1677.
- 659 [42] Stévant I, Kühne F, Greenfield A, Chaboissier MC, Dermitzakis ET, Nef S.
660 Dissecting Cell Lineage Specification and Sex Fate Determination in Gonadal Somatic
661 Cells Using Single-Cell Transcriptomics. *Cell reports* 2019; 26(12): 3272–3283.
- 662 [43] Han X, Wang R, Zhou Y, Fei L, Sun H, Lai S, et al. Mapping the mouse cell atlas
663 by microwell-seq. *Cell* 2018; 172(5): 1091–1107.
- 664 [44] Wolf FA, Angerer P, Theis FJ. SCANPY: large-scale single-cell gene expression
665 data analysis. *Genome biology* 2018; 19(1): 15.
- 666 [45] Van Der Maaten L. Accelerating t-SNE using tree-based algorithms. *The Journal*
667 of Machine Learning Research 2014; 15(1): 3221–3245.
- 668 [46] Becht E, McInnes L, Healy J, Dutertre CA, Kwok IW, Ng LG, et al.
669 Dimensionality reduction for visualizing single-cell data using UMAP. *Nature*
670 *biotechnology* 2019; 37(1): 38.

671 [47]Jin W, Jung J, Kang U. Supervised and extended restart in random walks for
672 ranking and link prediction in networks. PloS one 2019; 14(3): e0213857.
673

674 **Figure legends**

675 **Figure 1. Overview of DTFLOW algorithm.**

676 A. Step 1: Pre-process of a single-cell dataset into a gene expression matrix $\mathbf{X}_{N \times D}$; B.
677 Step 2: Compute the nearest neighbours for each cell, get a nearest neighbour graph
678 structure, and then transform the dataset into a Markov transition matrix \mathbf{T} ; C. Step 3:
679 Using the random walk with restart method to get a diffusion matrix \mathbf{S} in which each
680 cell is represented by a discrete distribution vector; \mathbf{v}_i . Step 4: Construct a
681 Bhattacharyya kernel matrix \mathbf{K} and a matrix \mathbf{P} based on the properties of kernel
682 method; E. Step 5: Perform singular value decomposition on \mathbf{K} to get
683 low-dimensional embedding \mathbf{z} ; F. Step 6: Calculate the new distance metric
684 based on the row of the matrix \mathbf{P} corresponding to the root cell \mathbf{v}_0 , and unitize it to
685 get the pseudo-time distances \mathbf{t} ; G. Step 7: Identify the multi-branches of cellular
686 differentiation by reverse searching based on the nearest neighbour graph structure.

687 **Figure 2. Developmental trajectories inferred by DTFLOW for the mouse embryo**
688 **single-cell dataset.**

689 A. Visualization of different time stages. B. Visualization of inferred temporal
690 trajectory. C. Visualization of inferred 5 sub-branches when the minimal cell number
691 required for forming a sub-branch satisfies $\text{min}(\mathbf{t}) > 0$. D. Visualization of inferred 3
692 sub-branches when the minimal cell number required for forming a sub-branch is larger
693 ($\text{min}(\mathbf{t}) > 1$).

694 **Figure 3. Performance of three inference methods for the mouse embryo**
695 **single-cell dataset.**

696 A. The accuracy of DTFLOW that is measured by Kendall rank correlation coefficient
697 for the networks determined by different values of restart probability α . B. Robustness
698 properties of three inference methods. The properties are obtained by randomly
699 sampling 90% of single cells from the whole dataset and comparing the ordering results
700 of the subset with those of the whole dataset.

701 **Figure 4. Developmental trajectories inferred by DTFLOW for the mouse**
702 **myeloid progenitor MARS-seq dataset.**

703 A. Visualization of different cell types. B. Visualization of inferred temporal trajectory.
704 C. Visualization of calculated 4 sub-branches when the minimal cell number required
705 for forming a sub-branch satisfies $n = 6 \sim 22$. D. Visualization of calculated 3
706 sub-branches when the minimal cell number required for forming a sub-branch is larger
707 ($n = 23 \sim 124$).

708 **Figure 5 Robustness properties of the three inference methods**

709 A. Robustness property for the mouse myeloid progenitor MARS-seq dataset. B.
710 Robustness property for the mouse female gonad scRNA-seq dataset.

711 **Figure 6 Developmental trajectories inferred by DTFLLOW for the mouse female
712 gonad scRNA-seq dataset.**

713 A. Visualization of different cell types. B. Visualization of inferred temporal trajectory.
714 C. Visualization of the calculated three sub-branches.

715 **Figure 7 Dimensional reduction results of the dataset [40] by using four different
716 methods.**

717 A. BFKD; B.PCA; C. tSNE; D.UMAP.

718 **Figure 8 Dimensional reduction results of the dataset [43] by using four different
719 methods.**

720 A. BFKD; B.PCA; C. tSNE; D.UMAP.

721

722

723

724

725

726

727

728

729

730

731

732

733 **Supplementary materials**

734

735 **File S1: Supplementary Information**

736

737 **Legends of Supplementary Figures**

738 **Figure S1** Branch Detection Algorithm by reverse-searching in the kNN Graph

739 **Figure S2: Visualization of the mouse embryo single-cell dataset by using Scanpy.**

740 A. Visualization of 7 time stages. B. Visualization of Pseudo-times. C. Visualization of
741 3 groups/subbranches.

742 **Figure S3: Visualization of the mouse embryo single-cell dataset by using
743 Monocle2.**

744 A. Visualization of 7 time stages. B. Visualization of Pseudotimes. C. Visualization of
745 4 states/sub-branches.

746 **Figure S4: Expression values of genes *Gata3* and *Sox2*.**

747 The ordered values from the mouse embryo single-cell dataset are plotted along
748 pseudo-time obtained by using the three inference methods. The lines correspond to the
749 results by using the Gaussian process regression for each branch.

750 **Figure S5: Robustness properties of three inference methods for the mouse
751 embryo single-cell dataset.**

752 The properties are obtained by randomly sampling 80% of single-cells from the whole
753 dataset and comparing the ordering results of the subset with those of the whole dataset.

754 **Figure S6: Visualization of the mouse myeloid progenitor MARS-seq dataset by
755 using Scanpy.**

756 A. Visualization of 19 cell clusters. B. Visualization of Pseudotimes. C. Visualization
757 of 4 groups/sub-branches.

758 **Figure S7: Visualization of the mouse myeloid progenitor MARS-seq dataset by
759 using Monocle2.**

760 A. Visualization of 10 cell types. B. Visualization of Pseudo-times. C. Visualization of
761 12 states/sub-branches.

762 **Figure S8:** Visualization of gene Elane and Klf1 of the mouse myeloid progenitor
763 MARS-seq dataset by using different methods:

764 A. DTFLOW, B. Scanpy and C. Monocle2.

765 **Figure S9: Visualization of the mouse female gonad scRNA-seq dataset by using**
766 **Scanpy.**

767 A. Visualization of 6 cell clusters. B. Visualization of Pseudo-times. C. Visualization of
768 4 groups/sub-branches.

769 **Figure S9: Visualization of the mouse female gonad scRNA-seq dataset by using**
770 **Monocle2.**

771 A. Visualization of 6 cell types. B. Visualization of Pseudo-times. C. Visualization of 3
772 states/sub-branches.

773

774

775

776

777

778

779

780

781

782

783

784

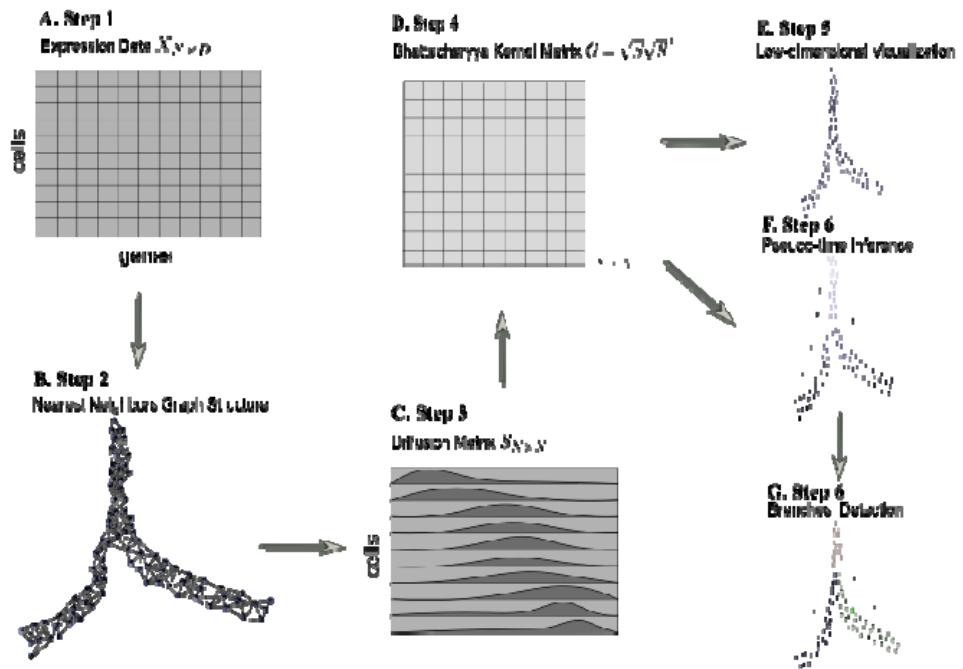
785

786

787

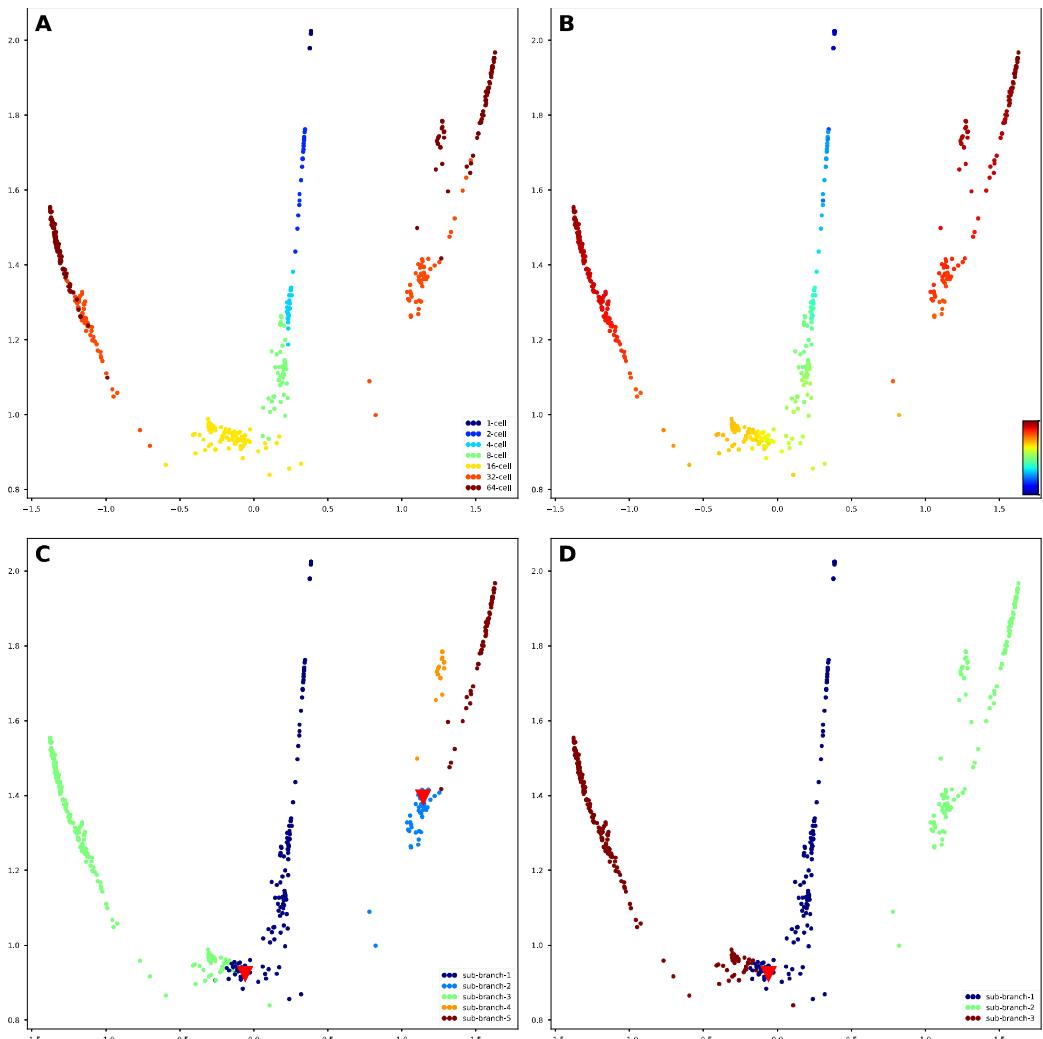
788

789



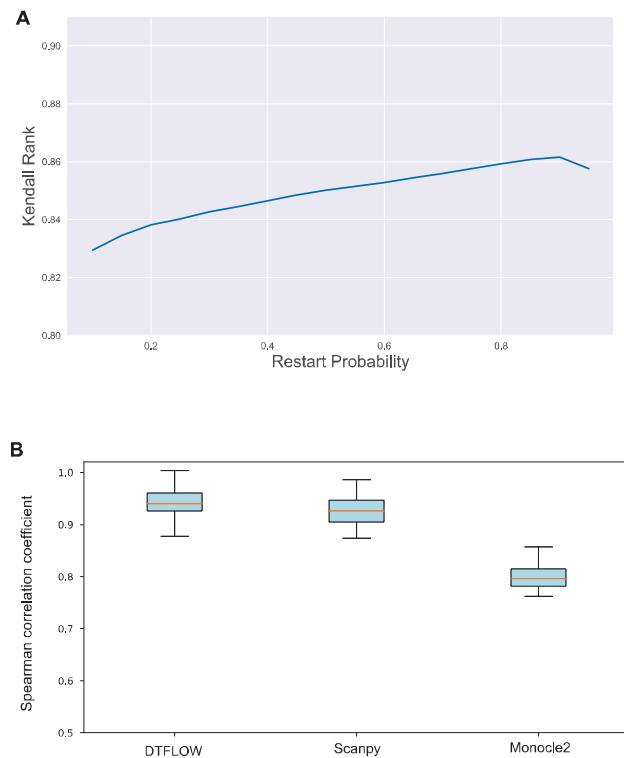
790

791 Figure 1.



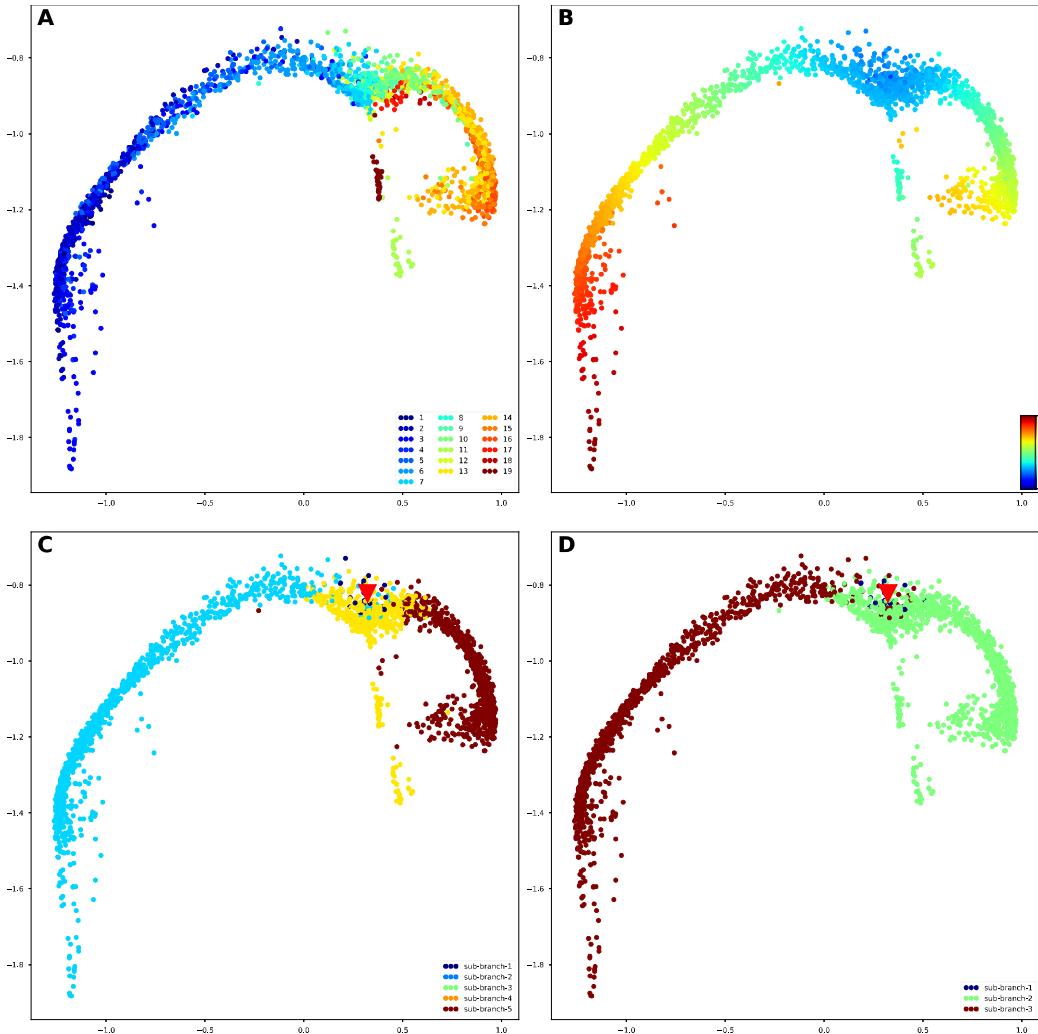
792

793 Figure 2



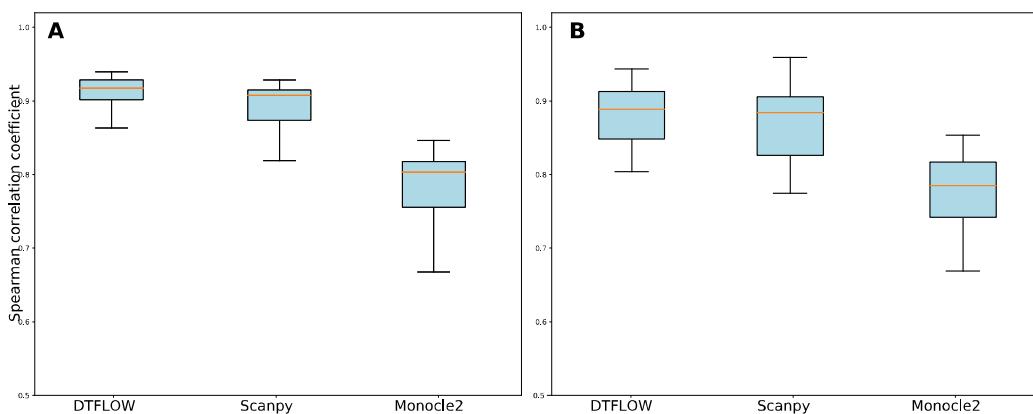
794

795 Figure 3.



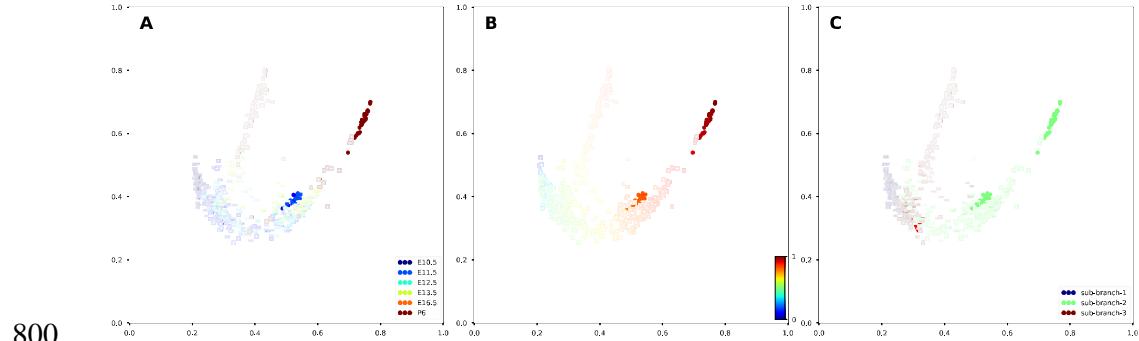
796

797 Figure 4.



798

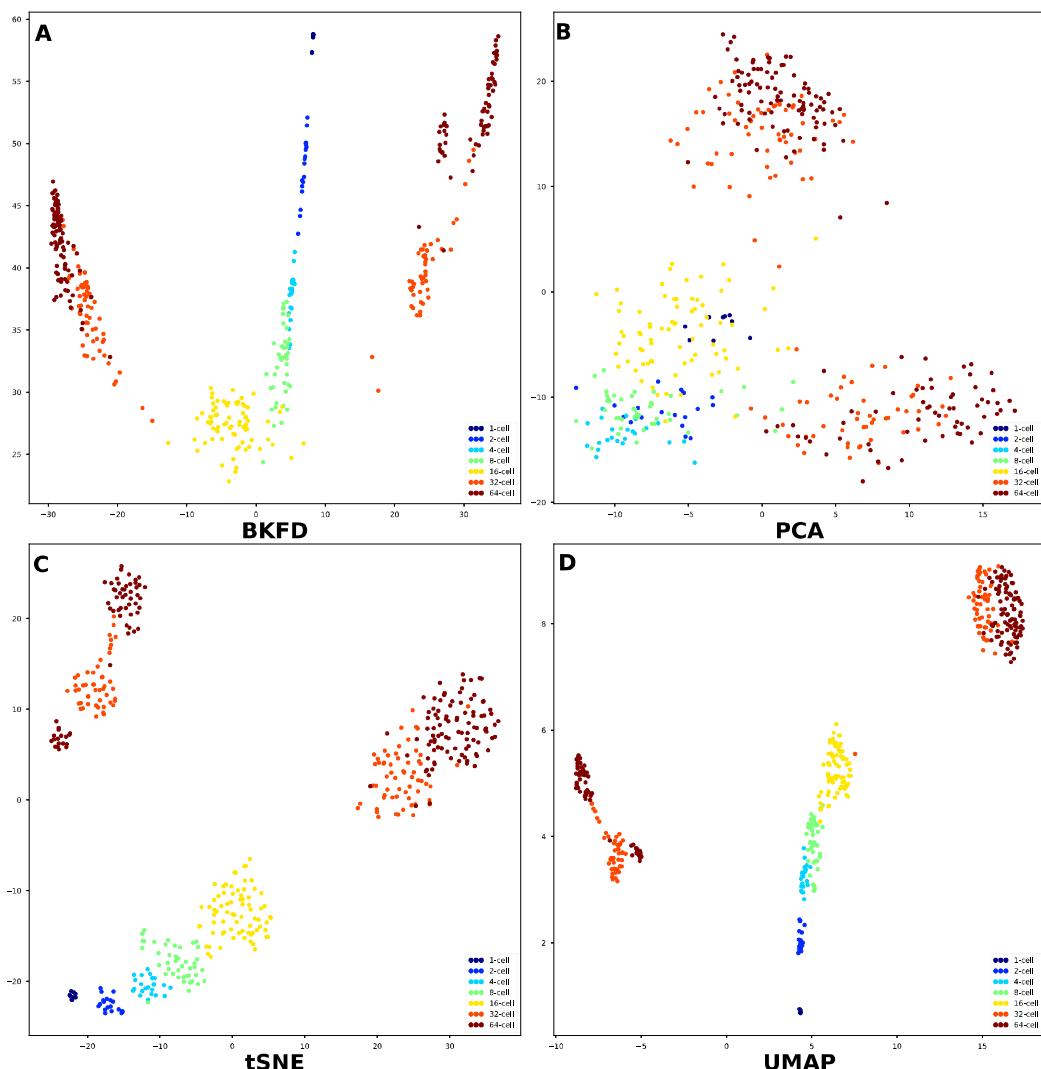
799 Figure 5.



800

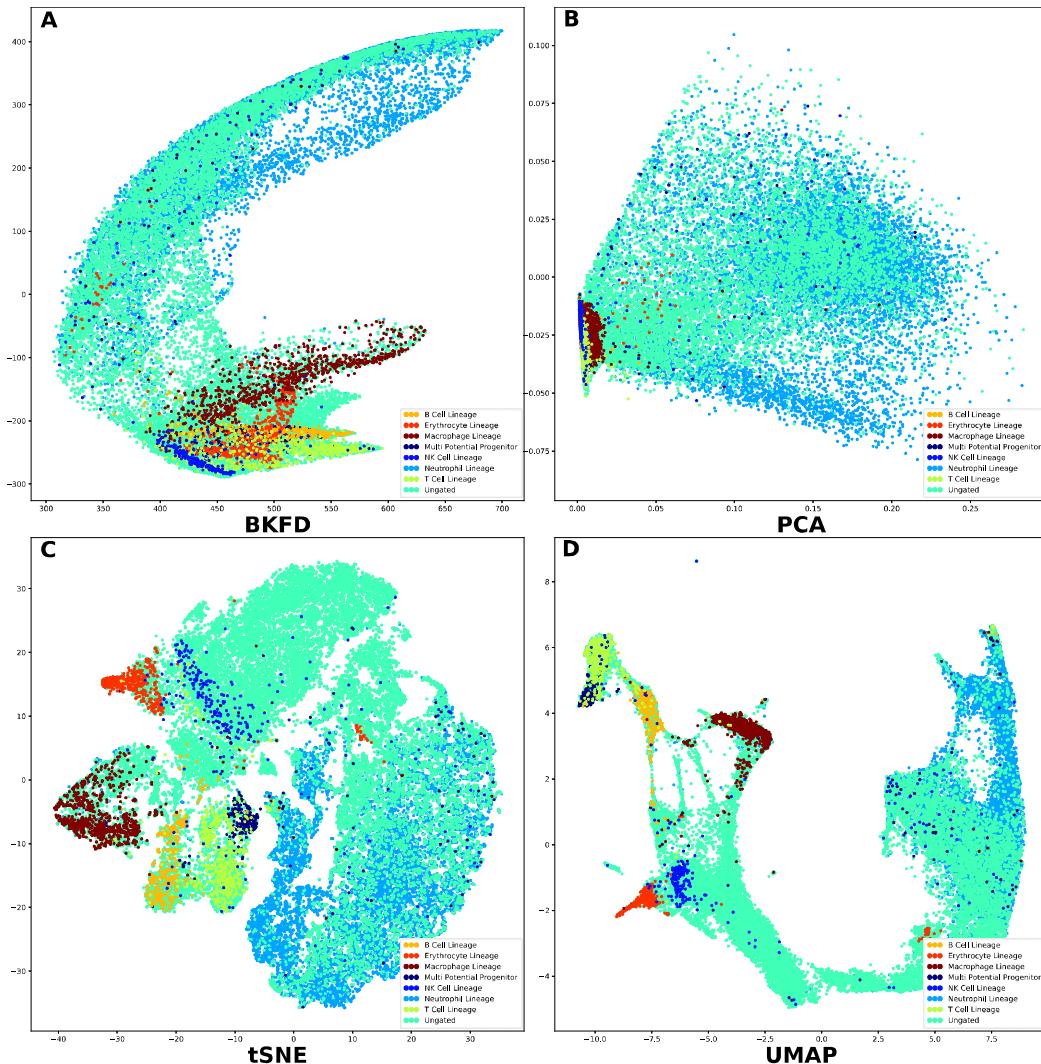
801 Figure 6.

802



803

804 Figure 7.



805

806

807

808

809

810

811

812

813

814

815

816

817

Supplementary Information

818

DTFLOW: Inference and Visualization of Single-cell Pseudo-temporal Trajectories Using Diffusion Propagation

820

Jiangyong Wei, Tianshou Zhou, Xinan Zhang and Tianhai Tian

821

S1. Comparison between DPT and DTFLOW

823 We note that DPT and DTFLOW use the similar propagation process. DPT provides the

824 (time independent) "path integral" by

825

$$\sum_{t=1}^{\infty} f(t) = f(0) \sum_{t=1}^{\infty} M^t,$$

826 where

is a probability density. Then the diffusion process of DPT is given

827 by

828

where the eigenvector π is corresponding to the eigenvalue 1 of the transition matrix

829

. The diffusion pseudo-time distance metric dpt is given by

830

which implies that it is also a kernel distance based on one kernel matrix $\pi^\top \pi$.

831

For the DTFLOW algorithm, the equation (7) can also be written as

832

Thus, DTFLOW can use parameter α to control the propagation procedure, which may

833

be the reason that the accuracy of DTFLOW is better than DPT.

834

835 **S2. The DTFLOW algorithm**

836 Algorithm 1 below gives the detailed steps of DTFLOW.

837 **Algorithm 1. DTFLOW**

838 **Input:** Single-cell data matrix $X_{N \times D}$, number of neighborhoods k , restart probability
839 $1 - p$, root cell r .

840 Part 1: *The Bhattacharyya Kernel Feature Decomposition Algorithm(BKFD)*. Output:
841 Low-dimensional structure Y_d and pseudo-time distances T .

- 842 1. Compute the k -nearest neighbors for each cell $X_{i,:}$ to get a nearest neighbor
843 graph structure.
- 844 2. $K \leftarrow$ transform the cell-cell. nearest neighbor distances of graph structure into a
845 symmetric Gaussian kernel weight matrix.
- 846 3. $M \leftarrow$ normalize K to get a Markov transition matrix.
- 847 4. Using the random walk with restart method to get one diffusion matrix $S \leftarrow (1 - p)(I - pM)^{-1}$.
- 848 5. Construct a Bhattacharyya kernel matrix $G \leftarrow \sqrt{S}\sqrt{S}^T$.
- 849 6. $\log G \leftarrow$ construct a new kernel matrix based on the property of kernel method.
- 850 7. $Y_d \leftarrow$ perform singular value decomposition on $\log G$.
- 851 8. $T \leftarrow$ for each cell i and the root cell r , calculate new distance metric $D_{ri} =$
852 $\sqrt{-2 \log G_{ij}}$, then unitize the distance metric.

853
854 Part 2. Branch Detection Algorithm by reverse-searching in the kNN Graph (see
855 Algorithm 2).

856

857 **S3. Branch Detection Algorithm by reverse-searching in the kNN Graph**

858 **Algorithm 2. RSKG**

859 **Input:** Indices array A of kNN graph, pseudo-times of single cells T , the minimum
860 number n of cells for a sub-branch.

861 **Output:** the sub-branch classification B_{vec} of single cells

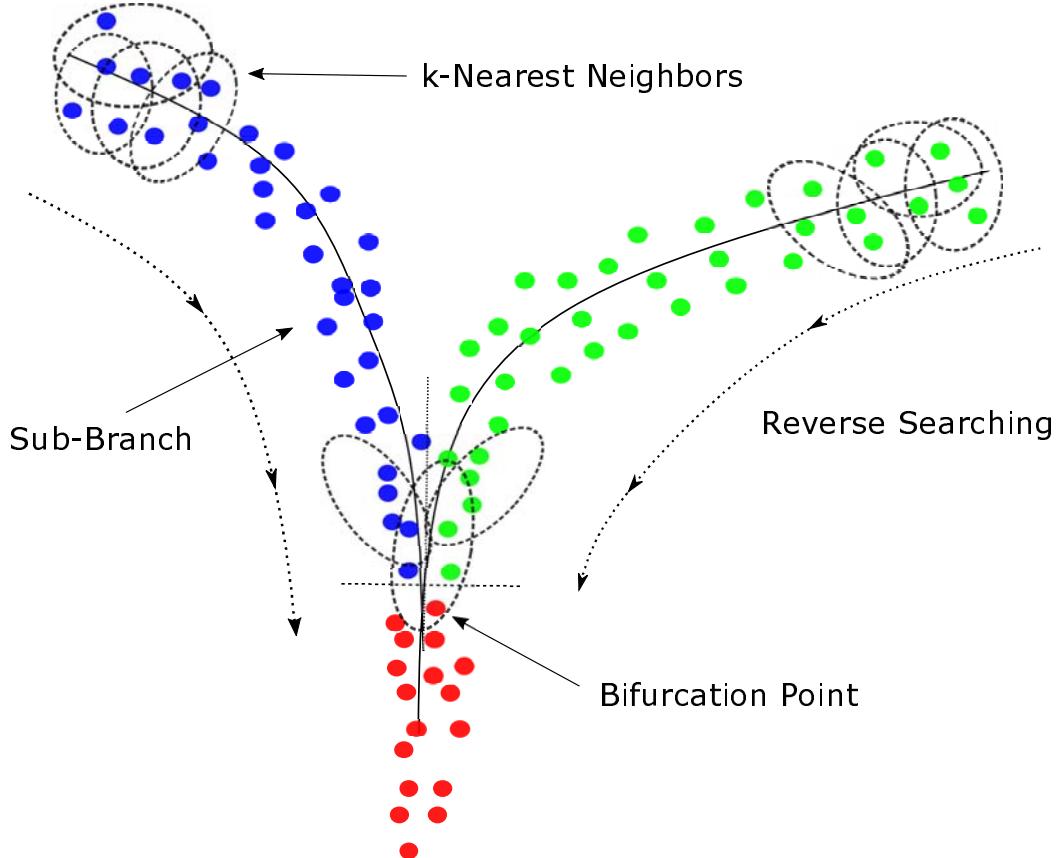
862 $R_{seq} \leftarrow$ the reverse indices ordering deduced by T .

863 For id in R_{seq} :

864 A[id] \leftarrow k-nearest neighbours of the \$id\$-th cell as one list group
865 If A[id] has no intersection with any list in prop-groups:
866 Append A[id] to prop-groups
867 Else
868 If A[id] has intersections with only one list in prop-groups:
869 Extend A[id] to that list of prop-groups
870 Else
871 If A[id] has intersection with two or more list of prop-groups:
872 Pop these intersected lists in prop-groups.
873 If At most one list has length $\geq n$:
874 Merge these intersected list together as one new list and append it to
875 prop-groups
876 Else
877 Let cell b_r with min pseudo-time in the intersection as bifurcation point.
878 Append the lists whose lengths $\geq n$ and only the cells whose pseudo-time
879 $> b_r$ to sub-branches.
880 Extend cells whose pseudo-time $> b_r$ in the other list to its nearest
881 sub-branch in sub-branches
882 Merge b_r and the cells whose pseudo-time $< b_r$ to one list and append
883 the list to prop-groups.
884 If A[id] has intersection with both prop-groups and sub-branches:
885 Let cell b_r with min pseudo-time in the intersection as bifurcation point.
886 Pop these intersected lists in prop-groups....
887 Append the lists whose lengths $\geq n$ and only the cells whose pseudo-time
888 $> b_r$ to sub-branches.
889 Extend cells whose pseudo-time $> b_r$ in other list of prop-groups to their
890 nearest sub-branches in the sub-branches.
891 Merge the cells whose pseudo-time $< b_r$ to one list and append the list to
892 prop-groups....

893 Assign the remaining lists in prop-groups to sub-branches, then get B_{vec} based on
894 sub-branches.

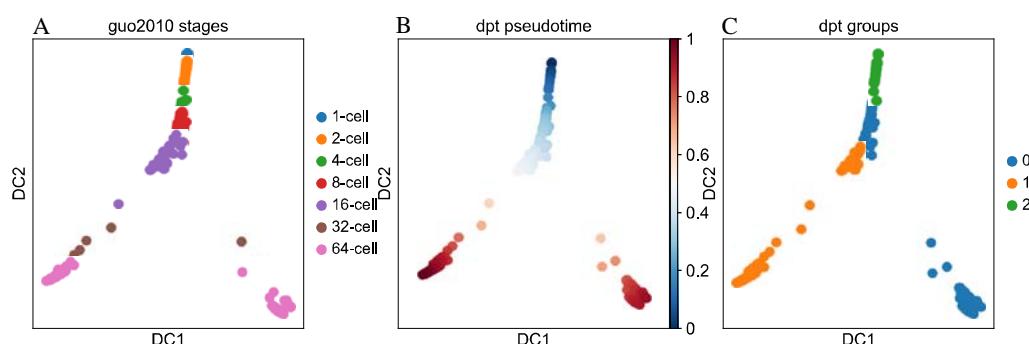
895



896

897 Supplementary Figure 1.

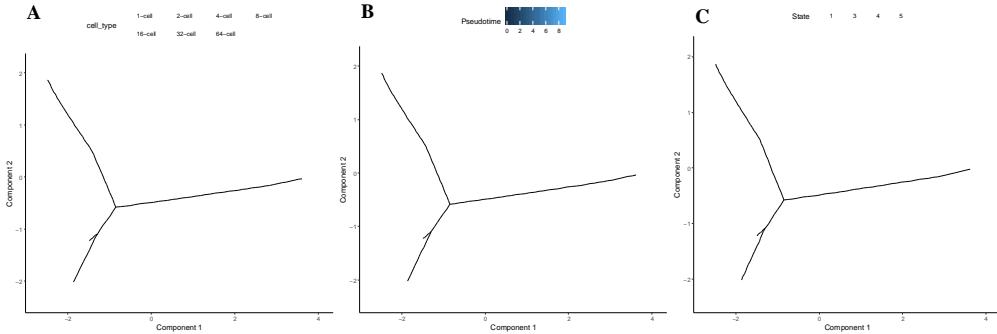
898



899

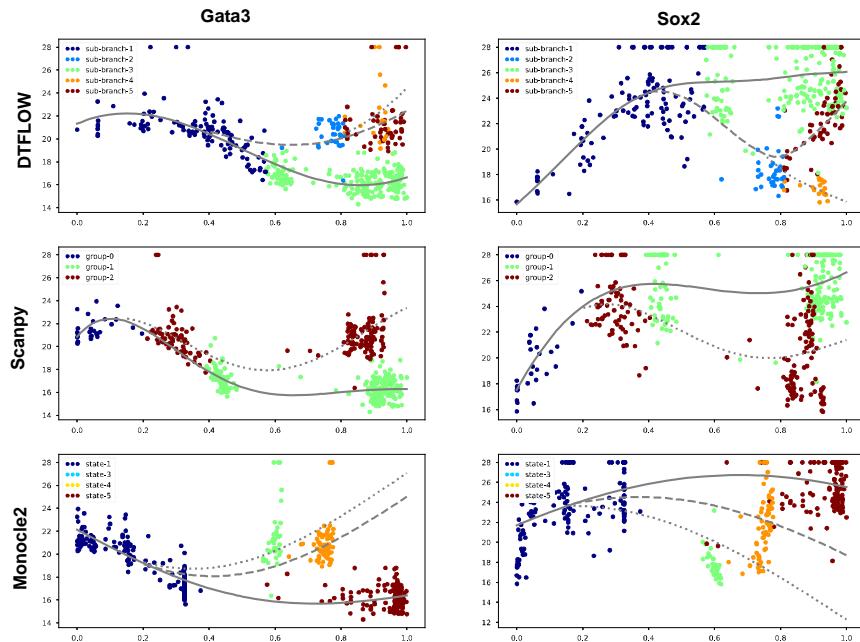
900 Supplementary Figure 2.

901



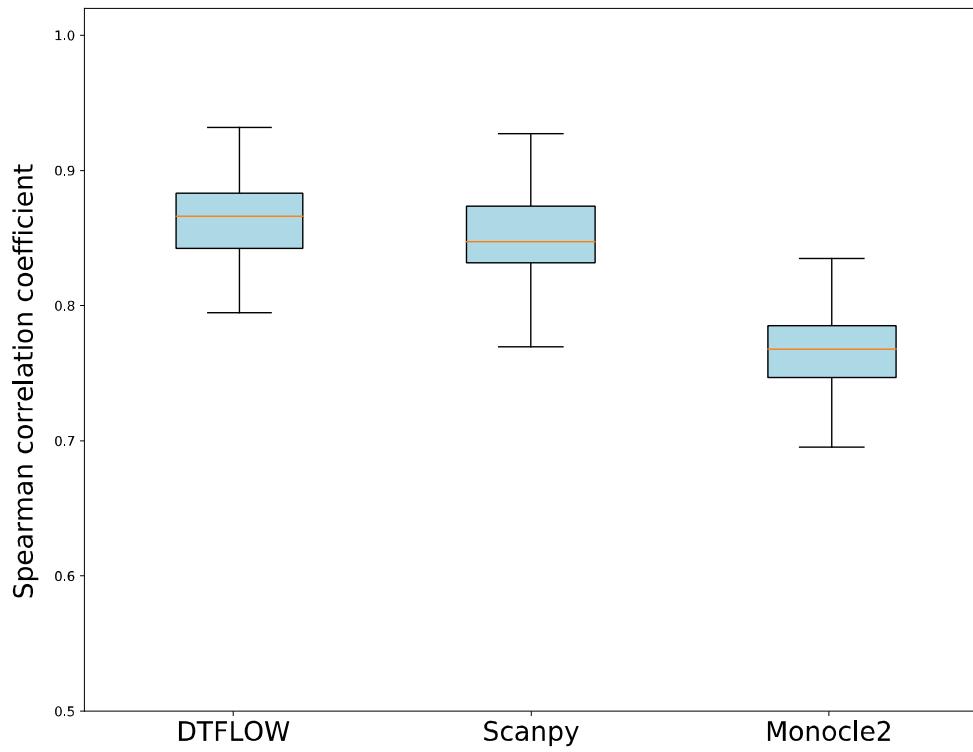
902

903 Supplementary Figure 3.



904

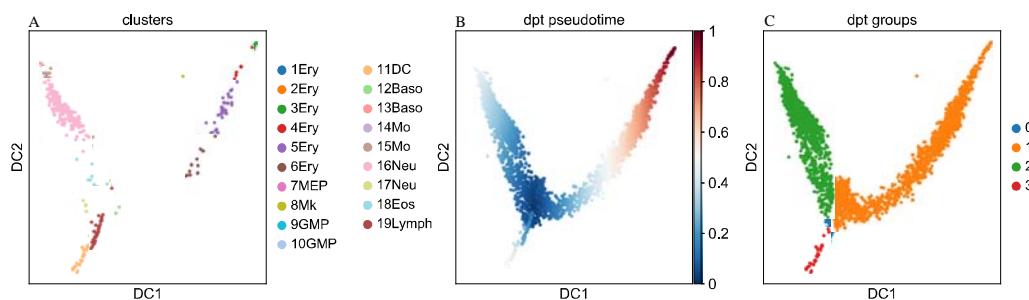
905 Supplementary Figure 4.



906

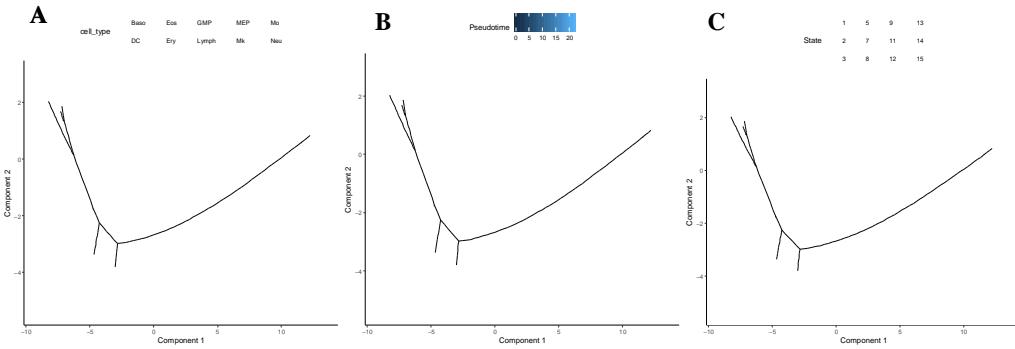
907 **Supplementary Figure 5.**

908



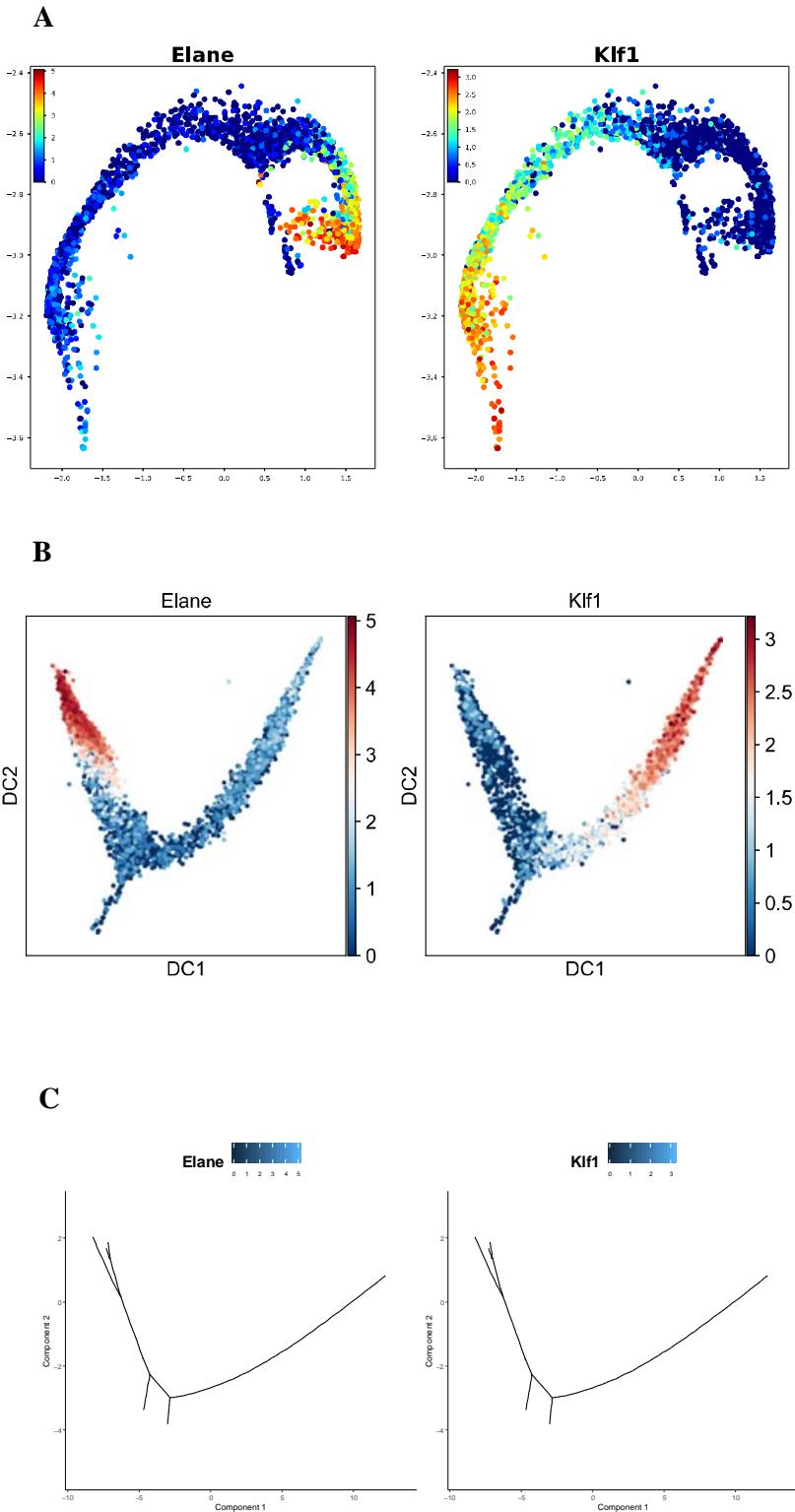
909

910 **Supplementary Figure 6.**



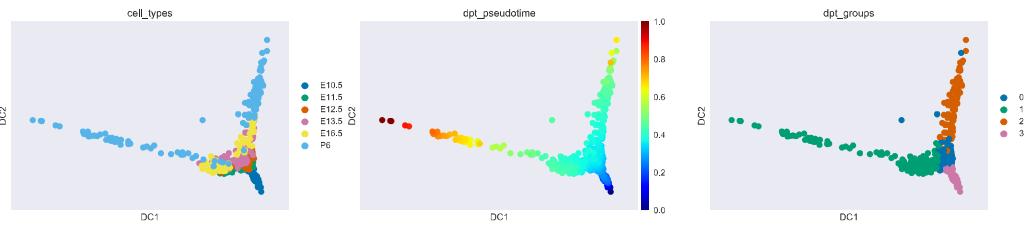
911

912 Supplementary Figure 7.



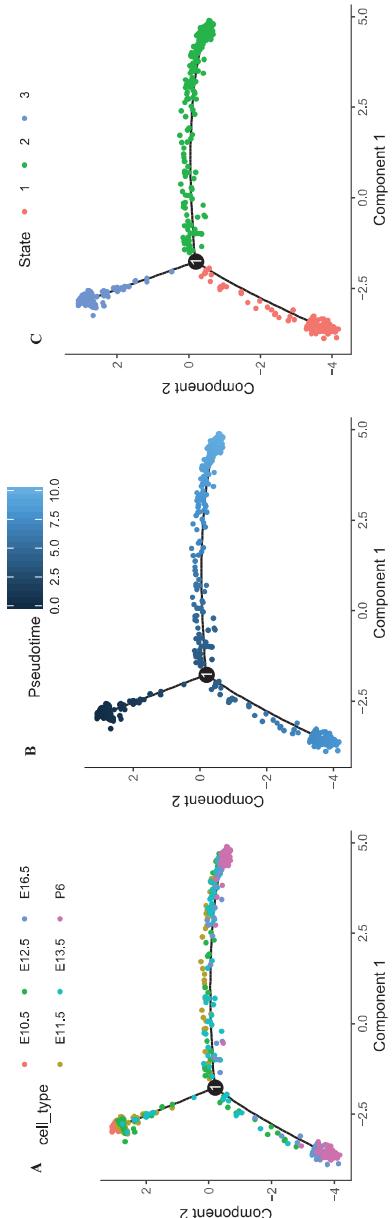
913

914 Supplementary Figure 8.



915

916 Supplementary Figure 9.



917

918 Supplementary Figure 10.

919