Elucidating Dynamic Cell Lineages and Gene Networks in Time-Course Single Cell Differentiation

Mengrui Zhang, Yongkai Chen, Dingyi Yu, Wenxuan Zhong, Jingyi Zhang, Ping Ma

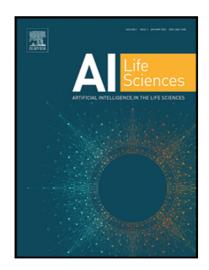
PII: S2667-3185(23)00012-0

DOI: https://doi.org/10.1016/j.ailsci.2023.100068

Reference: AILSCI 100068

To appear in: Artificial Intelligence in the Life Sciences

Received date: 21 December 2022 Revised date: 17 March 2023 Accepted date: 22 March 2023



Please cite this article as: Mengrui Zhang, Yongkai Chen, Dingyi Yu, Wenxuan Zhong, Jingyi Zhang, Ping Ma, Elucidating Dynamic Cell Lineages and Gene Networks in Time-Course Single Cell Differentiation, *Artificial Intelligence in the Life Sciences* (2023), doi: https://doi.org/10.1016/j.ailsci.2023.100068

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2023 Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (http://creativecommons.org/licenses/by-nc-nd/4.0/)

# Elucidating Dynamic Cell Lineages and Gene Networks in Time-Course Single Cell Differentiation

Mengrui Zhang\*, Yongkai Chen\*, Dingyi Yu, Wenxuan Zhong, Jingyi Zhang, Ping Ma

Dec. 2022

5 Abstract

Single cell RNA sequencing (scRNA-seq) technologies provide researchers with an unprecedented opportunity to exploit cell heterogeneity. For example, the sequenced cells belong to various cell lineages, which may have different cell fater in stein and progenitor cells. Those cells may differentiate into various mature cell types in a cell differentiation process. To trace the behavior of cell differentiation, researchers reconstruct cell lineages and predict cell fates by ordering cells chronologically into a traje tory with a pseudo-time. However, in scRNA-seq experiments, there are no cell-to-cell or respondences along with the time to reconstruct the cell lineages, which creates a significant of allenge for cell lineage tracing and cell fate prediction. Therefore, methods that can accurately reconstruct the dynamic cell lineages and predict cell fates are highly desirable.

In this article, we develop an innovative machine-learning framework called Cell Smoothing Transformation (Cen 'T) to elucidate the dynamic cell fate paths and construct gene networks in cell differentiation processes. Unlike the existing methods that construct one single bulk cell trajectory, ('ell'S') builds cell trajectories and tracks behaviors for each individual cell. Additionally, Cen'ST can predict cell fates even for less frequent cell types. Based on the individual cell fate trajectories, CellST can further construct dynamic gene networks to model gene-gene relationships along the cell differentiation process and discover critical genes that potentially regulate cells into various mature cell types.

Keywords: scRNA-seq, Optimal Transport, Smoothing Spline, Dynamic Gene Networks.

#### 1 Introduction

A comprehensive understanding of complex biological processes such as tissue development and regeneration requires the investigation of cell differentiation across a wide range of samples and experimental time points (Spiller et al., 2010). The cell differentiation process includes the differentiation of stem cells into different mature cell types (Guo et al., 2017; Burrows et al., 2020). Such a process is dynamic and continuous, including rapid changes in gene expressions and cell types over time. To profile such cell differentiation behaviors, single cell RNA-seq sequencing (scRNA-seq) technology has been developed rapidly (Nawy, 2013; Shapiro et al., 2013; Grün and Oudenaarden, 2015; Tanay and Regev, 2017). In particular, sch. A-seq enables researchers to observe the gene expressions of all cells simultaneously (Figure 1a) in both static or time-course 10 experiments (Figure 1b). The static scRNA-seq experiment takes a snapshot of all cells and their 11 gene expressions at one time (Lawson et al., 2015; H v. in et al., 2018), whereas the time-course 12 scRNA-seq experiments take snapshots at multiple t'm, points. Using scRNA-seq, researchers can 13 observe the behavior of individual cells in cell differentiation processes over time. Cell lineage tracing has been widely used to predic dy nanic cell fates by indicating the ancestor and posterity 15 cells in cell differentiation processes. For example, during a stem cell differentiation process, the 16 multipotent stem cells can deve's a noto multiple cell lineage endpoints (Figure 1c). Despite the 17 effectiveness, quantifying the virginic cellular changes of cell development is still challenging due 18 to the following technical limitations (Stegle et al., 2015). In time-course scRNA-seq experiments, 19 cells are sacrificed at d sequenced at each time point. Thus there is no cell-to-cell correspondence 20 information for cells between two time points, which creates a significant challenge in constructing 21 cell lineages and lucidating the dynamic cell behaviors in the differentiation process. Moreover, it 22 is very challenging to align different cells sequenced in two adjacent time points since expressions 23 of cells are high-dimensional and noisy, and the number of cells in each time point is large. Such 24 a large sample and high-dimensional and noisy data problem render many classical methods, such 25 as Euclidean distance or Pearson correlation, invalid (Alonge et al., 2020; Ren et al., 2017). 26 One natural approach to surmount the challenges is to order cells into a continuous cell trajec-27 tory. Many methods have been proposed to achieve this goal in static scRNA-seq experiments. In 28 these methods, researchers construct a pseudotime to order cells chronologically (Qiu et al., 2017;

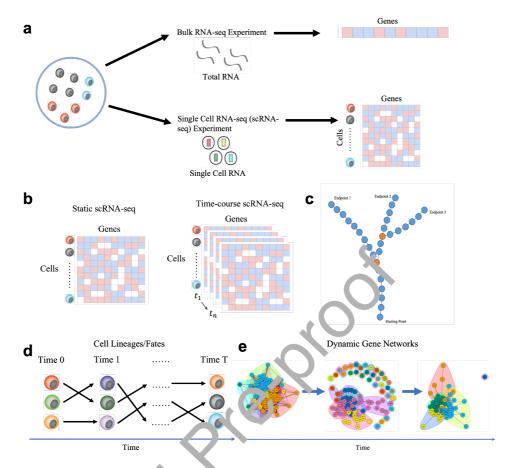


Figure 1: Single cell analysis and cell smoothing transformation (CellST) overview: **a**: The advantage of scRNA-seq analysis over pelk RNA-seq analysis. **b**: Data structures for static scRNA-seq experiments and time-course scR NA-seq experiments. Cells in time-course experiments have been marked with experimental time points. **c**: The multipotent stem cells can develop into multiple cell lineage endpoints. **c**' Ce'l lineages are constructed by connecting individual cells over time. Cell fate trajectories are constructed by smoothing the connected cell lineages. **e**: Dynamic gene networks are constructed based on the calculated dynamic relationship between genes.

- Cannoodt et al., 2016; Trapnell, 2015; Ji and Ji, 2016; Chen et al., 2019; Trapnell et al., 2014;
- <sup>2</sup> Liu et al., 2017). Despite their effectiveness, such methods may fail in the following circumstances
- 3 (Tritschler et al., 2019). First of all, most existing trajectory inference methods construct a bulk
- 4 cell trajectory, i.e., the mean trajectory of the population cells across time rather than that of
- 5 individual cells.
- 6 However, some individual cells' behaviors may oscillate up and down around their mean ex-
- 7 pressions or severely deviate from them. Cell differentiation behaviors are dominated by cells with
- 8 major cell types, and patterns with less frequent might be hidden in the dataset. Second, individual

cell developing trajectories may follow different complex topologies, including loops or alternative paths during the development. For example, analysis approaches in Moon et al. (2018), and Dai et al. (2020) used dimension reduction methods to identify a low-dimensional space of the gene expression space before constructing cell trajectories (Saelens et al., 2019; Wagner et al., 2018). Those methods may introduce a significant bias and are hard to validate, as cells are ordered based only on the selected reduced dimensions. Finally, the cells may not be synchronized at the same developing time points. Cells within the same time point can be expressed at different developing stages. In this situation, the bulk cell trajectory that takes the average pattern of cells at different stages might result in unreliable scientific discovery. In this article, we propose a novel analysis framework name? Cen Smoothing Transformation 10 (CellST) to overcome the aforementioned limitations. The CellST ramework elucidates dynamic 11 cell fates and constructs gene networks in the cell different ation process. In the CellST framework, 12 we propose a cell lineage tracing method, which alig is two individual cells between any adjacent 13 two time points via the optimal transport technique (Vihani, 2003; Meng et al., 2019), which is a powerful tool that can be used to model cel' dynamics (Schiebinger et al., 2019; Tong et al., 2020; 15 Zhang et al., 2020). Those aligned ce.'s can potentially represent individual cell lineages, tracing 16 cell differentiation behaviors by constructing cell-to-cell trajectories (Figure 1d). We then use a 17 smoothing spline model to precipe all fate trajectories and reduce both cell-cell variations. The smoothing spline method in odes the gene expression patterns in the aligned cell lineages from the 10 previous step and build, the estimated individual cell fate trajectories. Lastly, we narrow down our focus to utilize the concexpression patterns from those cell fate trajectories to construct dynamic gene networks (Figure 1e). The dynamic gene networks are constructed by estimating the dynamic relationship of pairwise gene expression patterns using the functional concurrent models (Wang et al., 2016) and smoothing spline models (Gu and Ma, 2005). The dynamic gene networks can be used to find critical genes by profiling genes with significantly different patterns from other genes. Our major contribution is developing the first analysis framework (CellST) to construct cell lineages and predict dynamic cell fates at the individual cell level, which can help researchers better observe cell behaviors in the differentiation process. In contrast, the existing methods only estimate the bulk trajectory in scRNA-seq experiments. Those analysis methods may overlook the hidden patterns in the cell differentiation process to create a spurious cell differentiation trajectory. We

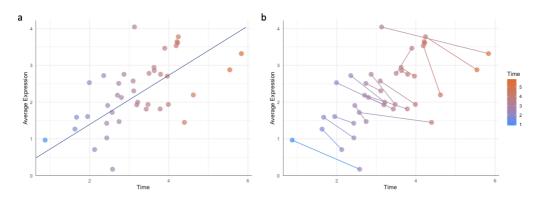


Figure 2: Example of cell-to-cell linking: **a**: Cell differentiation over time (x-axis) reflects an increasing trend in average cell expressions (y-axis). **b**: The individual cell correspondences at different time points reflect a decreasing trend in average cell expression (y-axis).

illustrate this problem by using a simulated time-course cell dat, set indicating the disadvantage of bulk cell trajectories (Figure 2). The cell-to-cell trajectories are ble to overcome the disadvantage and identify the real gene expression patterns in cell development. Under some cell development 13 and differentiation circumstances, cells' average expressions show an increasing pattern if we only construct one average cell trajectory to order elle (Figure 2a). However, when individual cells 15 are aligned at different time points, t'e ir lividual cell lineages' average expressions reflect unique 16 decreasing patterns, which are in contract to the bulk trajectory (Figure 2b). This means some 17 cells start at a lower expression le el, and the expression keeps going down over time. Those cell 18 development patterns can be easily misled by the average cell trajectory and thus reflect spurious cell differentiation behavers. Furthermore, we propose the dynamic gene networks based on the 20 individual cell fate to jectories to estimate the dynamic gene-gene relationship and critical genes 21 in the differentiation process. The empirical performance of the proposed framework is evaluated 22 by several simulated and real experiment studies. 23

#### 24 2 Method

- 25 In this section, we introduce the Cell Smooth Transformation (CellST) method, which constructs
- 1 the cell fate trajectories and dynamic gene networks for time-course scRNA-seq data.

#### 2 2.1 Cell lineage & Individual cell fate trajectories

- 3 To construct the cell fate trajectories, we first align the cells at different time points to construct
- 4 the cells' lineages information between time points. We then smooth the gene expression pattern
- 5 for each gene over time and extract the "mean curve" of all individual gene expression patterns in
- <sup>6</sup> single cell fate trajectories to obtain the general gene expression pattern.

#### <sup>7</sup> 2.1.1 Cell-to-cell lineages by optimal transport

8 Regarding cells at different time points as cells with genes of different domain spaces, we transform

9 the problem of aligning cells at different time points into a problem of domain adaptation. Specif-

ically, we denote the normalized gene expression for cell i at ti ne t by a d-dimensional vector  $\mathbf{x}_{i}^{t}$ ;

each dimension of  $\mathbf{x}_i^t$  represents a gene expression <sup>1</sup>. We write  $\mathbf{X}_t = \left\{\mathbf{x}_i^t\right\}_{i=1}^{n_t}$ , where  $n_t$  indicates

 $_{12}$  the number of cells at time t in single cell RNA-seq dataset. Our goal is to learn the transformation

between the domain spaces by aligning the distribution of  $X_t$  to  $X_{t+1}$ .

As a powerful tool to learn the transformation from one probability measure to another, optimal transport has been applied to solve the do nally adaptation problem (Courty et al., 2014). We thus apply optimal transport to obtain the comain adaptive coupling between  $\mathbf{X}_t$  and  $\mathbf{X}_{t+1}$ . In other words, we transform the cell alignment problem into an optimal transport problem. In particular, we formulate the problem as a Monge optimal transport by minimizing the cost for transporting a gene expression distribution  $\mu_t$  and  $\mu_{t+1}$  using a map  $\mathbf{T}_t$ :

$$\min_{\mathbf{T}_t} \int_t c(x, \mathbf{T}_t(x)) d\mu_t(x), \tag{1}$$

where  $\mathbf{T}_t \# \mu_t = \mu_{t+1}$ , # represents the push-forward operator, such that for any measurable x,  $\mathbf{T}_t \# \mu_t(x) = \mu_t(\mathbf{T}_t^{-1}(x))$ ,  $\mu_t$  and  $\mu_{t+1}$  are probability distribution of  $\mathbf{X}_t$  and  $\mathbf{X}_{t+1}$  in  $\mathbb{R}^d$ , where d is the dimension. We define the optimal transport map  $\mathbf{T}_t : \mathbb{R}^d \to \mathbb{R}^d$ , where  $\mathbb{R}^d$  can be interpreted as the domain space for  $\mathbf{x}_i^t$  or  $\mathbf{x}_i^{t+1}$ . In this optimal transport problem, one constraint for the transportation map  $\mathbf{T}_t$  from a measure  $\mu_t$  to a measure  $\mu_{t+1}$  is the so-called measurement-preserving, i.e.,  $\mathbf{T}_t \# \mu_t = \mu_{t+1}$ . Among all the measurement-preserving maps, the optimal  $\mathbf{T}_t$  is the one that

<sup>&</sup>lt;sup>1</sup>For  $\mathbf{x}_{i}^{t}$ , we (1) use all available genes, (2) select highly expressed genes, or (3) apply dimension reduction methods such as the principal component analysis (PCA). In the first two cases, each gene expression represents an individual gene; while in the last case, each gene feature represents a combination of all genes.

- 4 minimizes the transportation cost.
- 5 Since we can only observe gene expressions for sample cells at each time point, we focus on
- 6 the case where the probability distributions are discrete. The distributions  $\mu_t$  and  $\mu_{t+1}$  for gene
- <sup>7</sup> features at time points t and t+1 are defined as:

$$\mu_t = \frac{1}{n_t} \sum_{i=1}^{n_t} \delta_{t_i} \quad \text{and} \quad \mu_{t+1} = \frac{1}{n_{t+1}} \sum_{j=1}^{n_{t+1}} \delta_{t+1_j},$$
 (2)

where  $\delta_{t_i}$  and  $\delta_{t+1_j}$  are the Dirac measures at location  $\mathbf{x}_{t_i}$  and  $\mathbf{x}_{t+1_j}$  respectively. Denote the positions of the supporting points  $\mathbf{X}_t = (\mathbf{x}_{t_1}, \mathbf{x}_{t_2}, \dots, \mathbf{x}_{t_{n_t}})^T$ . In discrete cases, the transport  $\mathbf{T}_t$  from  $\mu_t$  to  $\mu_{t+1}$  can be denoted as  $\mathbf{T}_t(\mathbf{X}_t) = \Sigma \mathbf{X}_t$ , where  $\Sigma$  is an  $n_{t+1} \times n_t$  matrix. For simplicity, we first consider the equal-size mapping, i.e.,  $n_t = n_{t+1} = n$ . Notice that in this case, the transport between  $\mathbf{X}_t$  and  $\mathbf{X}_{t+1}$  is a one-to-one assignment with permutation,  $\Sigma$  then can be regarded as a permutation matrix with the (i,j)th element:

$$\Sigma_{i,j} = \begin{cases} 1 & \text{f } \Upsilon_t | \mathbf{x}_{t_j} ) = \mathbf{x}_{t+1_i}, \\ 0 & \text{otherwise,} \end{cases}$$
 (3)

Furthermore, the transportation cost  $C(\mathbf{T}_t)$  defined in (1) can be calculated as:

$$C(\mathbf{T}_t) = \sum_{i=1}^n \sum_{j=1}^n c\left(\mathbf{x}_{t_j}, \mathbf{x}_{t+1_i}\right) \Sigma_{i,j},\tag{4}$$

where  $c(\mathbf{x}_{t_j}, \mathbf{x}_{t+1_j})$  on be interpreted as the energy required to transform an individual cell from the stage as  $\mathbf{x}_j^t$  on the stage as  $\mathbf{x}_i^{t+1}$ . The optimal transport map  $\mathbf{T}_t$  then can be calculated through:

$$\min_{\Sigma} \sum_{i=1}^{n} \sum_{j=1}^{n} c(\mathbf{x}_{t_j}, \mathbf{x}_{t+1_i}) \Sigma_{i,j}.$$
 (5)

where  $c(\mathbf{x}_{t_j}, \mathbf{x}_{t+1_i}) = \|\mathbf{x}_{t_j} - \mathbf{x}_{t+1_i}\|^{\alpha}$  and  $\|\cdot\|$  is the Euclidean norm in  $\mathbb{R}^d$ . We set  $\alpha = 2$  in this paper. The minimum of the optimization problem (5) is called the  $L^{\alpha}$ -Wasserstein distance (to the power  $\alpha$ ) and is denoted by  $W_{\alpha}(\mu_t, \mu_{t+1})^{\alpha}$ . The  $W_{\alpha}$  defines a distance on the set of distributions (cells) that have moments of order  $\alpha$ . In general, the cell lineage construction by optimal transport can be summarized as three steps: Estimating empirical gene feature distributions  $\mu_t$  and  $\mu_{t+1}$  as

4 in (2). Finding an optimal transport map  $\mathbf{T}_t$  from  $\mu_t$  to  $\mu_{t+1}$  through (5). Applying  $\mathbf{T}_t$  to obtain the cell-to-cell coupling from  $\mathbf{X}_t$  to  $\mathbf{X}_{t+1}$ .

It's important to note that the optimal transport map discussed above could be unsuitable in some cases when the one-to-one cell differentiation assumption does not hold. However, our optimal transport framework can be modified to account for more general cell-to-cell relationships. Specifically, we consider the following two general scenarios. First, the number of cells may vary at 9 different time points, and as a result, some cells may need to be reused when constructing cell-to-10 cell lineages. This can lead to multiple lineages passing through a single cell at specific time points. Second, cells may exhibit different proliferation rates at the same time points, i.e., the numbers of 12 new cells produced by two cells at the same time point could be say in cantly different. This could 13 result in the varying proportions of cell groups across time. To valless these general assumptions, 14 we present a comprehensive discussion of our generalized methods and experimental results in the 15 supplementary materials.

#### 17 2.1.2 Individual cell fate trajectories by smoothing spline model

After we align the cells from different time points, we can obtain the individual cell lineages at time 18 points t and t+1. We then align central time points based on the cell couplings to construct each 19 cell's coarse cell fate trajectories access the timeline. Those cell fate trajectories are smoothed to 20 reduce the estimation varial ce in CellST by utilizing the smoothing spline models. The smoothing 21 spline model is a versa, emily of smoothing methods that are suitable for both univariate and 22 multivariate problems (Ju, 2013). To construct the proposed smoothed cell trajectories, we use 23 equation 6 to how, the behavior patterns of the gene expression along the cell fate trajectories. 24 Let t represent the time points in the time-course dataset, and  $g_i$  represent the gene expression for each gene within an aligned cell fate trajectory. For co-expressed genes, we model the gene expression patterns using a smoothing spline mix-effect model with  $\{g_i, t_i\}_{i=1}^n$  as the observations (Gu and Ma, 2005):

$$g_i = \eta(t_i) + \mathbf{z}_i^T \mathbf{b} + \varepsilon_i \tag{6}$$

i = 1, ..., n, where the regression function  $\eta(t_i)$  is assumed to be a smooth function on the genes domain space in a cell.  $\eta(t_i)$  are the fixed effects and  $\mathbf{z}_i^T \mathbf{b}$  are the random effects with  $\mathbf{b} \sim N(\mathbf{0}, B)$ 

- 5 and  $\varepsilon_i \sim N(0, \sigma^2)$ . The random effects are used to account for the co-expressed genes in one
- 6 individual cell trajectory. The model terms  $\eta(t)$  or  $\eta(t) + \mathbf{z}^T \mathbf{b}$  will be estimated using the penalized
- 7 (unweighted) least squares method through the minimization of

$$\frac{1}{n} \sum_{i=1}^{n} \left( g_i - \eta \left( t_i \right) - \mathbf{z}_i^T \mathbf{b} \right)^2 + \frac{1}{n} \mathbf{b}^T \Sigma \mathbf{b} + \lambda J(\eta), \tag{7}$$

where the first term measures the goodness-of-fit,  $J(\eta) = \int (\eta''(t))^2 dt$  quantifies the smoothness of  $\eta$ , and  $\lambda$  is the smoothing parameter controlling the trade-off between the goodness-of-fit and the smoothness of  $\eta$  (Wahba, 1990; Gu, 2013). Consider the minimization of the least squares

estimation (equation 7) in a space with basis  $\{\xi_1,\ldots,\xi_q\}$ , function  $\eta$  can be expressed as

$$\eta(t) = \sum_{j=1}^{d} c_j \xi_j(t) = \boldsymbol{\xi}^T(t) \mathbf{c}.$$
(8)

Plugging equation 8 into equation 7, thus  $\eta$  can be estimated by minimizing:

$$(\mathbf{g} - R\mathbf{c} - Z\mathbf{b})^{T}(\mathbf{g} - P\mathbf{c} - Z\mathbf{b}) + \mathbf{b}^{T}\Sigma\mathbf{b} + n\lambda\mathbf{c}^{T}Q\mathbf{c}.$$
 (9)

With the standard formulation of penalized least squares regression, the minimization of equation 7 is performed in a so-called reproducing kernel Hilbert space  $\mathcal{H} \subseteq \{\eta : J(\eta) < \infty\}$  in which  $J(\eta)$  is a square seminorm, and the solution resides in the space  $\mathcal{N}_J \oplus \operatorname{span} \{R_J(t_i, \cdot), i = 1, \dots, n\}$ , where  $\mathcal{N}_J = \{\eta : J(n) = 0\}$  is the null space of  $J(\eta)$  and  $R_J(\cdot, \cdot)$  is the so-called reproducing kernel in  $\mathcal{H} \oplus \mathcal{N}_J$ . The solution has an expression:

$$\eta(t) = \sum_{i=1}^{m} d_{\nu} \phi_{\nu}(t) + \sum_{i=1}^{n} \tilde{c}_{i} R_{J}(t_{i}, t)$$
(10)

where  $\{\phi_{\nu}\}_{\nu=1}^{m}$  is a basis of  $\mathcal{N}_{J}$ . It follows that  $R=(S,\tilde{Q})$ , where S is  $n\times m$  with the  $(i,\nu)$ th entry  $\phi_{\nu}(t_{i})$  and  $\tilde{Q}$  is  $n\times n$  with the (i,j) th entry  $R_{J}(t_{i},t_{j})$ . In the smoothing spline model, the estimation of  $\eta$  is highly related to the choosing of the smoothing parameter  $\lambda$ . We choose the smoothing parameter  $\lambda$  and estimate random effect  $\mathbf{b}$  by Generalized Cross-Validation (GCV) (Wahba, 1990; Gu and Ma, 2005). Since there are d gene expression patterns over t time points for the cell fate trajectories, the smoothing spline model estimates one expression pattern for individual

6 cells and smooths the expression patterns.

#### 7 2.1.3 Dynamic gene networks

- 8 We consider the connection of two genes to be dynamic and the relationship may smoothly change.
- Suppose we want to study the dynamic relationship of the lth gene and sth gene, where  $1 \le l, s \le l$
- $p,l \neq s$ . Denote  $X_i^{\langle l \rangle}(t)$  and  $X_i^{\langle s \rangle}(t)$  as the lth gene and sth gene's expression values of cell fate
- trajectories i, and  $i = 1, \dots, n$ . By taking lth gene as the response and sth gene as the covariate,
- we consider the functional concurrent linear model,

$$X_i^{\langle l \rangle}(t) = \beta^{\langle l, s \rangle}(t) X_i^{\langle s \rangle}(t) + \varepsilon_{i, t}^{\langle l, s \rangle} \tag{11}$$

where  $\beta^{\langle l,s\rangle}(t)$  models the dynamic linear relationship between two genes,  $\varepsilon_{i,t}^{\langle l,s\rangle}$ s are i.i.d. random errors with mean zero and constant variance. We extracte  $\beta^{\langle l,s\rangle}(t)$  by minimizing the following penalized least squares function,

$$\frac{1}{nK} \sum_{i=1}^{n} \sum_{k=1}^{K} \left( X_i^{\langle l \rangle} \left( z_{ik} \right) - \beta^{\langle l, s \rangle} \left( t_{ik} \right) X_i^{\langle s \rangle} \left( t_{ik} \right) \right)^2 + \lambda J(\beta^{\langle l, s \rangle}). \tag{12}$$

With the representer theorem ( $W_2h_{b,a}$ , 1990), the optimizer of 12 can be written as

$$\hat{\beta}_{x}^{(l,s)}(t) = \sum_{v=1}^{m} d_{v} \psi_{v}(t) + \sum_{i=1}^{n} \sum_{k=1}^{K} c_{ik} R_{1}(t_{ik}, t)$$
(13)

where  $\{\psi_v\}_{v=1}^m$  is the basis function of the m-dimensional null space  $\mathcal{H}_0$ , and  $R_J(\cdot, \cdot)$  is the reproducing kernel or  $\mathcal{H}_1$ . Moreover,  $d_v$  and  $c_{ik}$  are the coefficients to be estimated. By Plugging equation 13 to equation 12, we can yield the estimations of  $\mathbf{c} = (c_1, \cdots, c_{1K}, \cdots, c_{n1}, \cdots, c_{nK})^T$  and  $\mathbf{d} = (d_1, \cdots, d_{1K}, \cdots, d_{n1}, \cdots, d_{nk})^T$ , which follow

$$c = \left(\mathbf{M}^{-1} - \mathbf{M}^{-1}\mathbf{S}\left(\mathbf{S}^{T}\mathbf{M}^{-1}\mathbf{S}\right)^{-1}\mathbf{S}^{T}\mathbf{M}^{-1}\right)\mathbf{X}^{\langle s\rangle}\mathbf{X}^{\langle l\rangle}$$

$$d = \left(\mathbf{S}^{T}\mathbf{M}^{-1}\mathbf{S}\right)^{-1}\mathbf{S}^{T}\mathbf{M}^{-1}\mathbf{X}^{\langle l\rangle}$$
(14)

where  $\mathbf{X}^{\langle s \rangle} = diag((\mathbf{X}_1^{\langle s \rangle T}, \cdots, \mathbf{X}_n^{\langle s \rangle T}))$  with the vector  $\mathbf{X}_i^{\langle s \rangle T} = (X_i^{\langle s \rangle}(t_{i1}), \cdots, X_i^{\langle s \rangle}(t_{iK}))^T, \mathbf{X}^{\langle l \rangle} = (\mathbf{X}_1^{\langle l \rangle T}, \cdots, \mathbf{X}_n^{\langle l \rangle T}))^T$  with the vector  $\mathbf{X}_i^{\langle l \rangle} = (X_i^{\langle l \rangle}(t_{i1}), \cdots, X_i^{\langle l \rangle}(t_{iK}))^T, \mathbf{S} = (\mathbf{S}_1^T, \cdots, \mathbf{S}_n^T)^T$  with

- 5 the (k,v)th entry of the  $K \times m$  matrix  $\mathbf{S}_i$  equals to  $\psi_v(t_{ik})X_i^{\langle s \rangle}(t_{ik})$ ,  $\mathbf{M} = \mathbf{X}^{\langle s \rangle}\mathbf{Q}\mathbf{X}^{\langle s \rangle} + n\lambda \mathbf{I}$  and  $\mathbf{Q}$
- 6 is the  $nK \times nK$  block matrix with the (i,j)th block is the  $K \times K$  matrix with the (k,u)th entry
- 7 equals to  $R_1(t_{ik}, t_{ju})$ . Thus, the estimation of  $\beta^{\langle l, s \rangle}(t)$  can be written as

$$\hat{\beta}^{\langle l,s\rangle}(t) = \boldsymbol{\psi}^T \boldsymbol{d} + \boldsymbol{\xi}^T \boldsymbol{c} \tag{15}$$

- where  $\boldsymbol{\psi} = (\psi_1(t), \cdots, \psi_m(t))^T$  and  $\boldsymbol{\xi} = (R_1(t_{11}, t), \cdots, R_1(t_{1K}, t), \cdots, R_1(t_{n1}, t), \cdots, R_1(t_{nK}, t))^T$ .
- 9 Note that  $\beta^{(l,s)}(t)$  models the dynamic linear relationship between lth gene and sth gene, and
- $\beta^{\langle l,s\rangle}(t_0)=0$  means the correlation between gene l and gene s to be 0 at the time point  $t_0$ . We
- then derive the  $100(1-\alpha)\%$  confidence band of  $\beta^{\langle l,s\rangle}(t)$ . We adopt the Bayes model in (Gu, 2013)
- and get the posterior variance of  $\beta^{\langle l,s \rangle}(t)$  satisfies

$$\operatorname{Var}\left[\beta^{\langle l,s\rangle}(t) \mid \mathbf{X}, \mathbf{X}^{\langle l\rangle}\right] = \frac{\sigma^2}{nK\lambda} \left(R_1(t,t) + \boldsymbol{\psi}^T (\mathbf{S}^T \mathbf{M}^{-1} \mathbf{S})^{-1} \boldsymbol{\psi} - 2\boldsymbol{\psi}^T \boldsymbol{d}_{\xi} - \boldsymbol{\xi}^T \boldsymbol{c}_{\xi}\right)$$
(16)

13 where

$$c_{\xi} = \left(\mathbf{M}^{-1} - \mathbf{M}^{-1\zeta} \left(\mathbf{S}^{T} \mathbf{A}^{-1} \mathbf{S}\right)^{-1} \mathbf{S}^{T} \mathbf{M}^{-1}\right) \mathbf{X}^{\langle s \rangle} \boldsymbol{\xi}$$

$$d_{\xi} = \left(\mathbf{S}^{T} \mathbf{M}^{-1} \mathbf{S}^{\chi}\right)^{-1} \mathbf{S}^{T} \mathbf{M}^{-1} \mathbf{X}^{\langle s \rangle} \boldsymbol{\xi}$$
(17)

- Using equation (16), we can estimate the posterior variance of  $\beta^{(l,s)}(t_0)$  and write as  $\gamma^{(l,s)}(t_0)$ .
- We then construct the 100(1  $\alpha^{5/2}$ ) Bayesian confidence interval (BCI) of  $\beta^{\langle l,s\rangle}(t)$ :  $BCI_{\langle l,s\rangle}(t):=$
- $_{16}$   $\hat{\beta}^{\langle l,s\rangle}(t)\pm z_{\alpha/2}\sqrt{\gamma^{\langle l,s\rangle}(t)}, \text{ vhen } z_{\alpha/2} \text{ is the } 1-\alpha/2 \text{ quantile for standard normal distribution.}$
- We use Bayesian com lence intervals to construct the dynamic graph, where a node represents
- a gene, and are edge between two nodes exists if the two corresponding genes follow the model (11)
- with non-zero coefficient  $\beta^{\langle l,s\rangle}(t)$ .

#### 1 2.2 Test differentially expressed genes

- <sup>2</sup> We integrate a functional ANOVA test method (Górecki and Smaga, 2019) in our framework to
- 3 estimate deferentially expressed genes based on the constructed cell fate trajectories. For each gene
- 4 in those cell fate trajectories, we consider independent vectors of the random function  $\mathbf{X}_{ki}(t) =$
- 5  $(X_{ki1}(t), \dots, X_{kid}(t))^{\top}$ , where k indicates the number of trajectory groups, i indicates cells and d
- 6 indicates the number of genes in one individual cell trajectory, defined over the interval I. In the

- 7 multivariate analysis of variance problem for functional data (FMANOVA), we test the following
- 8 hypothesis

$$H_0: \boldsymbol{\mu}_1(t) = \dots = \boldsymbol{\mu}_k(t), t \in I,$$

$$H_A: \boldsymbol{\mu}_1(t) \neq \dots \neq \boldsymbol{\mu}_k(t), t \in I.$$
(18)

Wilk's lambda test statistics for testing significantly different genes are approximated using the fdANOVA method (Górecki and Smaga, 2019). The null distributions of test statistics are approximated by  $F_{(l-1)\kappa,(n-l)\kappa}$ -distribution,  $\kappa$  are estimated by the naive and biased-reduced methods (Zhang, 2014). The p-value is given by  $P\left(F_{(l-1)\kappa,(n-l)\kappa} > F_n\right)$ , where  $F_n$  denotes the test statistic. P-values for all genes tested are corrected by Benjamini & Y kutieli method (Benjamini and Yekutieli, 2001).

### 15 3 Results

We evaluated the performance of the CellST framev ork in cell lineage tracing and cell fate predic-16 tion in both simulated and real scRNA-seq ex eriments. The simulation analysis was conducted 17 in two scenarios: Firstly, we simulate a sc RNA-seq datasets with cells at two time points to only 18 investigate the accuracy of constructed cell-to-cell correspondence between time points. Secondly, 19 we simulated a time-course scRNA req dataset with multiple time points to examine individual cell 20 differentiation patterns in the collaste trajectories. For real scRNA-seq experiments, we conducted 21 cell lineage tracing in a sangle-cell mouse hematopoietic system experiment (Weinreb et al., 2018). 22 Moreover, we evaluated the entire proposed framework on a scRNA-seq experiment for zebrafish 23 cell embryoge. es. (N.acosko et al., 2015).

### 2 3.1 Simulation

#### 3 3.1.1 Reconstruct cell-to-cell correspondence along time points

To investigate the accuracy of cell aligning, we simulated scRNA-seq experiments with only two different time points. The simulated datasets, which contain the same number of cells and cell types, were generated independently for each time point. These simulation datasets contain five same cell types in both time points. In the simulation setting, the number of cells ranges from 200 to 600, and the number of genes in one cell ranges from 100 to 500. The cell alignment and



Figure 3: A simulation example of cell aligning process at two time points. **a**: Cell-to-cell alignment with five (right) cell types in both time points. **b**: Accuracy comparison of the cell aligning process (red) with other gene similarity measurements (Pearson correlation (blue) and Euclidean distance (green)).

cell-to-cell correspondences were constructed using the CellST has donly on the gene expression information of cells at each time point and no information on he penchmark labels of cell types. 10 Specifically, we estimated an empirical transportation cost for the individual cell alignment between 11 two time points using the gene expressions in cells. We aligned cells by selecting pairs with the 12 smallest transportation cost. Since cell dynaric is a continuous development process and cells 13 within the same cell type tend to have cimi ar gene expression profiles, the cell aligning accuracy 14 can be validated by counting the number of aligned cell pairs with the same cell type (Figure 3b). 15 We noticed that the accuracy of the cell aligning method has an increasing trend as we added 16 more genes in cells for the simal at a data. This observation is due to the fact that the CellST 17 gets more information to le rn the patterns of genes when more genes are simulated in each cell. Similarly, increasing cen numbers will also increase the aligning accuracy since cells can be treated 19 as information r pn ares to enhance the accuracy. We also compared the accuracy of coupled cells with the Lucadean distance and Pearson's correlation. Those two methods are the most commonly used distances or similarity measures for gene expression analysis. (Angermueller et al., 2016; Klimovskaia et al., 2020; Skinnider et al., 2019). The accuracy comparison results (Figure 3b) show the CellST method achieves the best cell aligning accuracy in the simulation settings. In summary, the cell aligning method achieves high accuracy and captures the significant gene expressions when aligning cells and constructing individual cell correspondences at two different time points. Accurate cell alignment is crucial for the down-streaming individual cell fate prediction when aligned cells are transformed into a trajectory over time.

#### 3.1.2 CellST estimate cell fates in two simulated pathways

To investigate the effectiveness in predicting the individual cell fates, we simulated a time-course scRNA-seq data with 13 experimental time stamps, and 160 cells were simulated at each time stamp. This simulation dataset has two pathways with distinct development expression patterns, and each pathway contains 100 genes. The first pathway was created using the contact inhibition genes that keep cells growing into only a layer one cell thick (mono-layer) (Pavel et al., 2018; Mendonsa et al., 2018). The growth of cells' average expression in this simulated pathway diminishes and approaches an equilibrium expression over time. We simulated the second pathway according to the cellular division process, which is more active in cells under mitosis and a sactive in cells in interphase (Tomasetti et al., 2017). Eighty cells contain only the contact inhibition pathway at each time point, and eighty cells contain only the cellular division pathway.

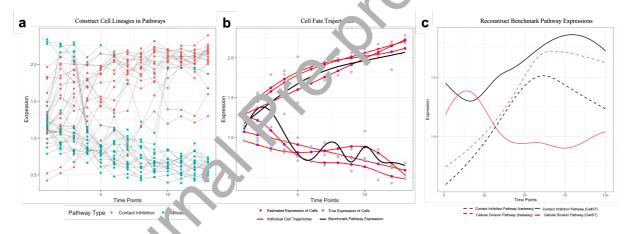


Figure 4: **a**: Cell couplines through all time points constructed by the CellST method. The cells are classified by in points through all time points constructed by the CellST method. The cells are classified by in points are classified by in points and be normally average expression cell trajectory(black curve). **c**: Development expression patterns for a small attended gene (m\_67). The red and black curves estimated by the CellST method indicate the general expression in two different pathways. The dotted two curves are constructed by the tradeseq method.

.

To observe and predict the dynamic cell fates, we utilized cells' experimental time information and built individual cell fate trajectories using CellST. The cell lineages between adjacent time points were constructed using cell lineage tracing in CellST (Figure 4a). Those connected cells were then smoothed using the smoothing spline technique in CellST to estimate the cell fate trajectories. Figure 4b illustrates the estimated individual cell fate trajectories (red curves). Based on the two distinct pathways, the two types of cells are automatically well separated by CellST. The expressions

of cell fate trajectories were compared with the benchmark pathway expression patterns (black curves). The expression of cell fate trajectories illustrates consistent patterns with the benchmark 11 expression of the two simulated pathways over time. In addition to the consistency, we observed that cells have unique behaviors over time from the cell fate trajectories. Some cells grow slower 13 and have lower expression values, while others grow faster and have higher expression values than the simulated average development patterns. The cell fate trajectories predict the unique cell 15 development behaviors by smoothing the constructed cell lineages to reduce cell-cell variance. 16 Next, we performed a comparative analysis of CellST with the existing trajectories analysis 17 method "tradeseq" (Van den Berge et al., 2020). The "tradeseq" is a trajectory-based method to 18 estimate the dynamic expressions of differentially expressed gences. By comparing the gene expres-19 sion patterns constructed by CellST and tradeseq (Figure 4c), remotice that the tradeseq method 20 constructed two similar expression patterns for a simular d gene expression (dotted curves), while 21 the CellST method built two distinct expression patterns ( lack and red curves). Those constructed 22 dynamic gene expression curves by CellST are also consistent with the simulated benchmark ex-23 pressions by showing distinct expression potents. When constructing the cell fate trajectories, the CellST method can automatically essify cells that contain different pathway expressions and 25 construct cell correspondences within the same pathway. 26

### 27 3.2 Real scRNA-Sec Experiments

#### 28 3.2.1 CellST construct accurate cell lineages

We applied the proposed method on a single cell mouse hematopoietic system experiment to evaluate the effectiven ss of constructing cell lineages (Weinreb et al., 2020). The dataset includes three experimental time points and the cells defined a continuous state map spanning from multipotent progenitors (MPPs) to nine mature cell types, including erythrocytes (Er), megakaryocytes (Mk), basophils (Ba), mast cells (Ma), eosinophils (Eos), neutrophils (Neu), monocytes (Mo), plasmacytoid dendritic cells (pDCs), Ccr7+ migratory DCs (migDCs), and lymphoid precursors (Ly) (Figure 5a). We constructed the cell developing lineages for neutrophils (Neu) cell type and compared the results with the benchmark cell pseudo-time from the original experiment (Figure 5a). CellST connected cells through the three experimental time points to represent the developing process's

- 1 penitential cell lineages. Specifically, we estimated an empirical transportation cost for the indi-
- 2 vidual cell correspondence between two time points using the gene expressions in cells and then
- 3 aligned cells by selecting pairs with the smallest transportation cost.

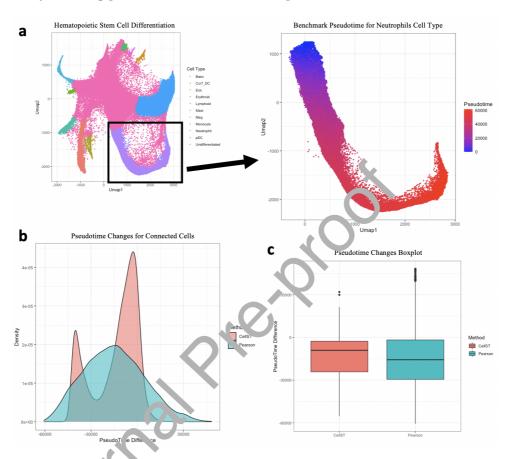


Figure 5: Constructing a dividual cell lineages with mouse hematopoietic system. **a**: scRNA-seq dataset for hematopoietic stem cell differentiation. We specifically focus on the cell differentiation of neutrophils (Nc.) and are cell types. **b**: The distribution pseudo time difference between aligned individual cells letween two time points. The distribution has been compared with the Pearson Correlation method, which measures the similarities between cells. **c**: Boxplot comparison for the pseudo time difference between CellST and Pearson correlation method.

- Since cell development is a gradual process and cells within the adjacent time points tend to have
- 5 similar gene expression profiles, we measured the differences from cell pseudo-time in all constructed
- 6 cell lineages. We compared the distribution of cell pseudo-time with Pearson's correlation method,
- <sup>7</sup> which has been widely used to measure similarity between two cells (Angermueller et al., 2016;
- 8 Klimovskaia et al., 2020; Skinnider et al., 2019). Comparing with Pearson's correlation method,
- 9 we observed that CellST constructs cell lineages with higher cell similarities, in which the changes

in pseudo-time of connected cells are gathering near zero (Figure 5b and c). The results indicate
that CellST can observe the gradually step-wise developing behaviors of cells at each time point.
CellST constructed cell lineages based on the cell-cell correspondence connection to represent the
cell-developing behaviors throughout the time points.

#### 14 3.3 Discover critical genes in zebrafish cell embryogenesis

To further investigate individual cell differentiation behaviors and gene-gene relationships, we per-15 formed CellST on a zebrafish embryogenesis scRNA-seq dataset. This dataset contains 38,731 cells 16 and 11,588 genes of early zebrafish development using Drop-seq (Macosko et al., 2015). Samples in the dataset are from the high blastula stage (3.3 hours post in the ation) when most cells are 18 pluripotent, to the six-somite stage (12 hours post-fertilize no.), when many cells have differentiated into different cell types. We observed that cells beer clustered together at the beginning 20 high blastula stage and differentiated into different all types in later development stages. Since 21 we are constructing cell trajectories for multiple ell types simultaneously, the proliferation rates 22 potentially vary across cells. To address this, we unlize our generalized CellST method to construct 23 cell fate trajectories (Figure 6a), which apture the unique individual cell development behaviors under more general scenarios. The method also enables us to estimate the proliferation rate for 25 each individual cell. The histogram of the normalized proliferation rates at different time points are presented in Figure 6b. Note that the mean of the normalized proliferation rates at each time 27 point equals 1, while the variance indicates the heterogeneity level of the proliferation rates among 28 cells at the corresponding time point. We notice that this heterogeneity is significantly high in the early stages of the cell differentiation process and decreases gradually over time.

Unlike the bulk cell trajectory, the CellST cell fate trajectories achieved full cell development coverage for all cells. The full coverage indicates that the cell fate trajectories can reveal less frequent cell development patterns overlooked by the bulk cell trajectory. The CellST constructed cell fate trajectories throughout the stages and illustrated the unique individual cell development behaviors. The cell fate trajectories return each cell's potential cell fate paths into different cell types throughout the 12 developmental stages.

Furthermore, as cells developed into multiple cell types at the 12.0-6-somite stage (last developmental stage), we built trajectory groups to those cell fate paths by CellST according to the cell

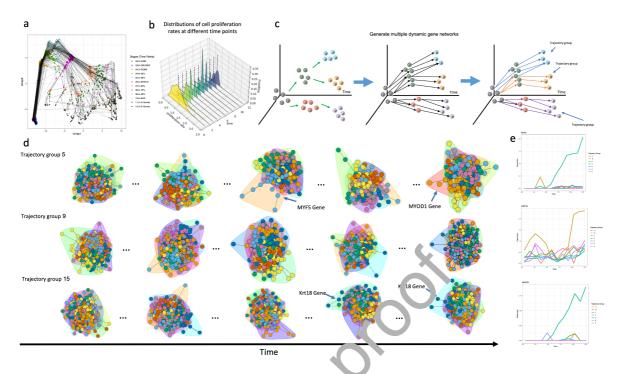


Figure 6: **a**: We constructed the individual cell fate trajectories by connecting cells through the 12 developing stages. **b**: Histograms of the distribution of the cell proliferation rates at different time points. **c**: Illustration of identifying multiplus cell trajectory groups based on the cell types at the 12.0-6-somite stage (last development a stage). **d**: Example of critical genes identified using CellST dynamic gene network. **e**: Dynamic gene networks constructed by CellST for the cell trajectory groups.

- types in the last developmer. Is 2 ge. We constructed the dynamic gene networks (Figure 6c) for each group of cell fate treecte ies. In those dynamic networks, we observed some genes that behave significantly differently firm other genes (Figure 6c) as cells developed into different cell types. For instance, MYF5, MY )D1, and KRT18 genes appeared to behave differently in two of the trajectory groups in later developmental stages. The MYF5 is a protein with a key role in regulating muscle differentiation or myogenesis, specifically the development of skeletal muscle (Esteves de Lima and Relaix, 2021; Agarwal et al., 2022), and MyoD1 is a key regulator that orchestrates skeletal muscle differentiation through the regulation of gene expression (Blum et al., 2012; Agaram et al., 2019). Moreover, KRT18 regulates the epithelial cell differentiation process (Jiang et al., 2020; Liu et al., 2021).
- We then visualized and validated the expressions of those critical genes (Figure 6d and Figure
- <sub>9</sub> 6e). The expressions of MYF5 and MYOD1 genes are significantly higher in trajectory group 5

versus in other trajectory groups, which is consistent with the discovery in the CellST dynamic networks. KRT18 is highly expressed in trajectory group 15, which is also consistent with the CellST dynamic network results. Additionally, we performed functional deferentially expressed gene tests based on the CellST cell fate trajectories. We discovered a total of 268 differentially expressed genes in this zebrafish cell development process dataset. We performed gene ontology annotations to those genes (Table 1), and the function of those genes is highly related to regulating the cell development/differentiation process.

Gene Ontology (GO) annotations	Count P-value
Multicellular organism development	133 2.100393e-45
Cell Differentiation	86 4.8e-25
Differentiation	36 8.6e-10
Cell fate specification	13 1.3e-5
Cell fate commitment	12 6.1e-5

Table 1: Top five gene functional annotation groups.

Those results proved that the cell fate trajectories and dynamic gene networks in the CellST method can be used to discover critical gener in a call differentiation process. We also demonstrated the CellST cell fate trajectories have from coverage on different cell lineages even in some rare cell types since the trajectories took individual cell behaviors. Lastly, those individual cell fate trajectories reflect unique gene expression patterns when cells develop into different mature cell types.

### 1 4 Discussion

16

- Understanding the dynamic of cell differentiation throughout a period is crucial for future research in scRNA-Seq analysis. We developed a novel machine learning analysis framework, CellST, to build cell fate trajectories and dynamic gene networks for time-course scRNA-seq datasets. The cell fate trajectories enabled researchers to observe the individual cell development behaviors and better use the benefit of the single cell sequencing technology. Compared to the existing bulk single-cell trajectory, we brought the cell development analysis into a more precise and unprecedented resolution. The dynamic gene networks estimated the dynamic relationship of genes and discovered potential critical genes during cell differentiation processes.
  - 19

There are three major advantages of the CellST analysis framework. Firstly, the cell lineages were constructed with high accuracy and provides unique individual cell differentiation behaviors between time points. Secondly, since the trajectory tracks individual cells, the cell fate trajectories will have full coverage on different cell lineages even in some rare cell types (Figure 6a). Thirdly, the dynamic gene networks analysis in the CellST framework can accurately estimate gene-gene relationships and discover critical genes in the cell differentiation process. Through the simulation and real dataset analysis, We constructed cell fate trajectories in single cell RNA-seq experiment and various dynamic cell differentiation behaviors were observed.

## 18 5 Data Availability

- 19 All data used in this paper are publicly available datagets. The mouse hematopoietic system
- experiment dataset can be found on GitHub (https://gith.b.com/AllonKleinLab/paper-data).
- 21 The zebrafish embryogenesis dataset from the original paper (Wagner et al., 2018) can be found in
- 1 the NCBI database with accession number CS 7.06 587.

## <sup>2</sup> 6 Code Availability

- 3 The CellST R package and example for constructing the cell fate trajectories and dynamic gene
- 4 networks analysis can be found on GitHub (https://github.com/zhanzmr/CellST).

## 5 7 Acknowledgements

- 6 This work was supported by National Science Foundation grants DMS-1903226, DMS-1925066,
- 7 DMS-2124493 and NIH grants R01GM1222080.

#### 8 Author information

- 9 Corresponding authors: Correspondence to Jingyi Zhang and Ping Ma.
- 10 Co-first author: Mengrui Zhang, Yongkai Chen

#### 11 References

- 12 Narasimhan P Agaram, Michael P LaQuaglia, Rita Alaggio, Lei Zhang, Yumi Fujisawa, Marc
- Ladanyi, Leonard H Wexler, and Cristina R Antonescu. Myod1-mutant spindle cell and scle-
- rosing rhabdomyosarcoma: an aggressive subtype irrespective of age. a reappraisal for molecular
- classification and risk stratification. *Modern pathology*, 32(1):27–36, 2019.
- 16 Megha Agarwal, Anushree Bharadwaj, and Sam J Mathew. Tle4 regulates muscle stem cell quies-
- cence and skeletal muscle differentiation. Journal of Cell Science, 135(4):jcs256008, 2022.
- Michael Alonge, Xingang Wang, Matthias Benoit, Sebastian Soyk, Lara Pereira, Lei Zhang, Ham-
- sini Suresh, Srividya Ramakrishnan, Florian Maumus, Danie le Ciren, et al. Major impacts of
- widespread structural variation on gene expression and c op improvement in tomato. Cell, 182
- (1):145–161, 2020.
- <sup>22</sup> Christof Angermueller, Stephen J Clark, Heather J Lee, tain C Macaulay, Mabel J Teng, Tim Xi-
- 23 aoming Hu, Felix Krueger, Sébastien A Smal voo l, Chris P Ponting, Thierry Voet, et al. Parallel
- single-cell sequencing links transcriptional and epigenetic heterogeneity. Nature Methods, 13(3):
- 25 229–232, 2016.
- Yoav Benjamini and Daniel Ye sv. i. The control of the false discovery rate in multiple testing
- under dependency. Ann. s of Statistics, pages 1165–1188, 2001.
- 1 Roy Blum, Vasupradha Vetnantham, Christopher Bowman, Michael Rudnicki, and Brian D Dyn-
- lacht. Gerom -wic  $\varepsilon$  identification of enhancers in skeletal muscle: the role of myod1. Genes  $\mathcal B$
- $development, \ 23(24):2763-2779, \ 2012.$
- <sup>4</sup> Natalie Burrows, Rachael JM Bashford-Rogers, Vijesh J Bhute, Ana Peñalver, John R Ferdinand,
- 5 Benjamin J Stewart, Joscelin EG Smith, Mukta Deobagkar-Lele, Girolamo Giudice, Thomas M
- 6 Connor, et al. Dynamic regulation of hypoxia-inducible factor- $1\alpha$  activity is essential for normal
- <sup>7</sup> b cell development. *Nature Immunology*, 21(11):1408–1420, 2020.
- 8 Robrecht Cannoodt, Wouter Saelens, and Yvan Saeys. Computational methods for trajectory
- 9 inference from single-cell transcriptomics. European Journal of Immunology, 46(11):2496–2506,
- 10 2016.

- 11 Huidong Chen, Luca Albergante, Jonathan Y Hsu, Caleb A Lareau, Giosuè Lo Bosco, Jihong
- Guan, Shuigeng Zhou, Alexander N Gorban, Daniel E Bauer, Martin J Aryee, et al. Single-
- cell trajectories reconstruction, exploration and mapping of omics data with stream. Nature
- Communications, 10(1):1-14, 2019.
- 15 Nicolas Courty, Rémi Flamary, and Devis Tuia. Domain adaptation with regularized optimal
- transport. In Joint European Conference on Machine Learning and Knowledge Discovery in
- 17 Databases, pages 274–289. Springer, 2014.
- 18 Karren Dai, Karthik Damodaran, Saradha Venkatachalapathy, Ali C Soylemezoglu, GV Shiv-
- ashankar, and Caroline Uhler. Predicting cell lineages using autoencoders and optimal transport.
- PLoS Computational Biology, 16(4):e1007828, 2020.
- 21 Joana Esteves de Lima and Frédéric Relaix. Master regul tors of skeletal muscle lineage develop-
- ment and pluripotent stem cells differentiation. Cel' R generation, 10:1–13, 2021.
- 23 Tomasz Górecki and Łukasz Smaga. fdanova: an software package for analysis of variance for
- univariate and multivariate functio al data. Computational Statistics, 34(2):571–597, 2019.
- Dominic Grün and Alexander Oug naarden. Design and analysis of single-cell sequencing experi-
- ments. Cell, 163(4):799–810, 2613.
- 1 Chong Gu. Smoothing vline ANOVA models, volume 297. Springer Science & Business Media,
- 2 2013.
- 3 Chong Gu and Ph. g Ma. Optimal smoothing in nonparametric mixed-effect models. The Annals
- of Statistics, 23(3):1357–1379, 2005.
- <sup>5</sup> Jingtao Guo, Edward J Grow, Chongil Yi, Hana Mlcochova, Geoffrey J Maher, Cecilia Lindskog,
- 6 Patrick J Murphy, Candice L Wike, Douglas T Carrell, Anne Goriely, et al. Chromatin and
- <sup>7</sup> single-cell rna-seq profiling reveal dynamic signaling and metabolic transitions during human
- spermatogonial stem cell development. Cell Stem Cell, 21(4):533–546, 2017.
- 9 Sinisa Hrvatin, Daniel R Hochbaum, M Aurel Nagy, Marcelo Cicconet, Keiramarie Robertson,
- Lucas Cheadle, Rapolas Zilionis, Alex Ratner, Rebeca Borges-Monroy, Allon M Klein, et al.

- Single-cell analysis of experience-dependent transcriptomic states in the mouse visual cortex.
- Nature Neuroscience, 21(1):120–129, 2018.
- <sup>13</sup> Zhicheng Ji and Hongkai Ji. Tscan: Pseudo-time reconstruction and evaluation in single-cell rna-seq
- analysis. Nucleic Acids Research, 44(13):e117-e117, 2016.
- 15 Peng Jiang, Rafael Gil de Rubio, Steven M Hrycaj, Stephen J Gurczynski, Kent A Riemondy,
- Bethany B Moore, M Bishr Omary, Karen M Ridge, and Rachel L Zemans. Ineffectual type
- 2-to-type 1 alveolar epithelial cell differentiation in idiopathic pulmonary fibrosis: persistence
- of the krt8hi transitional state. American journal of respiratory and critical care medicine, 201
- 19 (11):1443–1447, 2020.
- 20 Anna Klimovskaia, David Lopez-Paz, Léon Bottou, and Maximilian Nickel. Poincaré maps for
- 21 analyzing complex hierarchies in single-cell data. Natur Communications, 11(1):1–9, 2020.
- Devon A Lawson, Nirav R Bhakta, Kai Kessenbrock, Karin D Prummel, Ying Yu, Ken Takai,
- Alicia Zhou, Henok Eyob, Sanjeev Balakris nar, Chih-Yang Wang, et al. Single-cell analysis
- reveals a stem-cell program in hum in m stastatic breast cancer cells. *Nature*, 526(7571):131–135,
- 25 2015.
- <sup>26</sup> Yufan Liu, Jianjun Li, Bin Yao, Yin ii Wang, Rui Wang, Siming Yang, Zhao Li, Yijie Zhang, Sha
- Huang, and Xiaobing Fu. The stiffness of hydrogel-based bioink impacts mesenchymal stem cells
- differentiation toward sw at glands in 3d-bioprinted matrix. Materials Science and Engineering:
- C, 118:111387  $\angle 0$ . 1.
- <sup>4</sup> Zehua Liu, Huazi e Lou, Kaikun Xie, Hao Wang, Ning Chen, Oscar M Aparicio, Michael Q Zhang,
- 8 Rui Jiang, and Ting Chen. Reconstructing cell cycle pseudo time-series via single-cell transcrip-
- tome data. Nature Communications, 8(1):1-9, 2017.
- <sup>7</sup> Evan Z Macosko, Anindita Basu, Rahul Satija, James Nemesh, Karthik Shekhar, Melissa Goldman,
- 8 Itay Tirosh, Allison R Bialas, Nolan Kamitaki, Emily M Martersteck, et al. Highly parallel
- genome-wide expression profiling of individual cells using nanoliter droplets. Cell, 161(5):1202-
- 10 1214, 2015.

- <sup>11</sup> Alisha M Mendonsa, Tae-Young Na, and Barry M Gumbiner. E-cadherin in contact inhibition and
- cancer. Oncogene, 37(35):4769-4780, 2018.
- 13 Cheng Meng, Yuan Ke, Jingyi Zhang, Mengrui Zhang, Wenxuan Zhong, and Ping Ma. Large-scale
- optimal transport map estimation using projection pursuit. In Advances in Neural Information
- 15 Processing Systems, pages 8116–8127, 2019.
- 16 Kevin R Moon, Jay S Stanley III, Daniel Burkhardt, David van Dijk, Guy Wolf, and Smita Krish-
- naswamy. Manifold learning-based methods for analyzing single-cell rna-sequencing data. Current
- Opinion in Systems Biology, 7:36–46, 2018.
- <sup>19</sup> Tal Nawy. Single-cell sequencing. Nature Methods, 11(1):18, 20 3.
- Mariana Pavel, Maurizio Renna, So Jung Park, Fiona M Menzies Thomas Ricketts, Jens Füllgrabe,
- Avraham Ashkenazi, Rebecca A Frake, Alejandro Cornicer Lombarte, Carla F Bento, et al.
- 22 Contact inhibition controls cell survival and proliferation via yap/taz-autophagy axis. Nature
- 23 Communications, 9(1):1–18, 2018.
- <sup>24</sup> Xiaojie Qiu, Andrew Hill, Jonathan F. c'.er, Dejun Lin, Yi-An Ma, and Cole Trapnell. Single-cell
- mrna quantification and differential analysis with census. Nature Methods, 14(3):309, 2017.
- Gang Ren, Wenfei Jin, Kai ang Cui, Joseph Rodrigez, Gangqing Hu, Zhiying Zhang, Daniel R
- Larson, and Keji Zha Cof-mediated enhancer-promoter interaction is a critical regulator of
- cell-to-cell variation of tene expression. Molecular Cell, 67(6):1049–1058, 2017.
- Wouter Saelen, A. Lecht Cannoodt, Helena Todorov, and Yvan Saeys. A comparison of single-cell
- trajectory in Lence methods. Nature Biotechnology, 37(5):547–554, 2019.
- <sup>5</sup> Geoffrey Schiebinger, Jian Shu, Marcin Tabaka, Brian Cleary, Vidya Subramanian, Aryeh Solomon,
- Joshua Gould, Siyan Liu, Stacie Lin, Peter Berube, et al. Optimal-transport analysis of single-
- cell gene expression identifies developmental trajectories in reprogramming. Cell, 176(4):928–943,
- s 2019.
- 9 Ehud Shapiro, Tamir Biezuner, and Sten Linnarsson. Single-cell sequencing-based technologies will
- revolutionize whole-organism science. Nature Reviews Genetics, 14(9):618–630, 2013.

- <sup>11</sup> Michael A Skinnider, Jordan W Squair, and Leonard J Foster. Evaluating measures of association
- for single-cell transcriptomics. *Nature Methods*, 16(5):381–386, 2019.
- David G Spiller, Christopher D Wood, David A Rand, and Michael RH White. Measurement of
- single-cell dynamics. *Nature*, 465(7299):736–745, 2010.
- 15 Oliver Stegle, Sarah A Teichmann, and John C Marioni. Computational and analytical challenges
- in single-cell transcriptomics. Nature Reviews Genetics, 16(3):133–145, 2015.
- 17 Amos Tanay and Aviv Regev. Scaling single-cell genomics from phenomenology to mechanism.
- Nature, 541(7637):331–338, 2017.
- 19 Cristian Tomasetti, Rick Durrett, Marek Kimmel, Amaury Lanbett, Giovanni Parmigiani, Ann
- Zauber, and Bert Vogelstein. Role of stem-cell divisions in cancer risk. *Nature*, 548(7666):
- E13-E14, 2017.
- 22 Alexander Tong, Jessie Huang, Guy Wolf, Da 10, van Dijk, and Smita Krishnaswamy. Trajecto-
- rynet: A dynamic optimal transport pet ork for modeling cellular dynamics. arXiv Preprint
- 24 arXiv:2002.04461, 2020.
- <sup>25</sup> Cole Trapnell. Defining cell type: and states with single-cell genomics. Genome Research, 25(10):
- 1491–1498, 2015.
- Cole Trapnell, Davide Cachiarelli, Jonna Grimsby, Prapti Pokharel, Shuqiang Li, Michael Morse,
- Niall J Lennon Kennth J Livak, Tarjei S Mikkelsen, and John L Rinn. The dynamics and
- regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. Nature
- 4 Biotechnolog, 32(4):381, 2014.
- 5 Sophie Tritschler, Maren Büttner, David S Fischer, Marius Lange, Volker Bergen, Heiko Lickert,
- and Fabian J Theis. Concepts and limitations for learning developmental trajectories from single
- 7 cell genomics. Development, 146(12), 2019.
- 8 Koen Van den Berge, Hector Roux De Bezieux, Kelly Street, Wouter Saelens, Robrecht Cannoodt,
- 9 Yvan Saeys, Sandrine Dudoit, and Lieven Clement. Trajectory-based differential expression
- analysis for single-cell sequencing data. Nature Communications, 11(1):1–13, 2020.

- 11 Cédric Villani. Topics in optimal transportation. American Mathematical Soc., 2003.
- Daniel E Wagner, Caleb Weinreb, Zach M Collins, James A Briggs, Sean G Megason, and Allon M
- Klein. Single-cell mapping of gene expression landscapes and lineage in the zebrafish embryo.
- 14 Science, 360(6392):981–987, 2018.
- 15 Grace Wahba. Spline Models for Observational Data, volume 59 of CBMS-NSF Regional Conference
- Series in Applied Mathematics. SIAM, Philadelphia, 1990.
- Jane-Ling Wang, Jeng-Min Chiou, and Hans-Georg Müller. Functional data analysis. Annual
- 18 Review of Statistics and Its Application, 3:257–295, 2016.
- <sup>1</sup> Caleb Weinreb, Samuel Wolock, Betsabeh K Tusi, Merav Soccovs y, and Allon M Klein. Fun-
- damental limits on dynamic inference from single-cell staps of the National
- 3 Academy of Sciences, 115(10):E2467–E2476, 2018.
- 560 Caleb Weinreb, Alejo Rodriguez-Fraticelli, Fe n. 1do D Camargo, and Allon M Klein. Lineage
- tracing on transcriptional landscapes lin's state to fate during differentiation. Science, 367
- 562 (6479):eaaw3381, 2020.
- J Zhang. Analysis of variance for functional data. Monographs on Statistics and Applied Probability,
- 127:127, 2014.
- Jingyi Zhang, Wenxuar Zhong, and Ping Ma. A review on modern computational optimal transport
- methods with appliest one in biomedical research. arXiv preprint arXiv:2008.02995, 2020.

## Declaration of Competing Interest

- The authors declare that they have no known competing financial interests or personal relationships
- that could have appeared to influence the work reported in this paper.

