

FRONT MATTER

scDVF: Data-driven Single-cell Transcriptomic Deep Velocity Field Learning with Neural Ordinary Differential Equations

Authors

Zhanlin Chen¹, William C. King², Mark Gerstein^{1,3,4,*}, Jing Zhang^{5,*}

Affiliations

- ¹. Department of Statistics and Data Science, Yale University, New Haven, CT 06520, USA.
- ². Healthcare and Life Sciences, Microsoft, Redmond, WA 98052, USA.
- ³. Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT 06520, USA.
- ⁴. Department of Computer Science, Yale University, New Haven, CT 06520, USA.
- ⁵. Department of Computer Science, Yale University, New Haven, CT 06520, USA.

*Corresponding authors. Email: pi@gersteinlab.org, jing.zhang@uci.edu

Abstract

Recent advances in single-cell RNA sequencing technology provided unprecedented opportunities to simultaneously measure the gene expression profile and the transcriptional velocity of individual cells, enabling us to sample gene regulatory network dynamics along developmental trajectories. However, traditional methods have been challenged in offering a fundamental and quantitative explanation of the dynamics as differential equations due to the high dimensionality, sparsity, and complex gene interactions. Here, we present scDVF, a neural-network-based ordinary differential equation that can learn to model single-cell transcriptome dynamics and describe gene expression changes across time at a single-cell resolution. We applied scDVF on multiple published datasets from different technical platforms and demonstrate its utility to 1) formulate transcriptome dynamics of different timescales; 2) measure the instability of individual cell states; and 3) identify developmental driver genes upstream of the signaling cascade. Benchmarking with state-of-the-art vector-field learning methods shows that scDVF can improve representation accuracy by at least 50%. Further, our perturbation studies revealed that single-cell dynamical systems may exhibit properties similar to chaotic systems. In summary, scDVF allows for the data-driven discovery of differential equations that delineate single-cell transcriptome dynamics.

Teaser

Using neural networks to derive the ordinary differential equations behind single-cell transcriptome dynamics.

MAIN TEXT

Introduction

Single-cell RNA-sequencing (scRNA-seq) captures a transcriptomic snapshot of a dynamic biological process. However, many current analysis methods view scRNA-seq as

47 a static dataset. For example, Monocle constructs minimum spanning trees in the cellular
48 manifold as bifurcation trajectories (1). Palantir uses Markov transition matrices to model
49 neighboring cell transitions (2). More generally, diffusion pseudotime simulates diffusion
50 to create pseudo-temporal ordering of cells in the data manifold (3). Although these
51 computational methods have been effective in highlighting the dynamics behind single-
52 cell transcriptomes, a fundamental question remains: can we derive quantitative equations
53 that accurately explain the gene expression dynamics of transitioning single cells?
54 Discovering these equations as a function of time could answer questions about the cell
55 fates and the driving forces behind developmental trajectories.

56 Recovering the dynamics from sparse and noisy scRNA-seq data is a difficult task because
57 the cells are destroyed during data collection. With the development of RNA velocity, we
58 can compute the time derivative of the expression state using the ratio of unspliced versus
59 spliced transcripts (4). However, RNA velocity only predicts the future state of cells on
60 the timescale of hours. We reasoned that it might be possible to extrapolate farther into the
61 future by piecing together information from cells at different developmental times.
62 Nevertheless, it is challenging to explicitly derive differential equations that model all
63 gene interactions. Further, evaluating the generalizability of differential equations is still
64 an open question. Previous approaches have relied on time-resolved scRNA-seq and linear
65 ordinary differential equations (ODEs) to model the dynamics of regulatory networks (5,
66 6). However, linear systems may fail to capture the non-linearity of single-cell dynamics.
67 Moreover, single-cell dynamical systems have a high degrees-of-freedom due to the high
68 dimensionality of the data, which could lead to errors in any dimension (7).

69 Inspired by recent developments in neural ODEs and data-driven dynamical systems (8,
70 9), we present a computational framework called scDVF that learns to formulate the
71 dynamics underlying scRNA-seq experiments by modeling the gene expression changes
72 of single cells across time. With a deep-learning architecture, our approach can model
73 non-linear, high-dimensional gene interactions in single-cell dynamical systems. Further,
74 we can perform *in silico* studies to explore the behavior of biological processes over time.
75 In this regard, scDVF differs substantially from most single-cell methods, in that the
76 objective of our framework is to derive neural-network-based differential equations
77 describing single-cell gene expression dynamics. To illustrate the robustness and general
78 validity of our approach, we performed analyses on developmental mouse neocortex and
79 dentate gyrus, representing scRNA-seq experiments from different tissues, technical
80 platforms, and developmental time scales (10, 11). With three additional data sources
81 (mouse pancreatic endocrinogenesis, gastrulation, developing human forebrain), we
82 demonstrate the ability for scDVF to deconvolve gene co-expression networks and
83 benchmarked our method against a state-of-the-art vector-field learning approach (4, 12,
84 13).

85 **Results**

86 **Neural ODEs for Modeling Single-cell Transcriptome Dynamics**

87 In a gene regulatory network, the expression of certain genes can increase or decrease the
88 expression of other genes. In a broader biological context, a cell transitioning along its
89 developmental trajectory can signal a cascade of gene expression changes. These gene-to-
90 gene interactions can be formulated as a function of time using differential equations.
91 More specifically, each cell represents an instance of the dynamics sampled from the
92
93

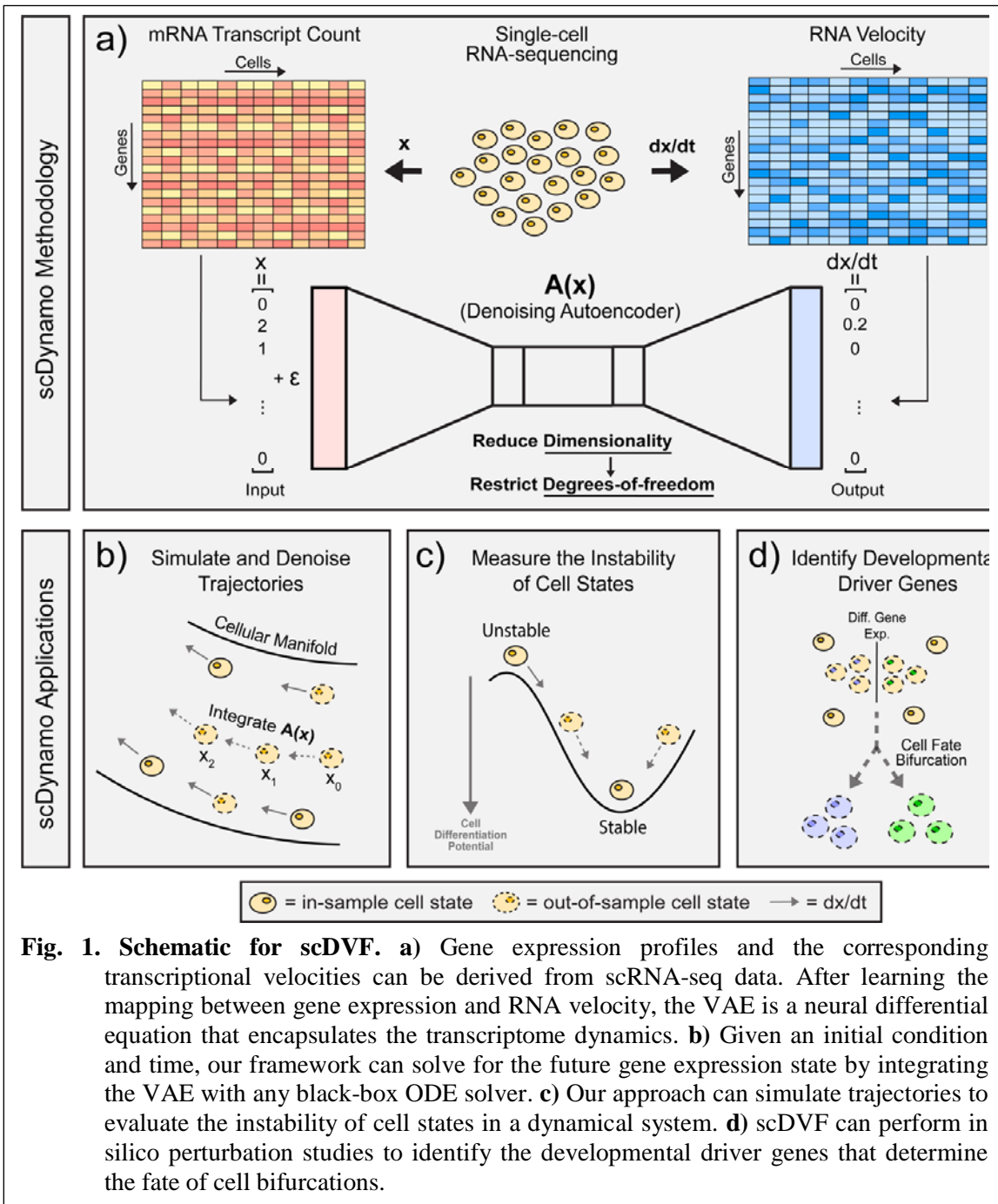


Fig. 1. Schematic for scDVF. **a)** Gene expression profiles and the corresponding transcriptional velocities can be derived from scRNA-seq data. After learning the mapping between gene expression and RNA velocity, the VAE is a neural differential equation that encapsulates the transcriptome dynamics. **b)** Given an initial condition and time, our framework can solve for the future gene expression state by integrating the VAE with any black-box ODE solver. **c)** Our approach can simulate trajectories to evaluate the instability of cell states in a dynamical system. **d)** scDVF can perform in silico perturbation studies to identify the developmental driver genes that determine the fate of cell bifurcations.

single-cell dynamical system. If the gene expression state of a cell is the vector \vec{x} , then the increase or decrease in the gene expression with respect to time is the RNA velocity vector $\frac{\partial \vec{x}}{\partial t}$. Rather than deriving a system of linear ODEs $\frac{\partial \vec{x}}{\partial t} = A\vec{x}$ with matrix A , we train a variational autoencoder (VAE) $A(\vec{x})$ to learn the mapping from the gene expression state \vec{x} to the RNA velocity $\frac{\partial \vec{x}}{\partial t}$ using data from each cell (Equation 1, Fig. 1a). Therefore, this VAE is a non-linear ODE and encapsulates the gene expression dynamics of individual cells in scRNA-seq. Then, given some initial gene expression state close to the data, we can numerically compute the future (or past) gene expression states with any black-box ODE solver. For example, given gene expression state \vec{x}_0 at time $t = 0$, we can use the Euler's method to find the gene expression state at \vec{x}_1 , and iteratively for $\vec{x}_2, \dots, \vec{x}_n$ (Equation 2, 3).

105
$$\frac{\partial \vec{x}_t}{\partial t} = A(\vec{x}_t) \quad (\text{Equation 1})$$

106
$$\vec{x}_{t+1} = \vec{x}_t + \frac{\partial \vec{x}_t}{\partial t}$$
 (Equation 2)

107
$$= \vec{x}_t + A(\vec{x}_t) \quad (\text{Equation 3})$$

108 By sequentially computing the next gene expression state, scDVF can outline the
109 developmental trajectory of single cells through time. Further, with different initial
110 conditions \vec{x}_0 , our framework can derive detailed insights into the future (or past) of
111 different cell states. Here, we explored three applications of scDVF. First, we simulated
112 and denoised developmental trajectories by extrapolating the dynamics to out-of-sample
113 cells (Fig. 1b). Second, we evaluated the instability of cell states by tracking gene
114 expression changes along simulated trajectories (Fig. 1c). Third, we performed *in silico*
115 perturbation studies to investigate how initial gene expression conditions impact the fate
116 of cell bifurcations (Fig. 1d).

117 **Deriving the Neural Equations Underlying Mouse Neocortex Development**

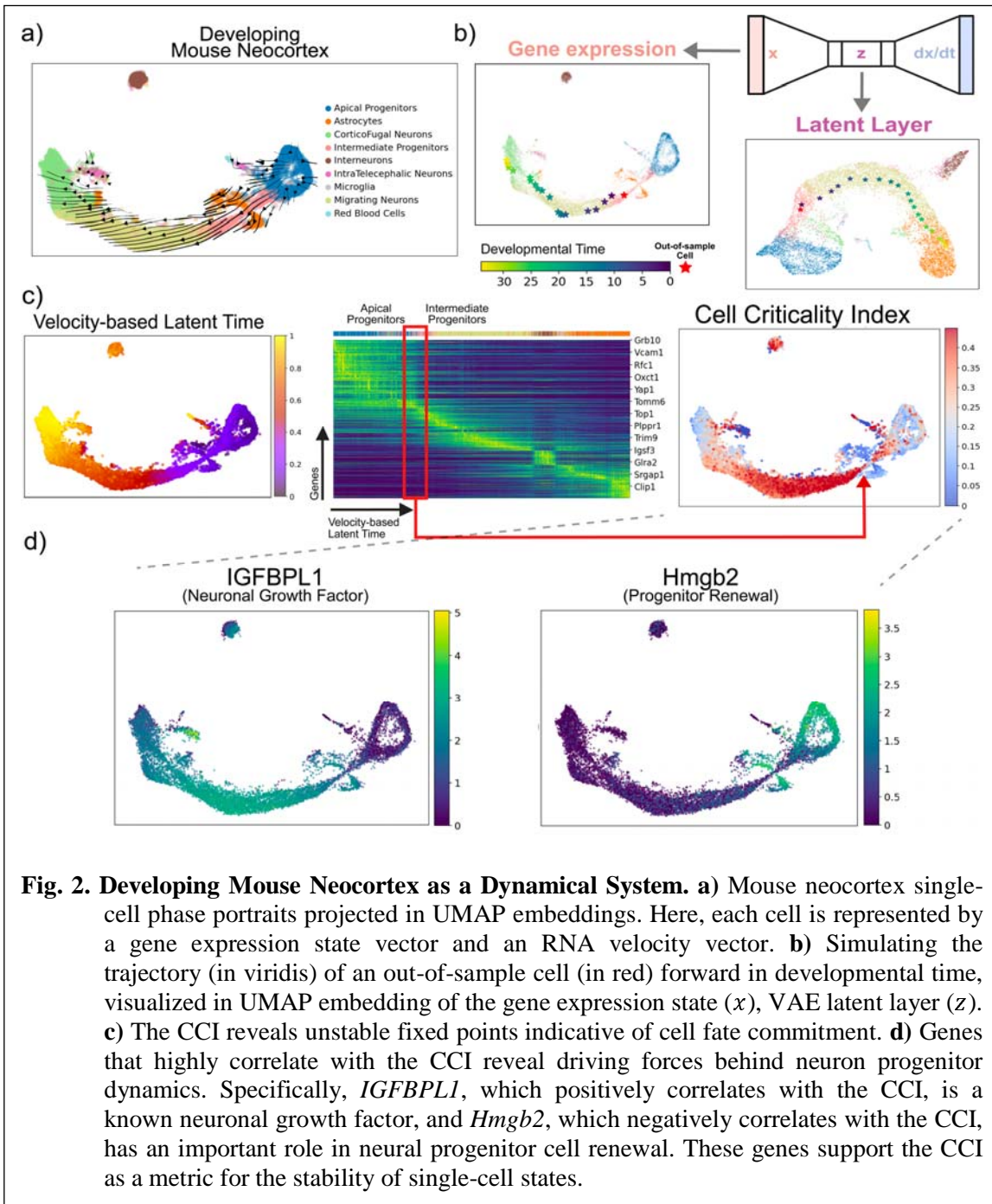
118 To evaluate whether scDVF can uncover the dynamics from sparse and noisy scRNA-seq
119 experiments, we considered a dataset of developing mouse neocortex with transcriptomes
120 profiled at E15.5 with the Chromium Single-cell 3' Library from 10x Genomics (Fig. 2a)
121 (12). Here, we show that summarizing the dynamics as neural ODEs can derive new
122 insights from the data.

123 First, we examined a hypothetical trajectory simulated from scDVF. After training the
124 VAE on neocortex cell states and velocities, we generated an out-of-sample cell as the
125 initial condition. The out-of-sample cell is simulated by adding noise to the gene
126 expression state of an existing cell, thereby representing a cell state that did not previously
127 exist in the data. Then, we incrementally solved for the future gene expression states of the
128 out-of-sample cell using scDVF. The simulated developmental path shows that our
129 predicted gene expression states moved along existing trajectories in the data manifold
130 (Fig. 2b). In the mouse neocortex, the out-of-sample cell started as an intermediate
131 progenitor, developed into a migrating neuron, and ultimately became a corticofugal
132 neuron (CFN). Further, when the VAE is solved with evenly distributed time increments,
133 the distances between intermediate states reflect the magnitude of the RNA velocity
134 vectors. Faster rates of change in gene expression generated more separated intermediate
135 states. Conversely, slower rates of change produced a denser collection of intermediate
136 points along the manifold.

137 When the VAE represents gene expression dynamics, we can visualize the latent layer
138 embeddings to gain insights into the low-dimensional dynamic manifold. Similar to gene
139 expression embeddings, the chronological and hierarchical order of developmental
140 trajectories in the latent layer are properly encoded (Fig. 2b). In the neocortex, apical
141 progenitors represent a major starting state, and CFNs represent a major terminal state.
142 The simulated cell migrates along existing trajectories in the low-dimensional dynamic
143 manifold.

144 **Characterizing Cell State Instability with the Cell Criticality Index**

145 Next, we aimed to characterize the stable and unstable fixed points of this single-cell
146 dynamical system. By looking forward in time, we can numerically approximate the
147 instability of single-cell states, which we call the cell criticality index (CCI). For a cell, the



148 CCI is defined as the cumulative information change, or the cumulative Kullback–Leibler
 149 (KL)-divergence, between gene expression distributions at each time step in the
 150 developmental trajectory. In order words, cell states that undergo large changes across
 151 time will have a high CCI, whereas cell states that only go through small changes will
 152 have a low CCI.

153 For each cell, we used scDVF to compute a developmental path such that the cell arrived
 154 at a steady terminal state. Then, we calculated the CCI along each path (Fig. 2c). The
 155 resulting developmental topology is similar to the classical Waddington landscape (14). In
 156 particular, the CCI can reveal unique topological information in the developmental
 157 landscape not directly observed in latent or pseudo time. For example, the intermediate
 158 progenitor states exhibit a higher criticality, whereas the apical progenitors and

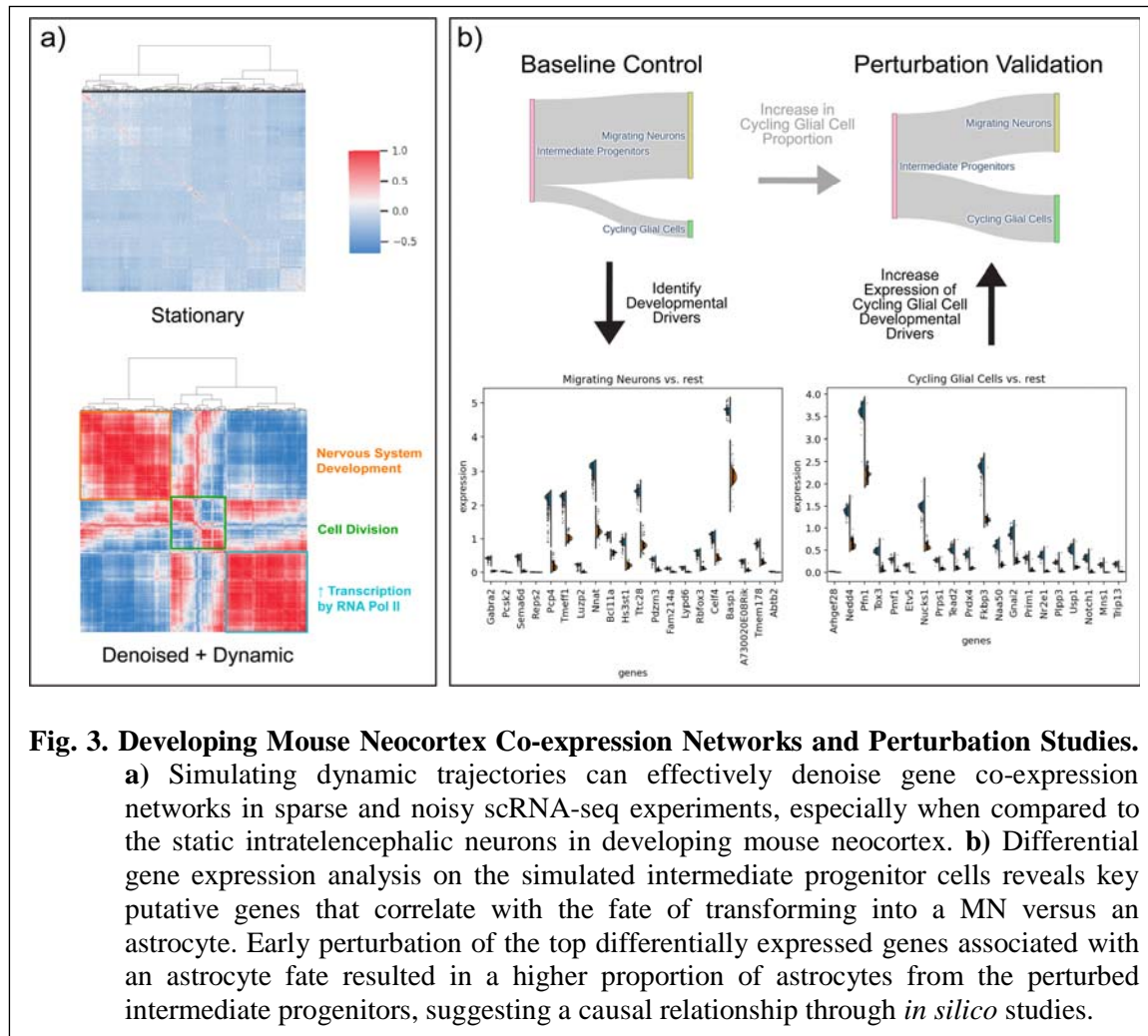
159 differentiated neuronal cell states experience a lower criticality. When progenitor cells are
160 differentiating into neuronal cell types, the heightened criticality at intermediate
161 progenitors represents fate commitment or a point of no return during development. In
162 dynamical systems, this suggests that the cell states with low criticality are located at a
163 stable fixed point, and the cell identity would remain stable even with small gene
164 expression perturbations. More interestingly, the intermediate progenitors are located at an
165 unstable fixed point with properties similar to a chaotic system in which a small
166 perturbation may result in large downstream changes. The instability of cell states can be
167 substantiated by examining the genes that best correlate with the CCI (Fig. 2d). For
168 example, previous experiments have shown that *IGFBPL1*, which positively correlates
169 with the CCI, is a known neuronal growth factor, and *Hmgb2*, which negatively correlates
170 with the CCI, has an important role in neural progenitor cell renewal (15, 16). The
171 expression of these genes supports the CCI as a metric for evaluating the instability of
172 single-cell states.

173 **Conducting *in Silico* Perturbation Studies with scDVF**

174 Lastly, we investigated the behavior of this dynamical system with similar perturbation
175 studies pioneered by (17). The goal of *in silico* perturbation studies is to computationally
176 identify which initial gene expression conditions impact the fate of cell bifurcations. In
177 short, we randomly sampled intermediate progenitors ($n = 1,000$) as the initial
178 conditions. By allowing these simulated intermediate progenitors to naturally evolve
179 according to the dynamics learned by scDVF, we observed a baseline 8:2 ratio of
180 migrating neuron (MN) versus astrocyte cell states. The ratio of future cell states indicates
181 that the MN cell state is a stronger attractive terminal state than the astrocyte cell state,
182 which corroborates with previous conclusions (11). Then, we performed differential gene
183 expression between initial conditions of different fates. The results suggest that early
184 expression perturbations in key upstream genes correlate with the fate of developmental
185 bifurcations (Fig. 3b).

186 Further, we formulated a way to perform hypothesis testing and to infer causal
187 relationships at developmental branching points (18). To investigate which developmental
188 driver genes cause progenitor cells to prefer one trajectory over another, we strategically
189 increased the expression of astrocyte-related developmental driver genes in another set of
190 simulated intermediate progenitors. We hypothesized that this perturbation would lead to a
191 larger proportion of astrocyte s as terminal cell states. Indeed, we observed a statistically
192 significant increase in the proportion of astrocytes (53%) compared with the baseline
193 (16.5%; binomial test $p < 10^{-8}$) under the dynamics learned by scDVF. Thus, *in silico*
194 perturbation studies can be used to efficiently and comprehensively identify
195 developmental driver genes upstream of the signaling cascade. More interestingly,
196 simulation results suggest that mouse neocortex development exhibits properties similar to
197 chaotic systems, where small perturbations in key upstream genes determine the fate of
198 cell bifurcations. In other words, small variations in the initial conditions of a cell may
199 result in large downstream changes.

200 **Exploring the Neural Equations Behind the Developing Mouse Dentate Gyrus**



201 Further, we evaluated whether scDVF can uncover the dynamics of a dataset from a
 202 different tissue, developmental timescale, and technical platform. We considered an
 203 scRNA-seq experiment of the developing mouse dentate gyrus with transcriptomes
 204 profiled using droplet-based scRNA-seq (Fig. 4a) (10). After obtaining a neural network
 205 representation of the dentate gyrus dynamics, an out-of-sample cell was simulated by
 206 perturbing the gene expression state of an *Nbl2* cell. With the out-of-sample cell as the
 207 initial condition, we used scDVF to simulate an out-of-sample cell trajectory, which
 208 moved along the existing granule cell trajectory in the data (Fig. 4b). Further, the VAE
 209 embeddings properly encoded the developmental hierarchy of cell types in the low-
 210 dimensional dynamic manifold (Fig. 4c).

211 When examining critical cell states in the dentate gyrus, we observed an abrupt gene
 212 expression change in the developmental manifold, which can be visualized when ordering
 213 cells in latent time (Fig. 4d). Specifically, the abrupt change in gene expression marks the
 214 transition from *nIPC* to *Nbli* cells and suggests fate commitment during the transition.
 215 After calculating the CCI, we found that cells experiencing this abrupt change also have a
 216 high criticality, which substantiates the CCI as a metric for quantifying the instability of
 217 cell states. The robustness of the CCI as an instability measure is also highlighted by its
 218 most strongly correlated genes in the dentate gyrus. For example, *IGFBPL-1*, which most
 219 positively correlates with the CCI, drives neuron differentiation in progenitor cells, and

220 *GRM5*, which most negatively correlates with the CCI, encodes glutamate receptors in
221 stable and differentiated neurons (Fig. 4e) (19, 20).

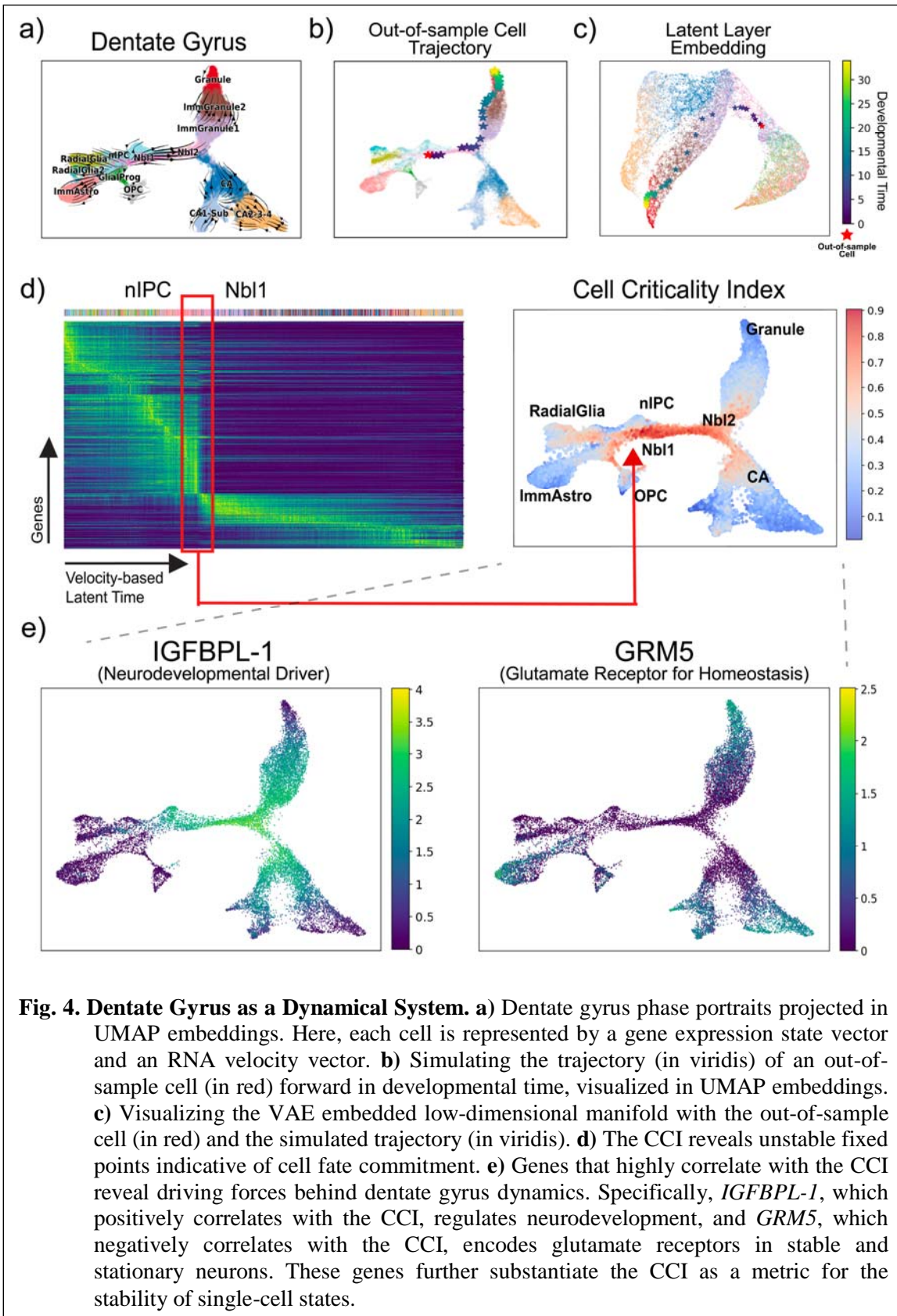


Fig. 4. Dentate Gyrus as a Dynamical System. **a)** Dentate gyrus phase portraits projected in UMAP embeddings. Here, each cell is represented by a gene expression state vector and an RNA velocity vector. **b)** Simulating the trajectory (in viridis) of an out-of-sample cell (in red) forward in developmental time, visualized in UMAP embeddings. **c)** Visualizing the VAE embedded low-dimensional manifold with the out-of-sample cell (in red) and the simulated trajectory (in viridis). **d)** The CCI reveals unstable fixed points indicative of cell fate commitment. **e)** Genes that highly correlate with the CCI reveal driving forces behind dentate gyrus dynamics. Specifically, *IGFBPL-1*, which positively correlates with the CCI, regulates neurodevelopment, and *GRM5*, which negatively correlates with the CCI, encodes glutamate receptors in stable and stationary neurons. These genes further substantiate the CCI as a metric for the stability of single-cell states.

222 Lastly, we conducted *in silico* perturbation studies to determine the genetic drivers behind
 223 dentate gyrus cell fate decisions. We randomly sampled upstream *Nbl2* cells ($n = 1,000$)
 224 as the initial conditions and allowed the simulated *Nbl2* cells to naturally evolve according
 225 to the dynamics captured by scDVF, which resulted in either terminal granule or

pyramidal cell states. Then, we performed differential expression analysis on the initial conditions (i.e., the simulated *Nbl2* cell states) of different fates (Fig. 5b). The top differentially expressed gene associated with a granule cell fate was *Prox1*. This gene has also been previously identified by RNA velocity and experimentally validated as being necessary for granule cell formation; moreover, the deletion of *Prox1* leads to the adoption of the pyramidal neuron fate (19). In addition, scDVF identified the top pyramidal neuron developmental driver gene as *Runx1t1*, which was recently shown to induce pyramidal neuron formation, with its deletion resulting in reduced neuron differentiation *in vitro* (21). As further validation, we increased the expression of pyramidal neuron developmental driver genes in simulated *Nbl2* cells and observed an elevated proportion of pyramidal neurons as terminal states (from 10% to 30%; binomial test $p < 10^{-7}$) under the dynamics captured by scDVF. In summary, *in silico* perturbation studies can be a low-cost alternative for identifying developmental driver genes. Further, the results show that scDVF is robust on scRNA-seq from different tissues, developmental timescales, and technical platforms.

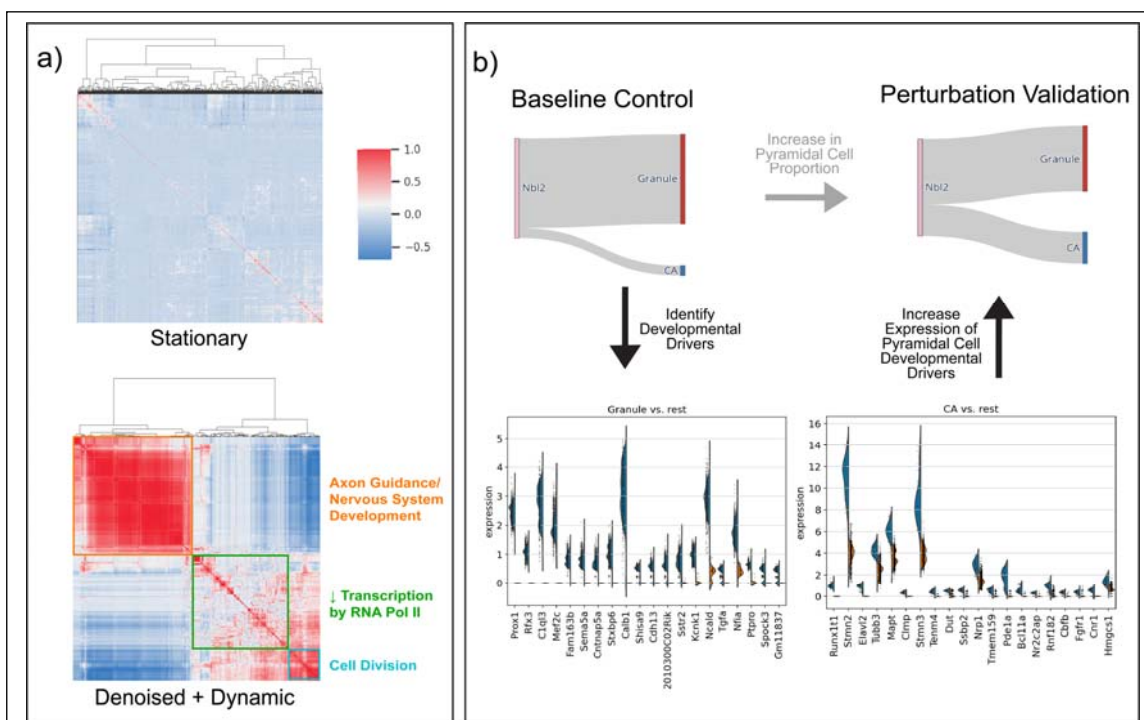


Fig. 5. Dentate Gyrus Co-expression Networks and Perturbation Studies. a) Simulating dynamic trajectories can effectively denoise gene co-expression networks in sparse and noisy scRNA-seq experiments, especially when compared to the static cells in dentate gyrus. b) Differential gene expression analysis on the simulated *Nbl2* cells reveals key putative genes that correlate with the fate of transforming into a pyramidal versus granule cells. Early perturbation of the top differentially expressed genes associated with a pyramidal cell fate resulted in a higher proportion of pyramidal cells from the perturbed *Nbl2* cells, suggesting a causal relationship through *in silico*

241 Comparing scDVF with Existing Methods

242 RNA velocity predicts gene expression change of individual cells in the timescale of
 243 hours. Previous simulations in this study used hypothetical progenitor cells as the initial
 244 conditions and computed trajectories into the future resulting in differentiated cells as
 245 terminal states. Conversely, we can use differentiated (or terminal) cells as the initial

246 conditions and rewind time with scDVF. Then, the retrograde developmental trajectory
247 represents the gene dynamics that would have resulted in the terminal cell types.

248 Due to sparse and noisy measurements, it is often challenging to detect strong correlation
249 between genes in scRNA-seq, thereby making it difficult to find coherent functional
250 modules in gene co-expression networks (22–24). However, denoising VAEs in scDVF
251 can reduce the variability along a developmental trajectory due to the sparsity and noise
252 associated with scRNA-seq (Fig. 6a). We hypothesize that cells in denoised trajectories
253 simulated from scDVF (with a representative initial condition) could amplify the
254 correlations within functional gene modules (Fig. 6b). Indeed, the gene co-expression
255 network of cells in retrograde trajectories has more significant gene correlations compared
256 to co-expression networks from static cells (Fig. 3a, Fig. 5a). Further, we biclustered the
257 co-expression matrix into gene clusters. By benchmarking our approach on four datasets,
258 we demonstrate that the gene clusters discovered from our method are more coherent by
259 comparing the gene ontology (GO) enrichments. The benchmarks show that functional
260 gene modules found from denoised and dynamic cells in retrograde trajectories have at
261 least two orders of magnitude higher enrichment for cell-type-specific GO terms
262 compared to static cell clusters (Fig. 6c). Therefore, the retrograde trajectories computed
263 by scDVF can effectively disentangle trajectory-specific gene regulatory networks and
264 serve as a computational solution for boosting signal-to-noise ratios in single-cell gene co-
265 expression networks.

266 Further, scDVF qualitatively differs from existing ODE-based regulatory networks (25).
267 First, explicitly deriving differential equations for biological processes is only feasible for
268 examining small-scale systems (26–29). In contrast, scDVF can capture high-dimensional
269 interactions and can scale to a large number of variables. Second, scDVF uses a neural
270 network to learn potentially non-linear gene interactions, which is more suitable for
271 modeling complex biological processes compared to linear ODEs and other kernel-based
272 sparse approximation methods (30, 31). In particular, we compared scDVF with state-of-
273 the-art vector field learning approach SparseVFC (6). Benchmarking results show that our
274 method has at least 50% reduction in out-of-sample velocity prediction loss across all
275 datasets, indicating that scDVF can learn a more accurate representation of the velocity
276 vector fields and can compute future cell states with better numerical precision (Fig. 6d).
277 Lastly, many previous ODE-based methods used pseudo-time as a substitute for time. In
278 comparison, scDVF uses RNA velocity, which reflects developmental time (5).

279 Discussion

280 Although many effective tools have been developed to illuminate the dynamics of single-
281 cell data, existing methods have mostly viewed single-cell datasets as a static manifold
282 (e.g., minimum spanning trees, Markov matrices, diffusion etc.). In reality, many
283 underlying biological processes captured by single-cell sequencing are dynamical systems,
284 where individual cells are transitioning from one state to another. Hence, deriving accurate
285 differential equations that quantify gene expression dynamics of single cells can answer
286 many questions about the cell fates and the genetic drivers behind developmental
287 trajectories.

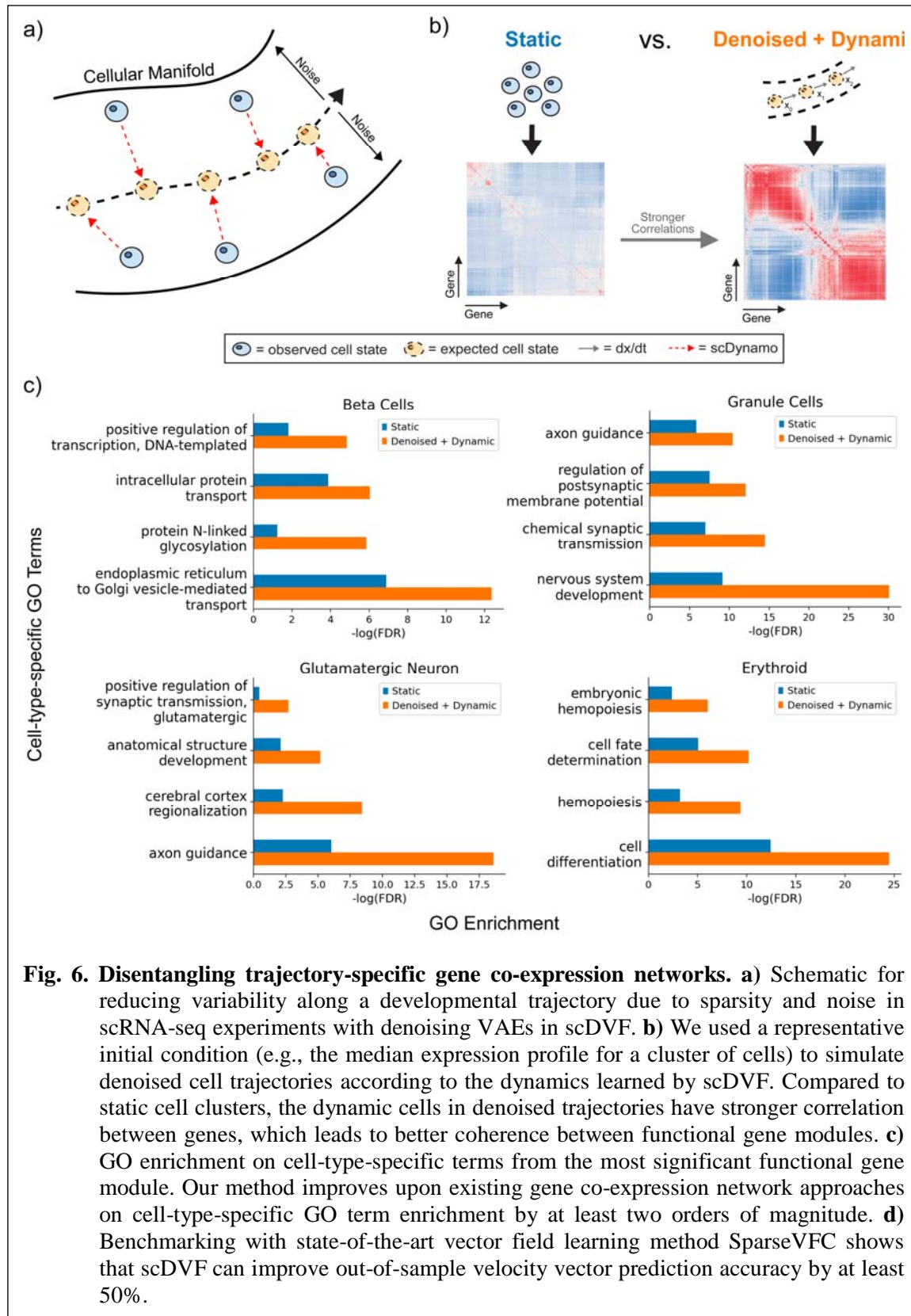


Fig. 6. Disentangling trajectory-specific gene co-expression networks. **a)** Schematic for reducing variability along a developmental trajectory due to sparsity and noise in scRNA-seq experiments with denoising VAEs in scDVF. **b)** We used a representative initial condition (e.g., the median expression profile for a cluster of cells) to simulate denoised cell trajectories according to the dynamics learned by scDVF. Compared to static cell clusters, the dynamic cells in denoised trajectories have stronger correlation between genes, which leads to better coherence between functional gene modules. **c)** GO enrichment on cell-type-specific terms from the most significant functional gene module. Our method improves upon existing gene co-expression network approaches on cell-type-specific GO term enrichment by at least two orders of magnitude. **d)** Benchmarking with state-of-the-art vector field learning method SparseVFC shows that scDVF can improve out-of-sample velocity vector prediction accuracy by at least 50%.

288
289
290

Explicitly deriving differential equations for all gene interactions is a challenging task. Therefore, we tackled the problem with a data-driven approach. We considered each cell in scRNA-seq as an instance sampled from a dynamic system, composed of a state vector

291 (gene expression, \vec{x}) and a velocity vector (RNA velocity, $\frac{\partial \vec{x}}{\partial t}$). Then, we trained a neural
292 network $A(\vec{x})$ to learn potentially non-linear mappings from the state \vec{x} to the velocity $\frac{\partial \vec{x}}{\partial t}$
293 of each cell. With a trained VAE $A(\vec{x})$ that takes part in the differential equation $\frac{\partial \vec{x}}{\partial t} =$
294 $A(\vec{x})$, we can integrate the VAE with any black box ODE solver to compute the future (or
295 past) gene expression states.

296 Overall, our scDVF framework allows hypothetical cells to evolve according to the
297 dynamics learned from existing cells in the data. Using the ability to simulate future gene
298 expression trajectories, we devised a metric to quantify the instability of individual cells
299 called the CCI. Through perturbing cell states with high criticality, *in silico* gene
300 perturbation studies can computationally identify key upstream driver genes that
301 determine the fate of cell bifurcations. Lastly, by rewinding the developmental time of
302 differentiated cells, retrograde trajectories can deconvolute trajectory-specific gene co-
303 expression networks and discover more coherent cell-type-specific gene modules.

304 Previous approaches utilize pseudo-time to construct a temporal-ordering of cells and
305 pluripotency metrics to measure the differentiation potential of a cell, similar to
306 quantifying the “potential energy” of Waddington landscapes. However, these “potential
307 energy” metrics are limited in describing dynamical systems. Theoretically, the potential
308 energy is converted into conservative forces, where the total work done by a cell becomes
309 independent of the developmental path taken. In order to more accurately capture the
310 expression changes along specific developmental trajectories, we designed a new metric
311 called the CCI, analogous to the “kinetic energy” of Waddington landscapes. In our
312 analysis, we demonstrated that this metric could highlight fixed points in single-cell
313 dynamical systems. Moreover, previous studies have formulated cell fate decisions as
314 high-dimensional critical state transitions (32, 33). Therefore, we further bring awareness
315 to the dynamical perspective of single-cell data and advocate for new metrics that quantify
316 the kinetics of single-cell experiments.

317 More interestingly, single-cell processes have long been hypothesized to exhibit properties
318 similar to chaotic systems (7, 34, 35). By recovering single-cell gene expression dynamics
319 with scDVF, we observed chaotic behaviors in the *in silico* gene perturbation studies,
320 where a small change in the initial gene expression state may result in a large difference in
321 the future states, also known as the butterfly effect. Specifically, small perturbations in
322 developmental driver genes of progenitor cells can alter the cell fate at developmental
323 branching points both *in vitro* and *in silico*. If single-cell dynamics exhibit chaotic
324 properties, under the right biological conditions, the chaos can spontaneously evolve into
325 lockstep patterns according to the Kuramoto model of synchronization (36). Hence,
326 synchronization models could be a possible explanation for emergent tissue-level
327 behaviors from single cells. These effects could be explored by incorporating the gene
328 interaction dynamics between cells. Currently, scDVF only models the gene expression
329 dynamics within a single cell. A future direction could expand the state space of scDVF
330 and incorporate gene interactions between spatially neighboring cells with spatial
331 transcriptomics (37). Another future direction includes incorporating concurrently
332 resolved protein and chromatin accessibilities and their velocities into the dynamical
333 model as a multi-modal representation of the cell state (38, 39).

335 **Materials and Methods**

337 **Data Collection and Preprocessing**

scRNA-seq data (pancreatic endocrinogenesis, dentate gyrus, mouse gastrulation, and human forebrain) were downloaded. After computing the gene expression count matrix. The top 3,000 variable genes and cells with a minimum of 20 transcripts were selected. Velocity genes were found using log-transformed data, and the moments were estimated using the top 30 principal components and the top 30 nearest neighbors. Dynamical velocity vectors were computed using the raw counts.

Variational Autoencoder Architecture

High-dimensional single-cell dynamical systems are difficult to model due to high degrees of freedom. For example, the number of features can sometimes be larger than the number of data points. Consequently, gene expression would only vary in a small portion of dimensions. Therefore, modeling the gene expression dynamics of a low-dimensional manifold embedded in high-dimensional data is a challenging task. Fortunately, autoencoders can reduce the dimensionality of the data by introducing an information bottleneck. Accordingly, when used to represent dynamical systems, autoencoders can restrict cell transitions to only movements along the low-dimensional manifold.

A variational autoencoder with four dense layers (size 64 as the intermediate layer and size 16 as the latent layer) was constructed using the Tensorflow and Keras packages (40, 41). The VAE takes the gene expression state as input, and outputs the RNA velocity. In the VAE, the encoder layers with weights W_e and biases b_e produces the hidden layer $h(x)$, which parametrizes the location and scale of i gaussian distributions. Then, a sample from each reparametrized gaussian distribution z_i is used as input for the decoder layer with weights W_d and biases b_d . The architecture can be expressed as:

$$\text{EncoderLayer}(x) = h(x) \quad (\text{Equation 4})$$

$$= \text{Relu}(b_e + W_e * x) \quad (\text{Equation 5})$$

$$\mu_i(x) = \text{EncoderLayer}(h(x)) \quad (\text{Equation 6})$$

$$\sigma^2_i(x) = \text{EncoderLayer}(h(x)) \quad (\text{Equation 7})$$

$$\epsilon_i \sim N(0, I) \quad (\text{Equation 8})$$

$$z_i \sim \mu_i\left(\frac{\partial x}{\partial t}\right) + \epsilon_i * \sigma^2_i\left(\frac{\partial x}{\partial t}\right) \quad (\text{Equation 9})$$

$$\text{DecoderLayer}(z_i) = \text{Relu}(b_d + W_d * z_i) \quad (\text{Equation 10})$$

where the $\text{Relu}(z)$ activation function is:

$$\text{Relu}(z) = \max(0, z) \quad (\text{Equation 11})$$

We used the mean squared error reconstruction loss with the Adam optimizer. To prevent overfitting and encourage a sparse representation of latent embeddings, L1 regularization was added to all layers. The evidence lower bound loss function with L1 regularization where $\lambda = 1 \times 10^{-6}$, $q(z | \frac{\partial x}{\partial t}) = N(\mu_i(\frac{\partial x}{\partial t}), \text{diag}(\sigma^2_i(\frac{\partial x}{\partial t})))$ $p(z) = N(0, I)$, can be described as:

$$L\left(\frac{\partial x}{\partial t}, \widehat{\frac{\partial x}{\partial t}}\right) = KL - Divergence + Reconstruction Loss + Regularization \quad (\text{Equation 12})$$

$$= KL(q\left(z \middle| \frac{\partial x}{\partial t}\right) || p(z)) - \sum_{i=1}^D \left(\frac{\partial x}{\partial t} - \widehat{\frac{\partial x}{\partial t}}\right)^2 + \lambda \sum_{i=1}^D |W_e| + \lambda \sum_{i=1}^D |W_d|$$

(Equation 13)

Because the input and output vectors are sparse, a small learning rate of 0.00001 was used. Early stopping was added once the validation loss did not improve for three consecutive epochs.

Initial Value Problems and ODE Solvers for Integration

Our framework can be used to predict gene expression profiles across time. Given t_0 and $\vec{x}(0) = x_0$, this is an initial value problem with the goal of solving $\vec{x}(t) = \vec{x}_t$ for any t :

$$\frac{\partial \vec{x}(t)}{\partial t} = f(t, \vec{x}(t)) \quad (\text{Equation 14})$$

Here, f is only a function of the state \vec{x}_t such that $f = A(\vec{x}_t)$. Then the equation becomes:

$$\frac{\partial \vec{x}_t}{\partial t} = A(\vec{x}_t) \quad (\text{Equation 15})$$

The first-order Euler's method for finding the state \vec{x}_{t+1} is:

$$\vec{x}_{t+1} = \vec{x}_t + \frac{\partial \vec{x}_t}{\partial t} \quad (\text{Equation 16})$$

$$= \vec{x}_t + A(\vec{x}_t) \quad (\text{Equation 17})$$

However, we can utilize higher-order ODE solvers from the SciPy package to find a more accurate solution (42). The explicit Runge-Kutta method of order 8 (DOP853) was used to obtain the most accurate solutions, but it has a slow runtime. Explicit Runge-Kutta method of order 5 (RK53) can be used to trade off accuracy for a faster runtime. In practice, cells in this study were integrated to a maximum of 35 discrete steps (each with 5 intermediate steps) forward in time, which should be experimentally derived for each dataset.

Addressing Drift Effects

In control theory, using only the previous state and the velocity vectors to predict the next state can result in a phenomenon called “dead reckoning,” where the errors accumulate after each step (43). To mitigate this effect, we employed two strategies:

1. Instead of a traditional VAE, we trained a denoising VAE to reduce the variance of predicted RNA velocity. By adding a small Gaussian noise ϵ to the gene expression input during training, we could increase the generalizability of the input space and improve extrapolations to out-of-sample cells.

$$\frac{\partial \vec{x}}{\partial t} = A(\vec{x} + \epsilon) \quad (\text{Equation 18})$$

2. As we integrated the VAE over time, we found reference cells in the data manifold every few steps and continued integration from the reference cell, as a form of high-gain Kalman filter. We designated the intermediate step size as a hyperparameter relative to the step size. For example, after integrating for five intermediate steps, we projected the predicted (or extrapolated) gene expression state to the original dataset using the top 30 principal components. Then, we identified the K -nearest neighbors ($K = 30$) within the PCA embeddings. The reference cell is defined as the median

429 expression profile among those K-nearest neighbor cells from the dataset, and ODE
430 integration continued from this reference cell. This allowed our prediction to adhere
431 closely to the data manifold and reduced the degree-of-freedom due to numerical
432 errors. Further, finding reference cells in the data also constructed boundary conditions
433 when integrating a dynamical system. For example, once the extrapolated state went
434 beyond the cellular manifold, there were no cells in the data to serve as a reference,
435 but the nearest neighboring cells from the dataset could still construct a reference cell
436 from where integration could continue.

437 **Measuring Instability with the Cell Criticality Index**

438 By solving for the developmental path of a single cell, we can measure the amount of gene
439 expression change along a trajectory, rather than comparing only the difference between
440 the start and end states. Previously, StemID used the entropy of the gene expression
441 distribution to heuristically identify stems cells in single-cell transcriptome data, where
442 pluripotent cells tend to have a more uniform gene expression distribution with a higher
443 entropy and differentiated cells tend to have a lower entropy (44). If \vec{x}^g denotes the
444 expression state of the genes g , then the StemID of the gene expression state is defined as:

$$445 \text{StemID}(\vec{x}) = -\sum_{i \in g} x^i \log\left(\frac{x^i}{\sum_{i \in g} x^i}\right) \quad (\text{Equation 19})$$

447 We reasoned that a change in the gene expression distribution (e.g., from high to low
448 entropy) can be captured using the relative entropy (or the KL-divergence). Based on this
449 idea, we devised a measure to quantify the capacity for any cell to undergo gene
450 expression change in the dynamical system. The CCI is calculated as the cumulative
451 information change, or the cumulative KL-divergence, between gene expression
452 distributions at each step in the developmental trajectory. Different from StemID, the CCI
453 can quantify the gene expression change of a cell regardless of the pluripotency. As an
454 analogy, StemID measures the “potential energy” of a cell's ability to differentiate,
455 whereas the CCI measures the “kinetic energy” of a cell's ability to change. If \vec{x}_t^g denotes
456 the expression state of the genes g at time t , then the cumulative KL-divergence for
457 $T = 35$ steps can be defined as:

$$458 \text{CCI}(\vec{x}) = \sum_{t=0}^T KL(\vec{x}_{t+1}^g || \vec{x}_t^g) \quad (\text{Equation 20})$$

$$459 = \sum_{t=0}^T \sum_{i \in g} x_{t+1}^i \log\left(\frac{x_{t+1}^i}{x_t^i}\right) \quad (\text{Equation 21})$$

462 ***In Silico* Perturbation Studies**

463 We divide an *in silico* perturbation study into three steps:

- 464 1. A sample of initial gene expression states ($n = 1,000$) was randomly generated. First,
465 we solved the random initial gene expression states over time to establish a
466 developmental baseline. Specifically, we aimed to observe the natural proportion of
467 terminal cell types that could arise from the dynamical system without any
468 intervention.
- 469 2. Then, we identified differentially expressed genes in the initial gene expression states
470 that correlate with development into a particular terminal cell type later in time.
471 Differential gene expression was performed using the *scanpy* package with Wilcoxon
472 test and Bonferroni corrections (45).

- 474 3. Lastly, we perturbed only the differentially expressed genes in another set of randomly
475 sampled initial gene expression states to test whether the perturbation increases the
476 proportion of cells for the terminal cell type.
477

478 To sample initial gene expression states, we computed the median expression profile of
479 progenitor cells and added Laplace distributed noise using the variance of those genes in
480 progenitor cells to randomly increase or decrease gene expression. For the perturbation,
481 exponentially distributed noise was added only to the top 100 differentially expressed
482 genes for the randomly sampled cells to specifically increase the expression of the top
483 differentially expressed genes. Terminal cell identity was determined by projecting the
484 data onto the top 30 principal components and by using K-nearest neighbor classification
485 (with $K = 30$). With the scVelo package, the dynamical mode estimates a variance for
486 each gene over all cells, whereas the stochastic mode estimates a variance for each cell.
487 Note that to model stochasticity in the stochastic mode, our framework could be easily
488 adapted to also learn the variance of the velocity vectors (as neural stochastic ODEs). All
489 initial gene expression states were integrated for 35 timesteps each with 5 intermediate
490 steps.
491

492 Retrograde Trajectory Simulation

493 Similar to *in silico* perturbation studies, we computed the median expression profile of a
494 terminal cell type (beta cells, granule cells, glutamatergic neurons, and erythroid) in each
495 scRNA-seq experiment (mouse pancreatic endocrinogenesis, dentate gyrus, human
496 forebrain, and mouse gastrulation) as the representative initial condition. A set of cells
497 ($n = 50$) were sampled from each representative initial condition by adding Laplace
498 distributed noise using the variance of the terminal cell type gene expression. The
499 retrograde trajectory for each cell was simulated by subtracting the predicted RNA
500 velocities from the gene expression state during integration:
501

$$\vec{x}_{t-1} = \vec{x}_t - \frac{\partial \vec{x}_t}{\partial t} \quad (\text{Equation 22})$$

$$= \vec{x}_t - A(\vec{x}_t) \quad (\text{Equation 23})$$

504
505 After integrating for 15 discrete steps each with 5 intermediate steps, a gene correlation
506 matrix of the cells in retrograde trajectories was calculated.
507

508 Gene Ontology Enrichment Analysis

509 Hierarchical biclustering was performed on the co-expression matrices, and three gene
510 clusters were identified from each co-expression matrix, representing three functional
511 modules. We performed GO enrichment analysis on each functional module using
512 GOATOOLS with Fisher's exact test (46). Further, we calculated the Benjamin-Hochberg
513 false discovery rates to correct for multiple testing. To compare between two co-
514 expression matrices, we considered the most significant enrichment out of the three
515 clusters for each GO term. In Fig. 3a and Fig. 5a, the most significantly enriched GO
516 terms associated with biological processes are listed next to each gene cluster.
517

518 References

- 519 1. C. Trapnell, D. Cacchiarelli, J. Grimsby, P. Pokharel, S. Li, M. Morse, N. J. Lennon, K. J.
520 Livak, T. S. Mikkelsen, J. L. Rinn, The dynamics and regulators of cell fate decisions are revealed
521 by pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32**, 381–386 (2014).
522 2. M. Setty, V. Kiseliouvas, J. Levine, A. Gayoso, L. Mazutis, D. Pe'er, Characterization of
523 cell fate probabilities in single-cell data with Palantir. *Nat. Biotechnol.* **37** (2019), pp. 451–460.

- 524 3. L. Haghverdi, M. Büttner, F. A. Wolf, F. Buettner, F. J. Theis, Diffusion pseudotime
525 robustly reconstructs lineage branching. *Nat. Methods*. **13**, 845–848 (2016).
- 526 4. G. La Manno, R. Soldatov, A. Zeisel, E. Braun, H. Hochgerner, V. Petukhov, K.
527 Lidschreiber, M. E. Kastrioti, P. Lönnerberg, A. Furlan, others, RNA velocity of single cells.
528 *Nature*. **560**, 494–498 (2018).
- 529 5. H. Matsumoto, H. Kiryu, C. Furusawa, M. S. Ko, S. B. Ko, N. Gouda, T. Hayashi, I.
530 Nikaido, SCODE: an efficient regulatory network inference algorithm from single-cell RNA-Seq
531 during differentiation. *Bioinformatics*. **33**, 2314–2321 (2017).
- 532 6. X. Qiu, Y. Zhang, J. D. Martin-Rufino, C. Weng, S. Hosseinzadeh, D. Yang, A. N.
533 Pogson, M. Y. Hein, K. H. J. Min, L. Wang, others, Mapping transcriptomic vector fields of
534 single cells. *Cell* (2022).
- 535 7. A. Uthamacumaran, A review of dynamical systems approaches for the detection of
536 chaotic attractors in cancer networks. *Patterns*. **2**, 100226 (2021).
- 537 8. R. T. Chen, Y. Rubanova, J. Bettencourt, D. Duvenaud, Neural ordinary differential
538 equations. *ArXiv Prepr. ArXiv180607366* (2018).
- 539 9. S. L. Brunton, J. L. Proctor, J. N. Kutz, Discovering governing equations from data by
540 sparse identification of nonlinear dynamical systems. *Proc. Natl. Acad. Sci.* **113**, 3932–3937
541 (2016).
- 542 10. H. Hochgerner, A. Zeisel, P. Lönnerberg, S. Linnarsson, Conserved properties of dentate
543 gyrus neurogenesis across postnatal development revealed by single-cell RNA sequencing. *Nat.*
544 *Neurosci.* **21**, 290–299 (2018).
- 545 11. D. J. Di Bella, E. Habibi, R. R. Stickels, G. Scalia, J. Brown, P. Yadollahpour, S. M.
546 Yang, C. Abbate, T. Biancalani, E. Z. Macosko, others, Molecular logic of cellular diversification
547 in the mouse cerebral cortex. *Nature*. **595**, 554–559 (2021).
- 548 12. A. Bastidas-Ponce, S. Tritschler, L. Dony, K. Scheibner, M. Tarquis-Medina, C. Salinno,
549 S. Schirge, I. Burtscher, A. Böttcher, F. J. Theis, others, Comprehensive single cell mRNA
550 profiling reveals a detailed roadmap for pancreatic endocrinogenesis. *Development*. **146**,
551 dev173849 (2019).
- 552 13. B. Pijuan-Sala, J. A. Griffiths, C. Guibentif, T. W. Hiscock, W. Jawaid, F. J. Calero-Nieto,
553 C. Mulas, X. Ibarra-Soria, R. C. Tyser, D. L. L. Ho, others, A single-cell molecular map of mouse
554 gastrulation and early organogenesis. *Nature*. **566**, 490–495 (2019).
- 555 14. E. Ferrell James, Bistability. *Bifurc. Waddingtons Epigenetic Landsc. Curr. Biol.* **22**,
556 R458–R466 (2012).
- 557 15. C. Guo, K.-S. Cho, Y. Li, K. Tchedre, C. Antolik, J. Ma, J. Chew, T. P. Utheim, X. A.
558 Huang, H. Yu, others, IGF1 regulates axon growth through IGF-1-mediated signaling
559 cascades. *Sci. Rep.* **8**, 1–13 (2018).
- 560 16. R. Bronstein, J. Kyle, A. B. Abraham, S. E. Tsirka, Neurogenic to gliogenic fate transition
561 perturbed by loss of HMGB2. *Front. Mol. Neurosci.* **10**, 153 (2017).
- 562 17. E. N. Lorenz, Deterministic nonperiodic flow. *J. Atmospheric Sci.* **20**, 130–141 (1963).
- 563 18. D. E. Wagner, A. M. Klein, Lineage tracing meets single-cell omics: opportunities and
564 challenges. *Nat. Rev. Genet.* **21** (2020), pp. 410–427.
- 565 19. Y. Gonda, H. Sakurai, Y. Hirata, H. Tabata, I. Ajioka, K. Nakajima, Expression profiles of
566 Insulin-like growth factor binding protein-like 1 in the developing mouse forebrain. *Gene Expr.*
567 *Patterns*. **7**, 431–440 (2007).
- 568 20. E. M. Batista, J. G. Doria, T. H. Ferreira-Vieira, J. Alves-Silva, S. S. Ferguson, F. A.
569 Moreira, F. M. Ribeiro, Orchestrated activation of mGluR5 and CB 1 promotes neuroprotection.
570 *Mol. Brain*. **9**, 1–17 (2016).
- 571 21. L. Zou, H. Li, X. Han, J. Qin, G. Song, Runx1t1 promotes the neuronal differentiation in
572 rat hippocampus. *Stem Cell Res. Ther.* **11**, 1–10 (2020).

- 573 22. S. G. Galfrè, F. Morandin, M. Pietrosanto, F. Cremisi, M. Helmer-Citterich, COTAN:
574 scRNA-seq data analysis based on gene co-expression. *NAR Genomics Bioinforma.* **3**, lqab072
575 (2021).
- 576 23. D. Mercatelli, F. Ray, F. M. Giorgi, Pan-Cancer and Single-Cell Modeling of Genomic
577 Alterations Through Gene Expression. *Front. Genet.* **10**, 671 (2019).
- 578 24. B. D. Harris, M. Crow, S. Fischer, J. Gillis, Single-cell co-expression analysis reveals that
579 transcriptional modules are shared across cell types in the brain. *Cell Syst.* (2021).
- 580 25. G. Karlebach, R. Shamir, Modelling and analysis of gene regulatory networks. *Nat. Rev.*
581 *Mol. Cell Biol.* **9**, 770–780 (2008).
- 582 26. S. Li, P. Brazhnik, B. Sobral, J. J. Tyson, A quantitative study of the division cycle of
583 *Caulobacter crescentus* stalked cells. *PLoS Comput. Biol.* **4**, e9 (2008).
- 584 27. H. H. McAdams, L. Shapiro, A bacterial cell-cycle regulatory network operating in time
585 and space. *Science.* **301**, 1874–1877 (2003).
- 586 28. J. Holtzendorff, D. Hung, P. Brende, A. Reisenauer, P. H. Viollier, H. H. McAdams, L.
587 Shapiro, Oscillating global regulators control the genetic circuit driving a bacterial cell cycle.
588 *Science.* **304**, 983–987 (2004).
- 589 29. E. G. Biondi, S. J. Reisinger, J. M. Skerker, M. Arif, B. S. Perchuk, K. R. Ryan, M. T.
590 Laub, Regulation of the bacterial cell cycle by an integrated genetic circuit. *Nature.* **444**, 899–904
591 (2006).
- 592 30. M. S. Yeung, J. Tegnér, J. J. Collins, Reverse engineering gene networks using singular
593 value decomposition and robust regression. *Proc. Natl. Acad. Sci.* **99**, 6163–6168 (2002).
- 594 31. D. C. Weaver, C. T. Workman, G. D. Stormo, in *Biocomputing '99* (World Scientific,
595 1999), pp. 112–123.
- 596 32. M. Mojtahedi, A. Skupin, J. Zhou, I. G. Castaño, R. Y. Leong-Quong, H. Chang, K.
597 Trachana, A. Giuliani, S. Huang, Cell fate decision as high-dimensional critical state transition.
598 *PLoS Biol.* **14**, e2000640 (2016).
- 599 33. A. E. Teschendorff, A. P. Feinberg, Statistical mechanics meets single-cell biology. *Nat.*
600 *Rev. Genet.* **22** (2021), pp. 459–476.
- 601 34. C. Furusawa, K. Kaneko, Chaotic expression dynamics implies pluripotency: when theory
602 and experiment meet. *Biol Direct.* **4** (2009), p. 17.
- 603 35. M. L. Heltberg, S. Krishna, M. H. Jensen, On chaotic dynamics in transcription factors
604 and the associated effects in differential gene regulation. *Nat. Commun.* **10**, 1–10 (2019).
- 605 36. Y. Kuramoto, *Chemical oscillations, waves, and turbulence* (Courier Corporation, 2003).
- 606 37. P. L. Ståhl, F. Salmén, S. Vickovic, A. Lundmark, J. F. Navarro, J. Magnusson, S.
607 Giacomello, M. Asp, J. O. Westholm, M. Huss, others, Visualization and analysis of gene
608 expression in tissue sections by spatial transcriptomics. *Science.* **353**, 78–82 (2016).
- 609 38. G. Gorin, V. Svensson, L. Pachter, Protein velocity and acceleration from single-cell
610 multiomics experiments. *Genome Biol.* **21**, 1–6 (2020).
- 611 39. M. Tedesco, F. Giannese, D. Lazarević, V. Giansanti, D. Rosano, S. Monzani, I. Catalano,
612 E. Grassi, E. R. Zanella, O. A. Botrugno, others, Chromatin Velocity reveals epigenetic dynamics
613 by single-cell profiling of heterochromatin and euchromatin. *Nat. Biotechnol.*, 1–10 (2021).
- 614 40. Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro,
615 Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow,
616 Andrew Harp, Geoffrey Irving, Michael Isard, Y. Jia, Rafal Jozefowicz, Lukasz Kaiser,
617 Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek
618 Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal
619 Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals,
620 Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, Xiaoqiang Zheng, TensorFlow:
621 Large-Scale Machine Learning on Heterogeneous Systems (2015), (available at
622 <https://www.tensorflow.org/>).

- 523 41. F. Chollet, Building autoencoders in keras. *Keras Blog*. **14** (2016).
524 42. P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E.
525 Burovski, P. Peterson, W. Weckesser, J. Bright, others, SciPy 1.0: fundamental algorithms for
526 scientific computing in Python. *Nat. Methods*. **17**, 261–272 (2020).
527 43. S. Beauregard, H. Haas, in *Proceedings of the 3rd Workshop on Positioning, Navigation*
528 *and Communication* (2006), pp. 27–35.
529 44. D. Grün, M. J. Muraro, J.-C. Boisset, K. Wiebrands, A. Lyubimova, G. Dharmadhikari,
530 M. van den Born, J. Van Es, E. Jansen, H. Clevers, others, De novo prediction of stem cell
531 identity using single-cell transcriptome data. *Cell Stem Cell*. **19**, 266–277 (2016).
532 45. F. A. Wolf, P. Angerer, F. J. Theis, SCANPY: large-scale single-cell gene expression data
533 analysis. *Genome Biol*. **19**, 1–5 (2018).
534 46. D. Klopfenstein, L. Zhang, B. S. Pedersen, F. Ramírez, A. W. Vesztrocy, A. Naldi, C. J.
535 Mungall, J. M. Yunes, O. Botvinnik, M. Weigel, others, GOATOOLS: A Python library for Gene
536 Ontology analyses. *Sci. Rep.* **8**, 1–17 (2018).
537

538 **Acknowledgments**

539
540 **Funding:** Research reported in this publication was supported by the National Institutes of
541 Health under award numbers [insert here].
542

543 **Author contributions:**

544 Conceptualization: ZC
545 Methodology: WK, ZC
546 Investigation: ZC, WK, MB, JZ
547 Visualization: ZC
548 Supervision: MB, JZ
549 Writing—original draft: ZC
550 Writing—review & editing: MB, JZ
551

552 **Competing interests:** Authors declare that they have no competing interests.
553

554 **Data and materials availability:** The mouse neocortex, pancreatic endocrinogenesis,
555 dentate gyrus, gastrulation, and human forebrain datasets used for this study can be found
556 in the NCBI Gene Expression Omnibus (GEO) repository with accession numbers
557 GSE153164, GSE132188, GSE95753, GSE87038, and in Sequence Read Archive (SRA)
558 under accession code SRP129388. All source code to reproduce this study can be found on
559 Github at <https://github.com/gersteinlab/scDVF>.