

# scMultiSim: simulation of multi-modality single cell data guided by cell-cell interactions and gene regulatory networks

Hechen Li (✉ [hli691@gatech.edu](mailto:hli691@gatech.edu))

Georgia Institute of Technology

Ziqi Zhang (✉ [ziqi.zhang@gatech.edu](mailto:ziqi.zhang@gatech.edu))

Georgia Institute of Technology <https://orcid.org/0000-0002-8198-0260>

Michael Squires (✉ [squiresmf@gatech.edu](mailto:squiresmf@gatech.edu))

Georgia Institute of Technology <https://orcid.org/0000-0002-3876-161X>

Xi Chen (✉ [chenx9@sustech.edu.cn](mailto:chenx9@sustech.edu.cn))

Southern University of Science and Technology <https://orcid.org/0000-0003-2648-3146>

Xiuwei Zhang (✉ [xiuwei.zhang@gatech.edu](mailto:xiuwei.zhang@gatech.edu))

Georgia Institute of Technology <https://orcid.org/0000-0002-1713-772X>

---

## Article

### Keywords:

DOI: <https://doi.org/>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

**Additional Declarations:** There is **NO** Competing Interest.

---

# 1 scMultiSim: simulation of multi-modality single 2 cell data guided by cell-cell interactions and gene 3 regulatory networks

4 Hechen Li<sup>1</sup>, Ziqi Zhang<sup>1</sup>, Michael Squires<sup>1</sup>, Xi Chen<sup>2</sup>, and Xiuwei Zhang<sup>1</sup>

5 <sup>1</sup>Georgia Institute of Technology, Atlanta, USA

6 <sup>2</sup>Southern University of Science and Technology, China

7 **Simulated single-cell data is essential for designing and evaluating computational methods in the**  
8 **absence of experimental ground truth. Existing simulators typically focus on modeling one or two**  
9 **specific biological factors or mechanisms that affect the output data, which limits their capacity to**  
10 **simulate the complexity and multi-modality in real data. Here, we present scMultiSim, an *in silico***  
11 **simulator that generates multi-modal single-cell data, including gene expression, chromatin accessibility,**  
12 **RNA velocity, and spatial cell locations while accounting for the relationships between modalities.**  
13 **scMultiSim jointly models various biological factors that affect the output data, including cell identity,**  
14 **within-cell gene regulatory networks (GRNs), cell-cell interactions (CCIs), and chromatin accessibility,**  
15 **while also incorporating technical noises. Moreover, it allows users to adjust each factor's effect**  
16 **easily. We validated scMultiSim's simulated biological effects and demonstrated its applications by**  
17 **benchmarking a wide range of computational tasks, including cell clustering and trajectory inference,**  
18 **multi-modal and multi-batch data integration, RNA velocity estimation, GRN inference and CCI inference**  
19 **using spatially resolved gene expression data. Compared to existing simulators, scMultiSim can**  
20 **benchmark a much broader range of existing computational problems and even new potential tasks.**

## 21 Introduction

22 In recent years, technologies that profile the transcriptome and other modalities (multi-omics) of single cell have  
23 brought remarkable advances in our understanding of cellular mechanisms [61]. For example, technologies  
24 have enabled the joint profiling of chromatin accessibility and gene expression data [10; 8; 41], as well  
25 as the measurement of surface protein abundance alongside transcriptome [56; 47]. Additionally, spatial  
26 locations of cells can be measured together with transcriptome profiles using imaging-based [52; 19; 63] or  
27 sequencing-based [55; 50] technologies.

28 The advent of single-cell multi-omics data has facilitated a more comprehensive understanding of cellular  
29 states, and more importantly, allowed researchers to explore the relationships between modalities and the  
30 causality across hierarchies [18]. Prior to the availability of single cell multi-omics data, gene regulatory  
31 network (GRN) inference methods were developed using only single-cell RNA sequencing (scRNA-seq) data [48].  
32 However, these methods mainly focused on transcription factors (TFs) as the sole factor affecting gene  
33 expressions. In reality, the observed gene-expression data is affected by multiple factors, such as the chromatin  
34 accessibility of corresponding regions. Consequently, newer methods utilizing both scRNA-seq and scATAC-seq  
35 data have been developed to infer GRNs [30; 62; 68]. Similarly, there has been a surge in the development  
36 of other computational tools that harness multi-modality information. For instance, Cell-Cell Interaction (CCI)  
37 inference methods seek to utilize both the gene expression and the spatial location modalities [16; 53; 5; 6] to  
38 learn the interactions with a lower false-positive rate than those using only scRNA-seq data [4; 26; 29]. Data  
39 integration methods combine multi-omics data to obtain a wholistic view of cells [58; 64; 1; 70; 36]. Moreover,  
40 RNA velocity can be inferred from unspliced and spliced counts using scRNA-seq data to indicate the near-future  
41 state of each cell [35; 3]. Recently, methods have also been proposed to infer RNA velocity from jointly profiled  
42 chromatin accessibility and transcriptomics data [38].

43 Overall, a large number of computational methods have been developed using scRNA-seq data or single cell  
44 multi- and spatial-omics data [66]. However, the scarcity of *ground truth* in experimental data makes it difficult  
45 to evaluate the performance of proposed computational methods. To address this, *de novo* simulators have  
46 been widely used to evaluate the accuracy of computational methods by generating data that models biological  
47 mechanisms and provides ground truth for benchmarking. SymSim [69], for example, provides ground truth cell  
48 identity and gene identity and thus can benchmark clustering, trajectory inference and differential expression  
49 detection. SERGIO [15], BEELINE [48] and dyngen [7] can simulate scRNA-seq data with given ground truth  
50 GRNs for testing GRN inference methods; while SERGIO, dyngen and VeloSim [71] can provide ground truth  
51 RNA velocity for testing RNA velocity inference methods. mistyR [60] generates single cell gene expression data  
52 from a given CCI network and can test CCI inference methods. With the *de novo* simulators, users can easily  
53 control the input parameters and obtain the exact ground truth. In addition to *de novo* simulators, Crowell *et*  
54 *al* [12] discussed another category of single cell data simulators, namely the reference-based methods, which

55 learn a generative model from a given real dataset and generate synthetic data [13; 59; 54; 2]. By design, these  
56 methods can output datasets that mimic the input reference data, but their flexibility can be limited by the specific  
57 reference dataset. Although they can provide ground truth cluster labels using annotations in the reference  
58 dataset or pre-determined labels during the simulation, none of the reference-based methods provides ground  
59 truth that is rarely available via domain knowledge, like GRNs, CCIs, or RNA velocity.

60 We consider that a desirable single cell simulator should meet several criteria: (1) it should generate as many  
61 modalities as possible to best represent a cell; (2) it should model as many biological factors and mechanisms that  
62 affect the output data as possible so that the output data has realistic complexity; and (3) it should provide ground  
63 truth of the biological factors to benchmark various computational methods. Most existing simulators generate  
64 only scRNA-seq data, and some generate only scATAC-seq data [44; 37]. Among the few ones that can generate  
65 multiple modalities, dyngen and SERGIO output unspliced and spliced counts with ground truth RNA velocity,  
66 while a reference-based simulator scDesign3 [54] can generate two modalities each with high dimensionality (*eg.*  
67 scRNA-seq and DNA methylation data), or one high-dimensional modality (*eg.* scRNA-seq) and spatial location  
68 data depending on the input reference dataset (Table S1).

69 In terms of the biological factors modeled in the simulator, existing *de novo* simulators model only one or  
70 a small subset of the following biological factors that affect gene expression in a cell: cell identity (cluster  
71 labels or positions on cell trajectories), chromatin accessibility, GRNs, and CCIs (Table S1). Data generated  
72 by reference-based simulators can inherently have these effects but it is challenging to obtain the ground truth of  
73 the biological factors, thus unable to measure the accuracy of a computational method.

74 In this paper, we present scMultiSim, a unified framework that models *all* the above biological factors as well  
75 as technical variations including sequencing noise and batch effect (Fig. 1a). For each single cell, it outputs  
76 the following modalities: unspliced and spliced mRNA counts, chromatin accessibility, and spatial location, while  
77 considering the cross-modality relationships. “Chromatin accessibility” is both an output modality (also called  
78 the scATAC-seq modality) and a biological factor that affects other output data (it affects the gene expression  
79 modality). scMultiSim provides ground truth information on cell identity (in terms of cell populations), RNA velocity,  
80 GRNs and CCIs, as well as relationships between chromatin accessibility and transcriptome data. Therefore, with  
81 one dataset, it can be used to evaluate methods for various computational tasks including clustering or trajectory  
82 inference, multi-modal and multi-batch data integration, RNA velocity estimation, GRN inference and CCI  
83 inference. Moreover, scMultiSim allows the users to adjust the effect of each biological factor on the output data,  
84 enabling them to investigate how the methods’ performance is affected by each factor when evaluating methods  
85 for a specific task. We present a comparison between scMultiSim and existing multi-modal simulators in Table S1.  
86 To our knowledge, scMultiSim is the most versatile simulator to date in terms of its benchmarking applications.

## 87 Results

88 In the following sections, we will provide a brief overview of the core concepts and the simulation process of  
89 scMultiSim. We will then demonstrate its capability to simulate multiple biological factors simultaneously by  
90 validating the effects of each factor on the output data. Furthermore, we will showcase the applications of  
91 scMultiSim by using it to benchmark a wide variety of computational tools.

### 92 scMultiSim overview

93 **The kinetic model and control of intrinsic noise.** In general, scMultiSim runs the simulation in two phases  
94 (Fig. 1b). In the first phase, scMultiSim employs the widely-accepted kinetic model [46] to generate the true  
95 gene expression levels in cells ("true counts"). In the second phase, scMultiSim introduces technical variations  
96 (library preparation noise, batch effects, etc) and generate scRNA-seq and scATAC-seq data that are statistically  
97 comparable to real data ("observed counts"). To model cellular heterogeneity and gene regulation effects,  
98 scMultiSim introduces two main concepts: *Cell Identity Factors* (CIFs) and *Gene Identity Vectors* (GIVs) (Fig. 1b  
99 (i, ii)). Biological factors, including cell population (cell identity), GRNs, and CCIs, are encoded in CIFs and GIVs  
100 (Fig. 2a). Additionally, to model single-cell chromatin accessibility, we also introduce Region Identity Vectors  
101 (RIVs, Fig. 1b(iii)). Further details on CIF, GIV and RIVs are provided in the next section.

102 When simulating single cell gene expression data, scMultiSim extends the idea of SymSim [69], where a kinetic  
103 model with three major parameters  $k_{on}$ ,  $k_{off}$ ,  $s$  was used to determine the expression pattern of a gene in a cell  
104 (Fig. 1b (vi)). In the kinetic model, a gene can switch between *on* and *off* states, with  $k_{on}$  and  $k_{off}$  be the rates of  
105 becoming *on* and *off*. When a gene is in the *on* state (which can be interpreted as promoter activation), mRNAs  
106 are synthesized at a rate  $s$  and degrade at a rate  $d$ . It is common to fix  $d$  at 1 and use the relative values for the  
107 other three parameters [43]. The kinetic parameters  $k_{on}$ ,  $k_{off}$ ,  $s$  are calculated from the CIF and GIV, as well as  
108 the corresponding scATAC-seq data (because chromatin accessibility is considered to affect gene expression).  
109 Since GIVs and CIFs encode information on cell identity, GRNs, and CCIs, the kinetic parameters thus capture  
110 the four biological factors that affect gene expression: cell identity, chromatin accessibility, GRNs, and CCIs.

111 The kinetic model used in scMultiSim provides two modes for generating true counts from the parameters,  
112 as shown in Fig. 1b (vii). The first mode is the full kinetic model, where genes undergo several cell cycles over  
113 time with *on/off* state changes, and the spliced/unspliced RNA counts are calculated. This mode provides the  
114 ground truth for RNA velocity since the RNA synthesis rate is known. The second mode is the Beta-Poisson  
115 model, which is equivalent to the kinetic model's master equation [31], and is faster to run than the full kinetic  
116 model. The Beta-Poisson model is recommended when RNA velocity is not needed. In the Beta-Poisson model,  
117 scMultiSim also introduces an intrinsic noise parameter  $\sigma_i$  that controls the amount of intrinsic noise caused by  
118 the transcriptional burst and the snapshot nature of scRNA-seq data. This parameter allows users to examine  
119 the influence of intrinsic noise on the performance of the computational methods. The two modes and the  $\sigma_i$

120 parameter are described in Methods.

121 **Modeling cellular heterogeneity and various biological effects.** The design of *Cell Identity Factors (CIFs)* and *Gene*  
122 *Identity Vectors (GIVs)* allows scMultiSim to encode cell identities and gene-level mechanisms (such as GRNs  
123 and CCIs) into the kinetic parameters and thereby impact the gene expression levels. This design also provides  
124 easy ways to adjust the effect of each factor on the output gene expression data.

125 The CIF of a cell is a 1D vector representing various biological factors that contributes to cellular heterogeneity,  
126 such as the cell condition (e.g. treated or untreated), or the expression of key TFs. The GIV of a gene act as  
127 the weights of the corresponding factors in the CIF, representing how strongly the corresponding CIF affect the  
128 gene's expression (Fig. 2a, Methods). By multiplying the CIF and GIV matrices, scMultiSim therefore generates  
129 a  $n_{\text{cell}} \times n_{\text{gene}}$  matrix, which is the desired kinetic parameter matrix with the cell and gene factors encoded.

130 Each CIF vector and GIV vector consists of four segments, each representing one type of extrinsic variation.  
131 They encode biological factors including cell identity (cell population, *i.e.*, the underlying cell trajectories or  
132 clusters), GRNs, and CCIs (Figs. 2a, S1a-b). We introduce the four segments in the following.

133 (i) Non-differential CIFs (**non-diff-CIF**) model the inherent cellular heterogeneity. They represent various  
134 environmental factors or conditions that are shared across all cells and are sampled from a Gaussian distribution  
135 with standard deviation  $\sigma_{\text{cif}}$ .

136 (ii) Differential CIFs (**diff-CIF**) control the user-desired cell population. These are the biological conditions that  
137 are unique to certain cell types. These factors lead to different cell types in the data. For a heterogeneous cell  
138 population, cells have different development statuses and types. Values for diff-CIFs are used to represent these  
139 cell differential factors, which are generated based on the user-input cell differential tree. When generating data  
140 for cells from more than one cell type, the minimal user input of scMultiSim is the cell differential tree, which  
141 controls the cell types (for discrete populations) or trajectories (for continuous populations) in the output. The  
142 process of generating diff-CIFs is described in Methods.

143 (iii) CIFs corresponding to Transcription Factors (**tf-CIF**) control the effects of GRNs. This segment, together  
144 with the TF segment in the GIV, model how a TF can affect expression of genes in the cell (Methods). Its length  
145 equals to the number of TFs. In other words, the GRN is encoded in the tf-CIFs and GIVs.

146 (iv) CIFs corresponding to ligands from neighboring cells (**lig-CIF**) control the effect of CCI. If CCI simulation  
147 is enabled, this segment together with the ligand segment in the GIV of the receptor gene encodes the ground  
148 truth CCI between two cells. This encoding ensures that a ligand and its interacting receptor have correlated  
149 gene expression. A receptor can also interact with ligands of multiple neighbors (Fig. 2a (viii)). The GIV matrices  
150 are generated carefully considering the nature of the kinetic model (Methods).

151 **The simulation process.** Fig. 1b shows an overview of the simulation process. The scATAC-seq data is generated  
152 at first (Fig. 2b(iv)), because we consider that the chromatin accessibility of a cell affects its gene expression.  
153 The scATAC-seq data also follows a pre-defined clustering or trajectory structure represented by the input cell  
154 differentiation tree. Similar to the gene expression, we multiply the CIF with a Region Identity Vector (RIV)  
155 matrix, which represents the effect of each CIF on the accessibility of chromatin regions. Details on generating  
156 the scATAC-seq data are included in Methods. The scATAC-seq data affects scRNA-seq data through the  $k_{on}$   
157 parameter, because chromatin accessibility controls the activated status of genes (Methods).

158 After obtaining all the kinetic parameters, scRNA-seq data can be generated in different modes: with or without  
159 CCIs and spatial locations, and with or without outputting RNA velocity data (Fig. 1b (vii, viii)). If the user specifies  
160 to generate RNA velocity, the full kinetic model is used, where cells undergo several cycles before the spliced  
161 and unspliced counts are outputted (Methods). Otherwise, if the Beta-Poisson model is used, and the true counts  
162 are sampled from the Beta-Poisson distribution. In this mode, RNA velocity and unspliced count data are not  
163 outputted.

164 **Simulating cell-cell interaction.** If specified to generate spatial-aware single cell gene expression data including  
165 cell spatial locations and CCI effects, scMultiSim uses a multiple-step approach that considers both time and  
166 space (Fig. 1b (viii), Fig. S1c). The simulation consists of a series of steps, with each step representing a time  
167 point. Cells are placed in a grid (Fig. 2a (ix), Fig. S1d), and one cell is added to the grid at each step, representing  
168 a newborn cell. Users can use the parameter  $p_n$  to control the probability for the newborn cell to locate with cells  
169 of the same type (Methods). As experimental data cannot measure cells at previous time points, scMultiSim  
170 outputs data only for cells at the final time point, which contains the accumulated CCI effects during the cells'  
171 developmental process.

172 To simulate CCI, scMultiSim requires a user-inputted list of ligand-receptor gene pairs that can potentially  
173 interact, which is called a ligand-receptor database. Users can input cell-type-level or single cell level CCI ground  
174 truth. If users do not provide ground truth CCIs, scMultiSim can randomly generate the ground truth from the  
175 ligand-receptor database.

176 **Technical variations and batch effects.** The steps described above belong to the first phase, which generates the  
177 “true” mRNA counts (and unspliced counts if RNA velocity mode is enabled) in the cells. In the second phase,  
178 scMultiSim simulates key experimental steps in wet labs that lead to technical noises in the data and output the  
179 observed scRNA-seq data. Batch effects can also be added to simulate datasets from a user-specified number  
180 of batches. Users can also control the amount of technical noise and batch effects between batches. These  
181 procedures are described in Methods. Next, we show the various output of scMultiSim and validate the effects  
182 present in the simulated data.

## 183 Design of simulation and datasets

184 We have generated a comprehensive set of datasets using scMultiSim to demonstrate the effects of different  
185 parameter configurations and to benchmark computational methods. These datasets contain both *main* and  
186 *auxiliary* datasets. The main datasets consists of 144 datasets with varying configurations of important  
187 parameters, including  $\sigma_{\text{cif}} \in \{0.1, 0.5\}$ ,  $n_{\text{cell}} \in \{500, 800\}$ ,  $n_{\text{gene}} \in \{110, 200, 500\}$ , and three different cell  
188 trajectories. The  $\sigma_{\text{cif}}$  parameter controls the standard deviation of the CIF and affects the within-cluster or  
189 within-neighborhood heterogeneity between cells. These main datasets contain all effects scMultiSim can  
190 simulate: GRN, chromatin accessibility, cell-cell interaction, technical noise and batch effect. Thus, the 144  
191 main datasets cover a wide range of variety, including different numbers of cells, genes, and trajectory shapes,  
192 to minimize potential bias and provide a more comprehensive benchmark of the computational methods.

193 As presented in Table 1, we label the main datasets with the following format:  $M\{p\}\{c\}\{s\}$ . The first letter M  
194 denotes the main dataset, followed by a letter  $p \in \{L, T, D\}$  that specifies the cell population as linear trajectory,  
195 tree trajectory or discrete, respectively. The number  $c \in [1, 12]$  denotes a particular configuration of  $\sigma_{\text{cif}}$ ,  $n_{\text{cell}}$ , and  
196  $n_{\text{gene}}$ , while the last lowercase letter  $s \in \{a, b, c, d\}$  represents random seed 1-4. For instance, the dataset MD5c  
197 has a discrete cell population,  $\sigma_{\text{cif}} = 0.1$ , 800 cells, 200 genes and random seed 3.

198 We have also generated auxiliary datasets with fewer types of effects and presented them in Table 2. These  
199 datasets allow us to explore the effect of other parameters and are compatible with computational methods  
200 that impose additional constraints on the input. In the remaining, we will primarily use the main datasets M for  
201 benchmarking and demonstration, while the auxiliary datasets will serve as additional and supplementary results.

## 202 scMultiSim generates multi-batch and multi-modality data from pre-defined clusters or trajectories

203 scMultiSim offers a key advantage in its ability to generate coupled scRNA-seq and scATAC-seq data while  
204 allowing users to control the shape of trajectories or clusters. It is accomplished by offering various parameters  
205 to control the structure of cell populations. First, the user can choose to generate “continuous” or “discrete”  
206 populations, and input a tree that represents the cell trajectories (in the case of “continuous” populations) or  
207 relationship between clusters (in the case of “discrete” populations). We name the tree “differentiation tree”.  
208 scMultiSim provides three example differentiation trees: Phyla1, Phyla3, and Phyla5, each having 1, 3, and  
209 5 leaves, as illustrated in Fig. 2b. The main datasets were simulated using these trees (Table 1). From a  
210 differentiation tree, scMultiSim is able to generate both discrete and continuous cell populations (Fig. 2c). Then,  
211 users can use these three parameters: intrinsic noise  $\sigma_i$ , CIF sigma  $\sigma_{\text{cif}}$  and Diff-to-nonDiff CIF ratio  $r_d$ , to  
212 control how clean or noisy the population structure is in the data (Fig. 2c-e).

213 For the continuous population, we visualize a dataset MT3a generated using tree Phyla3 in Fig. 2c. We  
214 can observe that the trajectories corresponding to the input differentiation tree are clearly visible for both the  
215 scRNA-seq and the scATAC-seq modality. For the discrete population, we visualize dataset MD3a and MD9a



216 generated with tree Phyla5 in Fig. 2d. The parameter  $\sigma_{\text{cif}}$  controls the standard deviation of the CIF, therefore  
217 with a smaller  $\sigma_{\text{cif}}$ , the clusters are tighter and better separated from each other. We then used the auxiliary  
218 dataset A (Table 2) to explore the effect of the intrinsic noise parameter  $\sigma_i$  and  $r_d$ , the ratio of number of diff-CIF  
219 to non-diff-CIFs. In Fig. 2e, we visualize the scRNA-seq modality generated using Phyla5 continuous mode  
220 with the same  $\sigma_{\text{cif}}$ . With a smaller Diff-to-nonDiff CIF ratio  $r_d$ , the trajectory is vague and more randomness is  
221 introduced, as the tree structure is encoded in the differential CIFs. With a smaller intrinsic noise  $\sigma_i$ , a fraction  
222 of the expression value is directly calculated from kinetic parameters without sampling from the Poisson model;  
223 As a result, the trajectory is more prominent. These patterns are much cleaner than real data because real data  
224 always has technical noise. We will show more results with technical noise in later sections and in Fig. S2.

225 **Coupling between scATAC-seq and scRNA-seq data.** In paired scATAC-seq and scRNA-seq data, these two  
226 data modalities are not independent of each other, as it is commonly considered that a gene's expression  
227 level is affected by the chromatin accessibility of the corresponding regions. If a gene's associated regions  
228 are accessible, this gene is more likely to be expressed. This mechanism can be naturally modeled in scMultiSim  
229 through the kinetic parameter  $k_{\text{on}}$  (Methods).

230 We provide a user-adjustable parameter, the ATAC-effect  $E_a$ , to control the extent of scATAC-seq data's effect  
231 on  $k_{\text{on}}$  (ranging from 0 and 1). In order to validate the connection between the scATAC-seq and scRNA-seq data,  
232 we calculate the mean Spearman correlation between these two modalities for genes that are controlled by one  
233 region in the scATAC-seq data. In Fig. 2f, we present the correlations under different  $E_a$  values. An averaged  
234 0.2-0.3 correlation can be observed using the default value (0.5), and the correlation increases with higher values  
235 of  $E_a$ . These results demonstrate that scMultiSim successfully models the connection between scATAC-seq and  
236 scRNA-seq data, enabling the generation of more realistic multi-omics datasets.

237 **scMultiSim simulates technical noise and batch effect.** The single cell gene expression data shown in Figs. 2c-f  
238 are "true" mRNA counts which do not have technical noise. scMultiSim can add technical noise including batch  
239 effects to the true counts to obtain observed counts (Methods). The amount of technical noise and batch effects  
240 can be adjusted through parameters, for example, the parameter  $E_{\text{batch}}$  can be used to control the amount of  
241 batch effect. Users can also specify the number of batches.

242 Fig. 2g shows the observed mRNA counts of dataset MD9a (true counts shown in Fig. 2d). The left plot shows  
243 data with one batch, and the right plot shows two batches. Technical noise and batch effects are also added to  
244 the scATAC-seq matrix. We further use the auxiliary dataset A to demonstrate the ability of scMultiSim to adjust  
245 the amount of technical noise and batch effect in both scRNA-seq and scATAC-seq modalities, in both continuous  
246 and discrete populations (Fig. S2). Here, we vary a main parameter for technical noise,  $\alpha$ , which denotes the  
247 capture efficiency that affects the detection ability of the dataset. Lower  $\alpha$  values correspond to poorer data  
248 quality.

## 249 **scMultiSim generates spliced and unspliced mRNA counts with ground truth RNA velocity**

250 If RNA velocity simulation is enabled, the kinetic model outputs the velocity ground truth using the RNA splicing  
251 and degradation rates. The Phyla5 tree in Fig. 2b is used to generate the results in Fig. 2h. The figure shows both  
252 the true spliced and unspliced counts, as well as the ground truth RNA velocity averaged by  $k$  nearest neighbor  
253 ( $k$ NN), which can be used to benchmark RNA velocity estimation methods. The RNA velocity vectors follow the  
254 cell trajectory (backbone and directions shown in red), which is specified by the user-inputted differentiation tree.

## 255 **scMultiSim generates single cell gene expression data driven by GRNs and cell-cell interactions**

256 The strength of scMultiSim also resides in its ability to incorporate the effect of GRN and CCI while preserving  
257 the pre-defined trajectory structures. In this section, we show that the GRN and CCI effects both exist in the  
258 simulated expression data. The main datasets (Table 1) used the 100-gene GRN from [15] as the ground truth  
259 GRN, which is visualized in Fig. 3a. We also incorporate CCIs by adding cross-cell ligand-receptor pairs to the  
260 within-cell GRNs. Specifically, we connect each cell's gene 99,101-104 to a neighbor cell's gene 91, 2, 6, 10  
261 (TFs), and 8 (non-TF) in the GRN (green edges in Fig. 3a). Next, we use one dataset (MT3a with a tree trajectory,  
262 500 genes, 500 cells, and  $\sigma_{\text{cif}} = 0.1$ ) to inspect the simulated effects in detail (Fig. 3b-e).

263 **GRN guided expression data.** We illustrate the gene regulation effects for dataset MT3a using a gene module  
264 correlation heatmap as shown in Fig.3b. The clustered heatmap is constructed by computing pairwise Spearman  
265 correlations between the expression levels of all genes. Each color on the top or left side of the heatmap  
266 represents a TF in the GRN. The figure shows that gene modules regulated by the same TF (genes with the  
267 same color) tend to be clustered together and have higher correlations with each other. These results suggest  
268 the presence of GRN effects in the expression data. To further illustrate the regulatory effects, we plot the  
269 expression of a specific regulator-target pair (gene 19-20) along one lineage (4-5-3 in Phyla3) in Fig. 3c. The  
270 plot clearly shows a correlation between the expression levels of the regulator and target genes. Moreover, we  
271 plot the accessibility levels for the corresponding chromatin region of gene 19 in Fig.3c. The plot indicates that  
272 significant drops in gene 19's expression occur when the related chromatin region is closed, providing further  
273 evidence for the regulatory effects of chromatin accessibility.

274 **Cell spatial locations.** scMultiSim provides convenient helper methods to visualize the cell spatial locations, as  
275 shown in Fig. 3d (dataset MT3a). For each ligand-receptor pair, arrows can be displayed between cells to show  
276 the direction of cell-cell interactions. We consider various biological scenarios when assigning the spatial location  
277 to each cell (Methods), for example, a newborn cell has a probability  $p_n$  of staying with a cell of the same type.  
278 Changing  $p_n$  allows us to generate different tissue layouts. In real data, how likely cells from the same cell type  
279 locate together depends on the tissue type, and scMultiSim provides  $p_n$  to tune this pattern. Fig. 3f shows the  
280 effect of varying  $p_n$ . The left figure in the panel was generated with  $p_n = 1$ , showing strong spatial clustering of

281 cells from the same cell type. The right figure in the panel was generated with  $p_n = 0.8$ , where cells from the  
282 same cell type are more spread out to enable more interactions across cell types.

283 **Correlations between interacting ligands and receptors.** scMultiSim simulates CCIs between single cells as well  
284 as between cell types. We validate the simulated CCI effects by comparing the correlations of expression levels  
285 between (i) neighboring cells with CCI, (ii) neighboring cells without CCI, and (iii) non-neighbor cells (Methods).  
286 As shown in Fig. 3e (using dataset MT3a), cells with CCI have an average pairwise correlation of 0.1, whereas  
287 cells without CCI exhibit approximately zero correlation, which is expected. We noticed that neighboring cells  
288 without CCI still have a slightly higher correlation compared to non-neighbor cells, which may be attributed to  
289 the dynamic nature of cell differentiation, where cells are evolving into new cell types over time, and CCI effects  
290 involved in an earlier cell type may remain in the final step.

### 291 **scMultiSim simulated datasets match real data**

292 We show that scMultiSim's output single cell gene expression data can statistically resemble real data. We  
293 used a spatially resolved single cell gene expression dataset measured with seqFISH+ technology [19; 16], and  
294 generated simulated data to match this real dataset (Methods). We used dyngen [7] as a baseline simulator to  
295 compare with, as it is also a *de novo* multi-modality simulator that shares a few functionalities with scMultiSim  
296 (Table S1). We compare the simulated data with real data in terms of the following properties: library size, zero  
297 counts per cell, zero counts per gene, mean count per gene, variation per gene, and the ratio between zero count  
298 and mean count per gene (Fig. 3g).

299 Fig. 3g shows that the library size, zero counts per cell, zero counts per gene and mean counts per gene  
300 simulated by scMultiSim are closer to that of real data than the dyngen simulated data, and both scMultiSim and  
301 dyngen are able to simulate data with realistic variation per gene. There is also usually a negative correlation  
302 between zero counts and mean counts in real data, and scMultiSim is able to simulate this relationship, matching  
303 well with the reference data.

### 304 **Benchmarking computational methods using scMultiSim**

305 We next show that scMultiSim can be used to benchmark a board range of computational tasks in single cell  
306 genomics, including clustering, trajectory inference, multi-modal data integration, RNA velocity estimation, GRN  
307 inference and CCI inference using spatially resolved single cell gene expression data. Using scMultiSim, we  
308 studied the performance of several recent methods on each task, and also investigated the effect of particular  
309 parameters for some of the benchmarks. As far as we know, scMultiSim is the only simulator that can benchmark  
310 all these tasks. It is noteworthy that our intention is not to perform a comprehensive benchmarking analysis, but  
311 rather to show evidence of scMultiSim's broad applications. We anticipate that these benchmarks can encourage  
312 forthcoming researchers to discover more use cases of scMultiSim.

### 313 **Benchmarking clustering and trajectory inference methods**

314 We first applied scMultiSim to test methods for two classic problems: cell clustering and trajectory inference,  
315 using the scRNA-seq modality in our discrete main datasets (MD, Table 1). We tested five clustering methods,  
316 PCA-KMeans, CIDR [39], SC3 [33], TSCAN [28], and Seurat [23] (Fig. 4a). For each method and each dataset in  
317 the main datasets, we vary the parameter “number of clusters”. Since Seurat does not provide direct control over  
318 the number of clusters, we varied the resolution parameter instead and plotted using the number of clusters  
319 in the results. From Fig. 4a, all methods have the best performance when the cluster number is the true  
320 value. In general, Seurat and SC3 have better performance than the others, which is consistent with previous  
321 benchmarking [17]. TSCAN performs better than PCA-KMeans in our results which is not the case in [17]. We  
322 also show the comparison separately for  $\sigma_{\text{cif}} = 0.1$  and  $\sigma_{\text{cif}} = 0.5$  in Fig. S3a-b. Comparing Fig. S3a with Fig. S3b,  
323 the methods generally have higher ARI with a smaller  $\sigma_{\text{cif}}$ , which is expected. Additionally, Seurat’s recommended  
324 resolution range (0.4-1.2) provides an accurate estimation of the number of clusters (Fig. S3c).

325 We evaluated the performance of five trajectory inference methods (PAGA [65], Monocle [49], Slingshot [57],  
326 MST [51], pCreode [25]) on tree-structured trajectories using the MT datasets. The result is shown in Fig. 4b,  
327 where we calculated the  $R^2$  and  $k$ NN purity (Methods) for each separate lineage in each dataset. Overall, PAGA  
328 Tree and Slingshot have the best performance, which is in line with results shown in previous benchmarking  
329 efforts [51; 69]. When comparing results on datasets with  $\sigma_{\text{cif}} = 0.1$  and  $\sigma_{\text{cif}} = 0.5$  (Fig. S4a-b), we again see  
330 that smaller  $\sigma_{\text{cif}}$  corresponds to better results. Furthermore, we tested on a simpler linear trajectory dataset ML1a  
331 (Fig. S4d), and the result was in line with a previous result shown in scDesign3 [54], which used a similar linear  
332 trajectory.

### 333 **Benchmarking multi-modal and multi-batch data integration methods**

334 A number of computational methods have been proposed to integrate single cell genomics data from multiple  
335 modalities and multiple batches [1]. We benchmarked three recently proposed multi-modal integration methods:  
336 Seurat bridge integration (Seurat-bridge) [24], UINMF [34] and Cobolt [21] that can integrate data matrices from  
337 multiple batches and modalities. We use all 144 main datasets to test their performance under various types of  
338 cell population. Each main dataset is divided into three batches (with batch effect 3), then the scRNA-seq data  
339 from batch 2 and scATAC-seq data from batch 3 are dropped intentionally to mimic a real scenario where some  
340 modalities are missing in certain batches (Fig. 4d). We use the following metrics to evaluate the performance  
341 of the integration methods: Adjusted Rand Index (ARI) and Normalized Mutual Information (NMI) as the metrics  
342 for cluster identity preservation, and Graph Connectivity and Average Silhouette Width (ASW) as metrics for  
343 batch mixing (Methods). These metrics were used in a recent paper on benchmarking single cell data integration  
344 methods [40].

345 The result is shown in Fig. 4c. Since Seurat-bridge does not output the latent embedding for the “bridge”

346 dataset (batch 1 in Fig 4d), only the two matrices from batches 2 and 3 (colored in Fig. 4d) were used for  
347 evaluation. We observe that UINMF has the best performance in terms of all four measurements. Seurat-bridge  
348 and Cobolt have comparable ARI and NMI but Cobolt has better batch mixing scores. When comparing the  
349 ARI and NMI scores for  $\sigma_{\text{cif}} = 0.1$  and  $\sigma_{\text{cif}} = 0.5$ , one can observe that these cell identity preservation scores  
350 are higher with smaller  $\sigma_{\text{cif}}$ . Comparing different cell population structures, we see that continuous populations  
351 (“Linear” and “Tree”) have lower ARI and NMI scores than discrete populations, potentially because that metrics  
352 like ARI and NMI are better suited for discrete populations.

353 We then ran the integration methods on a large dataset with 3000 cells and visualized the integrated  
354 latent embedding in Fig. S5, which helped us to understand each method’s behavior. We noticed that while  
355 Seurat-bridge has lower graph connectivity and ASW scores, different batches are located closely (but do not  
356 overlap) in the visualized latent space. That the reference and query data in the latent space do not overlap can  
357 cause the low batch mixing scores, but may not affect the ability of label transfer.

### 358 **Benchmarking RNA velocity estimation methods**

359 We demonstrate scMultiSim’s ability of benchmarking RNA velocity estimation methods by running two  
360 representative RNA velocity inference methods, scVelo [3] and VeloCyto [35], on the simulated data. We compare  
361 all three models in scVelo package, including the deterministic, stochastic, and dynamical models. The auxiliary  
362 dataset  $V$  (Table 2) was used, which contains 72 datasets of different numbers of cells and genes, with or without  
363 GRN. We evaluate the accuracy of inferred RNA velocity using cosine similarity score. The score measures the  
364 degree of mismatch between the direction of inferred and ground truth velocity, where a higher score shows a  
365 better inference result (Methods).

366 From the result shown in Fig. 4f, scVelo’s deterministic model has the highest cosine similarity score on all  
367 datasets. On the other hand, the dynamical model of scVelo, being considered a generalized version of VeloCyto,  
368 does not produce the best result. Interestingly, Gorin et al. also discussed a similar performance issue of the  
369 dynamical scVelo. They mentioned that the mismatch between the implicit assumption of dynamical scVelo and  
370 the true biological dynamics could be the cause of the performance issue [22]. In spite of the performance  
371 differences, the similarity scores are shown to be only around 0.2 for all methods. We suspect that it is the  
372 intrinsic noise within the simulated dataset that affect the inference accuracy of all methods. We further conduct  
373 experiments comparing the accuracy of inferred RNA velocity with and without  $k$ NN smoothing (Methods). By  
374 using  $k$ NN smoothing, the inferred RNA velocity of each cell is further averaged with the velocity of all its  
375 neighboring cells. Since  $k$ NN smoothing helps to reduce the noise effect on the inferred velocity, we expect  
376 that the overall performance should improve after the smoothing. The experiment results validate our assumption  
377 (Fig. 4e), where the average performance of all methods increases to 0.63. The experiments show that the  
378 intrinsic noise within the sequencing dataset heavily affects the accuracy of RNA velocity inference methods, and  
379 it is still a challenging task to infer RNA velocity from noisy scRNA-seq datasets.

## 380 **Benchmarking GRN inference methods**

381 Using scMultiSim, we benchmarked 11 GRN inference methods which were compared in a previous  
382 benchmarking paper [48]. Using the predicted networks, we calculate the AUROC (area under receiver operating  
383 characteristic curve) as well as the AUPRC (area under precision-recall curve) ratio, which is the AUPRC divided  
384 by the baseline value (a random predictor). These metrics were also used in previous benchmarking work [48].

385 We show results on the 144 main datasets in Fig. 5a. To further inspect the performance in a less-noisy  
386 scenario, we also generated auxiliary datasets G (Table 2) with a linear trajectory and without CCI effect.  
387 We benchmarked the methods using true counts and observed counts in G, respectively. The result of G is  
388 shown in Fig. 5b. All datasets use the same 100-gene GRN from [15]. We observed that PIDC [9] has the  
389 best overall performance, especially on true counts. Other methods like GENIE3 [27] and GRNBOOST2 [42]  
390 also have noteworthy precision. We then examined the effect of technical noise on the performance of GRN  
391 inference methods. On observed counts, both the AUPRC ratio and AUROC value suffer from a decline,  
392 indicating that it is significantly harder to infer the GRN from noisy data. However, PIDC continues to have  
393 the highest AUPRC and AUROC values, showing that its performance is more resistant to technical noises.  
394 SINCERITIES [45], PPCOR [32] and SINGE [14] perform well and beat GENIE3 and GRNBOOST2 on observed  
395 counts. Nevertheless, the absolute AUPRC values of all methods, even on true counts, are still far from satisfying,  
396 indicating that GRN inference is still a challenging problem.

397 Notably, the ordering of the methods tested using true counts is generally consistent with the ordering reported  
398 in [48] even though a different ground truth GRN was used. This fact not only validates the previous results  
399 but also suggests that scMultiSim can generate GRN-incorporated gene expression data comparable to other  
400 simulators. It indicates scMultiSim's practicality in benchmarking computational methods that involve GRNs.

## 401 **Benchmarking CCI inference methods**

402 Spatially resolved single cell gene expression data provides a powerful tool for understanding cellular processes,  
403 tissue organization, and disease mechanisms at the single cell level. Multiple methods have been proposed  
404 recently to infer CCIs based on spatial cell locations. However, these methods have yet to be compared in this  
405 relatively new field due to the scarcity of biological ground truth and spatial transcriptomics simulators.

406 We benchmarked three CCI inference methods based on spatially resolved single cell gene expression data,  
407 namely Giotto [16], SpaOTsc [5] and SpaTalk [53]. We run Giotto and SpaOTsc on the main datasets and show  
408 the result in Fig. 5c. Since SpaTalk needs a minimum of 3 genes from the receptor to a downstream activated  
409 TF, we also generated an auxiliary dataset C (Table 2) using an artificial GRN with long pathways to satisfy such  
410 requirement (Fig. S6a). There are totally eight C datasets with 500 cells, 200 genes and a linear trajectory, and the  
411 result is shown in fig. S6b. Again, we used AUPRC and AUROC as the metrics. When calculating the PRC and  
412 ROC curves, we applied different thresholds on Giotto's significance score and SpaTalk's Bonferroni corrected

413 p-values. Considering both AUROC and AUPRC, Giotto has the best performance with an average AUROC of  
414 0.68 and AUPRC of 0.54 on the main datasets. SpaTalk outputs too many identical p-values for different datasets  
415 on dataset C, causing the ROC and PRC curves to look unusual. Nevertheless, it has noteworthy performance in  
416 terms of AUROC and AUPRC values but is less accurate and stable than Giotto. The benchmarking results show  
417 that Giotto could be a versatile yet robust choice for CCI inference.

## 418 Discussion

419 We presented scMultiSim, a simulator of single cell multi-omics data which is able to incorporate biological  
420 factors including cell population, chromatin accessibility, RNA velocity, GRN and spatial CCIs to the output  
421 data. We demonstrated the presence of these simulated factors in the generated data, verified the relationship  
422 across modalities, and showcased the versatility of scMultiSim through benchmarking on various computational  
423 problems. Furthermore, by obtaining consistent benchmarking results with previous works like BEELINE [48] and  
424 dyngen [7], the simulated biological effects are validated to be practical and ready for real-world use.

425 Compared to existing simulators that mainly model one or two biological factors, scMultiSim generates data  
426 with more biological complexity similar to real data. This additional complexity enables researchers to better  
427 estimate the real-world performance of their methods on noisy experimental data. Furthermore, with the coupled  
428 data modalities in the output, researchers can benchmark computational methods that use multiple modalities,  
429 which was previously impossible.

430 scMultiSim's extensibility and versatility are central to its design, making it easy to include more biological  
431 factors and modalities in its simulations. For example, the framework used to model chromatin regions (RIV)  
432 and genes (GIV) can also be extended to include other data modalities, such as the protein abundance data.  
433 Additionally, we have shown that our CIF/GIV model is versatile enough to mathematically represent the effects  
434 of various biological mechanisms like GRNs and CCIs. In addition to the standard functions of scMultiSim, the  
435 model can be expanded to consider more realistic scenarios. For instance, the GRN can be set to a cell-specific  
436 and cell-type-specific mode, allowing for a more precise simulation of regulatory interactions. Moreover, the  
437 scATAC-seq data and scRNA-seq data can follow different trajectories or clustering structures, while the cell  
438 clusters can form less regular shapes than the current convex shapes.

439 scMultiSim's usability is supported by several key features. First, it requires minimal and easy-to-construct  
440 input. For example, users do not need to prepare a backbone for the trajectory to control the cell population;  
441 instead, only a plain R phylogenetic tree or a text file with the Newick format tree is needed. Second, scMultiSim  
442 has transparent parameters that are self-explanatory and have a clear effect on the result. The user explicitly sets  
443 crucial metrics such as the number of cells and genes. Third, scMultiSim's separated biological effects provide  
444 great flexibility. For example, the GRN can affect cell population shapes, but obtaining the desired trajectory using  
445 GRN alone is difficult without explicit control of the cell population. scMultiSim's diff and non-diff CIF mechanism

446 allows users to set the trajectory to any shape without affecting the GRN effects. Users can also let the GRN  
447 control the trajectory by increasing the number of non-diff CIF.

448 We underline that scMultiSim's major advantage is its ability to encode various factors into a single versatile  
449 model, thus creating a comprehensive multi-modal simulator that can benchmark an unprecedented range  
450 of computational methods. More importantly, the coupled data modalities in the output jointly provide more  
451 information than a single modality alone, making it ideal for designing and benchmarking new methods on  
452 multi-omics data. We believe that scMultiSim has the potential to be a powerful tool for fostering the development  
453 of new computational methods for single-cell multi-omics data. Moreover, as more benchmarks are conducted, it  
454 can help researchers select the appropriate tool based on the type of data they are working with, leading to more  
455 accurate and reliable analyses.



## 456 Methods

### 457 A. The Beta-Poisson model and intrinsic noise

The master equation of the kinetic model represents the steady state distribution of a gene's expression level given its kinetic parameters,  $k_{on}$ ,  $k_{off}$ , and  $s$  [43]. The Beta-Poisson model was shown to be equivalent to the master equation [31] with faster calculation. The gene expression level  $x$  (which is also the mRNA count) can be sampled from the following distribution:

$$y = \text{Beta}(k_{on}, k_{off})$$

$$x = \text{Poisson}(y \cdot s)$$

458 Using the above Beta-Poisson distribution to generate the gene expression level is one mode to obtain mRNA  
459 count for a gene in a cell. This works if we only need to generate the spliced mRNA counts. If users also need to  
460 generate unspliced mRNA counts and RNA velocity, the other mode, called the "full kinetic model" is used. The  
461 Beta-Poisson model is used by default when only generating spliced counts for lower running time.

The sampling process from the Beta-Poisson distribution to obtain  $x$  introduces intrinsic noise to the data, which corresponds to the intrinsic noise in real data caused by transcription burst. The theoretical mean of the kinetic model, which is  $(\frac{k_{on}}{k_{on}+k_{off}} \cdot s)$ , corresponds to the gene expression level of the gene with no intrinsic noise. We introduced parameter  $\sigma_i$  which controls the intrinsic noise by adjusting the weight between the random samples from the Poisson distribution and the theoretical mean:

$$x_{\sigma_i} = \sigma_i \cdot x + (1 - \sigma_i) \cdot (\frac{k_{on}}{k_{on} + k_{off}} \cdot s)$$

462 The intrinsic noise in scRNA-seq data is hard to reduce in experiments due to the snapshot nature of scRNA-seq  
463 data. The parameter  $\sigma_i$  allows users to investigate the effect of intrinsic noise on the performance of the  
464 computational methods.

### 465 B. Cell Identity Factors (CIFs) and Gene Identity Vectors (GIVs)

466 The length of the CIF and GIV, denoted by  $n_{cif}$ , can be adjusted by the user. Overall, we have a  $n_{cell} \times n_{cif}$  CIF  
467 matrix for each kinetic parameter (Fig. S1a), where each row is the CIF vector of a cell. Correspondingly, we also  
468 have the  $n_{cif} \times n_{gene}$  Gene Identity Vectors (GIV) matrix, (Fig. S1b) where each column is linked to a gene, acting  
469 as the weight of the corresponding row in the CIF matrix, i.e. how strong the corresponding CIF can affect the  
470 gene. In short, CIF encodes the *cell identity*, while GIV encodes the *strength of biological effects*. Therefore, by  
471 multiplying the CIF and GIV matrix, we are able to get a  $n_{cell} \times n_{gene}$  matrix, which is the desired kinetic parameter  
472 matrix with the cell and gene effects encoded. Each cell has three CIF vectors corresponding to the three kinetic  
473 parameters  $k_{on}$ ,  $k_{off}$ , and  $s$ , and similarly for the GIV vectors (Fig. S1a-b).

### 474 C. diff-CIF generates user-controlled trajectories or clusters.

475 When generating data for cells from more than one cell type, the minimal user input of scMultiSim is the cell  
 476 differentiation tree, which controls the cell types (for discrete population) or trajectories (for continuous population)  
 477 in the output. The generated scRNA-seq and scATAC-seq data reflect the tree structure through the diff-CIF  
 478 vectors. The diff-CIF vectors are generated as follows: starting from the root of the tree, a Gaussian random walk  
 479 along the tree (Fig. 2a) is performed for each cell to generate the  $n_{\text{diff-CIF}}$  dimension diff-CIF vector. Parameter  
 480  $\sigma_{\text{cif}}$  controls the standard deviation of the random walk, therefore a larger  $\sigma_{\text{cif}}$  will produce looser and noisier  
 481 trajectory structures. Another parameter  $r_d$  is used to control the relative number of diff-CIF to non-diff-CIF. With  
 482 a larger  $r_d$ , trajectories are clear and crisp in the output; with a smaller  $r_d$ , the trajectory is vague, and the shape  
 483 of the cell population is more controlled by other factors like GRN. For a discrete population, only the cell types at  
 484 the tree tips are used; then cells of each type are shifted by a Gaussian distribution, controlled by the same  $\sigma_{\text{cif}}$   
 485 parameter. Therefore, a smaller  $\sigma_{\text{cif}}$  will produce clearer cluster boundaries.

486 For a heterogeneous cell population, cells have different development stages and types. Users should input a  
 487 cell differentiation tree where each node represents a cell type. The tree provides a backbone for the trajectory  
 488 in the cell population. Each dimension of the diff-CIF vector is sampled along the tree via brownian motion. First,  
 489 cells start at the root of the tree; then for each dimension, the diff-CIF value for all cells  $v$  is

$$v_i = \sum_{j=1}^i q_j \text{ where } q_j = \mathcal{N}(0, \sigma_j).$$

490  $\sigma_j$  is the distance along the tree between cell  $j$  and  $j - 1$ . Alternatively, users can use an impulse model (using  
 491 the implementation in SymSim). The lengths of the non-diff-CIF and diff-CIF vectors can be controlled by the  
 492 user. More diff-CIFs will result in a more clear trajectory pattern in the cell population, which corresponds to the  
 493 input tree. With very few diff-CIFs, the cell population is mainly controlled by the GRN.

### 494 D. tf-CIF and GIV encode the GRN effects

495 To encode GRN effect in the simulated single cell gene expression data, the GIVs and CIFs are designed to  
 496 include a “TF part” (Fig. S1a). Cells are generated one by one along the given cell differentiation tree, where the  
 497 expressions of the TFs in the  $t^{\text{th}}$  cell affect the gene expression of cell  $t + 1$ . Formally, the  $i^{\text{th}}$  position of the TF  
 498 part (corresponding to the  $i^{\text{th}}$  TF) of in the CIV of cell  $t + 1$  is calculated as:

$$\text{tf-CIF}_i^{(t+1)} = \frac{\mathbf{x}_i^{(t)}}{\mathbf{x}_i^{(t)} + \frac{1}{n} \sum_l \mathbf{x}_l^{(t)}} \quad \forall i \in \text{TFs} \quad (1)$$

499 where  $\mathbf{x}_i^{(t)}$  is the expression level of the  $i^{\text{th}}$  TF in the  $t^{\text{th}}$  cell. The corresponding tf-CIF for the root cell is sampled  
 500 randomly from the Gaussian distribution  $\mathcal{N}_{\text{cif}}$  supplied by the user.

501 The TF part of the GIV for a gene also has length of  $n_{\text{TF}}$  (Fig. S1b). Considering all genes, we have a  $n_{\text{gene}} \times$   
 502  $n_{\text{TF}}$  matrix, which we call the GRN effect matrix. This matrix encodes the ground truth GRN that is supplied by

the user. Naturally, the GRN effect matrix is included in the GIV when calculating the  $s$  parameter, where the value at  $(i, j)$  is the regulation strength of TF  $j$  on gene  $i$ . Therefore, a larger regulation strength will lead to higher  $s$ , and consequently, higher expressions for the target genes. For  $k_{on}$  and  $k_{off}$ , the tf-CIF vector is sampled using  $\mathcal{N}_{cif}$ , assuming that the GRN does not affect the active state of a gene. However, in the scenario where it is desired to model GRN effect also in  $k_{on}$  and  $k_{off}$ , similar GRN effect matrix for  $s$  can be used for  $k_{on}$  and  $k_{off}$ .

scMultiSim also allows the use of ground truth GRNs which are cell specific. In this mode, random GRN edges are generated or deleted gradually along the pseudotime at a user-controlled speed. When simulating each cell, the tf-GIV will be filled with the current GRN effect matrix. The cell-specific GRN ground truth is outputted in this mode.

## E. lig-CIF and GIV encode cell-cell interactions

When simulating spatial transcriptomics data with CCI effects, we used a 2-D  $k \times k$  grid to model the spatial locations of cells (Fig. S1d). The grid size  $k$  is large enough to accommodate the  $n$  cells (can be specified by the user; if not provided, use 250% of cell number by default). A cell can have at most  $n_{nbs}$  neighbors with CCI (within the blue circle's range in Fig. 2a, and this radius can be adjusted). Therefore, the ligand CIF and GIV are of length  $n_{lig} \cdot n_{nbs}$ , where  $n_{lig}$  is the number of ligands.

The lig-GIV vector contains the CCI strength values, for example, the "n2 lg3" entry in Fig. 2a indicates how strong the ligand 3 from the neighbor at position 2 can affect the receptor 2 of this cell. The lig-CIF of each cell will inherit from its previous cell during the simulation process, which is similar to the tf-CIF mentioned above. Each entry of the lig-CIF vector corresponds to a ligand from one neighbor. The same Gaussian distribution  $\mathcal{N}_{cif}$  is used for  $k_{on}$  and  $k_{off}$ . For  $s$ , due to the similarity of the ligand-receptor pairs and the TF-target pairs, we use a similar strategy as tf-CIF (shown in Eq. 1): cell  $i$ 's lig-CIF is the normalized vector of cell  $i - 1$ 's gene expression counts of the ligand genes (See Fig. 2a, Fig. S1).

At each step  $t$ , a new cell is considered to be born and added to the grid. When adding a new cell, it has a probability of  $p_n$  to be a neighbor of an existing cell with the same cell type. We also provide other strategies to place a new cell, including (1) all cells placed at a random location, and (2) only the first  $m$  cells are randomly placed, and the remaining follow  $p_n$ . A pre-defined cell differentiation tree is required as input to define the differentiation topology in the cells. A new cell will always be in the initial state at the root of the differential tree. At each step, an existing cell moves forwards along a random path in the cell differential tree, representing the cell development. The gene expressions in the final step are output as the observed data. Fig. S1 shows the structure for the CCI mode.

To generate ground truth CCIs both at the cell types level and single cell level, scMultiSim pre-defines a ligand-receptor database, represented by a user input  $m \times 3$  matrix  $S$ . There are  $m$  ligand-target pairs in total that correspond to each row of  $S$ . For each pair  $i$ , there are three parameters: ligand gene  $L_i$ , receptor gene

536  $T_i$ , and the effect  $E_i$ , representing how strongly the ligand can affect the expression of the receptor. For each  
 537 cell type pair, the ground truth CCI between these two cell types are sampled from the ligand-receptor database  
 538 (corresponds to the columns in  $S$ ). For each neighboring cell pair, the ground truth CCIs between them follow  
 539 the cell-type-level ground truth CCIs: if the two cells belong to two cell types  $C_1$  and  $C_2$  respectively (where  $C_1$   
 540 can be the same as  $C_2$ ), the CCIs between these two cells follow the CCIs defined in  $S$  corresponding to pair  
 541  $(C_1, C_2)$ . Users can have further fine-grained control for each cell pair by letting it use a subset of ligand-receptor  
 542 pairs sampled from the cell-type level ground truth.

543 Although we collect cells at the last time point as our output (which is the case for real data), different cell  
 544 types are guaranteed to present in the last step since the cells are added at different time steps, therefore having  
 545 different development stages. In addition, we let the same cell (at the same location) have the same diff-CIF  
 546 across different time steps, so the trajectory encoded in the diff-CIF is preserved in the final step. A cell's TF and  
 547 ligand CIF for the current step is inherited from the previous one to make sure other factors stay the same.

548 We use the following steps to calculate the correlation between the expressions of neighboring cells in Fig. 3e.  
 549 First, a specific ligand-receptor pair  $(l, r)$  is chosen. Let  $T(a, b) = \{\text{true}, \text{false}\}$  denote that there is CCI between  
 550 cell  $a$  and cell  $b$  for  $(l, r)$ . Then, for each cell  $i$ , we get its neighbor list  $n_i$ , which is a vector of 4 cells. A vector of 4  
 551 non-adjacent cells  $m_i$  is also randomly sampled for this cell. Thus, let  $x_c^g$  denote the gene expression of cell  $c$  and  
 552 gene  $g$ . we calculate the "neighbor cells with CCI" correlation using the pairs  $\{(x_i^l, x_j^r) | j \in n_i, T(i, j) = \text{true}\}$ , the  
 553 "neighbor cells without CCI" correlation using the pairs  $\{(x_i^l, x_j^r) | j \in n_i, T(i, j) = \text{false}\}$ , and the "non-neighbor  
 554 cells" correlation using the pairs  $\{(x_i^l, x_j^r) | j \in m_i\}$ . Cell pairs of the same type are ignored while calculating the  
 555 correlations because they tend to have similar expressions.

## 556 F. Generating the Gene Identity Vectors

557 A gene's GIV vector has the same length as the CIF vectors. The values in the GIV of a gene act as the weights  
 558 of the corresponding factors in the CIF, *i.e.*, how strong the corresponding CIF can affect the gene (Fig. 2a). If we  
 559 have  $n_{\text{gene}}$  genes, we obtain a GIV matrix of size  $n_{\text{cif}} \times n_{\text{gene}}$ .

560 It can be divided into four submatrices as shown in Fig. S1b. For  $k_{on}$  and  $k_{off}$ , the nd-CIF and diff-CIF are  
 561 sampled from distribution  $\mathcal{G}$  as shown below:

$$\begin{cases} \mathcal{N}_{\text{giv}} & \text{w.p. } 1 - p_0^{\mathcal{G}} \\ 0 & \text{w.p. } p_0^{\mathcal{G}} \end{cases}$$

562 where  $p_0^{\mathcal{G}}$  is a parameter specifying the probability of being zero, and  $\mathcal{N}_{\text{giv}}$  is a user-adjustable gaussian  
 563 distribution. tf-GIV and lig-GIV are all zeros since TF/ligands affect  $s$  only. For  $s$ , the tf-GIV submatrix is the  
 564 GRN effect matrix, *i.e.* a  $n_{\text{TF}} \times n_{\text{gene}}$  matrix where the entry at  $(i, j)$  is the regulation effect between TF  $i$  and  
 565 gene  $j$ . Similarly, the lig-GIV submatrix is the cell-cell interaction effect matrix. The nd-GIV submatrix is sampled  
 566 from  $\mathcal{G}$ . For diff-GIV, we do the following steps to incorporate the connection between TFs and regulated genes:

567 (1) Randomly select 2 GIV entries for each TF gene and give them a fixed small number. (2) For every target  
 568 gene, it should use the same GIV vector as its regulators. If a gene has multiple regulators, its gene effects will  
 569 be the combination of that of the regulators. This is achieved by multiplying the  $n_{\text{diff}} \times n_{\text{TF}}$  GIV matrix in (1) and  
 570 the  $n_{\text{TF}} \times n_{\text{gene}}$  effect matrix. If a gene is both a TF and target, its GIV will be  $0.5 \cdot ((1) + (2))$ .

## 571 G. Simulating scATAC-seq data and relationship between scATAC-seq and scRNA-seq

572 Since scMultiSim incorporates the effect of chromatin accessibility in the gene expressions, the scATAC-seq data  
 573 is simulated before the scRNA-seq data. The cell types in the scATAC-seq data can follow the same differentiation  
 574 tree as in the scRNA-seq data (the scATAC-seq and scRNA-seq data share the same cells) or can follow a  
 575 different tree (to reflect the difference between modalities).

576 Similar to GIV, we use a randomly sampled *Region Identity Vector (RIV)* matrix to represent the chromatin  
 577 regions. Following the same mechanism, we multiply the CIF and RIV matrix, and obtained a “non-realistic  
 578 scATAC-seq” data matrix. Next, the scATAC-seq data matrix is obtained by scaling the “non-realistic” scATAC-seq  
 579 data to match a real distribution learned from real data. This is a step to capture the intrinsic variation of  
 580 the chromatin accessibility pattern, which we will also apply to the kinetic parameters when generating gene  
 581 expressions.

582 The RIV matrix is sampled from a distribution  $\mathcal{R}$  similar to  $\mathcal{G}$ :

$$\begin{cases} \mathcal{N}_{\text{riv}} & \text{w.p. } 1 - p_0^{\mathcal{R}} \\ 0 & \text{w.p. } p_0^{\mathcal{R}} \end{cases}$$

583 where  $p_0^{\mathcal{R}}$  is the probability of being zero and  $\mathcal{N}_{\text{riv}}$  is a user-adjustable Gaussian distribution. With the CIF and  
 584 RIV matrices, the  $n_{\text{cell}} \times n_{\text{region}}$  scATAC-seq can be generated by (1) multiplying the CIF matrix by the RIV matrix,  
 585 (2) scale the matrix to match the real data distribution, and (3) adding intrinsic noise (sampled from a small  
 586 Gaussian) to the scATAC-seq data. In Step (2), we use the same rank-based scaling process as used for the  
 587 kinetic parameters as described in Section “Preparing the kinetic parameters” above, and the real scATAC-seq  
 588 data distribution is obtained from the dataset in [11].

589 To incorporate the relationship between scATAC-seq and scRNA-seq data, we use the scATAC-seq data to  
 590 adjust the  $k_{\text{on}}$  parameter that is used to generate the scRNA-seq data, considering that chromatin accessibility  
 591 affects the activated status of genes. First, a region-to-gene matrix (Fig. 1b) is generated to represent the mapping  
 592 between chromatin regions and genes, where a gene can be regulated by 1-3 consecutive regions. Users  
 593 can input a region distribution vector  $\mathbf{r}$ , for example,  $(0.1, 0.5, 0.4)$  means a gene can be regulated by three  
 594 regions, and the probability of it being regulated by one, two and three consecutive regions are 0.1, 0.5 and 0.4,  
 595 respectively. The scATAC-seq data is also used to adjust  $k_{\text{on}}$  as described in the following section.

## 596 H. Preparing the kinetic parameters

597 The kinetic parameters,  $k_{on}$ ,  $k_{off}$  and  $s$  are needed when generating single cell gene expression data (mRNA  
598 counts) using the kinetic model or Beta-Poisson distribution (Fig. 1b). While the basic idea is to get the parameter  
599 matrix using CIFs and GIVs (Fig. 1b), the three parameters go through different post-processing after the step of  
600 CIF  $\times$  GIV. We first denote the result of CIF  $\times$  GIV for  $k_{on}$ ,  $k_{off}$  and  $s$  as  $M_1$ ,  $M_2$  and  $M_3$ , respectively.

601 (i)  $k_{on}$ . Since chromatin accessibility controls the activation of the genes, the scATAC-seq data is expected  
602 to affect the  $k_{on}$  parameter. We first prepare a  $n_{region} \times n_{gene}$  0-1 region-to-gene matrix  $Z$  using  $\mathbf{r}$ , where  $Z_{ij}$   
603 indicates region  $i$  is associated with gene  $j$  ( $Z$  is outputted as the region-to-gene matrix). We multiply the  
604 scATAC-seq matrix with  $Z$  to get the  $n_{cell} \times n_{gene}$  parameter matrix  $M'_1$ . Since the scATAC-seq data is sparse,  
605 there are many zeros in  $M'_1$ . Thus, we replace the zero entries in  $M'_1$  with the corresponding entries in  $M_1$   
606 (scaled to be smaller than the smallest non-zero entry in  $M'_1$ ) to help differentiate the zero entries. Finally,  $M'_1$  is  
607 sampled to match the distribution of  $k_{on}$  inferred from real data.

608 (ii)  $k_{off}$ . The parameters are obtained by scaling  $M_2$  to match the real data distribution. For both  $k_{on}$  and  $k_{off}$ ,  
609 it is possible to adjust the bimodality of gene expressions [69] through an optional bimodal factor  $B$ . A larger  $B$   
610 will downscale both  $k_{on}$  and  $k_{off}$ , therefore increasing the bimodality.

611 (iii)  $s$ . The parameters are obtained by scaling  $M_3$  to match the distribution of  $s$  inferred from real data. Then,  
612 users can also use a “scale.s” parameter to linearly scale  $s$ . It allows us to adjust the size of cells – some datasets  
613 may tend to large cells and some tend to have small cells depending on the cell types being profiled.

614 When scaling a matrix ( $M'_1$ ,  $M_2$ , or  $M_3$ ) to match a reference distribution (eg. the distributions of  $k_{on}$ ,  $k_{off}$   
615 and  $s$  estimated from real data), the procedure is as follows: denoting the reference distribution by  $\mathcal{D}$ , the matrix  
616 to rescale by  $X$ , and the number of elements in  $X$  by  $n$ , we sample  $n$  ordered values from  $\mathcal{D}$ , then replace  
617 the data in  $X$  using the same order. scMultiSim uses the reference kinetic distribution parameters provided in  
618 SymSim [69], where the kinetic parameters are estimated from real data via an MCMC approach. The data  
619 used are the UMI-based dataset of 3005 cortex cells by Zeisel et al. [67], and a non-UMI-based dataset of 130  
620 IL17-expressing T helper cells (Th17) by Gaublotte et al [20].

## 621 I. Generating RNA velocity with the full kinetic model

When using the full kinetic model, scMultiSim can generate the spliced and unspliced counts for each cell from  
the kinetic parameters. The starting spliced count  $x_s$  and unspliced count  $x_u$  for a cell are the previous cell's  
counts on the differential tree. For the first cell, the spliced/unspliced counts are

$$x_s = \frac{s \cdot k_{on} \cdot \beta}{k_{on} + k_{off}} \quad x_u = \frac{s \cdot k_{on} \cdot d}{k_{on} + k_{off}}$$

622 where  $\beta$  and  $d$  respectively represent the splicing and degradation rate of genes. Both  $\gamma$  and  $d$  are sampled from  
623 a user-controlled normal distribution.

We set the cell cycle length to be  $L = \frac{1}{k_{on}} + \frac{1}{k_{off}}$ , and divide it into multiple steps. The number of steps follows  $m = \left\lceil \frac{L}{\min(1/k_{on}, 1/k_{off})} \right\rceil$ . We also provide an optional cell length factor  $\eta_L$  parameter to scale the cycle length. The probabilities of gene switching on or off are then calculated with  $p_{on} = \frac{k_{on}}{m \cdot L}$  and  $p_{off} = \frac{k_{off}}{m \cdot L}$ . In each simulation step, we update the cell's current on/off state based on  $p_{on}$  and  $p_{off}$ , and generate the spliced/unspliced counts  $x_s$  and  $x_u$ . The spliced counts at step  $t$  are obtained by:

$$x_s^t = x_s^{t-1} + \frac{L}{m}(\beta \cdot x_u^{t-1} - d \cdot x_s^{t-1})$$

624 and the unspliced counts are obtained by:

$$x_u^t = \begin{cases} x_u^{t-1} + \frac{L}{m}(s - \beta \cdot x_u^{t-1}) & \text{if state is on} \\ x_u^{t-1} - \frac{L}{m}(\beta \cdot x_u^{t-1}) & \text{if state is off} \end{cases}$$

The outputted  $x_s$  and  $x_u$  are the values at the final step  $t = m$ . The ground truth RNA velocity is calculated as:

$$v = \beta \cdot x_u - d \cdot x_s$$

625 We obtain the KNN averaged RNA velocity by applying a Gaussian Kernel KNN on the raw velocity data, with  
626  $k = \lceil n_{\text{cell}}/50 \rceil$ . Then we normalize the velocity by calculating each cell's normalization factor  $s_i = |v_i|$ , where  $v_i$   
627 is the velocity vector for cell  $i$ .

## 628 J. Adding technical noise and batch effects to data

629 Technical noise is added to the true mRNA counts to generate observed counts (observed scRNA-seq data)  
630 (Fig. 1b). The workflow follows SymSim's approach [69]: we simulate multiple rounds of mRNA capture and PCR  
631 amplification, then sequencing and profiling with UMI or non-UMI protocols. The parameter  $\alpha$  controls the capture  
632 efficiency, that is, the rate of subsampling of transcripts during the capture step, which can vary in different cells,  
633 and user can specify it using a Normal distribution  $\alpha \sim \mathcal{N}(\alpha_\mu, \alpha_\sigma)$ . The sequencing depth  $d \sim \mathcal{N}(d_\mu, d_\sigma)$  is  
634 another parameter that controls the quality of the observed data.

635 Batch effects are added by first dividing the cells into batches, then adding gene-specific and batch-specific  
636 Gaussian noise based on shift factors. For each gene  $j$  in batch  $i$ , the shift factor is sampled from  $\text{Unif}(\mu_j -$   
637  $e_b, \mu_j + e_b)$ , where  $\mu_j \sim \mathcal{N}(0, 1)$ , and  $e_b$  is the parameter controlling the strength of batch effects. We provide  
638 several settings for adding highly expressed genes to help researchers fit the housekeeping genes in real data.  
639 scMultiSim also supports cell- and gene-wise tuning of the mRNA capture efficiency during the PCR process;  
640 therefore per-cell and per-gene metrics (such as zero count proportion and count variance) in the observed data  
641 can be controlled separately.

642 For scATAC-seq data, as the data is sampled from real data we do not explicitly simulate the experimental  
643 steps. We do provide methods to add batch effects to obtain multiple batches of scATAC-seq data.

## 644 K. Comparing statistical properties of simulated data with experimental data

645 To measure scMultiSim’s ability to generate realistic data while incorporating all the effects, we compare the  
 646 statistical properties of a real mouse somatosensory cortex seqFISH+ [19] dataset with simulated data generated  
 647 using selected parameters. The dataset, with 10000 genes and spatial locations of 523 cells, is featured in  
 648 Giotto [16]’s tutorial.

649 The scMultiSim simulated data has both GRN and CCI effects. The GRN used as input to scMultiSim is  
 650 obtained as follows: GENIE3 [27] was used to obtain an inferred GRN from the dataset, then after looking at the  
 651 output edge importance values, the top 200 edges were utilized to form a reference GRN. We used this GRN  
 652 (96 genes) and another randomly sampled 104 genes to generate a subsample of the data. We then simulated  
 653 a dataset with 200 genes and 523 cells using scMultiSim. After observing the dimension reduction of the real  
 654 dataset, a discrete cell population is assumed. We specify the cluster ground truth using the exact cell type labels  
 655 in the dataset. There are 10 cell types in total. We also used Giotto [16] to infer the cell-cell interactions between  
 656 cells. We chose the top-seven most significant ligand-receptor pairs from Giotto’s output, with p-value  $\leq 0.01$ ,  
 657 more than 10 ligand and 10 receptor cells, and the largest  $\log_2fc$  values.

658 We used dyngen [7] as a baseline simulator to compare with scMultiSim. We generated a simulated dataset  
 659 with dyngen, using the same GRN and number of cells. The cell types and cell-cell interaction ground truth were  
 660 not provided since dyngen does not support them. Yet, we supplied the raw mouse SS cortex count matrix to  
 661 dyngen’s `experiment_params` as a reference dataset.

662 We used the following metrics to compare the distribution of simulated and experimental datasets, which is  
 663 also used in [15]: library size (per cell), zero counts proportion (per cell), zero counts proportion (per gene), mean  
 664 counts (per gene), counts variance (per gene), and the relationship between zero counts and mean counts per  
 665 gene.

## 666 L. Evaluation metrics for benchmarking computational methods

When evaluating the trajectory inference methods, we calculate the coefficient of determination  $R^2$  and the  $k$ NN  
 purity for all cells on each lineage. Given the cells’ ground truth pseudotime vector  $t$  and the inferred pseudotime  
 $\hat{t}$ , the  $R^2$  is equal to the square of the Pearson correlation coefficient:

$$R^2 = 1 - \frac{\sum_i (t_i - \hat{t}_i)^2}{\sum_i (t_i - \bar{t})^2} = \rho^2(t, \hat{t})$$

667 where  $\bar{t}$  is the mean of  $t$ . Given a cell  $i$ ’s  $k$ NN neighborhood  $N_i^{\hat{t}}$  in  $\hat{t}$  and its  $k$ NN neighborhood  $N_i^t$  in  $t$ , the  $k$ NN  
 668 purity  $K_p$  for the cell is the Jaccard Index of  $N_i^t$  and  $N_i^{\hat{t}}$ .

669 The evaluation metrics used for multi-model data integration methods, Graph Connectivity and ASW, are  
 670 described as following.



671 Graph Connectivity is defined as:

$$GC = \frac{1}{|C|} \sum_{c \in C} \frac{|LCC(c)|}{|c|}$$

672 where  $C$  is all cell types,  $LCC(c)$  is in the largest connected component for cells of type  $c$ .

673 For the ASW:

$$batch\ ASW = \frac{1}{|M|} \sum_{k \in M} \frac{1}{|C_j|} \sum_{i \in C_j} 1 - |silhouette(i)|$$

674 where  $M$  is the set of all cell types, and  $C_j$  is all the cells of type  $j$ . We used the implementation in [40].

When evaluating RNA velocity inference methods, we used the *cosine similarity* between the averaged estimated velocity and the ground truth. Calculating the average of estimated velocity vectors is commonly used to reduce local noise [7]. In dyngen [7], averaged RNA velocities were calculated across cells at trajectory waypoints weighted through a Gaussian kernel using ground truth trajectory; while in scMultiSim, we averaged the raw velocity values by  $k$ NN with a Gaussian kernel and  $k = n_{\text{cells}}/50$  to achieve a similar averaging effect. Finally, cosine similarity is calculated as:

$$\frac{1}{n_{\text{cells}}} \sum_i \frac{v_i \cdot u_i}{\|v_i\| \|u_i\|}$$

675 where  $v_i$  is the ground truth velocity vector for cell  $i$ , and  $u_i$  is the predicted velocity vector.

## 676 M. Details on running clustering methods

677 We used CIDR 0.1.5, SC3 1.24.0, Seurat 4.1, and TSCAN 2.0. The parameters we specified are (1) SC3: `pct_dropout = [0,100]`, (2) Seurat: `dims.use = 30`. For PCA-Kmeans, we simply ran Kmeans clustering on  
678 the first 20 principle components using the default R implementation `prcomp` and `kmeans`. ARI is calculated by  
679 `adjustedRandIndex` from the R package `mclust`. Some code was adapted from [17].  
680

## 681 N. Details on running trajectory inference methods

682 We used the latest `dynverse` [51] package to run the trajectory inference methods. When running them, we  
683 provide the correct root cell ID, number of starting clusters and number of ending clusters. The  $R^2$  values are  
684 calculated between the inferred pseudotime and the ground truth for each separate lineage. The  $k$ NN purity value  
685 is calculated for each lineage as: for cell  $i$ , we obtain its  $k$  Nearest Neighbors  $N_i$  on the pseudotime with  $k = 50$ .  
686 Then the  $k$ NN purity for  $i$  is the Jaccard Index of  $N_i$  on the inferred pseudotime and  $N_i$  on the true pseudotime.  
687  $R^2$  measures the correctness of inferred pseudotime, but when there are multiple branches in the trajectory,  $R^2$   
688 does not distinguish cells with similar pseudotime but are on different branches. In this case, the  $k$ NN purity  
689 serves as a complementary measurement that measures the correctness of inferred trajectory backbone.

## 690 O. Details on running data integration methods

691 We use all 144 main datasets. Technical noise and batch effects were added using default parameters (non-UMI,  
 692  $\alpha \sim \mathcal{N}(0.1, 0.02)$ ,  $\text{depth} \sim \mathcal{N}(10^5, 3000)$ , ATAC observation probability 0.3). All integration methods were run on  
 693 the scRNA and scATAC data with technical noise and batch effect. For Seurat-bridge, we followed the vignette  
 694 “Dictionary Learning for cross-modality integration” in Seurat 4.1.0 using the default parameters. For UINMF,  
 695 we used the latest GitHub release. We followed the “UINMF integration of Dual-omics data” tutorial and ran the  
 696 `optimizeALS` method using  $k = 12$ . For Cobolt, we used the GitHub version `cd8015b`, with 10 latent dimensions,  
 697 learning rate 0.005. If the loss diverged, we automatically retry with learning rate 0.001. The metrics, including  
 698 ARI, NMI, Graph Connectivity, and ASW were computed using the `scib` [40] package.

## 699 P. Details on running RNA velocity estimation methods

700 We use the datasets `V` to benchmark RNA velocity inference methods as shown in Table 2. We used `scVelo`  
 701 0.2.4 and `VeloCyto` 0.17.17. We benchmarked `scVelo` with three modes: `deterministic`, `stochastic`, and  
 702 `dynamical`. For `VeloCyto`, we used the default options.

## 703 Q. Details on running GRN inference methods

704 We use the BEELINE [48] framework to benchmark GRN inference methods. Apart from the main datasets, The  
 705 dataset `G` (Table 2) was generated using the following configurations: The 100-gene GRN in Fig. 3, 1000 cells,  
 706 50 CIFs,  $r_d = 0.2$ ,  $\sigma_i = 1$ , with other default parameters. Eight datasets were generated for random seed 1 to 8,  
 707 and technical noise and batch effect was added using default parameters. We ran the BEELINE GitHub version  
 708 `79775f0`. In order to resolve runtime errors, all docker images were built locally, except that we used the provided  
 709 images on Docker Hub for `PIDC` and `Scribe`. We use BEELINE’s example workflow to infer GRN and calculate the  
 710 AUPRC ratio and AUROC for (a) true counts in the eight datasets, and (b) observed counts with batch effects in  
 711 the eight datasets. The AUPRC ratio is the AUPRC divided by the AUPRC of a random predictor, which equals to  
 712 the network density of the ground truth network. Eleven methods were benchmarked in total: `PIDC`, `GRNBoost2`,  
 713 `GENIE3`, `Sincerities`, `PPCOR`, `LEAP`, `GRISLI`, `SINGE`, `GRNVBEM`, `Scribe` and `SCODE`.

## 714 R. Details on running CCI inference methods

715 We generated 12 datasets using the following procedure. Apart from the main datasets, for each `C` dataset  
 716 (Table 2), we first construct the GRN (Fig. S6a): (1) let genes 1-6 be the transcription factors. Sample 70  
 717 edges from gene 1-6 to gene 7-53. (2) Connect gene 7-53 (regulator) to gene 54-100 (target) consecutively. (3)  
 718 Connect gene 54-100 to gene 110-156. In this way, we can generate a GRN with reasonable edge density and  
 719 make sure that there are three downstream genes for each TF, which is required by `SpaTalk`. Then we construct

720 the ligand-receptor pairs: let the ligands be gene 101-106 and receptors be gene 2, 6, 10, 8, 20, and 30. We  
721 divide a linear trajectory into 5 sections, corresponding to 5 cell types. Between each cell type pair (excluding  
722 same-type pairs), we sample 3-6 ligand-receptor pairs and enable cell-cell interactions with them for the two cell  
723 types. The dataset is then simulated using 160 genes in total, 500 cells, and 50 CIFs. We use the true counts to  
724 benchmark the methods.

To run SpaTalk, we modify the original `plot_lrpair_vln` method to return the p-value from the Wilcoxon rank sum test directly, rather than drawing a figure. Before using the p-values to calculate the precision and recall, we adjusted them using Bonferroni correction:

$$\hat{p}_i = \max(p_i \cdot |p|, 1)$$

725 where  $p$  is the p-value vector for all cell types and ligand-receptor pairs. For Giotto, we used the R package 1.1.2  
726 and followed the `mini_seqfish` vignette. For SpaOTsc, we used default parameters.

## 727 **Data and Code Availability**

728 The scMultiSim R package is available at <https://github.com/ZhangLabGT/scMultiSim>. The code  
729 for dataset generation and benchmarking is available at [https://github.com/ZhangLabGT/scMultiSim\\_](https://github.com/ZhangLabGT/scMultiSim_manuscript)  
730 [manuscript](https://github.com/ZhangLabGT/scMultiSim_manuscript).

## 731 **Acknowledgements**

732 This work was supported by the US National Science Foundation DBI-2019771 and National Institutes  
733 of Health grant R35GM143070 (HL, ZZ, XZ), the Guangdong Basic and Applied Basic Research  
734 Foundation (2022B1515120077 to XC) and the Shenzhen Innovation Committee of Science and Technology  
735 (20220815094330001 to XC).

## 736 **Competing Interests Statement**

737 The authors declared no competing interest.

## Bibliography

- 738 1. R. Argelaguet, A. S. E. Cuomo, O. Stegle, and J. C. Marioni. Computational principles and challenges in single-cell data  
739 integration. *Nat. Biotechnol.*, pages 1–14, May 2021.
- 740 2. G. Baruzzo, I. Patuzzi, and B. Di Camillo. SPARSim single cell: a count data simulator for scRNA-seq data.  
741 *Bioinformatics*, 36(5):1468–1475, Mar. 2020.
- 742 3. V. Bergen, M. Lange, S. Peidli, F. A. Wolf, and F. J. Theis. Generalizing RNA velocity to transient cell states through  
743 dynamical modeling. *Nat. Biotechnol.*, Aug. 2020.
- 744 4. R. Browaeys, W. Saelens, and Y. Saeys. NicheNet: modeling intercellular communication by linking ligands to target  
745 genes. *Nat. Methods*, 17(2):159–162, Feb. 2020.
- 746 5. Z. Cang and Q. Nie. Inferring spatial and signaling relationships between cells from single cell transcriptomic data. *Nat.*  
747 *Commun.*, 11(1):2084, Apr. 2020.
- 748 6. Z. Cang, Y. Zhao, A. A. Almet, A. Stabell, R. Ramos, M. V. Plikus, S. X. Atwood, and Q. Nie. Screening cell-cell  
749 communication in spatial transcriptomics via collective optimal transport. *Nat. Methods*, Jan. 2023.
- 750 7. R. Cannoodt, W. Saelens, L. Deconinck, and Y. Saeys. Spearheading future omics analyses using dyngen, a multi-modal  
751 simulator of single cells. *Nature Communications*, 12(1):1–9, 2021.
- 752 8. J. Cao, D. A. Cusanovich, V. Ramani, D. Aghamirzaie, H. A. Pliner, A. J. Hill, R. M. Daza, J. L. McFaline-Figueroa, J. S.  
753 Packer, L. Christiansen, F. J. Steemers, A. C. Adey, C. Trapnell, and J. Shendure. Joint profiling of chromatin accessibility  
754 and gene expression in thousands of single cells. *Science*, 361(6409):1380–1385, Sept. 2018.
- 755 9. T. E. Chan, M. P. H. Stumpf, and A. C. Babbie. Gene regulatory network inference from Single-Cell data using multivariate  
756 information measures. *Cell Syst*, 5(3):251–267.e3, Sept. 2017.
- 757 10. S. Chen, B. B. Lake, and K. Zhang. High-throughput sequencing of the transcriptome and chromatin accessibility in the  
758 same cell. *Nat. Biotechnol.*, 37(12):1452–1457, Dec. 2019.
- 759 11. X. Chen, R. J. Miragaia, K. N. Natarajan, and S. A. Teichmann. A rapid and robust method for single cell chromatin  
760 accessibility profiling. *Nature Communications*, 9(1):5345, 2018.
- 761 12. H. L. Crowell, S. X. M. Leonardo, C. Soneson, and M. D. Robinson. Built on sand: the shaky foundations of simulating  
762 single-cell RNA sequencing data. 2021.
- 763 13. H. L. Crowell, C. Soneson, P.-L. Germain, D. Calini, L. Collin, C. Raposo, D. Malhotra, and M. D. Robinson. muscat  
764 detects subpopulation-specific state transitions from multi-sample multi-condition single-cell transcriptomics data. *Nat.*  
765 *Commun.*, 11(1):6077, Nov. 2020.
- 766 14. A. Deshpande, L.-F. Chu, R. Stewart, and A. Gitter. Network inference with granger causality ensembles on single-cell  
767 transcriptomics. *Cell Rep.*, 38(6):110333, Feb. 2022.
- 768 15. P. Dibaeinia and S. Sinha. SERGIO: A Single-Cell expression simulator guided by gene regulatory networks. *Cell Syst*,  
769 Aug. 2020.
- 770 16. R. Dries, Q. Zhu, R. Dong, C. H. L. Eng, H. Li, K. Liu, Y. Fu, T. Zhao, A. Sarkar, F. Bao, R. E. George, N. Pierson, L. Cai,  
771 and G. C. Yuan. Giotto: a toolbox for integrative analysis and visualization of spatial expression data. *Genome Biology*,  
772 22, 2021.
- 773 17. A. Duò, M. Robinson, and C. Soneson. A systematic performance evaluation of clustering methods for single-cell rna-seq  
774 data. *F1000Research*, 7:1141, 11 2020.
- 775 18. M. Efremova and S. A. Teichmann. Computational methods for single-cell omics across modalities. *Nat. Methods*,  
776 17(1):14–17, Jan. 2020.
- 777 19. C.-H. L. Eng, M. Lawson, Q. Zhu, R. Dries, N. Kouloua, Y. Takei, J. Yun, C. Cronin, C. Karp, G.-C. Yuan, and L. Cai.  
778 Transcriptome-scale super-resolved imaging in tissues by RNA seqFISH. *Nature*, 568(7751):235–239, Apr. 2019.
- 779

- 780 20. J. T. Gaublot, N. Yosef, Y. Lee, R. S. Gertner, L. V. Yang, C. Wu, P. P. Pandolfi, T. Mak, R. Satija, A. K. Shalek, V. K.  
781 Kuchroo, H. Park, and A. Regev. Single-Cell Genomics Unveils Critical Regulators of Th17 Cell Pathogenicity. *Cell*,  
782 163(6):1400–1412, 2015.
- 783 21. B. Gong, Y. Zhou, and E. Purdom. Cobolt: integrative analysis of multimodal single-cell sequencing data. *Genome*  
784 *Biology*, 22(1):351, 2021.
- 785 22. G. Gorin, M. Fang, T. Chari, and L. Pachter. Rna velocity unraveled. *PLOS Computational Biology*, 18(9):e1010492,  
786 2022.
- 787 23. Y. Hao, S. Hao, E. Andersen-Nissen, W. M. M. III, S. Zheng, A. Butler, M. J. Lee, A. J. Wilk, C. Darby, M. Zagar,  
788 P. Hoffman, M. Stoeckius, E. Papalexi, E. P. Mimitou, J. Jain, A. Srivastava, T. Stuart, L. B. Fleming, B. Yeung, A. J.  
789 Rogers, J. M. McElrath, C. A. Blish, R. Gottardo, P. Smibert, and R. Satija. Integrated analysis of multimodal single-cell  
790 data. *Cell*, 2021.
- 791 24. Y. Hao, T. Stuart, M. Kowalski, S. Choudhary, P. Hoffman, A. Hartman, A. Srivastava, G. Molla, S. Madad,  
792 C. Fernandez-Granda, and R. Satija. Dictionary learning for integrative, multimodal, and scalable single-cell analysis.  
793 *bioRxiv*, 2022.
- 794 25. C. A. Herring, A. Banerjee, E. T. McKinley, A. J. Simmons, J. Ping, J. T. Roland, J. L. Franklin, Q. Liu, M. J. Gerdes, R. J.  
795 Coffey, and K. S. Lau. Unsupervised trajectory analysis of Single-Cell RNA-Seq and imaging data reveals alternative tuft  
796 cell origins in the gut. *Cell Syst*, 6(1):37–51.e9, Jan. 2018.
- 797 26. Y. Hu, T. Peng, L. Gao, and K. Tan. CytoTalk: De novo construction of signal transduction networks using single-cell  
798 transcriptomic data. *Sci Adv*, 7(16), Apr. 2021.
- 799 27. V. A. Huynh-Thu, A. Irrthum, L. Wehenkel, and P. Geurts. Inferring Regulatory Networks from Expression Data Using  
800 Tree-Based Methods. *PLOS ONE*, 5(9):1–10, 2010.
- 801 28. Z. Ji and H. Ji. TSCAN: Pseudo-time reconstruction and evaluation in single-cell RNA-seq analysis. *Nucleic Acids Res.*,  
802 44(13):e117, July 2016.
- 803 29. S. Jin, C. F. Guerrero-Juarez, L. Zhang, I. Chang, R. Ramos, C.-H. Kuan, P. Myung, M. V. Plikus, and Q. Nie. Inference  
804 and analysis of cell-cell communication using CellChat. *Nat. Commun.*, 12(1):1088, Feb. 2021.
- 805 30. K. Kamimoto, B. Stringa, C. M. Hoffmann, K. Jindal, L. Solnica-Krezel, and S. A. Morris. Dissecting cell identity via  
806 network inference and in silico gene perturbation. *Nature*, 614(7949):742–751, Feb. 2023.
- 807 31. J. Kim and J. C. Marionni. Inferring the kinetics of stochastic gene expression from single-cell RNA-sequencing data.  
808 14(1):R7, 2013.
- 809 32. S. Kim. ppcor: An R package for a fast calculation to semi-partial correlation coefficients. *Commun Stat Appl Methods*,  
810 22(6):665–674, Nov. 2015.
- 811 33. V. Y. Kiselev, K. Kirschner, M. T. Schaub, T. Andrews, A. Yiu, T. Chandra, K. N. Natarajan, W. Reik, M. Barahona, A. R.  
812 Green, and M. Hemberg. SC3: consensus clustering of single-cell RNA-seq data. *Nat. Methods*, 14(5):483–486, May  
813 2017.
- 814 34. A. R. Kriebel and J. D. Welch. UINMF performs mosaic integration of single-cell multi-omic datasets using nonnegative  
815 matrix factorization. *Nature Communications*, 13(1):780, 2022.
- 816 35. G. La Manno, R. Soldatov, A. Zeisel, E. Braun, H. Hochgerner, V. Petukhov, K. Lidschreiber, M. E. Kastrioti, P. Lönnberg,  
817 A. Furlan, J. Fan, L. E. Borm, Z. Liu, D. van Bruggen, J. Guo, X. He, R. Barker, E. Sundström, G. Castelo-Branco,  
818 P. Cramer, I. Adameyko, S. Linnarsson, and P. V. Kharchenko. RNA velocity of single cells. *Nature*, 560(7719):494–498,  
819 Aug. 2018.
- 820 36. B. Li, W. Zhang, C. Guo, H. Xu, L. Li, M. Fang, Y. Hu, X. Zhang, X. Yao, M. Tang, K. Liu, X. Zhao, J. Lin, L. Cheng,  
821 F. Chen, T. Xue, and K. Qu. Benchmarking spatial and single-cell transcriptomics integration methods for transcript  
822 distribution prediction and cell type deconvolution. *Nat. Methods*, May 2022.

- 823 37. C. Li, X. Chen, S. Chen, R. Jiang, and X. Zhang. simCAS: an embedding-based method for simulating single-cell  
824 chromatin accessibility sequencing data. Feb. 2023.
- 825 38. C. Li, M. Virgilio, K. L. Collins, and J. D. Welch. Single-cell multi-omic velocity infers dynamic and decoupled gene  
826 regulation. Dec. 2021.
- 827 39. P. Lin, M. Troup, and J. W. K. Ho. CIDR: Ultrafast and accurate clustering through imputation for single-cell RNA-Seq  
828 data. *bioRxiv*, page 068775, Aug. 2016.
- 829 40. M. D. Luecken, M. Büttner, K. Chaichoompu, A. Danese, M. Interlandi, M. F. Mueller, D. C. Strobl, L. Zappia, M. Dugas,  
830 M. Colomé-Tatché, and F. J. Theis. Benchmarking atlas-level data integration in single-cell genomics. *Nature Methods*,  
831 19(1):41–50, 2022.
- 832 41. S. Ma, B. Zhang, L. M. LaFave, A. S. Earl, Z. Chiang, Y. Hu, J. Ding, A. Brack, V. K. Kartha, T. Tay, T. Law, C. Lareau,  
833 Y.-C. Hsu, A. Regev, and J. D. Buenrostro. Chromatin potential identified by shared Single-Cell profiling of RNA and  
834 chromatin. *Cell*, 183(4):1103–1116.e20, Nov. 2020.
- 835 42. T. Moerman, S. Aibar Santos, C. Bravo González-Blas, J. Simm, Y. Moreau, J. Aerts, and S. Aerts. GRNBoost2 and  
836 arboreto: efficient and scalable inference of gene regulatory networks. *Bioinformatics*, 35(12):2159–2161, June 2019.
- 837 43. B. Munsky, G. Neuert, and A. van Oudenaarden. Using gene expression noise to understand gene regulation. *Science*,  
838 336(6078):183–187, Apr. 2012.
- 839 44. Z. Navidi, L. Zhang, and B. Wang. simATAC: a single-cell ATAC-seq simulation framework. *Genome Biol.*, 22(1):74, Mar.  
840 2021.
- 841 45. N. Papili Gao, S. M. M. Ud-Dean, O. Gandrillon, and R. Gunawan. SINCERITIES: inferring gene regulatory networks  
842 from time-stamped single cell transcriptional expression profiles. *Bioinformatics*, 34(2):258–266, Jan. 2018.
- 843 46. J. Peccoud and B. Ycart. Markovian modeling of gene-product synthesis. 48(2):222–234, Oct. 1995.
- 844 47. V. M. Peterson, K. X. Zhang, N. Kumar, J. Wong, L. Li, D. C. Wilson, R. Moore, T. K. McClanahan, S. Sadekova, and J. A.  
845 Klappenbach. Multiplexed quantification of proteins and transcripts in single cells. *Nat. Biotechnol.*, 35(10):936–939,  
846 Oct. 2017.
- 847 48. A. Pratapa, A. P. Jalihal, J. N. Law, A. Bharadwaj, and T. M. Murali. Benchmarking algorithms for gene regulatory network  
848 inference from single-cell transcriptomic data. *Nat. Methods*, Jan. 2020.
- 849 49. X. Qiu, Q. Mao, Y. Tang, L. Wang, R. Chawla, H. A. Pliner, and C. Trapnell. Reversed graph embedding resolves complex  
850 single-cell trajectories. *Nat. Methods*, 14(10):979–982, Oct. 2017.
- 851 50. S. G. Rodrigues, R. R. Stickels, A. Goeva, C. A. Martin, E. Murray, C. R. Vanderburg, J. Welch, L. M. Chen, F. Chen,  
852 and E. Z. Macosko. Slide-seq: A scalable technology for measuring genome-wide expression at high spatial resolution.  
853 *Science*, 363(6434), 2019.
- 854 51. W. Saelens, R. Cannoodt, H. Todorov, and Y. Saeys. A comparison of single-cell trajectory inference methods. *Nat.*  
855 *Biotechnol.*, Apr. 2019.
- 856 52. S. Shah, E. Lubeck, W. Zhou, and L. Cai. In situ transcription profiling of single cells reveals spatial organization of cells  
857 in the mouse hippocampus. *Neuron*, 92(2):342–357, Oct. 2016.
- 858 53. X. Shao, C. Li, H. Yang, X. Lu, J. Liao, J. Qian, K. Wang, J. Cheng, P. Yang, H. Chen, X. Xu, and X. Fan. Knowledge-graph-based cell-cell communication inference for spatially resolved transcriptomic data with SpaTalk. *Nature*  
859 *Communications*, 13(1):4429, 2022.
- 860 54. D. Song, Q. Wang, G. Yan, T. Liu, and J. J. Li. A unified framework of realistic in silico data generation and statistical  
861 model inference for single-cell and spatial omics. Sept. 2022.
- 862 55. P. L. Ståhl, F. Salmén, S. Vickovic, A. Lundmark, J. F. Navarro, J. Magnusson, S. Giacomello, M. Asp, J. O. Westholm,  
863 M. Huss, A. Mollbrink, S. Linnarsson, S. Codeluppi, Å. Borg, F. Pontén, P. I. Costea, P. Sahlén, J. Mulder, O. Bergmann,  
864 J. Lundeberg, and J. Frisén. Visualization and analysis of gene expression in tissue sections by spatial transcriptomics,  
865

- 866 2016.
- 867 56. M. Stoeckius, C. Hafemeister, W. Stephenson, B. Houck-Loomis, P. K. Chattopadhyay, H. Swerdlow, R. Satija, and  
868 P. Smibert. Simultaneous epitope and transcriptome measurement in single cells. *Nat. Methods*, 14(9):865–868, Sept.  
869 2017.
- 870 57. K. Street, D. Risso, R. B. Fletcher, D. Das, J. Ngai, N. Yosef, E. Purdom, and S. Dudoit. Slingshot: cell lineage and  
871 pseudotime inference for single-cell transcriptomics. *BMC Genomics*, 19(1):477, June 2018.
- 872 58. T. Stuart, A. Butler, P. Hoffman, C. Hafemeister, E. Papalexi, W. M. Mauck, 3rd, Y. Hao, M. Stoeckius, P. Smibert, and  
873 R. Satija. Comprehensive integration of Single-Cell data. *Cell*, 177(7):1888–1902.e21, June 2019.
- 874 59. T. Sun, D. Song, W. V. Li, and J. J. Li. scDesign2: a transparent simulator that generates high-fidelity single-cell gene  
875 expression count data with gene correlations captured. *Genome Biology*, 22(1):163, 2021.
- 876 60. J. Tanevski, R. O. Ramirez Flores, A. Gabor, D. Schapiro, and J. Saez-Rodriguez. Explainable multiview framework for  
877 dissecting spatial relationships from highly multiplexed data. *Genome Biology*, 23(97), 2022.
- 878 61. K. Vandereyken, A. Sifrim, B. Thienpont, and T. Voet. Methods and applications for single-cell and spatial multi-omics.  
879 *Nat. Rev. Genet.*, pages 1–22, Mar. 2023.
- 880 62. L. Wang, N. Trasanidis, T. Wu, G. Dong, M. Hu, D. E. Bauer, and L. Pinello. Dictys: dynamic gene regulatory network  
881 dissects developmental continuum with single-cell multi-omics. Sept. 2022.
- 882 63. X. Wang, W. E. Allen, M. A. Wright, E. L. Sylwestrak, N. Samusik, S. Vesuna, K. Evans, C. Liu, C. Ramakrishnan, J. Liu,  
883 G. P. Nolan, F.-A. Bava, and K. Deisseroth. Three-dimensional intact-tissue sequencing of single-cell transcriptional  
884 states. *Science*, 361(6400), July 2018.
- 885 64. J. D. Welch, V. Kozareva, A. Ferreira, C. Vanderburg, C. Martin, and E. Z. Macosko. Single-Cell multi-omic integration  
886 compares and contrasts features of brain cell identity. *Cell*, 177(7):1873–1887.e17, June 2019.
- 887 65. F. A. Wolf, F. K. Hamey, M. Plass, J. Solana, J. S. Dahlin, B. Göttgens, N. Rajewsky, L. Simon, and F. J. Theis.  
888 PAGA: graph abstraction reconciles clustering with trajectory inference through a topology preserving map of single  
889 cells. *Genome Biol.*, 20(1):59, Mar. 2019.
- 890 66. L. Zappia and F. J. Theis. Over 1000 tools reveal trends in the single-cell RNA-seq analysis landscape. *Genome Biol.*,  
891 22(1):301, Oct. 2021.
- 892 67. A. Zeisel, A. B. Muñoz-Manchado, S. Codeluppi, P. Lönnerberg, G. L. Manno, A. Juréus, S. Marques, H. Munguba,  
893 L. He, C. Betsholtz, C. Rolny, G. Castelo-Branco, J. Hjerling-Leffler, and S. Linnarsson. Cell types in the mouse cortex  
894 and hippocampus revealed by single-cell RNA-seq. *Science*, 347(6226):1138–1142, 2015.
- 895 68. S. Zhang, S. Pyne, S. Pietrzak, A. F. Siahpirani, R. Sridharan, and S. Roy. Inference of cell type-specific gene regulatory  
896 networks on cell lineages from single cell omic datasets. July 2022.
- 897 69. X. Zhang, C. Xu, and N. Yosef. Simulating multiple faceted variability in single cell RNA sequencing. *Nat. Commun.*,  
898 10(1):2611, June 2019.
- 899 70. Z. Zhang, C. Yang, and X. Zhang. scDART: integrating unmatched scRNA-seq and scATAC-seq data and learning  
900 cross-modality relationship simultaneously. *Genome Biology*, 23(1):139, 2022.
- 901 71. Z. Zhang and X. Zhang. VeloSim: Simulating single cell gene-expression and RNA velocity. *BioRxiv*, 2021.



## Tables

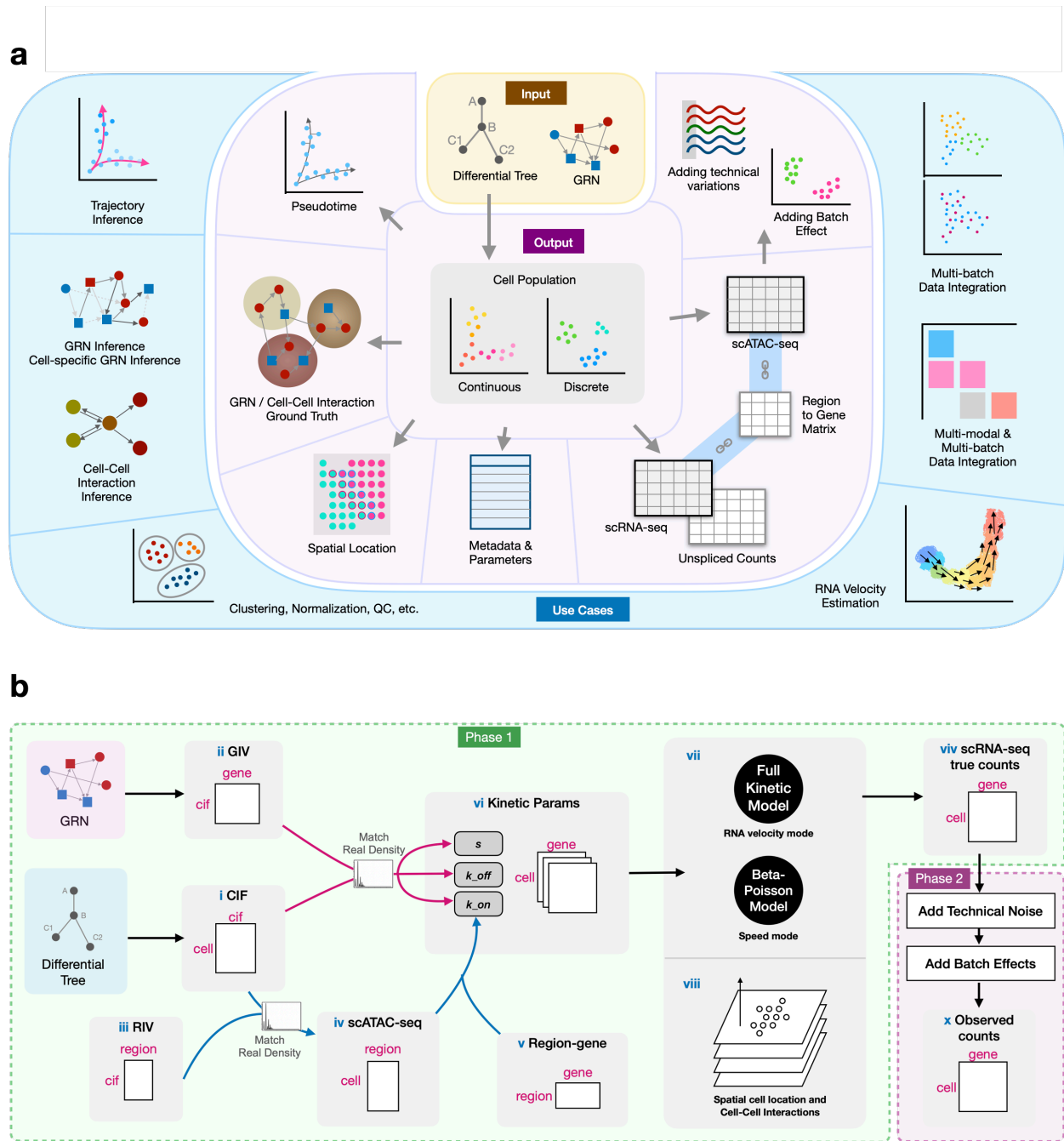
Name (#datasets)	Label	Trajectory	$\sigma_{cif}$	Cells	Genes	Velo	GRN	CCI	Seed
Main (144)	<b>M</b>	<b>L: Continuous Linear (Phyla1)</b>  <b>T: Continuous Tree (Phyla3)</b>  <b>D: Discrete (Phyla5)</b>	0.1	500	<b>1: 110 genes</b> <b>2: 200 genes</b> <b>3: 500 genes</b>	F	GRN_100	T	<b>a: 1</b> <b>b: 2</b> <b>c: 3</b> <b>d: 4</b>
				800	<b>4: 110 genes</b> <b>5: 200 genes</b> <b>6: 500 genes</b>				
			0.5	500	<b>7: 110 genes</b> <b>8: 200 genes</b> <b>9: 500 genes</b>				
				800	<b>10: 110 genes</b> <b>11: 200 genes</b> <b>12: 500 genes</b>				

**Table 1.** The main dataset contains 144 datasets with varying trajectory,  $\sigma_{cif}$ , number of cells and genes. For each parameter configuration, four datasets are generated using different random seeds. We number the datasets for easy referencing in the text: starting with the letter M, then a letter {L,T,D} specifying the trajectory; followed by a number 1-12 identifying the configuration of  $\sigma_{cif}$ , number of cells and genes; and last, a lowercase letter a-d indicating the random seed. For example, MD5c uses a discrete cell population,  $\sigma_{cif} = 0.1$ , 800 cells, 200 genes and random seed 3. Phyla1, Phyla3 and Phyla5 are the input tree structure used to generate the cell populations, and they are shown in Fig. 2b.

Name (#datasets)	Label	Trajectory	$\sigma_{cif}$	Cells	Genes	Velo	GRN	CCI	Seed	Other Params
Auxiliary (144)	<b>A</b>	Tree (Phyla5)	0.1, 0.5, 1	200	1000	F	GRN_100	F	1,2	$E_{atac} = 0.2, 0.5, 0.9$ $\sigma_i = 0.3, 1$ $r_d = 0.2, 0.8$
		Discrete (5 clusters)								
Velocity (72)	<b>V</b>	Tree (Phyla5)	0.1	500,750,1000	100,200,500	T	GRN_100 N/A	T	5-8	
Realistic (1)	<b>R</b>	Discrete (10 clusters)	0.1	523	200	F	Inferred	F	1	
Add. GRN (8)	<b>G</b>	Linear	0.1	110	1000	F	GRN_100	F	1-8	
Add. CCI (8)	<b>C</b>	Linear	0.1	200	500	F	Fig. S6	T	1-8	

**Table 2.** The auxiliary dataset and other datasets used in supplemental information.

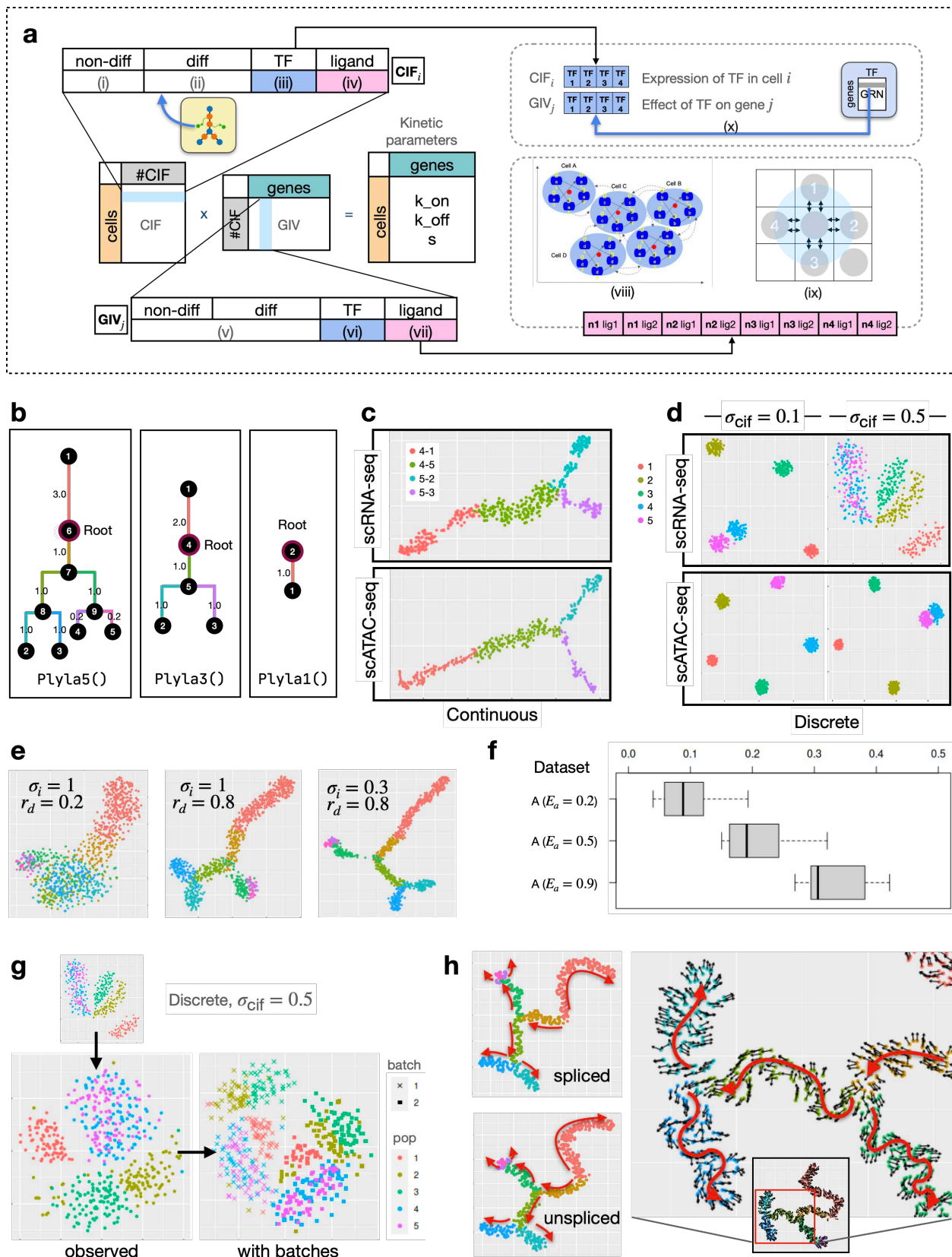
# Figures



**Figure 1. Overview of scMultiSim.** See the next page for descriptions.

. **Overview of scMultiSim.** **(a):** The input, output, and use cases. The minimal required input is a cell differential tree describing the differentiation relationship of cell types. It controls the cell trajectory or clusters in the output. A user-input ground truth GRN is recommended to guide the simulation. Users can also provide ground truth for cell-cell interaction and control each simulated biological effects using various parameters. **(b):** The overall structure of scMultiSim. The scATAC-seq data (iv) is firstly generated using CIF (i) and RIV (iii). The kinetic parameters used to generate scRNA-seq data (vi) is prepared using GIV (ii), CIF (i) and the scATAC-seq data with **(v)** a region-to-gene matrix. Using the parameters, either the full kinetic model (when RNA velocity is required), or the Beta-Poisson model (when running speed matters) will be used to generate the scRNA-seq data (vii). scMultiSim uses a multiple-step approach that considers both time and space when CCI is enabled (viii). With the simulated true counts (viv), technical noise and batch effects can be added to obtain the observed counts (x).

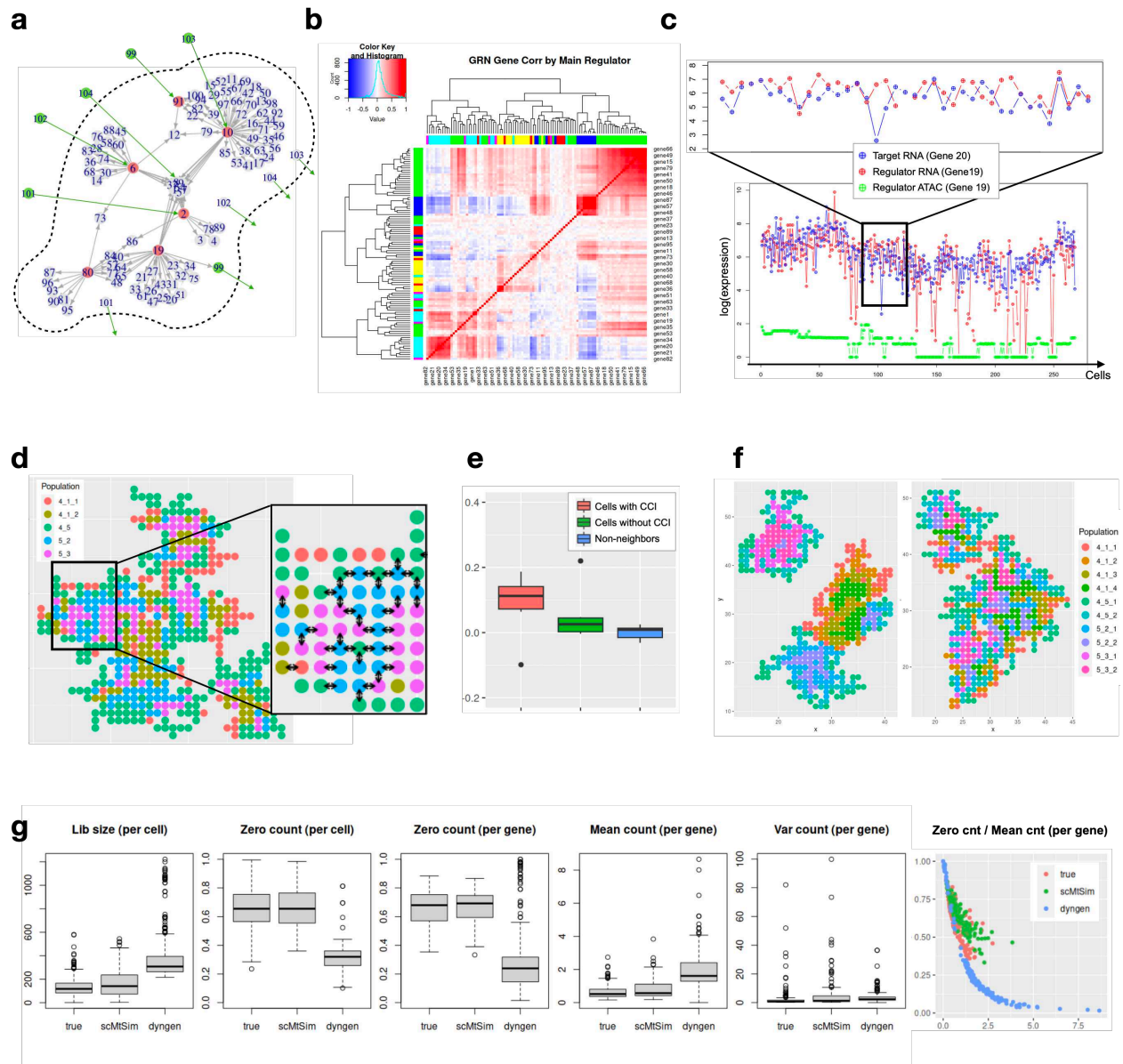
904



**Figure 2. scMultiSim generates multi-modal single cell data from pre-defined cell clustering structure or trajectories.** See the next page for descriptions.

. **scMultiSim generates multi-modal single cell data from pre-defined cell clustering structure or trajectories.** (a) The CIF and GIV matrix. We multiply the CIF and GIV matrix to get the cell $\times$ gene matrix for each kinetic parameter. CIFs and GIVs are divided into segments to encode different biological effects, where each segment encodes a certain type of biological factor. A cellular heterogeneity is modeled in the CIF, and regulation effects are encoded in the corresponding GIV vector. (viii) is the illustration of the cell-cell interactions and in-cell GRN in our model. (ix) is the grid system representing spatial locations of cells. A cell can have at most four neighbors (labeled 1-4) within a certain range (blue circle). The cell at the bottom right corner is not a neighbor of the center cell. (b) Three trees are provided by scMultiSim and used to produce the datasets. Phyla1 is a linear trajectory, while Phyla3 and Phyla5 has 3 and 5 leaves, respectively. (c) t-SNE visualization of the paired scRNA-seq and scATAC-seq data (without adding technical noise) from the main dataset MT3a (continuous populations following tree Phyla3), both having  $n_{\text{cell}} = n_{\text{gene}} = 500$ . (d) t-SNE visualization of the paired scRNA-seq and scATAC-seq data (without adding technical noise) from the main datasets MD3a and MD9a (discrete populations with five clusters, following tree Phyla5). (e) Additional results showing the effect of  $\sigma_i$  and  $r_d$  using datasets A. (f) Additional results exploring the ATAC effect parameter  $E_a$  using datasets A. Averaged Spearman correlation between scATAC-seq and scRNA-seq data for genes affected by one chromatin region, from 144 datasets using various parameters ( $\sigma_i$ ,  $\sigma_{\text{cif}}$ ,  $r_d$ , continuous/discrete). (g) The observed RNA counts in dataset MD9a with added technical noise and batch effects. (h) The spliced true counts, unspliced true counts, and the RNA velocity ground truth from dataset V. The velocity vectors point to the directions of differentiation indicated by red arrows, from the tree root to leaves.

905

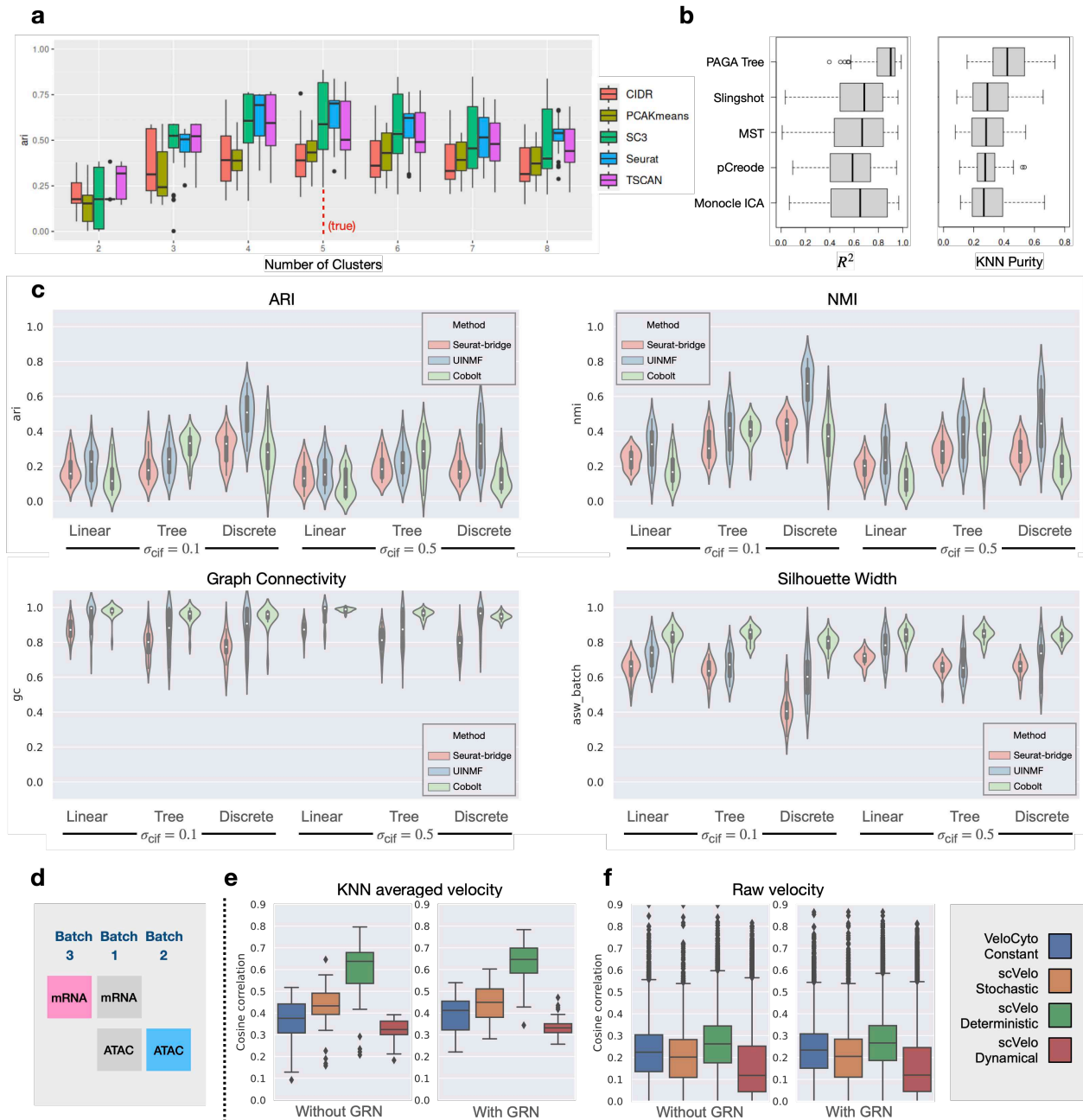


**Figure 3.** scMultiSim generates realistic single cell gene expression data driven by GRN and cell-cell interaction. See the next page for descriptions.

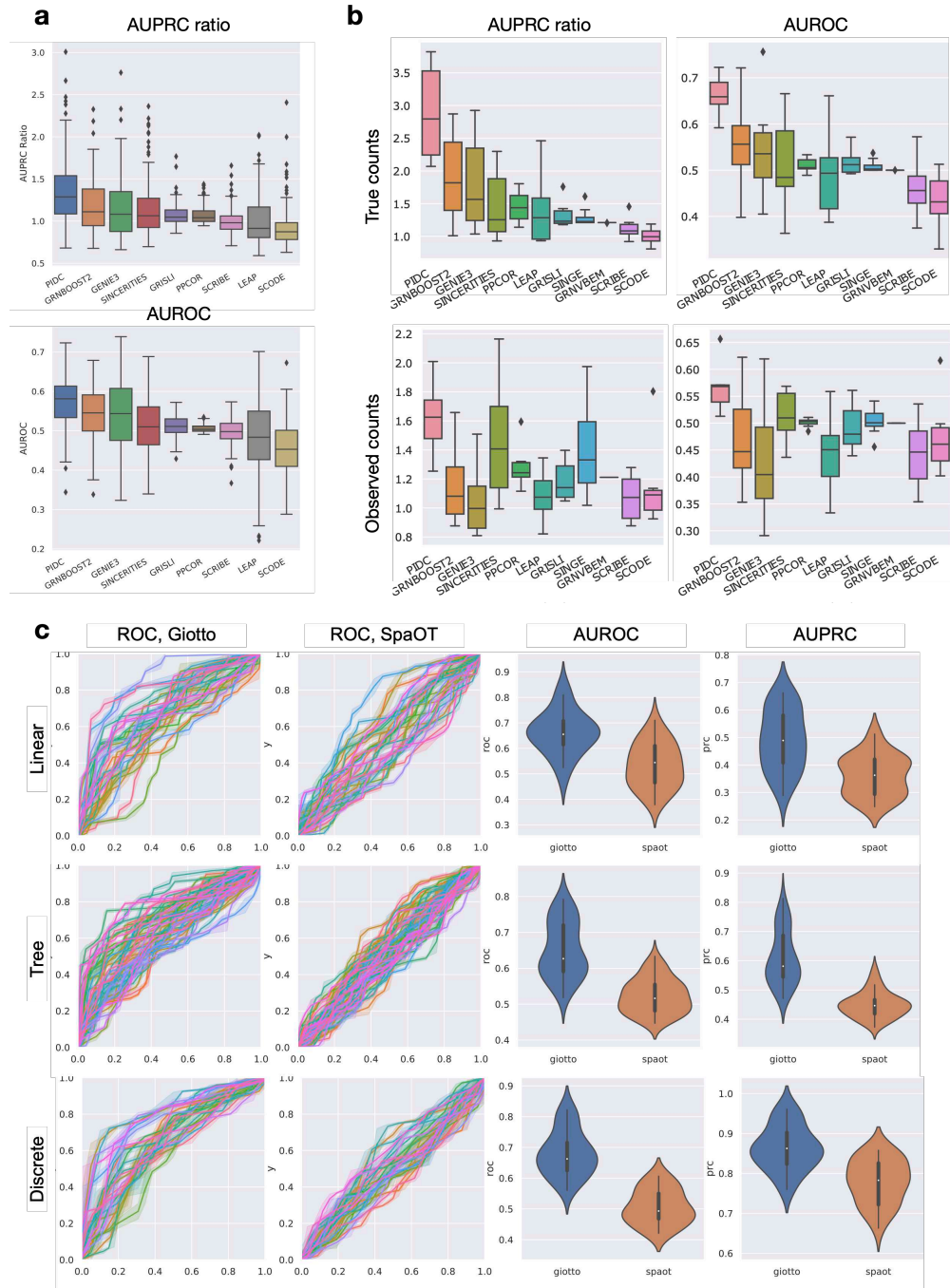
. **scMultiSim generates realistic single cell gene expression data driven by GRN and cell-cell interaction.** (a) The GRN and CCIs used to generate the main datasets. Red nodes are TF genes and green nodes are ligand genes. Green edges are the added ligand-receptor pairs when simulating cell-cell interactions. (b-e) Results from dataset MT3a, which uses Phyla3, 500 genes, 500 cells and  $\sigma_{\text{cif}} = 0.1$ . (b) The gene module correlation heatmap. The color at left or top represents the regulating TF of the gene. Genes regulated by the same TF have higher correlations and tend to be grouped together. (c) The log-transformed expression of a specific TF-target gene pair (gene19-gene20) for all cells on one lineage (4-5-3 in Phyla3). Correlation between the TF and target expressions can be observed. We also show the chromatin accessibility level for the TF gene 19, averaged from the two corresponding chromatin regions of the gene. Significant lower expression of gene 19 can be observed when the chromatin is closed. (d) The spatial location of cells, where each color represents a cell type. Arrows between two cells indicates that CCI exists between them for a specific ligand-receptor pair (gene101-gene2). By default, most cell-cell interactions occur between different cell types. (e) Gene expression correlation between (1) neighboring cells with CCI, (2) neighboring cells with CCI, and (3) non-neighbor cells. Cells with CCI have higher correlations. (f) scMultiSim provides options to control the the cell layout. We show the results of 1200 cells using same-type probability  $p_n = 1.0$  and 0.8, respectively. When  $p_n = 1.0$ , same-type cells tend to cluster together, while  $p_n = 0.8$  introduces more randomness. (g) Comparison between a real dataset and simulated data using multiple statistical measurements. Parameters were adjusted to match the real distribution as close as possible.

906





**Figure 4. Benchmarking clustering, trajectory inference, multi-modal data integration and RNA velocity estimation methods.** (a) Benchmarking clustering methods on dataset MD (discrete). Methods are grouped by number of clusters in the result. The vertical red dashed line shows the true number of clusters. A higher ARI indicates better clustering. (b) Benchmarking trajectory inference methods on dataset MT (continuous tree). Methods are evaluated based on their mean  $R^2$  and  $k$ NN purity on each lineage (higher is better). (c) Benchmarking multi-modal data integration methods. Metrics for the methods: ARI, NMI (higher = better at preserving cell identities), graph connectivity and average silhouette width of batch (higher = better merging batches). (d) The task illustration of multi-modal data integration. Only cells in batch 1 and 3 (pink and blue matrices) are used for evaluation. (e,f) Benchmarking RNA velocity estimation methods on auxiliary dataset V. The result is measured using cosine similarity.



**Figure 5. Benchmarking GRN inference and CCI inference methods.** (a) Benchmarking GRN inference methods. The upper figure shows AUPRC ratios (versus a random classifier), and the lower figure shows AUROC values. (b) Additional results on benchmarking GRN inference methods using datasets G that does not contain CCI effects. We also tested the performance on observed counts with technical noise. (c) Benchmarking cell-cell interaction inference methods. Each curve in the ROC/PRC plots correspondings to one dataset.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [scMultiSimsupp.pdf](#)