



scDesign3: in silico data generation for multimodal single-cell and spatial omics

Dongyuan Song, **Jingyi Jessica Li**

July 25, 2023

Department of Statistics
University of California, Los Angeles

<http://jsb.ucla.edu>

Introduction

Single-cell and spatial omics data: statistical characteristics

Processed data: a cell-by-feature matrix + cell covariates

Cell heterogeneity structures

- discrete cell types (known or latent)
- continuous trajectories (usually latent)
- spatial locations (known for spatial data)

Experimental designs

- batches (unwanted effects)
- conditions (biological signals)

Features

- gene expression (scRNA-seq, spatial transcriptomics, etc.)
- chromatin accessibility (scATAC-seq, SNARE-seq, etc.)
- protein abundance (CITE-seq, etc.)



Motivations of scDesign3

Computational benchmarking

- > 1000 computational tools at www.scrna-tools.org
- how to choose among competing computational tools?

Inference

Conditional on a cell covariate (type, pseudotime, or spatial location)

- every gene's distribution
- every gene pair's correlation

In silico controlled experiments

- negative control: to evaluate a pipeline's **false discoveries**
- positive control: to evaluate a pipeline's **discovery power**

A realistic simulator with interpretable parameters



Challenges in modeling single-cell and spatial multi-omics

- **High-dimensional:** 10^4 genes; for other features the number can be even larger
- **Correlation:** complex correlation structures between features
- **Diverse covariates:** cell types, continuous trajectories, spatial space
- **Multi-omics:** different omics may follow different distributions
- **Transparency:** not a black box



Brief Communication | [Published: 11 May 2023](#)

scDesign3 generates realistic in silico data for multimodal single-cell and spatial omics

[Dongyuan Song](#), [Qingyang Wang](#), [Guanao Yan](#), [Tianyang Liu](#), [Tianyi Sun](#) & [Jingyi Jessica Li](#) 

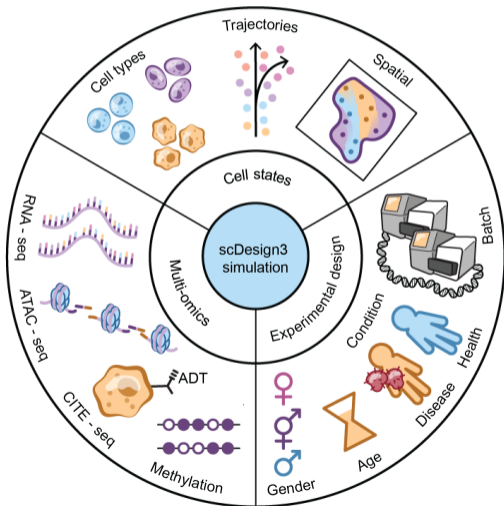
[Nature Biotechnology](#) (2023) | [Cite this article](#)

6740 Accesses | **1** Citations | **148** Altmetric | [Metrics](#)

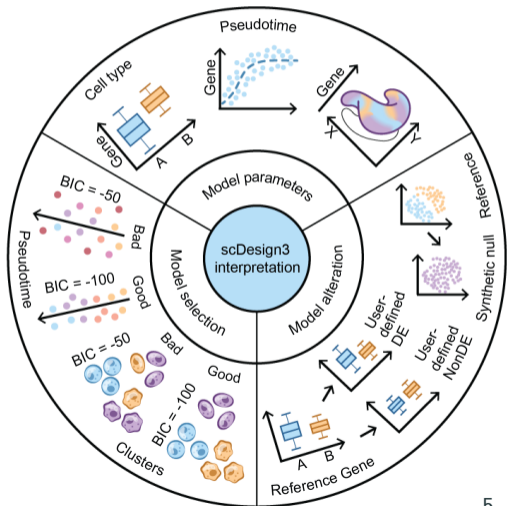


scDesign3's functionality

a



b



Methods

Mathematical notations of input data

- $\mathbf{Y} = [Y_{ij}] \in \mathbb{R}^{n \times m}$: cell-by-feature matrix
 - Y_{ij} : the measurement of feature j in cell i
 - \mathbf{Y} is usually a **count matrix** (i.e., $\mathbf{Y} \in \mathbb{N}^{n \times m}$)
- $\mathbf{X} = [\mathbf{x}_1; \dots; \mathbf{x}_n]^T \in \mathbb{R}^{n \times p}$: cell-by-state-covariate matrix, such as
 - **Cell type** ($p = 1$ categorical variable)
 - Cell **pseudotime** in p lineage trajectories (p continuous variables)
 - 2-dimensional **spatial coordinates** ($p = 2$ continuous variables)
- $\mathbf{Z} \in \mathbb{R}^{n \times q}$: cell-by-design-covariate matrix
 - $\mathbf{Z} = [\mathbf{b}; \mathbf{c}]$
 - $\mathbf{b} = (b_1; \dots; b_n)^T$ has $b_i \in \{1, \dots, B\}$; B representing cell i 's **batch**
 - $\mathbf{c} = (c_1; \dots; c_n)^T$ has $c_i \in \{1, \dots, C\}$; C representing cell i 's **condition**



Modeling features' marginal distributions

- First model the distribution of each feature j
- Use the generalized additive model for location, scale, and shape (**GAMLSS**) [Stasinopoulos and Rigby, 2008]
- The feature j 's regression model is:

$$\begin{aligned}
 Y_{ij} &\stackrel{\text{ind}}{\sim} F_j(\eta_{ij}; \eta_{ij}; p_{ij}) \\
 \eta_{ij} &= \eta_{j0} + \eta_{jb_i} + \eta_{jc_i} + f_{jc_i}(\mathbf{x}_i) \\
 \log(p_{ij}) &= \eta_{j0} + \eta_{jb_i} + \eta_{jc_i} + g_{jc_i}(\mathbf{x}_i) \\
 \text{logit}(p_{ij}) &= \eta_{j0} + \eta_{jb_i} + \eta_{jc_i} + h_{jc_i}(\mathbf{x}_i)
 \end{aligned} \tag{1}$$

where $\eta_j(\cdot)$ denotes feature j 's specific link function η_{ij} , depending on F_j

- The fitted distribution is denoted as $\hat{F}_j(\eta_j(\mathbf{x}_i; \mathbf{z}_i))$, $i = 1, \dots, n$; $j = 1, \dots, m$



Choices of marginal distributions

Distribution	PDF or PMF
Gaussian	$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}; x \in \mathbb{R}$
Bernoulli	$f(x) = \binom{1}{x} p^x (1-p)^{1-x}; x \in \{0, 1\}$
Poisson	$f(x) = \frac{\mu^x e^{-\mu}}{x!}; x \in \{0, 1, 2, \dots\}$
Negative Binomial	$f(x) = \frac{\Gamma(x + \frac{1}{\sigma})}{\Gamma(\frac{1}{\sigma})\Gamma(x+1)} \frac{1}{1+\sigma\mu} \left(\frac{\sigma\mu}{1+\sigma\mu}\right)^x; x \in \{0, 1, 2, \dots\}$
Zero-inflated Poisson	$f(x) = \begin{cases} p + (1-p)e^{-\mu}; & x = 0 \\ \frac{(1-p)\mu^x e^{-\mu}}{x!}; & x = 1, 2, 3, \dots \end{cases}$
Zero-inflated NB	$f(x) = \begin{cases} p + (1-p)\left(\frac{1}{1+\sigma\mu}\right)^{\frac{1}{\sigma}}; & x = 0 \\ \frac{(1-p)\Gamma(x + \frac{1}{\sigma})}{\Gamma(\frac{1}{\sigma})\Gamma(x+1)} \frac{1}{1+\sigma\mu} \left(\frac{\sigma\mu}{1+\sigma\mu}\right)^x; & x = 1, 2, 3, \dots \end{cases}$



Functions of modeling cell states

Covariate type	Covariate form	Function form ¹
Discrete cell type	$x_i \in \{1, \dots, K\}$	$f_{j_c i}(x_i) = \mathbb{1}_{j_c = x_i}$
One lineage	$x_i \in [0, 1)$	$f_{j_c i}(x_i) = \mathbb{1}_{j_c = \lfloor Kx_i \rfloor}$
p lineages	$\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^T \in [0, 1)^p$	$f_{j_c i}(\mathbf{x}_i) = \prod_{l=1}^p \mathbb{1}_{j_c = \lfloor Kx_{il} \rfloor}$
Spatial location	$\mathbf{x}_i = (x_{i1}, x_{i2})^T \in \mathbb{R}^2$	$f_{j_c i}(\mathbf{x}_i) = f_{j_c i}^{\text{GP}}(x_{i1}, x_{i2}; K)$

¹For simplicity, we only show the form of $f_{j_c i}(\cdot)$ because $g_{j_c i}(\cdot)$ and $h_{j_c i}(\cdot)$ have the same form.



Modeling features' joint distribution

- Denote cell i 's m features as a random vector $\mathbf{Y}_i = (Y_{i1}; \dots; Y_{im})^\top$
- Denote the joint CDF as: $F(\cdot | \mathbf{x}_i; \mathbf{z}_i) : \mathbb{R}^m \rightarrow [0; 1]$
- Modeling the joint CDF is challenging; thus we use **copula**
- Denote the conditional copula as $C(\cdot | \mathbf{x}_i; \mathbf{z}_i) : [0; 1]^m \rightarrow [0; 1]$:

$$F(\mathbf{y}_i | \mathbf{x}_i; \mathbf{z}_i) = C(F_1(y_{i1} | \mathbf{x}_i; \mathbf{z}_i); \dots; F_m(y_{im} | \mathbf{x}_i; \mathbf{z}_i) | \mathbf{x}_i; \mathbf{z}_i);$$

where $\mathbf{y}_i = (y_{i1}; \dots; y_{im})^\top$ is a realization of $\mathbf{Y}_i = (Y_{i1}; \dots; Y_{im})^\top$

- The simplest choice is a **Gaussian copula**:

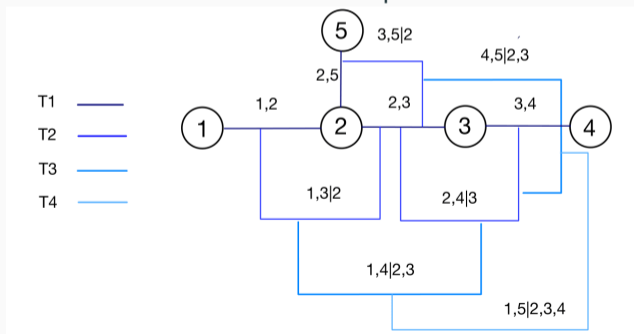
$$\begin{aligned} & C(F_1(y_{i1} | \mathbf{x}_i; \mathbf{z}_i); \dots; F_m(y_{im} | \mathbf{x}_i; \mathbf{z}_i) | \mathbf{x}_i; \mathbf{z}_i) \\ &= \Phi_m(\Phi^{-1}(F_1(y_{i1} | \mathbf{x}_i; \mathbf{z}_i)); \dots; \Phi^{-1}(F_m(y_{im} | \mathbf{x}_i; \mathbf{z}_i)); \mathbf{R}(\mathbf{x}_i; \mathbf{z}_i)) \end{aligned}$$

where $\Phi_m(\cdot | \mathbf{R})$ is a **joint Gaussian CDF** with a zero mean vector and a covariance matrix that is equal to the **correlation matrix \mathbf{R}**



From Gaussian copula to vine copula

- Since we often have $m \ll n$, Gaussian copula can be problematic
- One solution to model high-dimensional correlation is **Vine copula** [Czado et al., 2009]
- “Decompose” a high-dimensional copula into a sequence of bivariate copulas
- Graphic illustration of a 5-dimensional vine copula:



The plug-in estimation of copula

- To estimate $C(j \mathbf{x}_i; \mathbf{z}_i)$, we use the plug-in approach:

$$\hat{F}_1(j \mathbf{x}_i; \mathbf{z}_i); \dots; \hat{F}_m(j \mathbf{x}_i; \mathbf{z}_i)$$

- If $\hat{F}_j(j \mathbf{x}_i; \mathbf{z}_i)$ is a continuous distribution, each observed y_{ij} is transformed as:

$$u_{ij} = \hat{F}_j(y_{ij} | j \mathbf{x}_i; \mathbf{z}_i)$$

- If $\hat{F}_j(j \mathbf{x}_i; \mathbf{z}_i)$ is a discrete distribution, we use the **distributional transformation** to make it “continuous”:

$$u_{ij} = v_{ij} \hat{F}_j(y_{ij} | j \mathbf{x}_i; \mathbf{z}_i) + (1 - v_{ij}) \hat{F}_j(y_{ij} | j \mathbf{x}_i; \mathbf{z}_i); y_{ij} = 1; 2; \dots;$$

where v_{ij} 's are sampled independently from Uniform[0;1]

- $u_{ij} = \tilde{F}_j(y_{ij} | j \mathbf{x}_i; \mathbf{z}_i)$, where $\tilde{F}_j(j \mathbf{x}_i; \mathbf{z}_i)$ is the CDF of a continuous distribution
- Then $C(j \mathbf{x}_i; \mathbf{z}_i)$ is estimated from $\mathbf{u}_1; \dots; \mathbf{u}_n$, where $\mathbf{u}_i = (u_{i1}; \dots; u_{im})^T$



Generating the synthetic data

- Goal: generate $\mathbf{Y}^0 \in \mathbb{R}^{n^0 \times m}$ (n^0 synthetic cells and the same m features as \mathbf{Y})
- Given $\mathbf{X}^0 \in \mathbb{R}^{n^0 \times p}$ and $\mathbf{Z}^0 \in \mathbb{N}^{n^0 \times q}$,
 1. Sample a m -dimensional vector from the m -dimensional copula:

$$(U_{i'1}; \dots; U_{i'm})^T \sim \hat{C}(j; \mathbf{x}_{i'}; \mathbf{z}_{i'}); i^0 = 1; \dots; n^0$$

2. Calculate the marginal distribution:

$$Y_{i'j} \mid j; \mathbf{x}_{i'}; \mathbf{z}_{i'} \sim \hat{F}_j(j; \mathbf{x}_{i'}; \mathbf{z}_{i'}) = F_j(j; \mathbf{x}_{i'}; \mathbf{z}_{i'}; \hat{\mu}_{i'j}; \hat{\sigma}_{i'j}; \hat{\rho}_{i'j});$$

where

$$\begin{aligned} \hat{\mu}_{i'j} &= \hat{\mu}_{j0} + \hat{\mu}_{jb_{i'}} + \hat{\mu}_{jc_{i'}} + \hat{f}_{jc_{i'}}(\mathbf{x}_{i'}); \\ \hat{\sigma}_{i'j} &= \hat{\sigma}_{j0} + \hat{\sigma}_{jb_{i'}} + \hat{\sigma}_{jc_{i'}} + \hat{g}_{jc_{i'}}(\mathbf{x}_{i'}); \\ \hat{\rho}_{i'j} &= \hat{\rho}_{j0} + \hat{\rho}_{jb_{i'}} + \hat{\rho}_{jc_{i'}} + \hat{h}_{jc_{i'}}(\mathbf{x}_{i'}); \end{aligned}$$

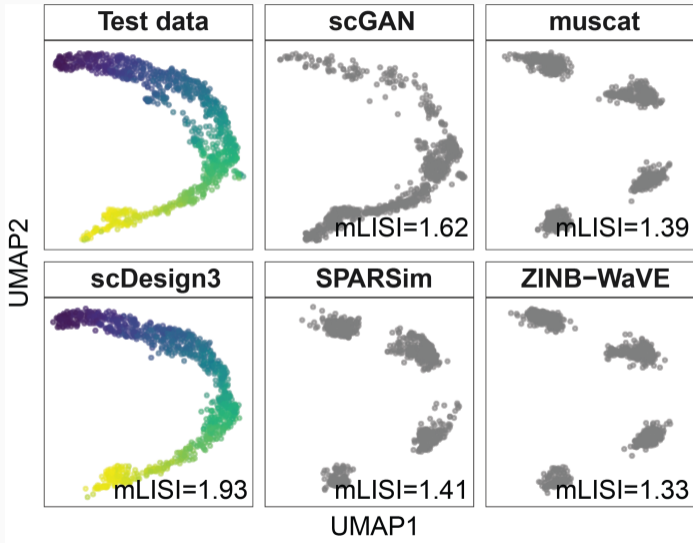
3. Get $(Y_{i'1}; \dots; Y_{i'm})^T$ by inverse CDF:

$$Y_{i'j} = \hat{F}_j^{-1}(U_{i'j} \mid j; \mathbf{x}_{i'}; \mathbf{z}_{i'}); j = 1; \dots; m$$

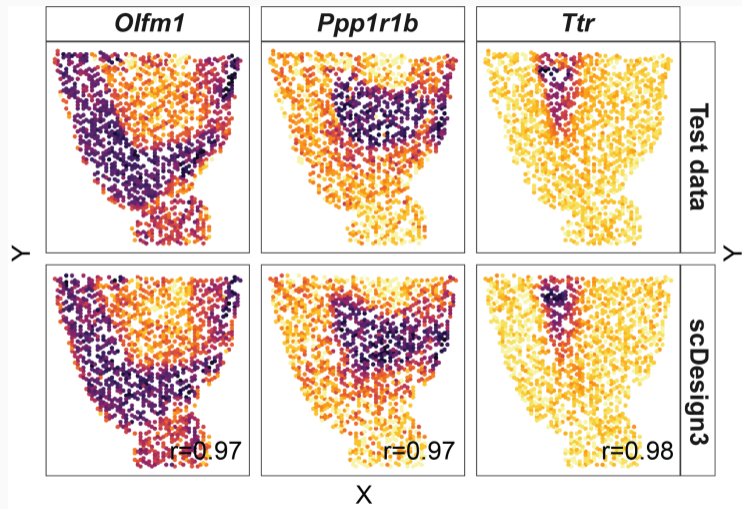


Results

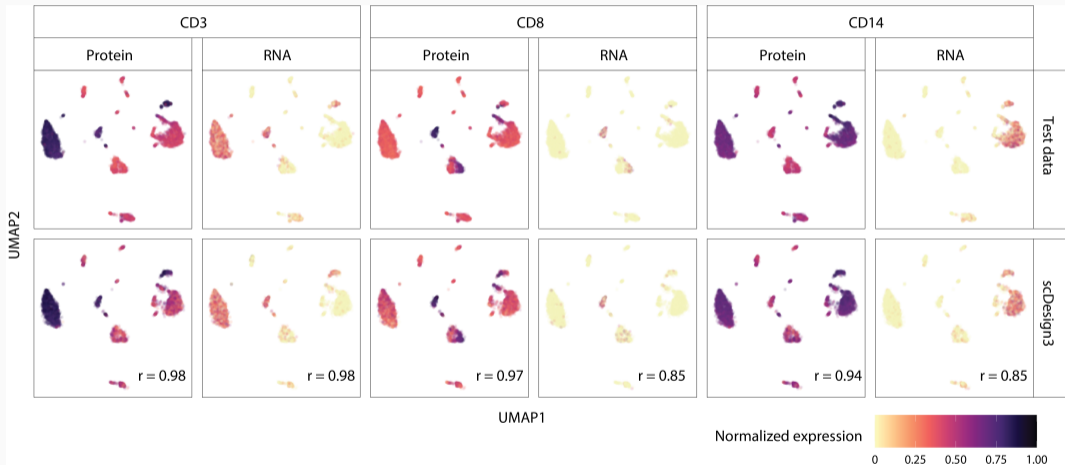
scDesign3 simulates continuous cell differentiation



scDesign3 simulates brain spatial patterns



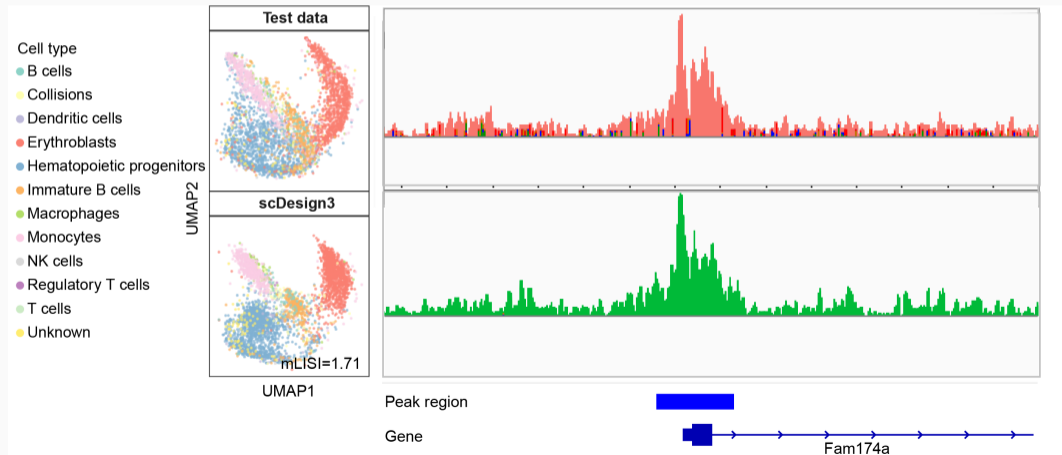
scDesign3 simulates RNA and protein co-expression in blood cells



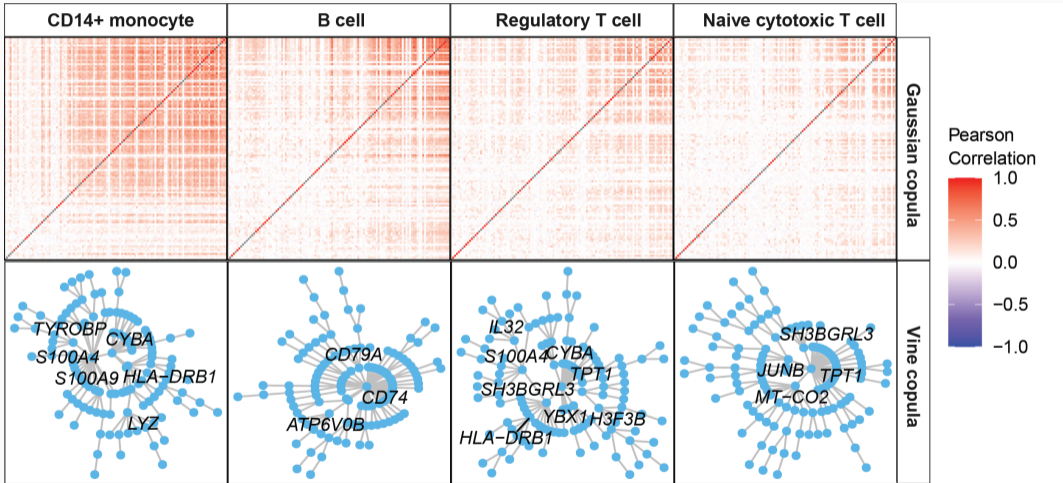
scDesign3 simulates counts and reads of bone marrow ATAC-seq data

scDesign3 +

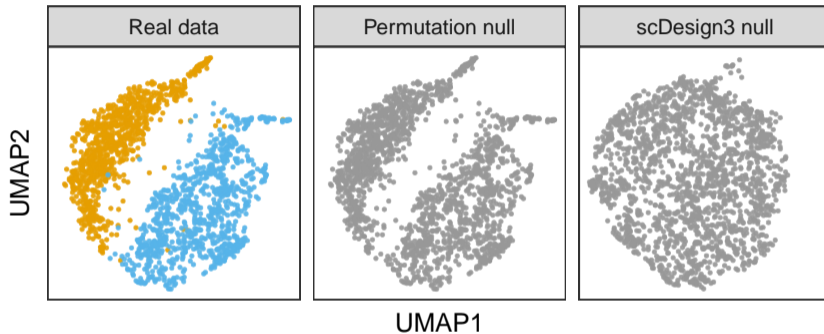
scReadSim(<https://www.biorxiv.org/content/10.1101/2022.05.29.493924v3>)



Copula reveals biological differences between immune cell types



scDesign3 generates in silico negative control



Cell type ● Naive cytotoxic T cell ● Regulatory T cell ● Null



