



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

Dongyuan Song

June 28, 2023

Bioinformatics Program

University of California, Los Angeles

Advisor: Dr. Jingyi Jessica Li

<http://jsb.ucla.edu>

About me

- Ph.D. candidate in Bioinformatics, UCLA
- Currently, visiting student at Harvard Stats
- My advisor, Jingyi Jessica Li, is a Professor in Statistics at UCLA

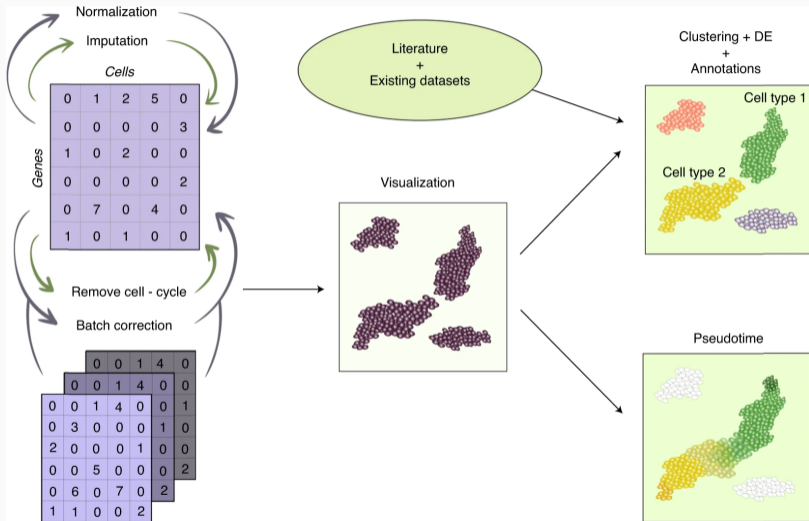


Dr. Jingyi Jessica Li

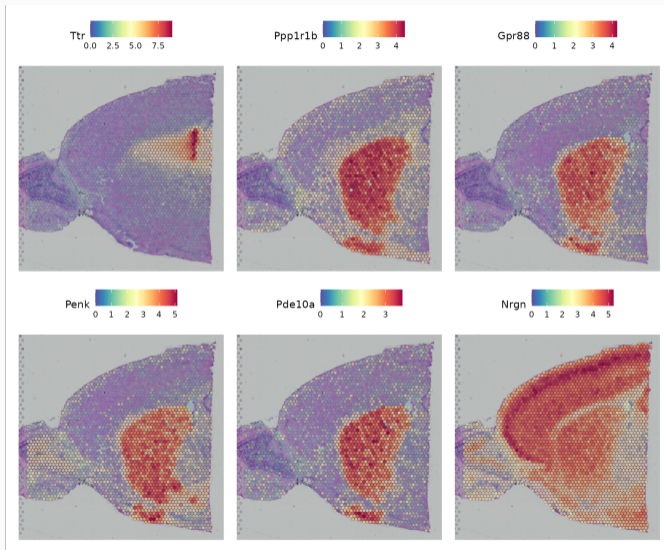


Introduction

Single-cell RNA sequencing (scRNA-seq)

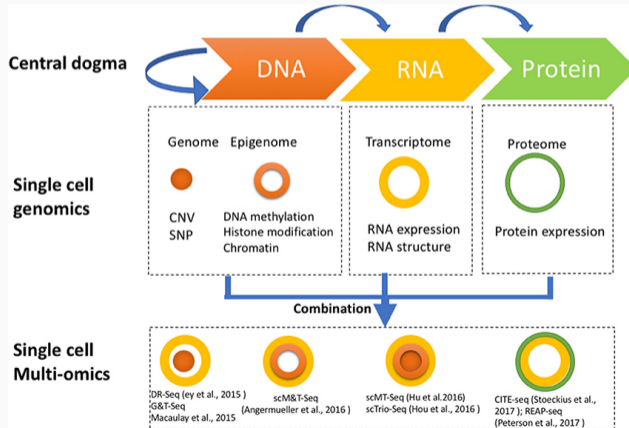


Spatial transcriptomics: gene expression in space



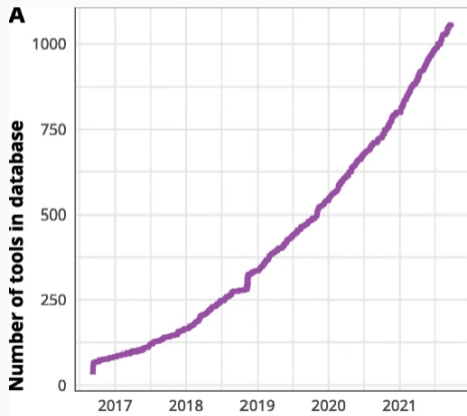
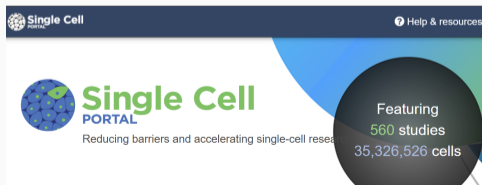
Not only genes: other omics and multi-omics

- Measure other types of features (omics) rather than genes
- Measure several types of features simultaneously (multi-omics)



Massive data and numerous tools

- So many datasets and computational tools! Therefore, people want:
 - **A unified probabilistic model** for interpreting single-cell and spatial data?
 - **An all-in-one simulator** for comparing various computational tools?



[Zippa et al., 2021]

Challenges in modeling single-cell and spatial multi-omics

- **High-dimensional:** $\sim 10^4$ 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](#) 

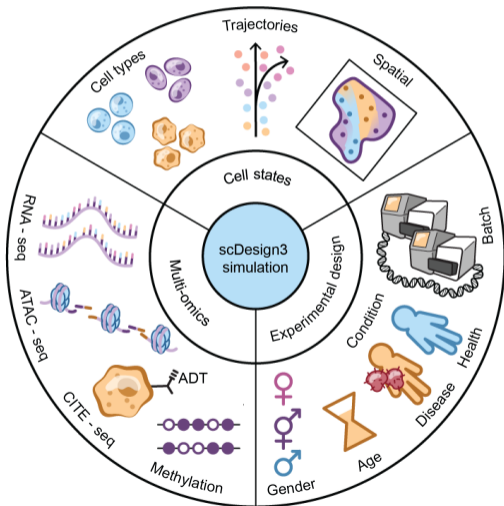
[Nature Biotechnology](#) (2023) | [Cite this article](#)

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

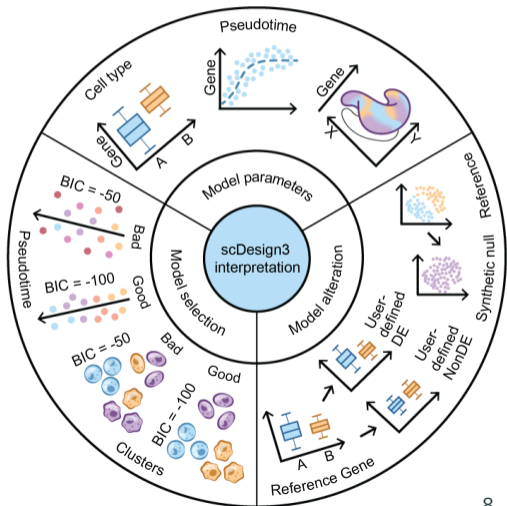


scDesign3's functionality

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 often 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 $b_i \in \{1, \dots, B\}$ representing cell i 's batch
 - $\mathbf{c} = (c_1, \dots, c_n)^T$ has $c_i \in \{1, \dots, C\}$ 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 regression model is:

$$\begin{cases} Y_{ij} \mid \mathbf{x}_i, \mathbf{z}_i & \stackrel{\text{ind}}{\sim} F_j(\cdot \mid \mathbf{x}_i, \mathbf{z}_i ; \mu_{ij}, \sigma_{ij}, p_{ij}) \\ \theta_j(\mu_{ij}) & = \alpha_{j0} + \alpha_{jb_i} + \alpha_{jc_i} + f_{jc_i}(\mathbf{x}_i) \\ \log(\sigma_{ij}) & = \beta_{j0} + \beta_{jb_i} + \beta_{jc_i} + g_{jc_i}(\mathbf{x}_i) \\ \text{logit}(p_{ij}) & = \gamma_{j0} + \gamma_{jb_i} + \gamma_{jc_i} + h_{jc_i}(\mathbf{x}_i) \end{cases}, \quad (1)$$

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

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



Choices of marginal distributions

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



Functions of modeling cell states

Covariate type	Covariate form	Function form ¹
Discrete cell type	$x_i \in \{1, \dots, K_C\}$	$f_{j_{c_i}}(x_i) = \alpha_{j_{c_i} x_i}$
One lineage	$x_i \in [0, \infty)$	$f_{j_{c_i}}(x_i) = \sum_{k=1}^K b_{j_{c_i} k}(x_i) \beta_{j_{c_i} k}$
p lineages	$\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^T \in [0, \infty)^p$	$f_{j_{c_i}}(\mathbf{x}_i) = \sum_{l=1}^p \sum_{k=1}^K b_{j_{c_i} kl}(x_{il}) \beta_{j_{c_i} lk}$
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 measurements of the m features as: a random vector $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{im})^T$
- Denote the joint CDF as: $F(\cdot \mid \mathbf{x}_i, \mathbf{z}_i) : \mathbb{R}^m \rightarrow [0, 1]$
- Modeling the joint CDF is challenging; thus we use **copula**
- Denote the conditional copula as $C(\cdot \mid \mathbf{x}_i, \mathbf{z}_i) : [0, 1]^m \rightarrow [0, 1]$:

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

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

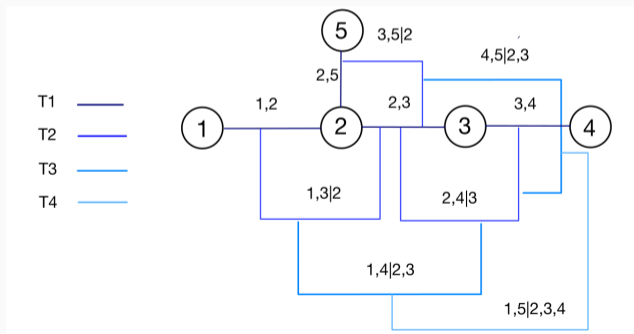
- The simplest choice is a **Gaussian copula**:

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



From Gaussian copula to vine copula

- Since we often have $m > n$, 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
- A graph can describe it:



The plug-in estimation of copula

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

$$\hat{F}_1(\cdot | \mathbf{x}_i, \mathbf{z}_i), \dots, \hat{F}_m(\cdot | \mathbf{x}_i, \mathbf{z}_i)$$

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

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

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

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

where v_{ij} 's are sampled independently from $\text{Uniform}[0, 1]$

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



Generating the synthetic data

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

$$(U_{i'1}, \dots, U_{i'm})^T \sim \hat{C}(\cdot \mid \mathbf{x}_{i'}, \mathbf{z}_{i'}), \quad i' = 1, \dots, n'$$

2. Calculate the marginal distribution:

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

where

$$\begin{cases} \theta(\hat{\mu}_{i'j}) &= \hat{\alpha}_{j0} + \hat{\alpha}_{jb_{i'}} + \hat{\alpha}_{jc_{i'}} + \hat{f}_{jc_{i'}}(\mathbf{x}_{i'}), \\ \log(\hat{\sigma}_{i'j}) &= \hat{\beta}_{j0} + \hat{\beta}_{jb_{i'}} + \hat{\beta}_{jc_{i'}} + \hat{g}_{jc_{i'}}(\mathbf{x}_{i'}), \\ \text{logit}(\hat{p}_{i'j}) &= \hat{\gamma}_{j0} + \hat{\gamma}_{jb_{i'}} + \hat{\gamma}_{jc_{i'}} + \hat{h}_{jc_{i'}}(\mathbf{x}_{i'}). \end{cases}$$

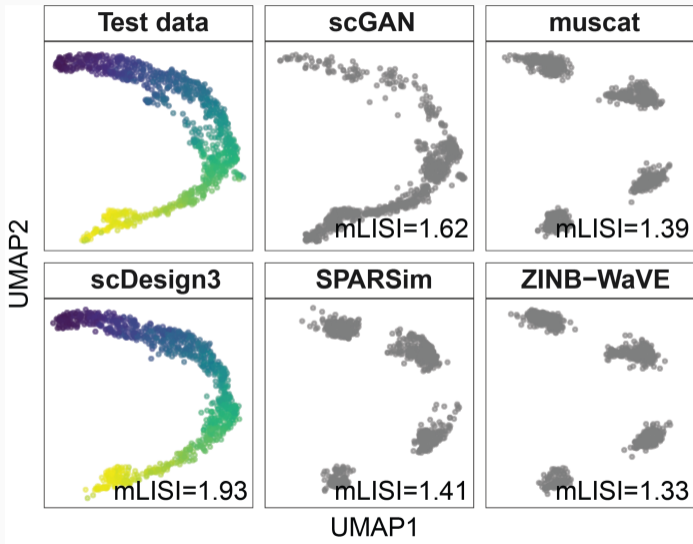
3. Get $(Y_{i'1}, \dots, Y_{i'm})^T$ by inverse CDF:

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

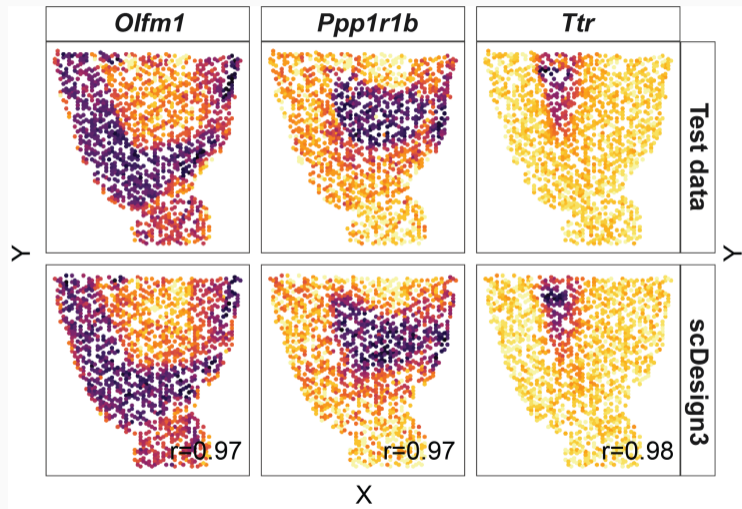


Results

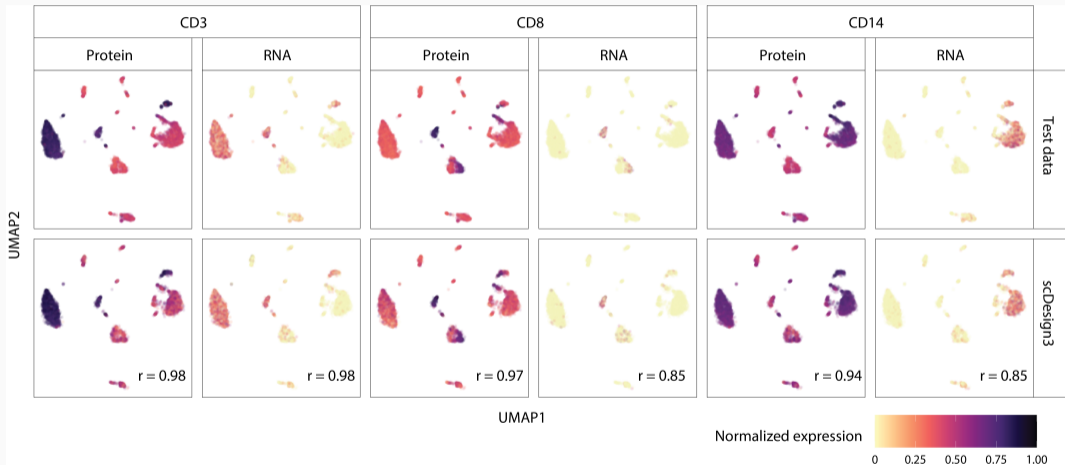
scDesign3 simulates continuous cell differentiation



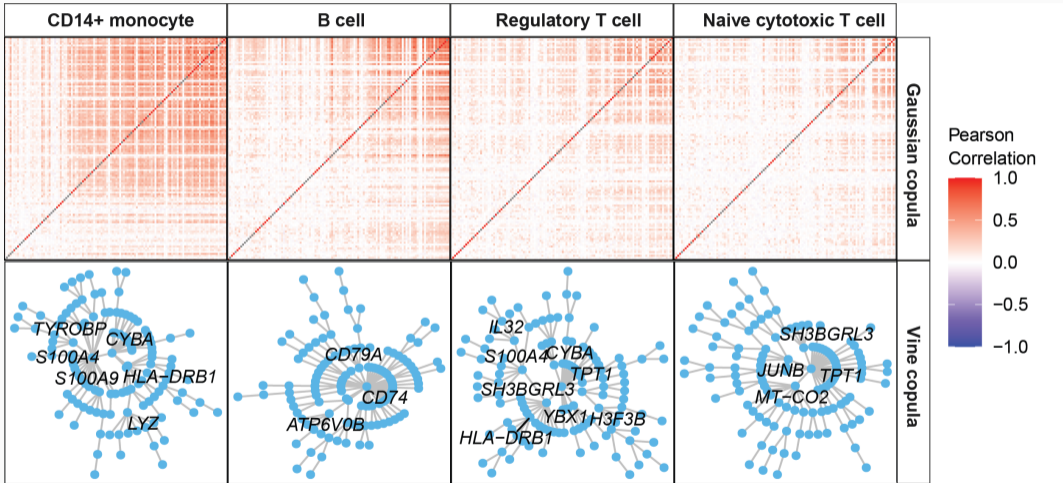
scDesign3 simulates brain spatial patterns



scDesign3 simulates RNA and protein co-expression in blood cells



Copula reveals biological differences between immune cell types



Links and Acknowledgements



Qingyang Wang



Dr. Jingyi Jessica Li

- Paper link: <https://www.nature.com/articles/s41587-023-01772-1>
- Software: <https://github.com/SONGDONGYUAN1994/scDesign3>
- My email: dongyuansong@ucla.edu
- Many thanks to the BKT workshop 2023!

