





Gene expression

scDiffusion: conditional generation of high-quality single-cell data using diffusion model

Erpai Luo ^{1,†}, Minsheng Hao ^{1,†}, Lei Wei ¹, Xuegong Zhang ^{1,2,*}

¹MOE Key Lab of Bioinformatics and Bioinformatics Division of BNRIST, Department of Automation, Tsinghua University, Beijing 100084, China

²School of Life Sciences and School of Medicine, Center for Synthetic and Systems Biology, Tsinghua University, Beijing 100084, China

*Corresponding author. E-mail: zhangxg@tsinghua.edu.cn

[†]= equal contribution.

Associate Editor: Anthony Mathelier

Abstract

Motivation: Single-cell RNA sequencing (scRNA-seq) data are important for studying the laws of life at single-cell level. However, it is still challenging to obtain enough high-quality scRNA-seq data. To mitigate the limited availability of data, generative models have been proposed to computationally generate synthetic scRNA-seq data. Nevertheless, the data generated with current models are not very realistic yet, especially when we need to generate data with controlled conditions. In the meantime, diffusion models have shown their power in generating data with high fidelity, providing a new opportunity for scRNA-seq generation.

Results: In this study, we developed scDiffusion, a generative model combining the diffusion model and foundation model to generate high-quality scRNA-seq data with controlled conditions. We designed multiple classifiers to guide the diffusion process simultaneously, enabling scDiffusion to generate data under multiple condition combinations. We also proposed a new control strategy called Gradient Interpolation. This strategy allows the model to generate continuous trajectories of cell development from a given cell state. Experiments showed that scDiffusion could generate single-cell gene expression data closely resembling real scRNA-seq data. Also, scDiffusion can conditionally produce data on specific cell types including rare cell types. Furthermore, we could use the multiple-condition generation of scDiffusion to generate cell type that was out of the training data. Leveraging the Gradient Interpolation strategy, we generated a continuous developmental trajectory of mouse embryonic cells. These experiments demonstrate that scDiffusion is a powerful tool for augmenting the real scRNA-seq data and can provide insights into cell fate research.

Availability and implementation: scDiffusion is openly available at the GitHub repository <https://github.com/EperLuo/scDiffusion> or Zenodo <https://zenodo.org/doi/10.5281/zenodo.13268742>.

1 Introduction

Single-cell RNA sequencing (scRNA-seq) data offer comprehensive depictions of the gene expression profile of every single cell, gaining a more systematic and precise understanding of the development and function of living organisms (Gohil *et al.* 2021, Jovic *et al.* 2022). Although current sequencing technologies have come a long way, the cost and difficulty of sequencing remain high. Besides, the biological samples are sometimes hard to obtain (Suvà and Tirosh 2019, Jiang *et al.* 2022, Ke *et al.* 2022), and certain cell types within a sample may be too rare to be analyzed. It is still challenging to obtain enough high-quality scRNA-seq data of interest, which may impede biological discovery as most tools for scRNA-seq analysis require a certain amount of high-quality data.

Some researchers have endeavored to generate *in silico* gene expression data that obviate the need for further biological samples, thus mitigating the limited availability of scRNA-seq data. This pseudo data is designed to meet specific criteria, thereby facilitating more effective downstream analysis. There are two main types of *in silico* data generation methods: statistical models and deep generative models.

Statistical models use some well-studied statistical distributions such as Zero-inflated Negative Binomial (ZINB) (Greene 1994) to model the expression data, and new data are generated by manually setting certain parameters of the distributions (Zappia *et al.* 2017, Lindenbaum *et al.* 2018, Li and Li 2019, Dibaenia and Sinha 2020). Some statistical models such as SPARSim (Baruzzo *et al.* 2020) and SCRIP (Qin *et al.* 2022) may oversimplify the complex patterns in real scenarios, making these methods hard to mimic real gene expression data well. Besides, statistical models such as scDesign3 (Song *et al.* 2024) need to carefully design the model to cope with the different statistical laws in different conditions and are hard to generate data under conditions that are not taken into account in the original design.

The recent prosperity of deep generative models brings great chances for the *in silico* transcriptomic data generation (Lopez *et al.* 2020). By principle, the current models can be summarized into two types: the variational autoencoder (VAE) based and the generative adversary network (GAN) based. Although the VAE-based models such as scVI (Lopez *et al.* 2018) are the most prominent in this field (Kingma and

Welling 2013), these models are mainly focused on downstream analysis tasks (e.g. batch correction and clustering), rather than generating gene expression profiles of cells. GAN-based models such as scGAN were proposed for generating new cells and accomplishing downstream tasks (Marouf et al. 2020, Xu et al. 2020, Lall et al. 2022). However, these GAN-based models can only generate data from a known distribution, unsatisfying the need to supplement unmeasured data. Besides, in their conditional generation experiments, they usually cluster the dataset first and then control the conditional generation based on the labels of clusters, which is not enough when we need to generate cells with more fine-grained conditions. Furthermore, GAN requires careful designing and tuning of model architectures, as well as optimization tricks, to achieve a stable training process (Saxena and Cao 2021, Yang et al. 2022), hindering the smooth application on new datasets and the generation of data under certain conditions.

Recently, the latent diffusion model (LDM) (Rombach et al. 2022) has demonstrated excellent performance in several areas such as images, audio, and videos (Cao et al. 2022, Yang et al. 2022, Croitoru et al. 2023). Compared with the GAN model, it has a stable training process and can easily generate samples conditioned on complex prompts (Dhariwal and Nichol 2021, Bond-Taylor et al. 2021, Zhang et al. 2023). However, few studies have deployed it in the single-cell area. One of the challenges lies in the fact that LDM needs a pre-trained autoencoder model to link the data in the latent and original space. There are no such models in the single-cell field until the recent emergence of the foundation models (Theodoris et al. 2023, Bian et al. 2024, Cui et al. 2024, Hao et al. 2024). By using massive parameters and being trained on extensive datasets, these models can learn the unified representation of the gene expression data, which facilitates a variety of downstream tasks and can be used as the autoencoder model in LDM.

Here we propose scDiffusion, an *in silico* scRNA-seq data generation model combining LDM with the foundation model, to generate single-cell gene expression data with given conditions. scDiffusion has three parts, an autoencoder, a denoising network, and a condition controller. We used the pre-trained model SCimilarity (Heimberg et al. 2023) as an autoencoder to rectify the raw distribution and reduce the dimensionality of scRNA-seq data, which can make the data amenable to diffusion modeling. The denoising network was redesigned based on a skip-connected multilayer perceptron (MLP) to learn the reversed diffusion process. The conditional controller is a cell type classifier, enabling scDiffusion to generate data specific to a particular cell or organ type according to diverse requirements.

We conducted a series of experiments to evaluate the performance of our model. First, we evaluated the single-conditional generation ability of scDiffusion. We generated new data on three different datasets and evaluated the generated data with different metrics. The results showed that the scDiffusion could generate realistic scRNA-seq data and had superior conditional generation ability. Then we assessed the multi-conditional generation capabilities of scDiffusion. We deployed two separate classifiers to guide the generation. The results showcased the model's ability to generate out-of-distribution data through the amalgamation of known conditions. We also proposed a new condition control strategy, Gradient Interpolation, to interpolate continuous cell trajectories from discrete cell states.

We used this strategy to reconstruct the intermediate states within an embryonic development process. The results showed that the scDiffusion could bridge the gaps between sequencing intervals and provide a more comprehensive developmental timeline. With the powerful generation ability, we believe scDiffusion has the potential to augment existing scRNA-seq data and contribute to the investigation of undersampled or even unseen cell states.

2 Materials and methods

The scDiffusion model consists of three parts, a pre-trained foundation model SCimilarity (Heimberg et al. 2023) as the autoencoder, a denoising network, and a conditional classifier, as depicted in Fig. 1. At the training stage, the SCimilarity model is first finetuned based on the pre-trained weight to embed the gene expression profile. After that, the diffusion process is applied to each embedding derived by the autoencoder and produces a series of noisy embeddings. These noisy embeddings serve as the training data for the backbone network. Meanwhile, the conditional classifier processes the embeddings to predict associated labels, such as cell types. At the inference stage, the denoising network denoises the input noise embeddings and generates new embeddings. The generation can be guided by the classifier or the Gradient Interpolation strategy. The generated embeddings are finally fed into the decoder to obtain full gene expression. Detailed descriptions of scDiffusion are provided below.

2.1 Finetuning the pre-trained foundation model as autoencoder

An important prerequisite for LDM is to have a suitable autoencoder to concatenate the data in latent space and original space. To satisfy this prerequisite, we used the pre-trained foundation model SCimilarity as the autoencoder, encoding the gene expression data of every single cell S_{ori} into a latent space embedding x_0 . SCimilarity is an encoder-decoder style network trained on a 22.7 million cell corpus assembled across 399 published scRNA-seq studies and could be used to extract unifying representation from cell expression profiles.

We finetuned the pre-trained model weights with our data and used SCimilarity's encoder and decoder separately. The input of the encoder is a gene expression profile that is normalized by $1e4$ total counts and then logarithmized, and the output is a 128D latent space embedding as it was set in the pre-trained model. The decoder subsequently accepts the latent space embedding and generates the corresponding expression profile S_{new} as the output. Leveraging the powerful representation and generalization capabilities of foundation model, we could extract complete information from cellular expression into latent representations and accurately rebuild the original gene expression.

During the finetuning process, we used the pre-trained weight except for the first layer and the last layer since the numbers of genes in our data differ from that in the pre-trained dataset. Compared to training a new autoencoder from scratch, this method gives faster and better access to the desired autoencoder (Fig. 2a). As shown in Supplementary Fig. S1, the distribution of gene expression is transformed into a Gaussian-like distribution by the autoencoder, which is in line with the Gaussian distribution used in the diffusion process and makes it much easier for the denoising network to learn the reverse process.

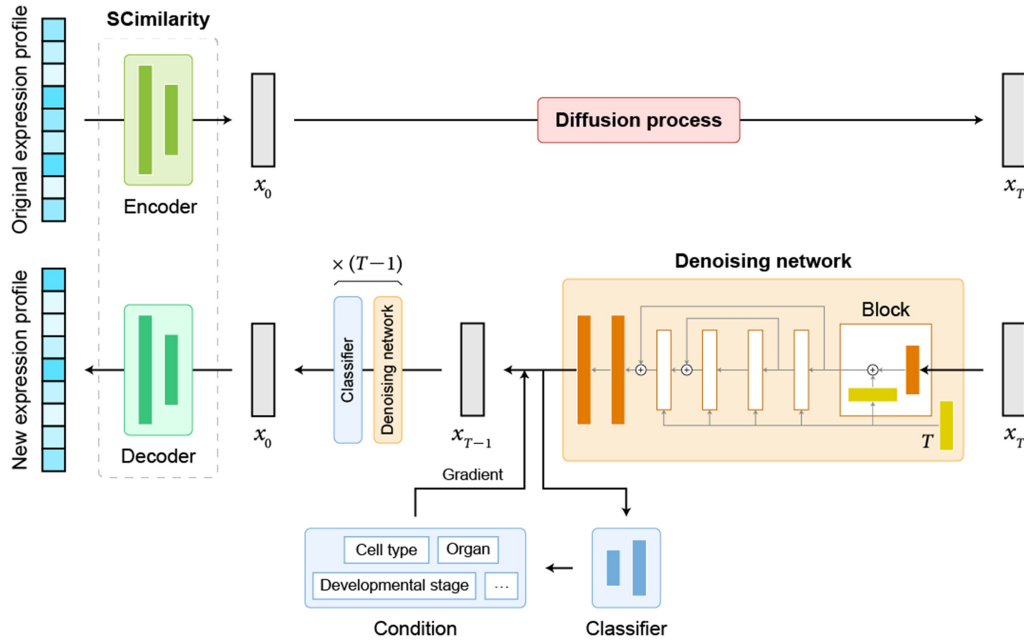


Figure 1. The overall structure of scDiffusion.

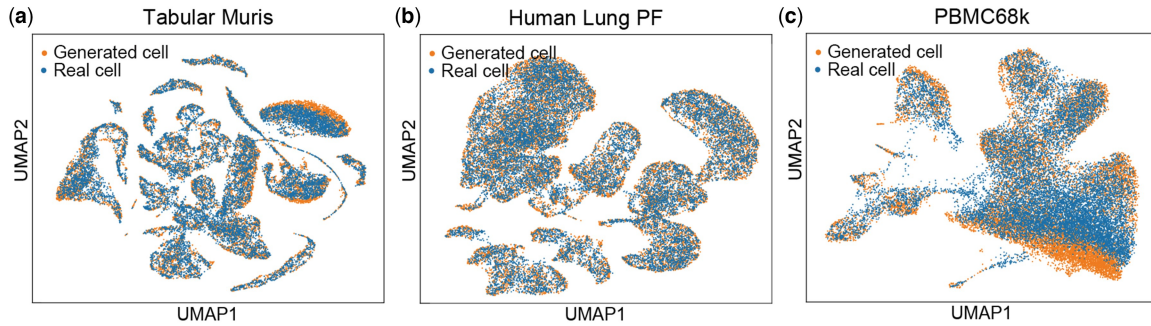


Figure 2. scDiffusion can generate realistic cell data. (a) UMAP of scDiffusion-generated Tabular Muris data and real Tabular Muris data. (b) UMAP of scDiffusion-generated Human Lung PF data and real Human Lung PF data. (c) UMAP of scDiffusion-generated PBMC68k data and real PBMC68k data.

2.2 Training the denoising network

After getting embeddings from the encoder, the diffusion process is applied to each embedding to obtain noisy data. The denoising network is trained on these noisy embeddings to learn the reversed process. Classical denoising network models such as convolutional neural networks are not applicable to gene expression, as a gene expression profile of scRNA-seq data is a long, sparse, and unordered vector. Thus, we developed a new architecture as the backbone, with fully connected layers and a skip-connected structure (Fig. 1). The skip-connected structure can help to maintain the characteristics of features at different levels and reduce the loss of information.

In the diffusion process, the original cell embedding x_0 becomes a noisy embedding x_T by iteratively adding noise through T steps. For the i -th step, the embedding x_i is sampled from the following distribution:

$$q(x_i | x_{i-1}) = N\left(x_i | \sqrt{1 - \beta_i} x_{i-1}, \beta_i I\right), \approx \text{where } \beta_i \in (0, 1) \quad (1)$$

where I stands for standard Gaussian noise. β_i is a coefficient that varies with time step, and β_{min} and β_{max} are two parameters that control the scale of β_i in the diffusion process:

$$\beta_i = \frac{\beta_{min}}{T} + \frac{i-1}{T-1} \left(\frac{\beta_{max}}{T} - \frac{\beta_{min}}{T} \right) \quad (2)$$

The training goal is to learn the reverse diffusion process $p(x_{i-1} | x_i)$. In each iteration, x_{i-1} at step $i-1$ is predicted, given an embedding x_i at step i . Such process also follows the Gaussian distribution. According to previous works (Ho et al. 2020, Dhariwal and Nichol 2021), the mean and variance are parameterized as:

$$p_\theta(x_{i-1} | x_i) = N(x_{i-1} | \mu_\theta(x_i, i), \exp(w\beta_i)I) \quad (3)$$

where w in the variance is an adjustable weight that controls the randomness of the reverse process. The mean $\mu_\theta(x_i, i)$ can be written as:

$$\mu_\theta(x_i, i) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \varepsilon_\theta(x_t, t) \right) \quad (4)$$

where $\alpha_t = 1 - \beta_t$ and $-\alpha_t = \prod_{s=1}^t \alpha_s$. $\varepsilon_\theta(x_t, t)$ is the added noise predicted by the backbone network. In other words, the backbone network takes the cell's latent space embedding x_i and the timestamp i as inputs to predict the noise.

In the inference process, the diffusion model takes the Gaussian noise as the initial input and denoises it iteratively through T steps. Eventually, we can get the new cellular latent space embedding x_0 and put it into the decoder to get the final gene expression data.

2.3 Conditional generation and the gradient interpolation strategy

We use the classifier guidance method to perform conditional generation. This method does not interfere with the training of the denoising network model. Instead, the classifier is first trained separately by using condition labels like cell types and then provides gradients to guide cell generation. The classifier for guidance is very flexible. Theoretically, all deep learning-based cell classifiers that can generate gradients can be used to guide the conditional generation process, such as scANVI (Lotfollahi et al. 2022), MARS (Brbic et al. 2020), and scFoundation (Hao et al. 2024). Here, we designed the cell classifier as a four-layer MLP. After generating a series of embeddings from cells with labels, the classifier takes both timestamp i and cell embedding x_i as inputs and predicts the cell labels y paired with x_0 . The cross-entropy loss is used for training. It is worth noting that only the embeddings between step 0 and step $T/2$ of the diffusion process are used for training the classifier, considering that the signal in the rest part is too noisy to be predicted.

As for inference, given each step i between the last part of the reverse process (between timestamp 0 and $T/2$), the classifier receives the intermediate state x_i and outputs the predicted probability for every cell type. By computing the cross entropy loss between the predicted and desired condition given by the user, the gradient derived from the classifier can guide the denoising network model to generate a designated endpoint. The new embedding with the guidance is now sampled from:

$$p_\theta(x_{i-1} | x_i, y) = \mathcal{N}(x_{i-1} | \mu_\theta(x_i, I) + \beta_i \gamma \nabla_{x_i} \log p_\phi(y | x_i), \exp(w \beta_i) I) \quad (5)$$

where $p_\phi(y | x_i)$ stands for the classifier's result, and γ is a weight that controls the effectiveness of the classifier to the reverse process. ϕ indicates the trainable parameters in the classifier. This guidance will affect every step of the reverse process and finally help the model's output reach a certain condition.

Since the classifier is trained aside from the diffusion model and is only used in the inference stage, we can train multiple classifiers $\{\phi_1, \phi_2, \dots\}$ to control different conditions separately. The gradient that guides the diffusion process is the summation of all the classifiers' gradients with different weights $\{\gamma_1, \gamma_2, \dots\}$.

We proposed the Gradient Interpolation strategy to generate continuous cell condition guidance. A classifier receives two different conditions such as the initial and end state of cell differentiation, and generates two gradients at the same time. These gradients are then integrated to guide the diffusion to an unseen intermediate state. Specifically speaking, the $\beta_i \gamma \nabla_{x_i} \log p_\phi(y | x_i)$ in Equation (5) is replaced by:

$$\begin{aligned} & \beta_i \gamma \nabla_{x_i} \log p_\phi(y | x_i) \\ \rightarrow & \beta_i (\gamma_1 \nabla_{x_i} \log p_\phi(y_1 | x_i) + \gamma_2 \nabla_{x_i} \log p_\phi(y_2 | x_i)) \end{aligned} \quad (6)$$

where γ_1 and γ_2 represent two adjustable coefficients that control the distance between the generated cells and the two

target cell states. By tuning these coefficients, scDiffusion can decide which cell state the generated cell is closer to, thus generating cells with continuous states. With this strategy, the initial state of the diffusion generation process is changed from pure Gaussian noise to the latent space embedding of cells of the initial condition, following a noise addition process:

$$x_{init} = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \quad (7)$$

where x_{init} is the initial state, and t is a parameter that is smaller than the total diffusion step. α_t is the same thing as in Equation (4). This modification preserves the general characteristics of the initial cells, allowing the model to generate a series of new cell states for each given initial state. These generated cells can constitute a continuous trajectory of cell states.

2.4 Evaluation metrics

To compare the similarity between generated and real cells, we evaluated the generated data with various metrics. The statistical indicators consist of Spearman Correlation Coefficient (SCC), Maximum Mean Discrepancy (MMD) (Gretton et al. 2012), local inverse Simpson's index (LISI) (Haghverdi et al. 2018), and quantile-quantile plot (QQ-plot). We log-norm the gene expression data of generated and real cells and calculated SCC between them. The LISI score was calculated on the data-integrated KNN graph by using the Python package scib (Luecken et al. 2022). The QQ-plot was drawn using both real and generated expression data of a specific gene for a given cell type.

The nonstatistical metrics include Uniform Manifold Approximation and Projection (UMAP) visualization (McInnes et al. 2018), marker gene expression, CellTypist classification (Domínguez Conde et al. 2022), random forest evaluation, and KNN evaluation. The UMAP plot was used to visualize the generated and real expression data on a 2D plane to provide a subjective judgment for the generated data.

CellTypist is used to judge whether the conditionally generated data can be classified into the right type. The random forest evaluation shares the same idea with scGAN, which uses a random forest model with 1000 trees and 5 maximum depths to distinguish cells from real and generated, and the more similar these two cells are, the closer the area under the receiver operating characteristic (ROC) curve (AUC) metric for random forests approaches 0.5. The KNN evaluation metric is the same as the random forests metrics except the classifier model is switched to the KNN model.

3 Results

We conducted four experiments to demonstrate the capability of scDiffusion. First, we investigated the data generation ability of scDiffusion and compared it with the deep learning-based method scGAN and statistical learning based methods scDesign3, SPARSim, and SCRIP. We then assessed scDiffusion on a single-conditional generation task to generate specific cell types. Furthermore, we applied scDiffusion in a multi-conditional generation case with both cell types and organs as conditions and used it to generate new cells under an unseen condition which is out of the distribution of the training data. Lastly, we used the Gradient Interpolation strategy to generate intermediate states in cell reprogramming.

We used five single-cell transcriptomic datasets in these experiments. The Human Lung Pulmonary fibrosis (PF) dataset (Habermann *et al.* 2020) is a large scRNA-seq dataset that contains >110 thousand cells of human lungs with PF. Tabular Muris (Schaum *et al.* 2018) is a large-scale single-cell transcriptomic database of mice across 12 organs. The Tabular Sapiens (TTS Consortium *et al.* 2022) is the first-draft human cell atlas of nearly 500 000 cells from 24 organs of 15 normal human subjects, we selected six cell types from five of these organs to do the experiments. The Waddington-OT dataset (Schiebinger *et al.* 2019) is a cell reprogramming dataset of mouse embryonic fibroblasts (MEFs), containing cells with different timestamps during an 18-day reprogramming process. The PBMC68k dataset (Zheng *et al.* 2017) is a classical scRNA-seq dataset that contains 11 different cell types of human peripheral blood mononuclear cells (PBMCs). As the CD4+ T helper 2 cells had an extremely low number and could not be classified by CellTypist, we removed them for downstream analysis. For all 5 datasets, we filtered out cells with less than 10 expression counts and genes that expressed in <3 cells. In the CellTypist training, the data were split into training and testing sets with a ratio of 0.8–0.2, whereas the random forest and KNN models utilized splits of 0.75:0.25 and 0.7:0.3, respectively.

In all experiments, we set the diffusion step to 1000. The parameter γ in Equation (5) was set to 2. The parameter w in Equations (3) and (5) was set to 0.5. The parameter t in Equation (7) was set to 600. Detailed strategies of training and inference as well as the ablation study of scDiffusion were provided in the Supplementary Material.

3.1 Realistic scRNA-seq data generation

We first examined the ability of the diffusion model to generate without the guidance of the classifier, such ability is the foundation of the conditional generation. We applied scDiffusion on the Tabular Muris dataset, Human Lung PF dataset, and PBMC68k dataset to generate new cells (Fig. 2a–c). For comparison, we also generated cells with scGAN and scDesign3 using their default parameter settings.

We evaluated the performance of scDiffusion, scGAN, scDesign3, SPARSim, and SCRIP, with various metrics, and the results indicated that scDiffusion could generate realistic scRNA-seq data that was comparable with the state of the art method. The average SCC of scDiffusion in three datasets was 0.984, the mean MMD score was 0.018, the mean LISI score was 0.887 and the mean AUC of the random forest was 0.697, which are similar to scGAN and scDesign3, and outperform SPARSim and SCRIP (Supplementary Fig. S2).

We then evaluated whether using the pre-trained foundation model as the autoencoder impacts the performance of scDiffusion. We trained three additional scDiffusion models from scratch with varying training steps. Results showed that the pre-trained foundation model can obtain higher generation quality with shorter training time (Supplementary Table S1), confirming the important role of the pre-trained foundation model in improving the efficiency of scDiffusion.

3.2 Conditionally generating specific cell types

We next trained a cell type classifier according to the annotations provided by the Tabular Muris dataset to guide the conditional generation of scDiffusion. For each cell type, we conditionally generated the same number of cells as the real data.

As shown in Fig. 3a and Supplementary Fig. S3, the conditionally generated cells visually overlapped with the real cells on the UMAP plot. In order to compare the quality of conditionally generated cells, we trained the cscGAN model (Marouf *et al.* 2020) with the same dataset and labels, conditionally generated each type of cell to make the comparison. Additionally, we added the cells generated by scDesign3 to the comparison since they have cell type labels, too.

We used CellTypist to classify these conditionally generated cells. As shown in Supplementary Table S2, the classification accuracies of the diffusion-generated cells were close to those of the real cells in the test set, with the diffusion-generated cells attaining an average accuracy of 0.93 across all cell types, compared to the 0.98 of the real cells. The cells generated by the scDesign3 also had high accuracies, achieving a mean accuracy of 0.99. In contrast, cells generated by the cscGAN could not be distinguished by the CellTypist, culminating in a mean accuracy of 0.04.

Considering that CellTypist was mainly used to distinguish between different cell types, we further used KNN model to distinguish between cells of each cell type and the corresponding real cells. The results in Fig. 3b showed that KNN models could not distinguish between diffusion-generated cells of a specified type and the real cells of that type, as the AUC scores in all kinds of cells are near 0.5, while it can easily distinguish the GAN-generated cells and scDesign3-generated cells, with basically all the AUC scores >0.7. We further examined the expressions of key genes in the Tabular Muris dataset. We selected five transcription factors (Klf13, Ybx1, Hnrnpk, Cnbp, Hmgb2) that have the highest mean Gini importance when making cell type classification in the original paper (Schaum *et al.* 2018), and drew the QQ-plots of these genes in different cell types between real and conditional generated data (Supplementary Fig. S4a), which indicated that the expression of the key genes generated by the scDiffusion was close to that of real.

Similarly, we applied this analytical procedure to the PBMC68k dataset and observed congruent outcomes. (Fig. 3c and d, Supplementary Figs S3 and S4b, Supplementary Table S3). It's worth mentioning that rare cell types, such as the Thymus cell in the Tabular Muris dataset (2.5% in the whole dataset) and the CD34+ cell in the PBMC68k dataset (0.4% in the whole dataset), can also be well generated.

3.3 Generating out-of-distribution cell data with multiple conditions

We then tried to generate cells with multiple conditions based on the Tabular Muris dataset. We trained two classifiers to separately control different conditions, one for organ type and the other for cell type. We selected three cell groups, mammary gland T cell, spleen T cell, and spleen B cell, from the dataset for training. We would like to generate cells with a new combination of conditions (mammary B cell) which was not seen in the training data, or in other words, out of the distribution of the training data.

To test the generated multi-conditional data, we trained a CellTypist model with all kinds of cells in the mammary gland and used it to classify the real and generated mammary gland B cells. The result showed that 98% of the generated cells and 92% of the real B cells in the test set were categorized into the B cell, which showed that the scDiffusion can generate mammary B cells comparable to the real one. We

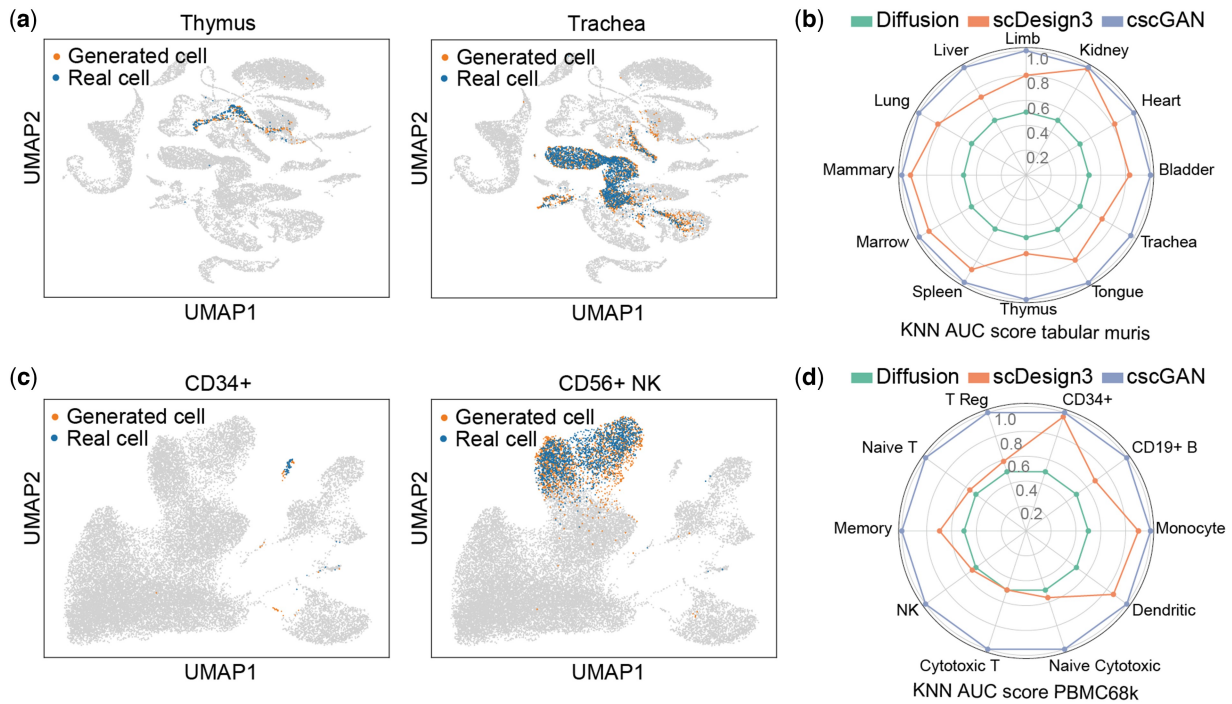


Figure 3. (a) UMAP of different cell types in the Tabular Muris dataset generated by conditional diffusion. The Thymus cell is a rare cell type. (b) The AUC score of KNN in different cell types in the Tabular Muris dataset. (c) UMAP of different cell types in the PBMC68k dataset generated by conditional diffusion. The CD34+ cell is a rare cell type. (d) The AUC score of KNN in different cell types in the PBMC68k dataset.

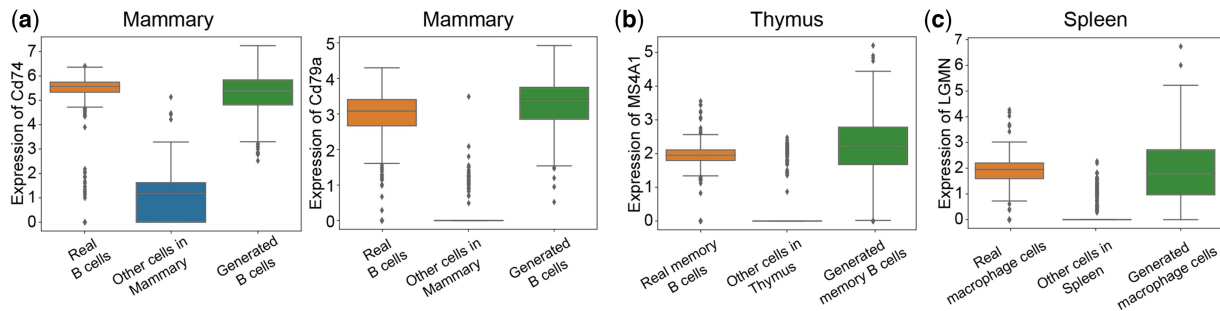


Figure 4. Marker genes' expression levels of real and multi-conditionally generated cells. (a) Marker genes of mammary B cells. (b) Marker gene of thymus memory B cells. (c) Marker gene of spleen macrophage cells.

then picked two marker genes of mammary B cells [CD74 (Bhatt et al. 2021) and CD79A (Hilton et al. 2019)] and drew the box plot for them. The results in Fig. 4a showed that the generated marker genes have similar expression levels as the real ones.

We further selected cells from five organs and six cell types from the Tabular Sapiens dataset and did the same training process as above. The targets were changed to generate spleen thymus memory B cells and macrophage cells, which were removed from the training data. The results of the CellTypist showed that 96.75% of the thymus memory B cells were categorized into memory B cells, and 96.63% of the generated spleen macrophage cells were categorized into macrophage cells. As a comparison, the accuracy of real Thymus memory B cells and spleen macrophage cells were 96.91% and 99.53%. The box plot of marker genes of thymus memory B cells [CD79A (Garman et al. 2020), SPIB, CD19 (de Masson et al. 2014), and MS4A1 (Nieto et al. 2021)] and spleen macrophage cells [C1QB, C1QC (Missarova et al. 2021), CD68

(Brown et al. 2005), and LGMN (Zhao et al. 2020)] also indicated similar levels of expressions (Fig. 4b and c, Supplementary Figs S5 and S6). These marker genes were selected using CellMarker2.0 (Hu et al. 2023).

All the results suggested that scDiffusion could effectively generate realistic out-of-distribution cells by learning the expression patterns of known cells.

3.4 Generating intermediate cell states during cell reprogramming

We used the Gradient Interpolation strategy to generate the intermediate cell states during cell reprogramming in the Waddington-OT dataset. We trained scDiffusion on the Waddington-OT dataset, which contains MEFs with the induction of reprogramming to induced pluripotent stem cells (iPSCs). The data were across 18 days since induction with a half-day interval, and a part of the cells were induced to redifferentiate at day 8.

We first trained scDiffusion with all integer days and generated cells in the middle of two integer days. We compared the results with linear interpolation. As some of the cells were induced to redifferentiate at day 8, we interpolated cells within each treatment group separately. The interpolation weights of both methods were set to 1:1. As shown in Fig. 5a and b, scDiffusion exhibited better performance than linear interpolation in MMD metrics and LISI metrics. The mean MMD and LISI of scDiffusion are 0.3217 and 0.4488, while the result of linear interpolation is 0.5206 and 0.3217, respectively. It is worth noting that scDiffusion was not trained with the information of different treatments, but its performance was still better than linear interpolation according to the treatment information, suggesting that the diffusion model can well capture the miscellaneous distribution of cells and well fit their intermediate states.

We then chose all samples from day 0 to day 8 with the exception of day 3.5 and day 4 to train scDiffusion, and trained the classifier with the same dataset using the timestamp as the label. We then sent two conditions, day 3 and day 4.5, to the classifier and used Gradient Interpolation to generate a series of cell states between day 3 and day 4.5 in the development trajectory (Fig. 5c). The initial noise was set to be the noised latent space embeddings of day-3 cells.

We generated 10 states between day 3 and day 4.5 (Fig. 5d). We calculated the diffusion pseudotime (Haghverdi et al. 2016) of different states (Supplementary Fig. S7) and found that state 6 and state 8 are the closest to real cells of day 3.5 and day 4, respectively. These cells were previously stripped out from the training data. We compared these two states with the linear interpolation of day 3 and day 4.5 whose weight was the same as Gradient Interpolation. The MMD scores of state 6 and state 8 are 0.047 and 0.0545, while the linear interpolation's scores are 0.1289 and 0.1167. The LISI scores of these two states with real day 3.5 and day 4 were 0.64 and 0.38, while the linear interpolations were 0.20 and 0.34. These results showed that scDiffusion can generate cells that are closer to the real intermediate state.

To explore whether scDiffusion can benefit biological discovery, we examined the expression values of several key Transcription factors (TF) and marker genes across different timestamps in real and interpolated cells. As shown in Supplementary Fig. S8a, gene *Shisa8*, an early marker of successful mesenchymal-to-epithelial transition (Pei and Grishin 2012) reported in the original paper (Schiebinger et al. 2019), reached its maximum expression at day 1.5 in the scDiffusion-interpolated data, aligning with the ground truth. In contrast, when using linear interpolation, the expression of *Shisa8* reached its maximum value at day 1, missing the critical time point in the real world. Similar results were shown in the cases of the key TF *Prrx1* (Charlier et al. 2022) and the

marker gene *Fut9* (Pei and Grishin 2012) (Supplementary Fig. S8b and c). These results showed that scDiffusion could better capture the nonlinear changing profile of key genes and may thus benefit the biological discovery in the intermediate states by precise interpolation.

4 Discussion

In this article, we presented a deep generative neural network scDiffusion based on LDM and the foundation model. scDiffusion can be stably trained on large datasets to learn complex data distributions and generate realistic data. Besides, scDiffusion can generate gene expression profiles under any customized single-cell level conditions, even for rare-type cells the amount of which is limited in training samples. Furthermore, using multi-conditional generation and the Gradient Interpolation strategy, scDiffusion is uniquely capable of generating out-of-distribution data as well as intermediate states between two known cell states. To the best of our knowledge, this functionality is unique to scDiffusion.

More efforts are needed to design better evaluation metrics to accommodate the increasing complexity of generated data. A comprehensive and universally applicable evaluation metric is essential for assessing the validity of generated cells. Although there are already many metrics that measure the quality of generated data in different aspects, different metrics may be needed in different data usage scenarios. For example, the accuracy in CellTypist classification is more important in evaluating key biological characteristics, and RF AUC or LISI may play vital roles when caring the global similarity of cells. Besides, there is still no method for assessing the confidence of generated data. The current evaluation metrics rely on known data, and may not be enough for evaluating the generated unseen data. It's hard to define a universal threshold to gauge the confidence of the new data for all scenarios. Though, we found that the MMD of unseen cells shows high correlation with the MMD of known cells during training (Supplementary Fig. S9), suggesting that it is possible to estimate the confidence of unseen cells by the results of known cells.

More downstream applications may be accomplished by scDiffusion. A very natural thing is multi-omics data generation. Theoretically, scDiffusion can generate any kind of single-cell data. Besides, scDiffusion can also be used in the quality improvement of single-cell data. For instance, by learning the overall expression paradigm in clean data, scDiffusion can perform denoising operations for contaminated data. In the future, we will try to replace the classifier with more powerful tools such as CLIP (Radford et al. 2021) in the stable diffusion (Rombach et al. 2022). In this way we may use more complex conditions to control the generating

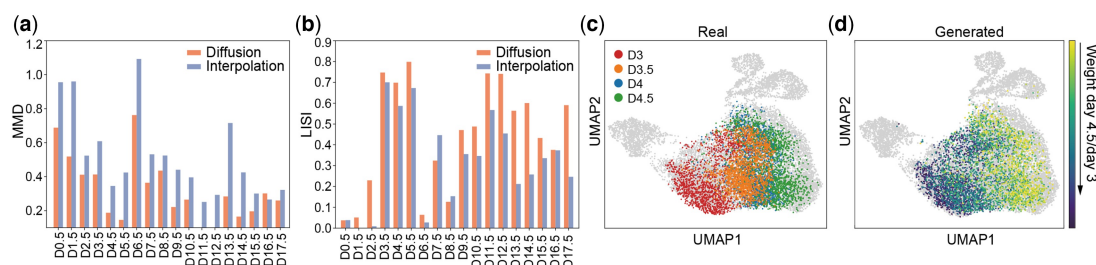


Figure 5. (a) The MMD score of different methods at different timestamps. (b) The LISI score of different methods at different timestamps. (c) UMAP of real cells. (d) UMAP of cells generated by Gradient Interpolation.

process and enable more complex tasks such as *in silico* cell perturbation, providing important help for drug selection and the control of cell state transition.

Supplementary data

Supplementary data are available at *Bioinformatics* online.

Conflict of interest

None declared.

Funding

The work was supported in part by National Key R&D Program of China [2021YFF1200901], and National Natural Science Foundation of China [62250005, 61721003, 62373210].

Data availability

The code and datasets of scDiffusion are available at <https://github.com/EperLuo/scDiffusion>. The code is also accessible via Zenodo at <https://zenodo.org/doi/10.5281/zenodo.13268742>.

References

- Baruzzo G, Patuzzi I, Di Camillo B. Sparsim single cell: a count data simulator for scRNA-seq data. *Bioinformatics* 2020;36:1468–75.
- Bhatt D, Kang B, Sawant D *et al.* STARTRAC analyses of scRNA-seq data from tumor models reveal T cell dynamics and therapeutic targets. *J Exp Med* 2021;218:20201329.
- Bian H, Chen Y, Dong X *et al.* scMulan: a multitask generative pre-trained language model for single-cell analysis. In: *International Conference on Research in Computational Molecular Biology*. Cham: Springer Nature Switzerland, 2024.
- Bond-Taylor S, Leach A, Long Y *et al.* Deep generative modelling: a comparative review of VAEs, GANs, normalizing flows, energy-based and autoregressive models. *IEEE Trans Pattern Anal Mach Intell* 2021;44:7327–47.
- Brbic M, Zitnik M, Wang S *et al.* Mars: discovering novel cell types across heterogeneous single-cell experiments. *Nat Methods* 2020;17:1200–6.
- Brown KE, Broadhurst KA, Mathahs MM *et al.* Immunodetection of aldehyde reductase in normal and diseased human liver. *Histol Histopathol* 2005;22:1133–42.
- Cao H, Tan C, Gao Z *et al.* A survey on generative diffusion models. *IEEE Trans Knowl Data Eng* 2024;36:2814–30.
- Charlier F, Weber M, Izak D *et al.* Statannotations. Zenodo 2022. <https://doi.org/10.5281/zenodo.7213391>
- Croitoru F-A, Hondru V, Ionescu RT *et al.* Diffusion models in vision: a survey. *IEEE Trans Pattern Anal Mach Intell* 2023;45:10850–69.
- Cui H, Wang C, Maan H *et al.* scGPT: toward building a foundation model for single-cell multi-omics using generative AI. *Nat Methods* 2024;21:1470–80.
- de Masson A, Le Buanec H, Bouaziz J-D. Purification and immunophenotypic characterization of human B cells with regulatory functions. In: Vitale G, Mion F (eds), *Regulatory B Cells. Methods in Molecular Biology*. New York, NY: Humana Press, 2014, Vol. 1190, 45–52.
- Dhariwal P, Nichol A. Diffusion models beat GANs on image synthesis. *Adv Neural Inf Process Syst* 2021;34:8780–94.
- Dibaenia P, Sinha S. Sergio: a single-cell expression simulator guided by gene regulatory networks. *Cell Syst* 2020;11:252–71.e11.
- Domínguez Conde C, Xu C, Jarvis LB *et al.* Cross-tissue immune cell analysis reveals tissue-specific features in humans. *Science* 2022;376:eabl5197.
- Garman L, Pelikan RC, Rasmussen A *et al.* Single cell transcriptomics implicate novel monocyte and T cell immune dysregulation in sarcoidosis. *Front Immunol* 2020;11:567342.
- Gohil SH, Iorgulescu JB, Braun DA *et al.* Applying high-dimensional single-cell technologies to the analysis of cancer immunotherapy. *Nat Rev Clin Oncol* 2021;18:244–56.
- Greene WH. Accounting for excess zeros and sample selection in Poisson and negative binomial regression models. NYU Working Paper 1994; No. EC-94-10.
- Gretton A, Borgwardt KM, Rasch MJ *et al.* A kernel two-sample test. *J Mach Learn Res* 2012;13:723–73.
- Habermann AC, Gutierrez AJ, Bui LT *et al.* Single-cell RNA sequencing reveals profibrotic roles of distinct epithelial and mesenchymal lineages in pulmonary fibrosis. *Sci Adv* 2020;6:eaba1972.
- Haghverdi L, Büttner M, Wolf FA *et al.* Diffusion pseudotime robustly reconstructs lineage branching. *Nat Methods* 2016;13:845–8.
- Haghverdi L, Lun AT, Morgan MD *et al.* Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nat Biotechnol* 2018;36:421–7.
- Hao M, Gong J, Zeng X *et al.* Large-scale foundation model on single-cell transcriptomics. *Nat Methods* 2024;21:1481–91.
- Heimberg G, Kuo T, DePianto D *et al.* Scalable querying of human cell atlases via a foundational model reveals commonalities across fibrosis-associated macrophages. *bioRxiv*, 2023, preprint: not peer reviewed. doi: <https://doi.org/10.1101/2023.07.18.549537>.
- Hilton HG, Rubinstein ND, Janki P *et al.* Single-cell transcriptomics of the naked mole-rat reveals unexpected features of mammalian immunity. *PLoS Biol* 2019;17:e3000528.
- Ho J, Jain A, Abbeel P. Denoising diffusion probabilistic models. *Adv Neural Inf Process Syst* 2020;33:6840–51.
- Hu C, Li T, Xu Y *et al.* Cellmarker 2.0: an updated database of manually curated cell markers in human/mouse and web tools based on scRNA-seq data. *Nucleic Acids Res* 2023;51:D870–6.
- Jiang P, Sinha S, Aldape K *et al.* Big data in basic and translational cancer research. *Nat Rev Cancer* 2022;22:625–39.
- Jovic D, Liang X, Zeng H *et al.* Single-cell RNA sequencing technologies and applications: a brief overview. *Clin Transl Med* 2022;12:e694.
- Ke M, Elshenawy B, Sheldon H *et al.* Single cell RNA-sequencing: a powerful yet still challenging technology to study cellular heterogeneity. *Bioessays* 2022;44:e2200084.
- Kingma DP, Welling M. Auto-encoding variational Bayes. arXiv, arXiv:1312.6114, 2013. preprint: not peer reviewed.
- Lall S, Ray S, Bandyopadhyay S. LSH-GAN enables in-silico generation of cells for small sample high dimensional scRNA-seq data. *Commun Biol* 2022;5:577.
- Li WV, Li JJ. A statistical simulator scDesign for rational scRNA-seq experimental design. *Bioinformatics* 2019;35:i41–50.
- Lindenbaum O, Stanley J, Wolf G *et al.* Geometry based data generation. In: Bengio S, Wallach H, Larochelle H *et al.* (eds), *Proceedings of the 32nd International Conference on Neural Information Processing Systems* 2018. Red Hook, NY, USA: Curran Associates Inc., 2018.
- Lopez R, Gayoso A, Yosef N. Enhancing scientific discoveries in molecular biology with deep generative models. *Mol Syst Biol* 2020;16:e9198.
- Lopez R, Regier J, Cole MB *et al.* Deep generative modeling for single-cell transcriptomics. *Nat Methods* 2018;15:1053–8.
- Lotfollahi M, Naghipourfar M, Luecken MD *et al.* Mapping single-cell data to reference atlases by transfer learning. *Nat Biotechnol* 2022;40:121–30.
- Luecken MD, Büttner M, Chaichoompu K *et al.* Benchmarking atlas-level data integration in single-cell genomics. *Nat Methods* 2022;19:41–50.
- Marouf M, Machart P, Bansal V *et al.* Realistic in silico generation and augmentation of single-cell RNA-seq data using generative adversarial networks. *Nat Commun* 2020;11:166.

- McInnes L, Healy J, Melville J. UMAP: uniform manifold approximation and projection for dimension reduction. arXiv:1802.03426, 2018, preprint: not peer reviewed.
- Missarova A, Jain J, Butler A *et al.* genebasis: an iterative approach for unsupervised selection of targeted gene panels from scRNA-seq. *Genome Biol* 2021;22:333.
- Nieto P, Elosua-Bayes M, Trincado JL *et al.* A single-cell tumor immune atlas for precision oncology. *Genome Res* 2021;31:1913–26.
- Pei J, Grishin NV. Unexpected diversity in shisa-like proteins suggests the importance of their roles as transmembrane adaptors. *Cell Signal* 2012;24:758–69.
- Qin F, Luo X, Xiao F *et al.* Scrip: an accurate simulator for single-cell RNA sequencing data. *Bioinformatics* 2022;38:1304–11.
- Radford A, Kim J W, Hallacy C *et al.* Learning transferable visual models from natural language supervision. In: *Proceedings of the International Conference on Machine Learning*. New York, NY: Online. PMLR, 2021.
- Rombach R, Blattmann A, Lorenz D *et al.* High-resolution image synthesis with latent diffusion models. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, New Orleans, LA. Piscataway, NJ: IEEE, 2022.
- Saxena D, Cao J. Generative adversarial networks (GANs) challenges, solutions, and future directions. *ACM Comput Surv* 2021;54:1–42.
- Schaum N *et al.* Single-cell transcriptomics of 20 mouse organs creates a tabula muris: the Tabula Muris Consortium. *Nature* 2018;562:367.
- Schiebinger G, Shu J, Tabaka M *et al.* Optimal-transport analysis of single-cell gene expression identifies developmental trajectories in reprogramming. *Cell* 2019;176:928–43.e22.
- Song D, Wang Q, Yan G *et al.* scDesign3 generates realistic in silico data for multimodal single-cell and spatial omics. *Nat Biotechnol* 2024;42:247–52.
- Suvà ML, Tirosh I. Single-cell RNA sequencing in cancer: lessons learned and emerging challenges. *Mol Cell* 2019;75:7–12.
- Theodoris CV, Xiao L, Chopra A *et al.* Transfer learning enables predictions in network biology. *Nature* 2023;618:616–24.
- TTS Consortium*, Jones RC, Karkanias J *et al.* The Tabula Sapiens: a multiple-organ. Single-cell transcriptomic atlas of humans. *Science* 2022;376:eabl4896.
- Xu Y, Zhang Z, You L *et al.* scIGANS: single-cell RNA-seq imputation using generative adversarial networks. *Nucleic Acids Res* 2020;48:e85.
- Yang L, Zhang Z, Song Y *et al.* Diffusion models: a comprehensive survey of methods and applications. *ACM Comput Surveys* 2022;56(4):1–39.
- Zappia L, Phipson B, Oshlack A. Splatter: simulation of single-cell RNA sequencing data. *Genome Biol* 2017;18:174.
- Zhang C, Zhang C, Zhang M *et al.* Text-to-image diffusion model in generative AI: a survey. arXiv, arXiv:2303.07909, 2023, preprint: not peer reviewed.
- Zhao J, Zhang S, Liu Y *et al.* Single-cell RNA sequencing reveals the heterogeneity of liver-resident immune cells in human. *Cell Discov* 2020;6:22.
- Zheng GXY, Terry JM, Belgrader P *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat Commun* 2017;8:14049.