Deep learning tackles single-cell analysis -

A survey of deep learning for scRNA-seq analysis

3

1

2

- 4 Mario Flores^{1§}, Zhentao Liu¹, Tinghe Zhang¹, Md Musaddaqui Hasib¹, Yu-Chiao Chiu²,
- 5 Zhenqing Ye^{2,3}, Karla Paniagua¹, Sumin Jo¹, Jianqiu Zhang¹, Shou-Jiang Gao^{4,6}, Yu-Fang
- 6 Jin¹, Yidong Chen^{2,3}§, and Yufei Huang^{5,6}§

7

- ¹Department of Electrical and Computer Engineering, the University of Texas at San Antonio, San
- 9 Antonio, TX 78249, USA
- ²Greehey Children's Cancer Research Institute, University of Texas Health San Antonio, San
- 11 Antonio, TX 78229, USA
- ³Department of Population Health Sciences, University of Texas Health San Antonio, San
- 13 Antonio, TX 78229, USA
- ⁴Department of Microbiology and Molecular Genetics, University of Pittsburgh, Pittsburgh,
- 15 Pennsylvania, PA 15232, USA
- ⁵Department of Medicine, School of Medicine, University of Pittsburgh, PA 15232, USA
- 17 ⁶UPMC Hillman Cancer Center, University of Pittsburgh, PA 15232, USA

18

19

- 20 §Correspondence should be addressed to Mario Flores (<u>mario.flores@utsa.edu</u>);
- 21 Yidong Chen (cheny8@uthscsa.edu); Yufei Huang (yuh119@pitt.edu)

22

23 **Running title**: Deep learning for single-cell RNA-seq analysis

- Mario Flores, Ph.D., is an Assistant Professor in the Department of Electrical and Computer
- 26 Engineering at the University of Texas at San Antonio. His research focuses on DNA and RNA

- 1 sequence methods, transcriptomics analysis (including scRNA-seq), epigenetics, comparative
- 2 genomics, and deep learning to study mechanisms of gene regulation.
- 3 **Zhentao Liu** is a Ph.D. student in the Department of Electrical and Computer Engineering, the
- 4 University of Texas at San Antonio. His research focuses on deep learning for cancer genomics
- 5 and drug response prediction.
- 6 Tinghe Zhang is a Ph.D. student in the Department of Electrical and Computer Engineering, the
- 7 University of Texas at San Antonio. His research focuses on deep learning for cancer genomics
- 8 and drug response prediction.
- 9 Md Musaddaqui Hasib is a Ph.D. student in the Department of Electrical and Computer
- 10 Engineering, the University of Texas at San Antonio. His research focuses on interpretable deep
- learning for cancer genomics.
- 12 **Zhenqing Ye**, Ph.D., is an Assistant Professor in the Department of Population Health Sciences
- and the bioinformatics facility manager at Greehey Children's Cancer Research Institute at the
- 14 University of Texas Health San Antonio. His research focuses on computational methods on next-
- generation sequencing and single-cell RNA-seq data analysis.
- 16 **Sumin Jo** is a Ph.D. student in the Department of Electrical and Computer Engineering, the
- 17 University of Texas at San Antonio. Her research focuses on m⁶A mRNA methylation and deep
- 18 learning for biomedical applications.
- 19 Karla Paniagua is a Ph.D. student in the Department of Electrical and Computer Engineering,
- the University of Texas at San Antonio. Her research focuses on applications of deep learning
- 21 algorithms.
- 22 **Yu-Chiao Chiu,** Ph.D., is a Postdoctoral Fellow at the Greehey Children's Cancer Research
- 23 Institute at the University of Texas Health San Antonio. His postdoctoral research is focused on
- 24 developing deep learning models for pharmacogenomic studies.

- Jianqiu Zhang, Ph.D., is an Associate Professor in the Department of Electrical and Computer
- 2 Engineering at the University of Texas at San Antonio. Her current research focuses on deep
- 3 learning for biomedical applications such as m⁶A mRNA methylation.
- 4 Shou-Jiang Gao, Ph.D., is a Professor in UPMC Hillman Cancer Center and Department of
- 5 Microbiology and Molecular Genetics, University of Pittsburgh. His current research interests
- 6 include Kaposi's sarcoma-associate herpesvirus (KSHV), AIDS-related malignancies,
- 7 translational and cancer therapeutics, and systems biology.
- 8 **Yu-Fang Jin**, Ph.D., is a Professor in the Department of Electrical and Computer Engineering at
- 9 the University of Texas at San Antonio. Her research focuses on mathematical modeling of
- 10 cellular responses in immune systems, data-driven modeling and analysis of macrophage
- activations, and deep learning applications.
- 12 **Yidong Chen**, Ph.D., is a Professor in the Department of Population Health Sciences and the
- director of Computational Biology and Bioinformatics at Greehey Children's Cancer Research
- 14 Institute at the University of Texas Health San Antonio. His research interests include
- 15 bioinformatics methods in next-generation sequencing technologies, integrative genomic data
- analysis, genetic data visualization and management, and machine learning in translational
- 17 cancer research.
- 18 **Yufei Huang**, Ph.D., is a Professor in the Department of Medicine, University of Pittsburgh
- 19 School of Medicine and Leader in Al Research at UPMC Hillman Cancer Center. His current
- 20 research interests include uncovering the functions of m⁶A mRNA methylation, cancer virology,
- 21 and medical Al & deep learning.

Abstract

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

Since its selection as the method of the year in 2013, single-cell technologies have become mature enough to provide answers to complex research questions. With the growth of single-cell profiling technologies, there has also been a significant increase in data collected from single-cell profilings, resulting in computational challenges to process these massive and complicated datasets. To address these challenges, deep learning (DL) is positioned as a competitive alternative for single-cell analyses besides the traditional machine learning approaches. Here we survey a total of 25 DL algorithms and their applicability for a specific step in the single cell RNA-seq processing pipeline. Specifically, we establish a unified mathematical representation of variational autoencoder, autoencoder, generative adversarial network, and supervised DL models, compare the training strategies and loss functions for these models, and relate the loss functions of these models to specific objectives of the data processing step. Such a presentation will allow readers to choose suitable algorithms for their particular objective at each step in the pipeline. We envision that this survey will serve as an important information portal for learning the application of DL for scRNA-seg analysis and inspire innovative uses of DL to address a broader range of new challenges in emerging multiomics and spatial single-cell sequencing.

Key points: 1

6

7

8

9

10

11

13

- 2 Single cell RNA sequencing technology generates a large collection of transcriptomic 3 profiles of up to millions of cells, enabling biological investigation of hidden expression 4 functional structures or cell types, predicting their effects or responses to treatment more 5 precisely, or utilizing subpopulations to address unanswered hypotheses.
 - Twenty-five deep learning-based approaches for single cell RNA seg data analysis are systematically reviewed in this paper according to the challenge they address and their roles in the analysis pipeline.
 - A unified mathematical description of the surveyed DL models is presented and the specific model features were discussed when reviewing each approach.
- A comprehensive summary of the evaluation metrics, comparison algorithms, and 12 datasets by each approach is presented.
- 14 **Keywords**: deep learning; single-cell RNA-seq; imputation; dimensionality reduction; clustering;
- 15 batch correction; cell type identification; functional prediction; visualization

1. Introduction

Single cell sequencing technology has been a rapidly developing area to study genomics, transcriptomics, proteomics, metabolomics, and cellular interactions at the single cell level for cell-type identification, tissue composition, and reprogramming [1, 2]. Specifically, sequencing of the transcriptome of single cells, or single-cell RNA-sequencing (scRNA-seq), has become the dominant technology in many frontier research areas such as disease progression and drug discovery [3, 4]. One particular area where scRNA-seq has made a tangible impact is cancer, where scRNA-seq is becoming a powerful tool for understanding invasion, intratumor heterogeneity, metastasis, epigenetic alterations, detecting rare cancer stem cells, and therapeutic response [5, 6]. Currently, scRNA-seg is applied to develop personalized therapeutic strategies that are potentially useful in cancer diagnosis, therapy resistance during cancer progression, and the survival of patients [5, 7]. The scRNA-seq has also been adopted to combat COVID-19 to elucidate how the innate and adaptive host immune system miscommunicates, worsening the immunopathology produced during the viral infection [8, 9].

These studies have led to a massive amount of scRNA-seq data deposited to public databases such as the 10X single-cell gene expression dataset, Human Cell Atlas, and Mouse Cell Atlas. Expressions of millions of cells from 18 species have been collected and deposited, waiting for further analysis (Single Cell Expression Atlas, EMBL-EBI, October 2021), . On the other hand, due to biological and technical factors, scRNA-seq data presents several analytical challenges related to its complex characteristics like missing expression values, high technical and biological variance, noise and sparse gene

coverage, and elusive cell identities [1]. These characteristics make it difficult to directly apply commonly used bulk RNA-seq data analysis techniques and have called for novel statistical approaches for scRNA-seq data cleaning and computational algorithms for data analysis and interpretation. To this end, specialized scRNA-seq analysis pipelines such as Seurat [10] and Scanpy [11], along with a large collection of task-specific tools, have been developed to address the intricate technical and biological complexity of scRNA-seq data.

Recently, deep learning has demonstrated its significant advantages in natural language processing and speech and facial recognition with massive data [12-14]. Such advantages have initiated the application of DL in scRNA-seq data analysis as a competitive alternative to conventional machine learning (ML) approaches for uncovering cell clustering [15, 16], cell type identification [15, 17], gene imputation [18-20], and batch correction [21] in scRNA-seq analysis. Compared to conventional ML approaches, DL is more powerful in capturing complex features of high-dimensional scRNA-seq data. It is also more versatile, where a single model can be trained to address multiple tasks or adapted and transferred to different tasks. Moreover, DL training scales more favorably with the number of cells in scRNA-seq data size, making it particularly attractive for handling the ever-increasing volume of single cell data. Indeed, the growing body of DL-based tools has demonstrated DL's exciting potential as a learning paradigm to significantly advance the tools we use to interrogate scRNA-seq data.

In this paper, we present a comprehensive review of the recent advances of DL methods for solving the challenges in scRNA-seq data analysis (Table 1) from the quality control, normalization/batch effect correction, dimensionality reduction, visualization, feature selection, and data interpretation by surveying deep learning papers published up to April 2021. In order to maintain high quality for this review, we choose not to include any (bio)archival papers, although a proportion of these manuscripts contain important new findings that would be published after completing their peer-reviewed process. Previous efforts to review the recent advances in ML methods focused on efficient integration of single cell data [22, 23]. A recent review of DL applications on single cell data has summarized 21 DL algorithms that might be deployed in single cell studies [24]. It also evaluated the clustering and data correction effect of these DL algorithms using 11 datasets.

In this review, we focus more on the DL algorithms with a much detailed explanation and comparison. Further, to better understand the relationship of each surveyed DL model with the overall scRNA-seq analysis pipeline, we organize the surveys according to the challenge they address and discuss these DL models following the analysis pipeline. A unified mathematical description of the surveyed DL models is presented and the specific model features are discussed when reviewing each method. This will also shed light on the modeling connections among the surveyed DL methods and the recognization of the uniqueness of each model. Besides the models, we also summarize the evaluation matrics used by these DL algorithms and methods that each DL algorithm was compared with. The online location of the code, the development platform, the used datasets for

each method are also cataloged to facilitate their utilization and additional effort to

2 improve them. Finally, we also created a companion online version of the paper at

3 https://huang-ai4medicine-lab.github.io/survey-of-DL-for-scRNA-seq-

4 analysis/gitbook/ book, which includes expanded discussion as well as a survey of

5 additional methods. We envision that this survey will serve as an important information

6 portal for learning the application of DL for scRNA-seq analysis and inspire innovative

use of DL to address a broader range of new challenges in emerging multi-omics and

8 spatial single-cell sequencing.

2. Overview of the scRNA-seq processing pipeline

Various scRNA-seq techniques (like SMART-seq, Drop-seq, and 10X genomics sequencing) [25, 26] are available nowadays with their sets of advantages and disadvantages. Despite the differences in the scRNA-seq techniques, the data content and processing steps of scRNA-seq data are quite standard and conventional. A typical scRNA-seq dataset consists of three files: genes quantified (gene IDs), cells quantified (cellular barcode), and a count matrix (number of cells x number of genes), irrespective of the technology or pipeline used. A series of essential steps in the scRNA-seq data processing pipeline and optional tools for each step with both ML and DL approaches are illustrated in **Fig. 1**.

With the advantage of identifying each cell and unique molecular identifiers (UMIs) for expressions of each gene in a single cell, scRNA-seq data are embedded with increased technical noise and biases [27]. **Quality control (QC)** is the first and the key step to filter out dead cells, double-cells, or cells with failed chemistry or other technical artifacts. The

- 1 most commonly adopted three QC covariates include the number of counts (count depth)
- 2 per barcode identifying each cell, the number of genes per barcode, and the fraction of
- 3 counts from mitochondrial genes per barcode [28].

- 5 **Normalization** is designed to eliminate imbalanced sampling, cell differentiation, viability,
- 6 and many other factors. Approaches tailored for scRNA-seg have been developed
- 7 including the Bayesian-based method coupled with spike-in, or BASiCS [29],
- 8 deconvolution approach, scran [30], and scTransfrom in Seurat where regularized
- 9 Negative Binomial Regression was proposed [31]. Two important steps, batch correction
- and imputation, will be carried out if required by the analysis.
- Batch Correction is a common source of technical variation in high-throughput sequencing
- experiments due to variant experimental conditions such as technicians and experimental time,
- imposing a major challenge in scRNA-seq data analysis. Batch effect correction algorithms
- include detection of mutual nearest neighbors (MNNs) [32], canonical correlation analysis
- 15 (CCA) with Seurat [33], and Harmony algorithm through cell-type representation [34].
- **Imputation** step is necessary to handle high sparsity data matrix, due to missing value or
- dropout in scRNA-seq data analysis. Several tools have been developed to "impute" zero
- values in scRNA-seq data, such as SCRABBLE [35], SAVER [36] and scImpute [37].
- 19 Dimensionality reduction and visualization are essential steps to represent
- 20 biologically meaningful variations and high dimensionality with significantly reduced
- 21 computational cost. Dimensionality reduction methods, such as principal component
- 22 analysis (PCA), are widely used in scRNA-seq data analysis to achieve that purpose.
- 23 More advanced nonlinear approaches that preserve the topological structure and avoid
- overcrowding in lower dimension representation, such as LLE [38] (used in SLICER [39]),

tSNE [40], and UMAP [41], have also been developed and adopted as a standard in

single-cell data visualization.

3

5

7

8

9

2

4 **Clustering analysis** is a key step to identify cell subpopulations or distinct cell types to

unravel the extent of heterogeneity and their associated cell-type-specific markers.

6 Unsupervised clustering is frequently used to categorize cells into clusters according to

their similarity often measured in the aforementioned dimensionality-reduced

representations. Some of popular algorithms include the community detection algorithm

Louvain [42] and Leiden [43], and data-driven dimensionality reduction followed with k-

Means cluster by SIMLR [44].

11

12

13

14

15

16

17

18

10

Feature selection is another important step in single-cell RNA-seg analysis to select a

subset of genes, or features, for cell-type identification and functional enrichment of each

cluster. This step is achieved by differential expression analysis designed for scRNA-seq,

such as MAST that used linear model fitting and likelihood ratio testing [45]; SCDE that

adopted a Bayesian approach with a Negative Binomial model for gene expression and

Poisson process for dropouts [46], or DEsingle that utilized a Zero-Inflated Negative

Binomial model to estimate the dropouts [47].

19

20

21

Besides these key steps, downstream analyses include cell type identification,

coexpression analysis, prediction of perturbation response, where DL has also been

22 applied. Other advanced analyses including trajectory inference and velocity and

1 pseudotime analysis are not discussed here because most of the approaches on these

topics are non-DL based.

3

2

- 4 3. Overview of common deep learning models for scRNA-seq analysis
- 5 We start our review by introducing the general formulations of widely used deep learning
- 6 models. As most of the tasks including batch correction, dimensionality reduction,
- 7 imputation, and clustering are unsupervised learning tasks, we will give special attention
- 8 to unsupervised models including variational autoencoder (VAE), the autoencoder (AE),
- 9 or generative adversarial networks (GAN). We will also discuss the general supervised
- and transfer learning formulations, which find their applications in cell type predictions
- and functional studies. We will discuss these models in the context of scRNA-seq,
- detailing the different features and training strategies of each model and bringing attention
- to their uniqueness.

14 15

3.1. Variational Autoencoder

- Let x_n represent a $G \times 1$ vector of expression levels (UMI counts or normalized, log-
- 17 transformed expression) of G genes in cell n, where $p(x_{gn}|v_{gn},\alpha_{gn})$ follows some
- distribution (e.g., zero-inflated negative binomial (ZINB) or Gaussian), where v_{gn} and α_{gn}
- are distribution parameters (e.g., mean, variance, or dispersion) (Fig. 2A). We consider
- v_{gn} to be of particular interest (e.g., the mean counts) and is thus further modeled by a
- decoder neural network D_{θ} (**Fig. 2A**) as

$$\mathbf{v}_n = D_{\boldsymbol{\theta}}(\mathbf{z}_n, s_n), \tag{1}$$

- 22 where the gth element of \mathbf{v}_n is \mathbf{v}_{gn} and $\boldsymbol{\theta}$ is a vector of decoder weights, $\mathbf{z}_n \in \mathbb{R}^d$
- 23 represents a latent representation of gene expression and is used for visualization and

- clustering and s_n is an observed variable (e.g., the batch ID). For VAE, z_n is commonly
- assumed to follow a multivariate standard Normal prior, i.e., $p(\mathbf{z}_n) = \mathcal{N}(0, \mathbf{I}_d)$ with \mathbf{I}_d
- being a $d \times d$ identity matrix. Further, α_{gn} of $p(x_{gn} | \nu_{gn}, \alpha_{gn})$ is a nuisance parameter,
- 4 which has a prior distribution $p(\alpha_{gn})$ and can be either estimated or marginalized in
- 5 variational inference. Now define $\Theta = \{\theta, \alpha_{ng} \ \forall n, g\}$. Then, $p(x_{gn} | v_{gn}, \alpha_{gn})$ and (1)
- 6 together define the likelihood $p(\mathbf{x}_n|\mathbf{z}_n, s_n, \mathbf{\Theta})$.

8 The goal of training is to compute the maximum likelihood estimate of **9**

$$\widehat{\mathbf{\Theta}}_{ML} = \operatorname{argmax}_{\mathbf{\Theta}} \ \sum_{n=1}^{N} \log p(\mathbf{x}_n | \mathbf{s}_n, \mathbf{\Theta}) \approx \operatorname{argmax}_{\mathbf{\Theta}} \ \sum_{n=1}^{N} \mathcal{L}(\mathbf{\Theta}),$$
 (2)

9 where $\mathcal{L}(\mathbf{\Theta})$ is the evidence lower bound (ELBO),

$$\mathcal{L}(\mathbf{\Theta}) = \mathbf{E}_{q(\mathbf{Z}_n | \mathbf{X}_n, S_n, \mathbf{\Theta})}[\log p(\mathbf{X}_n | \mathbf{Z}_n, S_n, \mathbf{\Theta})] - D_{KL}[q(\mathbf{Z}_n | \mathbf{X}_n, S_n, \mathbf{\Theta}) || p(\mathbf{Z}_n)], \tag{3}$$

and $q(\mathbf{z}_n|\mathbf{x}_n, s_n)$ is an approximate to $p(\mathbf{z}_n|\mathbf{x}_n, s_n)$ and assumed as

$$q(\mathbf{z}_n|\mathbf{x}_n,s_n) = \mathcal{N}\left(\boldsymbol{\mu}_{z_n},diag(\boldsymbol{\sigma}_{z_n}^2)\right),\tag{4}$$

with $\{\mu_{z_n}, \sigma_{Z_n}^2\}$ given by an encoder network E_{ϕ} (Fig. 2A) as

$$\left\{\boldsymbol{\mu}_{z_n}, \boldsymbol{\sigma}_{z_n}^2\right\} = E_{\boldsymbol{\phi}}(\boldsymbol{x}_n, s_n), \tag{5}$$

- where ϕ is the weights vector. Now, $\Theta = \{\theta, \phi, \alpha_{ng} \, \forall n, g\}$ and equation (2) is solved by
- 13 the stochastic gradient descent approach while a model is trained.
- All the surveyed papers that deploy VAE follow this general modeling process.
- 15 However, a more general formulation has a loss function defined as

$$L(\mathbf{\Theta}) = -\mathcal{L}(\mathbf{\Theta}) + \sum_{k=1}^{K} \lambda_k L_k(\mathbf{\Theta}), \tag{6}$$

- where $L_k \forall k=1,...,K$ are losses for different functions (clustering, cell type prediction,
- etc) and λ_k s are the Lagrange multipliers. With this general formulation, for each paper,
- we examined the specific choices of data distribution $p(x_{gn} | v_{gn}, \alpha_{gn})$ that define $\mathcal{L}(\Theta)$,

- different L_k designed for specific functions, and how the decoder and encoder were
- 2 applied to model different aspects of scRNA-seq data.

4

5 3.2. Autoencoders

- 6 AEs learn the low dimensional latent representation $\mathbf{z}_n \in \mathbb{R}^d$ of expression \mathbf{x}_n . The AE
- 7 includes an encoder E_{ϕ} and a decoder D_{θ} (Fig. 2B) such that

$$\mathbf{z}_n = E_{\boldsymbol{\phi}}(\mathbf{x}_n); \ \widehat{\mathbf{x}}_n = D_{\boldsymbol{\theta}}(\mathbf{z}_n), \tag{7}$$

- where $\mathbf{\Theta} = \{ \pmb{\theta}, \pmb{\phi} \}$ are encoder and decoder weight parameters and $\widehat{\pmb{x}}_n$ defines the
- 9 parameters (e.g. mean) of the likelihood $p(x_n|\Theta)$ (Fig. 2B) and is often considered as
- imputed and denoised expressions. Additional design can be included in an AE model
- 11 for batch correction, clustering, and other objectives.
- The training of an AE model is generally carried out by stochastic gradient descent
- algorithms to minimize the loss similar to Eq. (6) except $\mathcal{L}(\mathbf{\Theta}) = -\log p(\mathbf{x}_n|\mathbf{\Theta})$. When
- 14 $p(x_n|\Theta)$ is the Gaussian, $\mathcal{L}(\Theta)$ becomes the mean square error (MSE) loss

$$\mathcal{L}(\mathbf{\Theta}) = \sum_{n=1}^{N} ||\boldsymbol{x}_n - \widehat{\boldsymbol{x}}_n||_2^2.$$
 (8)

- 15 Because different AE models differ in their AE architectures and loss functions, we will
- discuss the specific architecture and loss functions for each reviewed DL model in Section
- 17 4.

18 19

3.3. Generative adversarial networks

- 20 GANs have been used for imputation, data generation, and augmentation of the scRNA-
- seq analysis. Without loss of generality, the GAN, when applied to scRNA-seq, is designed
- 22 to learn how to generate gene expression profiles from p_x , the distribution of x_n . The

vanilla GAN consists of two deep neural networks [48]. The first network is the generator $G_{\theta}(\mathbf{z}_n, y_n)$ with parameter θ , a noise vector \mathbf{z}_n from the distribution p_z and a class label y (e.g. cell type) and is trained to generate x_f , a "fake" gene expression (Fig. **2C**). The second network is the discriminator network D_{ϕ_D} with parameters ϕ_D , trained to distinguish the "real" x from fake x_f (Fig. 2C). Both networks, G_{θ} and D_{ϕ_D} are trained to outplay each other, resulting in a minimax game, in which G_{θ} is forced by $D_{\phi_{\mathbb{D}}}$ to produce better samples, which, when converge, can fool the discriminator $D_{m{\phi}_{\mathbb{D}}}$, thus becoming samples from $p_{\scriptscriptstyle \mathcal{X}}.$ The vanilla GAN suffers heavily from training instability and mode collapsing[49]. To that end, Wasserstein GAN (WGAN) [49] was developed with the WGAN loss [50]:

$$L(\boldsymbol{\theta}) = \max_{\emptyset_{D}} \sum_{n=1}^{N} D_{\emptyset_{D}}(\boldsymbol{x}_{n}) - \sum_{n=1}^{N} D_{\emptyset_{D}}(G_{\theta}(\boldsymbol{z}_{n}, y_{n})).$$
(9)

Additional terms can also be added to equation (9) to constrain the functions of the generator. Training based on the WGAN loss in Eq. (9) amounts to a min-max optimization, which iterates between the discriminator and the generator, where each optimization is achieved by a stochastic gradient descent algorithm through back-propagation. The WGAN requires $D_{\phi_{\mathbb{D}}}$ to be K-Lipschitz continuous [50], which can be satisfied by adding the gradient penalty to the WGAN loss [49]. Once the training is done, the generator G_{ϕ_G} can be used to generate gene expression profiles of new cells.

3.4. Supervised deep learning models

Supervised deep learning models, including deep neural networks (DNN), convolutional neural network (CNN), and capsule networks (CapsNet), have been used for cell type

- identifications [51-53] and functional predictions [54]. The general supervised deep
- learning model F takes x_n as an input and outputs $p(y_n|x_n)$, the probability of phenotype
- 3 label y_n (e.g. a cell type) as

$$p(y_n|\mathbf{x}_n) = F(\mathbf{x}_n), \tag{10}$$

- 4 where F can be DNN, CNN, or CapsNet. We omit the discussion of DNN and CNN as they
- 5 are widely used in different applications and there are many excellent surveys on them
- 6 [55]. We will focus our discussion on CasNet next.
- 7 A CasNet takes an expression x_n to first form a feature extraction network (consisting of
- 8 L parallel single-layer neural networks) followed by a classification capsule network. Each
- 9 of the *L* parallel feature extraction layers generates a primary capsule $\mathbf{u}_l \in \mathbb{R}^{d_p}$ as

$$\mathbf{u}_l = ReLU(\mathbf{W}_{P,l}\mathbf{x}_n) \,\forall l = 1, \dots, L, \tag{11}$$

- where $\mathbf{W}_{P,l} \in \mathbb{R}^{d_p \times G}$ is the weight matrix. Then, the primary capsules are fed into the
- capsule network to compute K label capsules $v_k \in \mathbb{R}^{d_t}$, one for each label, as

$$\boldsymbol{v}_k = squash\left(\sum_{l=1}^{L} c_{kl} \boldsymbol{W}_{kl} \boldsymbol{u}_l\right) \, \forall k = 1, ..., K,$$
 (12)

- where *squash* is the squashing function [56] to normalize the magnitude of its input vector
- to be less than one, W_{kl} is another trainable weight matrix, and $c_{kl} \forall l = 1, ..., L$, are the
- coupling coefficients that represent the probability distribution of each primary capsule's
- impact on the predicted label k. Parameters c_{kl} are not trained but computed through the
- dynamic routing process proposed in the original capsule networks [52]. The magnitude
- of each capsule v_k represents the probability of predicting label k for input x_n . Once
- trained, the important primary capsules for each label and then the most significant genes

- 1 for each important primary capsule can be used to interpret biological functions associated
- with the prediction.
- 3 The training of the supervised models for classification overwhelmingly minimizes the
- 4 cross-entropy loss by stochastic gradient descent computed by a back-propagation
- 5 algorithm.

7

4. Survey of deep learning models for scRNA-seq analysis

- 8 In this section, we survey applications of DL models for scRNA-seq analysis. To better
- 9 understand the relationship between the problems that each surveyed work addresses
- and the key challenges in the general scRNA-seq processing pipeline, we divide the
- survey into sections according to steps in the scRNA-seq processing pipeline illustrated
- in **Fig. 1**. For each DL model, we present the model details under the general model
- 13 framework introduced in Section 3 and discuss the specific loss functions. We also survey
- 14 the evaluation metrics and summarize the evaluation results. To facilitate cross-
- references of the information, we summarized all algorithms reviewed in this section in
- 16 **Table 1** and tabulated the datasets and evaluation metrics used in each paper in **Tables**
- 17 **2 & 3.** We also listed in **Fig. 3** all other algorithms against which each surveyed method
- evaluated, highlighting the extensiveness that these algorithms were assessed for their
- 19 performance.

20

4.1. Imputation

2223

21

4.1.1. DCA: deep count autoencoder

- DCA [18] is an AE for imputation (Figs. 2B, 4B) and has been integrated into the Scanpy
- 2 framework.
- 3 <u>Model.</u> DCA models UMI counts with missing values using the ZINB distribution

$$p(x_{gn}|\Theta) = \pi_{gn}\delta(0) + (1 - \pi_{gn})NB(\nu_{gn}, \alpha_{gn}), \text{ for } g = 1, ... G; n = 1, ... N,$$
 (13)

- 4 where $\delta(\cdot)$ is a Dirac delta function, $NB(\cdot, \cdot)$ denotes the negative binomial distribution,
- 5 and π_{gn} , ν_{gn} , α_{gn} , representing dropout rate, mean, and dispersion, respectively, are
- functions of the output (\hat{x}_n) of the decoder in the DCA as follows,

$$\pi_n = sigmoid(\mathbf{W}_n \widehat{\mathbf{x}}_n); \ \mathbf{v}_n = \exp(\mathbf{W}_v \widehat{\mathbf{x}}_n); \ \mathbf{\alpha}_n = \exp(\mathbf{W}_\alpha \widehat{\mathbf{x}}_n),$$
 (14)

- 7 where W_{π} , W_{v} , and W_{α} are additional weights to be estimated. The DCA encoder and
- 8 decoder follow the general AE formulation as in Eq. (7) but the encoder takes the
- 9 normalized, log-transformed expression as input. To train the model, DCA uses a
- 10 constrained log-likelihood as the loss function

$$L(\mathbf{\Theta}) = \sum_{n=1}^{N} \sum_{g=1}^{G} \left(-logp(x_{gn}|\mathbf{\Theta}) + \lambda \pi_{gn}^{2} \right), \tag{15}$$

- with $\Theta = \{ \theta, \phi, W_{\pi}, W_{v}, W_{\alpha} \}$. Once the DCA is trained, the mean counts v_n are used as
- the denoised and imputed counts for cell n.
- 13 Results. For evaluation, DCA was compared to other methods using simulated data
- 14 (using Splatter R package), and real bulk transcriptomics data from a developmental C.
- 15 elegans time-course experiment was used with added simulating single-cell specific
- 16 noise. Gene expression was measured from 206 developmentally synchronized young
- 17 adults over a twelve-hour period (C. elegans). Single-cell specific noise was added in
- 18 silico by genewise subtracting values drawn from the exponential distribution such that

- 1 80% of values were zeros. The paper analyzed the Bulk contains less noise than single-
- 2 cell transcriptomics data and can thus aid in evaluating single-cell denoising methods by
- 3 providing a good ground truth model. The authors also did a comparison of other methods
- 4 including SAVER [36], scImpute [37], and MAGIC[57]. DCA denoising recovered original
- 5 time-course gene expression pattern while removing single-cell specific noise. Overall,
- 6 DCA demonstrated the strongest recovery of the top 500 genes most strongly associated
- 7 with development in the original data without noise; DCA was shown to outperform other
- 8 existing methods in capturing cell population structure in real data using PBMC, CITE-
- 9 seq, runtime scales linearly with the number of cells.

4.1.2. SAVER-X: single-cell analysis via expression recovery harnessing external data

- 12 SAVER-X [58] is an AE model (Figs. 2B, 4B) developed to denoise and impute scRNA-
- seg data with transfer learning from other data resources.
- 14 Model. SAVER-X decomposes the variation in the observed counts x_n with missing
- 15 values into three components: i) predictable structured component representing the
- shared variation across genes, ii) unpredictable cell-level biological variation and gene-
- specific dispersions, and iii) technical noise. Specifically, x_{gn} is modeled as a Poisson-
- 18 Gamma hierarchical model,

$$p(x_{gn}|\Theta) = Poisson(l_n x'_{gn}), \qquad p(x'_{gn}|\nu_{gn},\alpha_g) = Gamma(\nu_{gn},\alpha_g \nu_{gn}^2), \tag{16}$$

- where l_n is the sequencing depth of cell n, v_{gn} is the mean, and α_g is the dispersion. This
- 20 Poisson-Gamma mixture is an equivalent expression to the NB distribution and thus, the
- 21 ZINB distribution as Eq. (13) is adopted to model missing values.

The loss is similar to Eq. (15). However, v_{qn} is initially learned by an AE pre-trained using external datasets from an identical or similar tissue and then transferred to x_n to be denoised. Such transfer learning can be applied to data between species (e.g., human and mouse in the study), cell types, batches, and single-cell profiling technologies. After v_{qn} is inferred, SAVER-X generates the final denoised data \hat{x}_{qn} by an empirical Bayesian shrinkage. Results. SAVER-X was applied to multiple human single-cell datasets of different scenarios: i) T-cell subtypes, ii) a cell type (CD4+ regulatory T cells) that was absent from the pretraining dataset, iii) gene-protein correlations of CITE-seq data, and iv) immune cells of primary breast cancer samples with a pretraining on normal immune cells. SAVER-X with pretraining on HCA and/or PBMCs outperformed the same model without pretraining and other denoising methods, including DCA [28], scVI[17], scImpute [37], and MAGIC [57]. The model achieved promising results even for genes with very low UMI counts. SAVER-X was also applied for a cross-species study in which the model was pretrained on a human or mouse dataset and transferred to denoise another. The results demonstrated the merit of transferring public data resources to denoise in-house scRNAseg data even when the study species, cell types, or single-cell profiling technologies are

19 20

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

different.

4.1.3. DeepImpute: Deep neural network Imputation

- 21 DeepImpute [20] imputes genes in a divide-and-conquer approach, using a bank of DNN
- 22 models (**Fig. 4A**) with 512 output, each to predict gene expression levels of a cell.
- 23 <u>Model.</u> For each dataset, DeepImpute selects to impute a list of genes or highly variable
- 24 genes (variance over mean ratio, default = 0.5). Each sub-neural network aims to

- 1 understand the relationship between the input genes and a subset of target genes. Genes
- 2 are first divided into *N* random subsets of 512 target genes. For each subset, a two-layer
- 3 DNN is trained where the input includes genes that are among the top 5 best-correlated
- 4 genes to target genes but not part of the target genes in the subset. The loss is defined
- 5 as the weighted MSE

$$\mathcal{L}(\mathbf{\Theta}) = \sum \mathbf{x}_n (\mathbf{x}_n - \widehat{\mathbf{x}}_n)^2, \tag{17}$$

- 6 which gives higher weights to genes with higher expression values.
- 7 Result. DeepImpute had the highest overall accuracy and offered shorter computation
- 8 time with less demand on computer memory than other methods like MAGIC, DrImpute,
- 9 ScImpute, SAVER, VIPER, and DCA. Using simulated and experimental datasets (**Table**
- 10 2), DeepImpute showed benefits in improving clustering results and identifying
- 11 significantly differentially expressed genes. DeepImpute and DCA, show overall
- 12 advantages over other methods and between which DeepImpute performs even better.
- 13 The properties of DeepImpute contribute to its superior performance include 1) a divide-
- 14 and-conquer approach which contrary to an autoencoder as implemented in DCA,
- resulting in a lower complexity in each sub-model and stabilizing neural networks, and 2)
- 16 the subnetworks are trained without using the target genes as the input which reduces
- overfitting while enforcing the network to understand true relationships between genes.
- 19 4.1.4. LATE: Learning with AuToEncoder
- 20 LATE [59] is an AE (**Figs. 2B, 4B**) whose encoder takes the log-transformed expression
- as input.

- 1 *Model.* LATE sets zeros for all missing values and generates the imputed expressions.
- 2 LATE minimizes the MSE loss as Eq. (8). One issue is that some zeros could be real and
- 3 reflect the actual lack of expressions.
- 4 Result. Using synthetic data generated from pre-imputed data followed by random
- 5 dropout selection at different degrees, LATE outperforms other existing methods like
- 6 MAGIC, SAVER, DCA, scVI, particularly when the ground truth contains only a few or no
- 7 zeros. However, when the data contain many zero expression values, DCA achieved a
- 8 lower MSE than LATE, although LATE still has a smaller MSE than scVI. This result
- 9 suggests that DCA likely does a better job identifying true zero expressions, partly
- 10 because LATE does not make assumptions on the statistical distributions of the single-
- cell data that potentially have inflated zero counts.

12 **4.1.5.** scGMAI

- 13 Technically, scGMAI [60] is a model for clustering but it includes an AE (**Figs. 2B, 4B**) in
- the first step to combat dropout.
- 15 <u>Model.</u> To impute the missing values, scGMAI applies an AE like LATE to reconstruct
- log-transformed expressions with dropout but chooses a smoother Softplus activation
- 17 function instead. The MSE loss as in Eq. (8) is adopted.
- After imputation, scGMAI uses fast independent component analysis (ICA) on the
- 19 AE reconstructed expressions to reduce the dimension and then applies a Gaussian
- 20 mixture model on the ICA reduced data to perform the clustering.
- 21 Results. To assess the performance, the AE in scGMAI was replaced by five other
- imputation methods including SAVER [36], MAGIC [57], DCA [28], scImpute [37], and
- 23 CIDR[61]. A scGMAI implementation without AE was also compared. Seventeen scRNA-

- seq data (part of them are listed in **Tables 2b & c** as marked) were used to evaluate cell
- 2 clustering performances. The results indicated that the AEs significantly improved the
- 3 clustering performance in eight of seventeen scRNA-seq datasets.

5 **4.1.6.** scIGANs

- 6 Imputation approaches based on information from cells with similar expressions suffer
- 7 from oversmoothing, especially for rare cell types. sclGANs [19] is a GAN-based
- 8 imputation algorithm (Figs. 2C, 4E), which overcomes this problem by training a GAN
- 9 model to generate samples with imputed expressions.
- 10 <u>Model.</u> scIGAN takes the image-like reshaped gene expression data x_n as input. The
- model follows a BEGAN [62] framework, which replaces the GAN discriminator D with a
- 12 function R_{ϕ_R} to compute the reconstruction MSE. Then, the Wasserstein distance loss
- between the reconstruction errors between the real and generated samples are computed

$$L(\boldsymbol{\theta}, \boldsymbol{\Phi}) = \max_{\emptyset_R} \sum_{n=1}^N R_{\emptyset_R}(\boldsymbol{x}_n) - \sum_{n=1}^N R_{\emptyset_R}(G_{\theta}(E_{\Phi}(\boldsymbol{x}_n), y),$$
(18)

- 14 This framework forces the model to meet two computing objectives, i.e. reconstructing the
- 15 real samples and discriminating between real and generated samples. Proportional
- 16 Control Theory was applied to balance these two goals during the training [63].
- 17 After training, the decoder G_{θ} is used to generate new samples of a specific cell
- type. Then, the k-nearest neighbors (KNN) approach is applied to the real and generated
- samples to impute the real samples' missing expressions.
- 20 Results. scIGANs was first tested on simulated samples with different dropout rates.
- 21 Performance of rescuing the correct clusters was compared with 11 existing imputation
- 22 approaches including DCA, DeepImpute, SAVER, scImpute, MAGIC, etc. scIGANs

reported the best performance for all metrics. scIGAN was next evaluated for its ability to correctly cluster cell types on the Human brain scRNA-seq data, which showed superior performance than existing methods again. scIGANs was then evaluated for identifying cell-cycle states using scRNA-seq datasets from mouse embryonic stem cells. The results showed that scIGANs outperformed competing existing approaches for recovering subcellular states of cell cycle dynamics. scIGANs were further shown to improve the identification of differentially expressed genes and enhance the inference of cellular trajectory using time-course scRNA-seq data from the differentiation from H1 ESC to definitive endoderm cells (DEC). Finally, scIGAN was also shown to scale to scRNA-seq methods and data sizes.

11 12

10

1

2

3

4

5

6

7

8

9

4.2. Batch effect correction

13

14 4.2.1. BERMUDA: Batch Effect ReMoval Using Deep Autoencoders

- BERMUDA [64] deploys a transfer-learning method (**Figs. 2B, 4B**) to remove the batch effect. It performs correction to the shared cell clusters among batches and therefore preserves batch-specific cell populations.
- 18 <u>Model.</u> BERMUDA is an AE that takes normalized, log-transformed expression as input.
- 19 Its consists of two parts as

$$L(\mathbf{\Theta}) = \mathcal{L}(\mathbf{\Theta}) + \lambda L_{MMD}(\mathbf{\Theta}), \tag{19}$$

- where $\mathcal{L}(\mathbf{\Theta})$ is the MSE loss and L_{MMD} is the maximum mean discrepancy (MMD) [65] loss that measures the differences in distributions between pairs of similar cell clusters
- shared among batches as:

$$L_{MMD}(\mathbf{\Theta}) = \sum_{i_a, i_b, j_a, j_b} M_{i_a, j_a, i_b, j_b} MMD(\mathbf{z}_{i_a, j_a}, \mathbf{z}_{i_b, j_b}), \tag{20}$$

- where $z_{i,j}$ is the latent variable of $x_{i,j}$, the input expression of a cell from cluster j of batch
- 2 i, M_{i_a,j_a,i_b,j_b} is 1 if cluster i_a of batch j_a and cluster i_b of batch j_b are determined to be
- 3 similar by MetaNeighbor [66] and 0, otherwise. The MMD equals zero when the
- 4 underlying distributions of the observed samples are the same.
- 5 Results. BERMUDA was shown to outperform other methods like mnnCorrect [32],
- 6 BBKNN[67], Seurat [10], and scVI [17] in removing batch effects on simulated and human
- 7 pancreas data while preserving batch-specific biological signals. BERMUDA provides
- 8 several improvements compared to existing methods: 1) capable of removing batch
- 9 effects even when the cell population compositions across different batches are vastly
- different; and 2) preserving batch-specific biological signals through transfer-learning
- which enables discovering new information that might be hard to extract by analyzing
- 12 each batch individually.

4.2.2. DESC: batch correction based on clustering

- DESC [68] is an AE model (Figs. 2B, 4B) that removes batch effect through clustering
- with the hypothesis that batch differences in expressions are smaller than true biological
- 17 variations between cell types, and, therefore, properly performing clustering on cells
- across multiple batches can remove batch effects without the need to define batches
- 19 explicitly.
- 20 <u>Model.</u> DESC has a conventional AE architecture. Its encoder takes normalized, log-
- 21 transformed expression and uses decoder output, \widehat{x}_n as the reconstructed gene
- expression, which is equivalent to a Gaussian data distribution with \hat{x}_n being the mean.
- The loss function is similar to Eq. (19) and except that the second loss L_c is the clustering
- loss that regularizes the learned feature representations to form clusters as in the deep

embedded clustering [69]. The model is first trained to minimize $\mathcal{L}(\mathbf{\Theta})$ only to obtain the initial weights before minimizing the combined loss. After the training, each cell is

3 assigned with a cluster ID.

4 Results. DESC was applied to the macaque retina dataset, which includes animal level,

5 region level, and sample-level batch effects. The results showed that DESC is effective

6 in removing the batch effect, whereas CCA [33], MNN [32], Seurat 3.0 [10], scVI [17],

7 BERMUDA [64], and scanorama [70] were all sensitive to batch definitions. DESC was

then applied to human pancreas datasets to test its ability to remove batch effects from

multiple scRNA-seq platforms and yielded the highest ARI among the comparing

approaches mentioned above. When applied to human PBMC data with interferon-beta

stimulation, where biological variations are compounded by batch effect, DESC was

shown to be the best in removing batch effect while preserving biological variations.

13 DESC was also shown to remove batch effect for the monocytes and mouse bone marrow

data and DESC was shown to preserve the pseudotemporal structure. Finally, DESC

scales linearly with the number of cells, and its running time is not affected by the

increasing number of batches.

17 18

19

21

16

8

9

10

11

12

14

15

4.2.3. iMAP: Integration of Multiple single-cell datasets by Adversarial Paired-style transfer networks

iMAP [71] combines AE (Figs. 2B, 4B) and GAN (Figs. 2C, 4E) for batch effect removal.

It is designed to remove batch biases while preserving dataset-specific biological

22 variations.

- 1 Model. iMAP consists of two processing stages, each including a separate DL model. In
- the first stage, a special AE, whose decoder combines the output of two separate
- 3 decoders D_{θ_1} and D_{θ_2} , is trained such that

$$\mathbf{z}_n = E_{\phi}(\mathbf{x}_n); \ \hat{\mathbf{x}}_n = D_{\theta}(\mathbf{z}_n, s_n) = ReLu(D_{\theta_1}(s_n) + D_{\theta_2}(\mathbf{z}_n, s_n)), \tag{21}$$

- 4 where s_n is the one-hot encoded batch number of cell n. D_{θ_1} can be understood as
- 5 decoding the batch noise, whereas D_{θ_2} reconstructs batch-removed expression from the
- 6 latent variable z_n . The training minimizes the loss in Eq. (19) except the 2^{nd} loss is the
- 7 content loss

$$L_t(\mathbf{\Theta}) = \sum_{n=1}^{N} \left\| \mathbf{z}_n - E_{\phi} \left(D_{\theta}(\mathbf{z}_n, \tilde{s}_n) \right) \right\|_2^2, \tag{22}$$

- 8 where \tilde{s}_n is a random batch number. Minimizing $L_t(\Theta)$ further ensures the reconstructed
- 9 expression \hat{x}_n would be batch agnostic and has the same content as x_n .
- However, due to the limitation of AE, this step is still insufficient for batch removal.
- 11 Therefore, a second stage is included to apply a GAN model to make expression
- distributions of the shared cell type across different baches indistinguishable. To identified
- the shared cell types, a mutual nearest neighbors (MNN) strategy adapted from [32] was
- developed to identify MNN pairs across batches using batch effect independent \mathbf{z}_n as
- opposed to x_n . Then, a mapping generator G_{θ_G} is trained using MNN pairs based on GAN
- such that $x_n^{(A)} = G_{\theta_G}(x_n^{(S)})$, where $x_n^{(S)}$ and $x_n^{(A)}$ are the MNN pairs from batch S and an
- anchor batch A. The WGAN-GP loss as in Eq. (9) was adopted for the GAN training. After
- training, G_{θ_G} is applied to all cells of a batch to generate batch-corrected expression.
- 19 Results: iMAP was first tested on benchmark datasets from human dendritic cells and
- 20 Jurkat and 293T cell lines and then human pancreas datasets from five different

platforms. All the datasets contain both batch-specific cells and batch-shared cell types.

2 iMAP was shown to separate the batch-specific cell types but mix batch shared cell types

and outperformed 9 other existing batch correction methods including Harmony, scVI,

fastMNN, Seurat, etc. iMAP was then applied to the large-scale Tabula Muris datasets

containing over 100K cells sequenced from two platforms. iMAP could not only reliably

integrate cells from the same tissues but identify cells from platform-specific tissues.

Finally, iMAP was applied to datasets of tumor-infiltrating immune cells and shown to

reduce the dropout ratio and the percentage of ribosomal genes and non-coding RNAs,

thus improving detection of rare cell types and ligand-receptor interactions. iMAP scales

with the number of cells, showing minimal time cost increase after the number of cells

exceeds thousands. Its performance is also robust against model hyperparameters.

4.3. Dimensionality reduction, latent representation, clustering, and data augmentation

Dimensionality reduction is indispensable for many type of scRNA-seq data analysis, considering the limited number of cell types in each biospecimen. Furthermore, biological processes of interests often involve the complex coordination of many genes, therefore, latent representation which capture biological variation in reduced dimensions are useful in interpreting experiment conditions and cell heterogeneity. Both AE- and VAE-based are capable of learning latent representations. VAE-based models have the benefit of regularity of the latent space and generative factors. The GAN-based models can produce augmented data that may in return to enhance the clustering, e.g., due to low representation of certain cell types.

4.3.1. Dimensionality reduction by AEs with gene-interaction constrained architecture

1 This study [72] considers AEs (Figs. 2B, 4B) for learning the low-dimensional representation and specifically explores the benefit of incorporating prior biological 2 knowledge of gene-gene interactions to regularize the AE network architecture. 3 4 Model. Several AE models with single or two hidden layers that incorporate gene 5 interactions reflecting transcription factor (TF) regulations and protein-protein interactions 6 (PPIs) are implemented. The models take normalized, log-transformed expressions and 7 follow the general AE structure, including dimension-reducing and reconstructing layers, 8 but the network architectures are not symmetrical. Specifically, gene interactions are 9 incorporated such that each node of the first hidden layer represented a TF or a protein 10 in the PPI; only genes that are targeted by TFs or involved in the PPI were connected to the node. Thus, the corresponding weights of $E_{\pmb{\phi}}$ and $D_{\pmb{\theta}}$ are set to be trainable and 11 12 otherwise fixed at zero throughout the training process. Both unsupervised (AE-like) and 13 supervised (cell-type label) learning were studied. 14 Results. Regularizing encoder connections with TF and PPI information considerably 15 reduced the model complexity by almost 90% (7.5-7.6M to 1.0-1.1M). The clusters formed on the data representations learned from the models with or without TF and PPI 16 17 information were compared to those from PCA, NMF, independent component analysis 18 (ICA), t-SNE, and SIMLR [44]. The model with TF/PPI information and 2 hidden layers 19 achieved the best performance by five of the six measures and the best average 20 performance. In terms of the cell-type retrieval of single cells, the encoder models with 21 and without TF/PPI information achieved the best performance in 4 and 3 cell types, 22 respectively. PCA yielded the best performance in only 2 cell types. The DNN model with TF/PPI information and 2 hidden layers again achieved the best average performance 23

- across all cell types. In summary, this study demonstrated a biologically meaningful way
- 2 to regularize AEs by the prior biological knowledge for learning the representation of
- 3 scRNA-seq data for cell clustering and retrieval.

- 5 4.3.2. Dhaka: a VAE-based dimension reduction model.
- 6 Dhaka [73] was proposed to reduce the dimension of scRNA-seq data for efficient
- 7 stratification of tumor subpopulations.
- 8 Model. Dhaka adopts a general VAE formulation (Figs. 2A, 4C). It takes the normalized,
- 9 log-transformed expressions of a cell as input and outputs the low-dimensional
- 10 representation.
- 11 Result. Dhaka was first tested on the simulated dataset. The simulated dataset contains
- 12 500 cells, each including 3K genes, clustered into 5 different clusters with 100 cells each.
- 13 The clustering performance was compared with other methods including t-SNE, PCA,
- 14 SIMLR, NMF, an autoencoder, MAGIC, and scVI. Dhaka was shown to have an ARI
- 15 higher than most other comparing methods. Dhaka was then applied to the
- 16 Oligodendroglioma data and could separate malignant cells from non-malignant
- microglia/macrophage cells. It also uncovered the shared glial lineage and differentially
- expressed genes for the lineages. Dhaka was also applied to the Glioblastoma data and
- 19 revealed an evolutionary trajectory of the malignant cells where cells gradually evolve
- 20 from a stemlike state to a more differentiated state. In contrast, other methods failed to
- 21 capture this underlying structure. Dhaka was next applied to the Melanoma cancer
- 22 dataset [74] and uncovered two distinct clusters that showed the intra-tumor
- 23 heterogeneity of the Melanoma samples. Dhaka was finally applied to copy number

- variation data [75] and shown to identify one major and one minor cell clusters, of which
- 2 other methods could not find.

- 4 4.3.3. scvis: a VAE for capturing low-dimensional structures
- 5 scvis [76] is a VAE network (Figs. 2A, 4C) that learns the low-dimensional
- 6 representations capture both local and global neighboring structures in scRNA-seg data.
- 7 *Model:* scvis adopts the generic VAE formulation described in section 3.1. However, it has
- 8 a unique loss function defined as

$$L(\mathbf{O}) = -\mathcal{L}(\mathbf{O}) + \lambda L_t(\mathbf{O}), \tag{23}$$

- 9 where $\mathcal{L}(\mathbf{\Theta})$ is ELBO as in Eq. (3) and L_t is a regularizer using non-symmetrized t-SNE
- objective function [76], which is defined as

$$L_t(\mathbf{\Theta}) = \sum_{i=1}^{N} \sum_{j=1, j \neq i}^{N} p_{j|i} \log \frac{p_{j|i}}{q_{j|i}},$$
(24)

where i and j are two different cells, $p_{i|j}$ measures the local cell relationship in the data 11 space, and $q_{j|i}$ measures such relationship in the latent space. Because t-SNE algorithm 12 preserves the local structure of high dimensional space, $\,L_t\,$ learns local structures of cells. 13 Results. scvis was tested on the simulated data and outperformed t-SNE in a nine-14 15 dimensional space task. scvis preserved both local structure and global structure. The relative positions of all clusters were well kept but outliers were scattered around clusters. 16 17 Using simulated data and comparing to t-SNE, scvis generally produced consistent and 18 better patterns among different runs while t-SNE could not. scvis also presented good 19 results on adding new data to an existing embedding, with median accuracy on new data 20 at 98.1% for K= 5 and 94.8% for K= 65, when train K cluster on original data then test the 21 classifier on new generated sample points. The scvis was subsequently tested on four 1 real datasets including metastatic melanoma, oligodendroglioma, mouse bipolar and

2 mouse retina datasets. In each dataset, scvis was showed to preserve both the global

and local structure of the data.

4

7

8

10

11

12

13

14

15

18

19

20

21

22

23

3

5 4.3.4. scVAE: VAE for single-cell gene expression data

6 scVAE [77] includes multiple VAE models (Figs. 2A, 4C) for denoising gene expression

levels and learning the low-dimensional latent representation of cells. It investigates

different choices of the likelihood functions in the VAE model to model different data sets.

9 <u>Model.</u> scVAE is a conventional fully connected network. However, different distributions

have been discussed for $p(x_{gn}|v_{gn},\alpha_{gn})$ to model different data behaviors. Specifically,

scVAE considers Poisson, constrained Poisson, and negative binomial distributions for

count data, piece-wise categorical Poisson for data including both high and low counts,

and zero-inflated version of these distributions to model missing values. To model

multiple modes in cell expressions, a Gaussian mixture is also considered for

 $q(\mathbf{z}_n|\mathbf{x}_n,s_n)$, resulting in a GMVAE. The inference process still follows that of a VAE as

discussed in section 3.1.

17 <u>Results.</u> scVAEs were evaluated on the PBMC data and compared with factor analysis

(FA) models. The results showed that GMVAE with negative binomial distribution

achieved the highest lower bound and ARI. Zero-inflated Poisson distribution performed

the second best. All scVAE models outperformed the baseline linear factor analysis

model, which suggested that a non-linear model is needed to capture single-cell genomic

features. GMVAE was also compared with Seurat and shown to perform better using the

withheld data. However, scVAE performed no better than scVI [17] or scvis [76], both are

24 VAE models.

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

4.3.5. VASC: VAE for scRNA-seq

3 VASC [78] is another VAE (Figs. 2A, 4C) for dimension reduction and latent

Model: VASC's input is the log-transformed expression but rescaled in the range [0,1]. A

4 representation but it models dropout.

dropout layer (dropout rate of 0.5) is added after the input layer to force subsequent layers to learn to avoid dropout noise. The encoder network has three layers fully connected and the first layer uses linear activation, which acts like an embedded PCA transformation. The next two layers use the ReLU activation, which ensures a sparse and stable output. This model's novelty is the zero-inflation layer (ZI layer), which is added after the decoder to model scRNA-seq dropout events. The probability of dropout event is defined as $e^{-\hat{x}^2}$ where \hat{x} is the recovered expression value obtained by the decoder network. Since back-propagation cannot deal with a stochastic network with categorical variables, a Gumbel-softmax distribution [79] is introduced to address the difficulty of the ZI layer. The loss function of the model takes the form $L = \mathcal{L}(\mathbf{0}) + \lambda L_{KL}(\mathbf{0})$, where \mathcal{L} is the binary entropy because the input is scaled to [0 1], and L_{KL} a loss performed using KL divergence on the latent variables. After the model is trained, the latent code can be used as the dimension-reduced feature for downstream tasks and visualization. Results. VASC was compared with PCA, t-SNE, ZIFA, and SIMLR on 20 datasets. In the study of embryonic development from zygote to blast cells, all methods roughly reestablished the development stages of different cell types in the dimension-reduced space. However, VASC showed the better performance to model embryo developmental progression. In the Goolam, Biase and Yan datasets, scRNA-seq data were generated

through embryonic development stages from zygote to blast, VASC re-established

1 development stage from 1, 2, 4, 8, 16 to blast, while other methods failed. In the Pollen, 2 Kolodziejczyk, and Baron dataset, VASC formed an appropriate cluster, either with homogeneous cell type, preserved proper relative positions, or minimal batch influence. 3 4 Interestingly, when tested on the PBMC dataset, VASC was shown to identify the major 5 global structure (B cells, CD4+, CD8+ T cells, NK cells, Dendritic cells), and detect subtle 6 differences within monocytes (FCGR3A+ vs CD14+ monocytes), indicating the capability 7 of VASC handling a large number of cells or cell types. Quantitative clustering performance in NMI, ARI, homogeneity and completeness was also performed. VASC 8 9 always ranked top two in all the datasets. In terms of NMI and ARI, VASC best performed 10 on 15 and 17 out of 20 datasets, respectively.

1112

4.3.6. scDeepCluster

scDeepCluster [80] is an AE network (Figs. 2B, 4B) that simultaneously learns feature 13 14 representation and performs clustering via explicit modeling of cell clusters as in DESC. <u>Model:</u> Similar to DCA, scDeepCluster adopts a ZINB distribution for x_n as in Eq. (13) 15 16 and (15). The loss is similar to Eq. (19) except that the first term is the negative loglikelihood of the ZINB data distribution as defined in Eq. (15) and the second L_c is a 17 clustering loss performed using KL divergence as in DESC algorithm. Compared to csvis, 18 19 scDeepcluster focuses more on clustering assignment due to the KL divergence. 20 Results. scDeepCluster was first tested on the simulation data and compared with other 21 seven methods including DCA [18], two multi-kernel spectral clustering methods MPSSC 22 [81] and SIMLR [44], CIDR [61], PCA + k-mean, scvis [76] and DEC[82]. In different dropout rate simulations, scDeepCluster significantly outperformed the other methods 23 24 consistently. In signal strength, imbalanced sample size, and scalability simulations,

scDeepcluster outperformed all other algorithms and scDeepCluster and most notably advantages for weak signals, robust against different data imbalance levels and scaled linearly with the number of cells. scDeepCluster was then tested on four real datasets (10X PBMC, Mouse ES cells, Mouse bladder cells, Worm neuron cells) and shown to outperform all other comparing algorithms. Particularly, MPSSC and SIMLR failed to process the full datasets due to quadratic complexity.

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

6

1

2

3

4

5

4.3.7. cscGAN: Conditional single-cell generative adversarial neural networks

cscGAN [83] is a GAN model (Figs. 2C, 4E) designed to augment the existing scRNAseg samples by generating expression profiles of specific cell types or subpopulations. Two models, csGAN and cscGAN, were developed following the general formulation of WGAN described in section 3.3. The difference between the two models is that cscGAN is a conditional GAN such that the input to the generator also includes a class label y or cell type, i.e. $\phi_G(\mathbf{z}, y)$. The projection-based conditioning (PCGAN) method [84] was adopted to obtain the conditional GAN. For both models, the generator (three layers of 1024, 512, and 256 neurons) and discriminator (three layers of 256, 512, and 1024 neurons) are fully connected DNNs. Results: The performance of scGAN was first evaluated using PBMC data. The generated samples were shown to capture the desired clusters and the real data's regulons. Additionally, the AUC performance for classifying real from generated samples by a Random Forest classifier only reached 0.65, performance close to 0.5. Finally, scGAN's generated samples had a smaller MMD than those of Splatter, a state-of-the-art scRNAseg data simulator [85]. Even though a large MMD was observed for scGAN when compared with that of SUGAR, another scRNA-seg simulator, SUGAR [86] was noted

- 1 for prohibitively high runtime and memory. scGAN was further trained and assessed on
- 2 the bigger mouse brain data and shown to model the expression dynamics across tissues.
- 3 Then, the performance of cscGAN for generating cell-type-specific samples was
- 4 evaluated using the PBMC data. cscGAN was shown to generate high-quality scRAN-
- 5 seq data for specific cell types. Finally, the real PBMC samples were augmented with the
- 6 generated samples. This augmentation improved the identification of rare cell types and
- 7 the ability to capture transitional cell states from trajectory analysis.

4.4. Multi-functional models

- 10 Given the versatility of AE and VAE in addressing different scRAN-seq analysis
- challenges, DL models possessing multiple analysis functions have been developed. We
- survey these models in this section.

13 4.4.1. scVI: single-cell variational inference

- scVI [17] is designed to address a range of fundamental analysis tasks, including batch
- 15 correction, visualization, clustering, and differential expression.
- 16 <u>Model.</u> scVI is a VAE (**Figs. 2A, 4C**) that models the counts of each cell from different
- batches. scVI adopts a ZINB distribution for x_{an}

$$p(x_{gn}|\pi_{gn}, L_n, \nu_{gn}, \alpha) = \pi_{gn}\delta(0) + (1 - \pi_{gn})NB(L_n\nu_{gn}, \alpha_g),$$
(25)

- which is defined similarly as Eq (14) in DCA, except that L_n denotes the scaling factor for
- cell n, which follows a log-Normal $(log\mathcal{N})$ prior as $p(L_n) = log\mathcal{N}(\mu_{L_n}, \sigma_{L_n}^2)$, therefore, v_{gn}
- represents the mean counts normalized by L_n . Now, let $s_n \in \{0,1\}^B$ be the batch ID of cell
- 21 n with B being the total number of batches. Then, v_{gn} and π_g are further modeled as

- 1 functions of the d-dimension latent variable $\mathbf{z}_n \in \mathbb{R}^d$ and the batch ID s_n by the decoder
- 2 networks $D_{\theta_{\nu}}$ and $D_{\theta_{\pi}}$ as

$$\mathbf{v}_n = D_{\boldsymbol{\theta}_{\nu}}(\mathbf{z}_n, s_n), \ \boldsymbol{\pi}_n = D_{\boldsymbol{\theta}_{\pi}}(\mathbf{z}_n, s_n), \tag{26}$$

- 3 where the gth element of \mathbf{v}_n and $\mathbf{\pi}_n$ are \mathbf{v}_{gn} and $\mathbf{\pi}_g$, respectively, and $\mathbf{\theta}_{v}$, and $\mathbf{\theta}_{\pi}$ are the
- 4 decoder weights. Note that the lower layers of the two decoders are shared. For inference,
- 5 both \mathbf{z}_n and L_n are considered as latent variables and therefore $q(x_n, s_n) =$
- 6 $q(\mathbf{z}_n|\mathbf{x}_n,s_n)q(L_n|\mathbf{x}_n,s_n)$ is a mean-field approximate to the intractable posterior
- 7 distribution $p(\mathbf{z}_n, L_n | \mathbf{x}_n, s_n)$ and

$$q(\mathbf{z}_{n}|\mathbf{x}_{n}, s_{n}) = \mathcal{N}\left(\boldsymbol{\mu}_{z_{n}}, diag(\boldsymbol{\sigma}_{Z_{n}}^{2})\right),$$

$$q(L_{n}|\mathbf{x}_{n}, s_{n}) = log\mathcal{N}\left(\boldsymbol{\mu}_{L_{n}}, diag(\boldsymbol{\sigma}_{L_{n}}^{2})\right),$$
(27)

- 8 whose means and variances $\{\pmb{\mu}_{^{Z}n},\pmb{\sigma}_{^{Z}n}^2\}$ and $\{\pmb{\mu}_{^{L}n},\pmb{\sigma}_{^{L}n}^2\}$ are defined by the encoder
- 9 networks E_Z and E_L applied to x_n and s_n as

$$\{\boldsymbol{\mu}_{z_n}, \boldsymbol{\sigma}_{z_n}^2\} = E_{\boldsymbol{\phi}_z}(\boldsymbol{x}_n, s_n),$$

$$\{\boldsymbol{\mu}_{L_n}, \boldsymbol{\sigma}_{L_n}^2\} = E_{\boldsymbol{\phi}_L}(\boldsymbol{z}_n, s_n)$$
(28)

where $\, {m \phi}_{\scriptscriptstyle Z} ,$ and ${m \phi}_{\scriptscriptstyle L}$ are the encoder weights. Note that, like the decoders, the lower layers 10 11 of the two encoders are also shared. Overall, the model parameters to be estimated by the variational optimization is $\Theta = \{ \theta_{\nu}, \theta_{\pi}, \phi_{z}, \phi_{L}, \alpha_{g} \}$. After inference, \mathbf{z}_{n} are used for 12 visualization and clustering. v_{gn} provides a batch-corrected, size-factor normalized 13 estimate of gene expression for each gene g in each cell n. An added advantage of the 14 15 probabilistic representation by scVI is that it provides a natural probabilistic treatment of 16 the subsequent differential analysis, resulting in lower variance in the adopted hypothesis 17 tests.

Results: scVI was evaluated for its scalability, the performance of imputation. For 2 scalability, ScVI was shown to be faster than most nonDL algorithms and scalable to handle twice as many cells as nonDL algorithms with a fixed memory. For imputation, ScVI, together with other ZINB-based models, performed better than methods using 5 alternative distributions. However, it underperformed for the dataset (HEMATO) with fewer cells. For the latent space, scVI was shown to provide a comparable stratification 7 of cells into previously annotated cell types. Although scVI failed to ravel SIMLR, it is among the best in capturing biological structures (hierarchical structure, dynamics, etc.) and recognizing noise in data. For batch correction, it outperforms ComBat. For 10 normalizing sequencing depth, the size factor inferred by scVI was shown to be strongly correlated with the sequencing depth. Interestingly, the negative binomial distribution in 12 the ZINB was found to explain the proportions of zero expressions in the cells, whereas the zero probability π_{gn} is found to be more correlated with alignment errors. For 14 differential expression analysis, scVI was shown to be among the best.

15

16

1

3

4

6

8

9

11

13

4.4.2. LDVAE: linearly decoded variational autoencoder

- LDVAE [87] is an adaption of scVI to improve the model interpretability but it still benefits 17 18 from the scalability and efficiency of scVI. Also, this formulation applies to general VAE 19 models and thus is not restricted to scRNA-seg analysis.
- <u>Model.</u> LDVAE follows scVI's formulation but replaces the decoder $D_{\theta_{\nu}}$ in Eq. (26) by a 20 linear model 21

$$\mathbf{v}_n = \mathbf{W}\mathbf{z}_n,\tag{29}$$

- where $\mathbf{W} \in \mathbb{R}^{d \times G}$ is the weight matrix. Being the linear decoder provides interpretability in
- 2 the sense that the relationship between latent representation \mathbf{z}_n and gene expression \mathbf{v}_n
- 3 can be readily identified. LDVAE still follows the same loss and non-linear inference
- 4 scheme as scVI.
- 5 Results. LDVAE's latent variable z_n could be used for clustering of cells with similar
- 6 accuracy as a VAE. Although LDVAE had a higher reconstruction error than VAE, due
- 7 to the linear decoder, the variations along the different axes of z_n establish direct linear
- 8 relationships with input genes. As an example from analyzing mouse embryo scRNA-seq,
- 9 $z_{1,n}$, the second element of z_n , is shown to relate to simultaneous variations in the
- 10 expression of gene Pou5f1 and Tdgf1. In contrast, such interpretability would be
- intractable without approximation for a VAE. LDVAE was also shown to induce fewer
- 12 correlations between latent variables and to improve the grouping of the regulatory
- programs. LDVAE is capable to scale to a large dataset with ~2M cells.

- 15 **4.4.3. SAUCIE**
- SAUCIE [15] is an AE (Figs. 2B, 4B) designed to perform multiple functions, including
- 17 clustering, batch correlation, imputation, and visualization. SAUCIE is applied to the
- 18 normalized data instead of count data.
- 19 <u>Model.</u> SAUCIE includes multiple model components designed for different functions.
- 1. Clustering: SAUCIE first introduced a "digital" binary encoding layer $h^c \in \{0,1\}^J$ in the
- decoder *D* that functions to encode the cluster ID. To learn this encoding, an entropy
- 22 loss is introduced

$$L_D = \sum_{k=1}^K p_k \log p_k, \tag{30}$$

- where p_k is the probability (proportion) of activation on neuron k by the previous layer.
- 2 Minimizing this entropy loss promotes sparse neurons, thus forcing a binary encoding.
- To encourage clustering behavior, SAUCIE also introduced an intracluster loss as

$$L_{C} = \sum_{i,j:h_{i}^{C} = h_{j}^{C}} \|\widehat{x}_{i} - \widehat{x}_{j}\|^{2},$$
(31)

- 4 which computes the distance L_c between the expressions of a pair of cells (\hat{x}_i, \hat{x}_i)
- 5 that have the same cluster ID $(h_i^c = h_i^c)$.

- 7 2. Batch correction: To correct the batch effect, an MMD loss is introduced to measure
- the differences in terms of the distribution between batches in the latent space

$$L_B = \sum_{l=1, l \neq ref}^{B} MMD(\mathbf{z}_{ref}, \mathbf{z}_l), \tag{32}$$

- where B is the total number of batches and \mathbf{z}_{ref} is the latent variable of an arbitrarily
- 10 chosen reference batch.
- 11 3. Imputation and visualization: The output of the decoder is taken by SAUCIE as an
- imputed version of the input gene expression. To visualize the data without performing
- an additional dimension reduction directly, the dimension of the latent variable \mathbf{z}_n is
- 14 forced to 2.
- 15 Training the model includes two sequential runs. In the first run, an autoencoder is trained
- to minimize the loss $L_0 + \lambda_B L_B$ with L_0 being the MSE reconstruction loss defined in (9)
- so that a batch-corrected, imputed input \tilde{x} can be obtained at the output of the decoder.
- In the second run, the bottleneck layer of the encoder from the first run is replaced by a
- 19 2-D latent code for visualization and a digital encoding layer is also introduced. This model
- 20 takes the cleaned \widetilde{x} as the input and is trained for clustering by minimizing the loss L_0 +

- 1 $\lambda_D L_D + \lambda_C L_C$. After the model is trained, \tilde{x} is the imputed, batch-corrected gene
- 2 expression. The 2-D latent code is used for visualization and the binary encoder encodes
- 3 the cluster ID.
- 4 Results. SAUCIE was evaluated for clustering, batch correction, imputation, and
- 5 visualization on both simulated and real scRNA-seq and scCyToF datasets. The
- 6 performance was compared to minibatch *kmeans*, Phenograph [88] and 22 for clustering;
- 7 MNN [32] and canonical correlation analysis (CCA) [33] for batch correction; PCA,
- 8 Monocle2 [89], diffusion maps, UMAP [90], tSNE [91] and PHATE [92] for visualization;
- 9 MAGIC [57], scImpute [37] and nearest neighbors completion (NN completion) for
- imputation. Results showed that SAUCIE had a better or comparable performance with
- other approaches. Also, SAUCIE has better scalability and faster runtimes than any of
- the other models. SAUCIE's results on the scCyToF dengue dataset were further
- analyzed in greater detail. SAUCIE was able to identify subtypes of the T cell populations
- and demonstrated distinct cell manifold between acute and healthy subjects.

- 16 **4.4.4.** scScope:
- scScope [93] is an AE (Figs. 2B, 4D) with recurrent steps designed for imputation and
- 18 batch correction.
- 19 *Model.* scScope has the following model design for batch correction and imputation.
- 1. Batch correction: A batch correction layer is applied to the input expression as

$$\mathbf{x}_n^c = ReLu(\mathbf{x}_n - \mathbf{B}\mathbf{u}_c), \tag{33}$$

- where and ReLU is the ReLu activation function, $\mathbf{B} \in \mathbb{R}^{G \times K}$ is the batch correction
- matrix, $u_c \in \{0,1\}^{K \times 1}$ is an indicator vector with entry 1 indicates the batch of the input,
- and *K* is the total number of batches.

2. Recursive imputation: Instead of using the reconstructed expression \hat{x}_n as the imputed expression like in SAUCIE, scScope adds an imputer to \hat{x}_n to recursively improve the imputation result. The imputer is a single-layer autoencoder, whose decoder performs the imputation as

$$\widehat{\widehat{\mathbf{x}}}_n = P_I \big[D_I \big(\widehat{\widehat{\mathbf{z}}}_n \big) \big], \tag{34}$$

where $\hat{\mathbf{z}}_n$ is the output of the imputer encoder, D_I is the imputer decoder, and P_I is a masking function that set the elements in $\hat{\mathbf{x}}_n$ that correspond to the non-missing values to zero. Then, $\hat{\mathbf{x}}_n$ will be fed back to fill the missing value in the batch corrected input as $\mathbf{x}_n^c + \hat{\mathbf{x}}_n$, which will be passed on to the main autoencoder. This recursive imputation can iterate multiple cycles as selected.

10 The loss function is defined as

$$\mathcal{L}(\mathbf{\Theta}) = \sum_{n=1}^{N} \sum_{t=1}^{T} ||P_{I}^{-}[\mathbf{x}_{n}^{c} - \widehat{\mathbf{x}}_{n}^{t}]||^{2}, \qquad (35)$$

recursion, P_l^- is another masking function that forces the loss to compute only the non-missing values in \boldsymbol{x}_n^c .

Results. scScope was evaluated for its scalability, clustering, imputation, and batch correction. It was compared with PCA, MAGIC, ZINB-WaVE, SIMLR, AE, scVI, and DCA. For scalability and training speed, scScope was shown to offer scalability (for >100K cells) with high efficiency (faster than most of the approaches). For clustering results, scScope outperformed most of the algorithms on small simulated datasets but offered similar performance on large simulated datasets. For batch correction, scScope performed comparably with other approaches but with faster runtime. For imputation, scScope

where T is the total number of recursion, \hat{x}_n^t is the reconstructed expression at tth

- 1 produced smaller errors consistently across a different range of expression. scScope was
- 2 further shown to be able to identify rear cell populations from a large mix of cells.

- 4 4.5. Automated cell type identification
- 5 scRNA-seq can catalog cell types in complex tissues under different conditions. However,
- 6 the commonly adopted manual cell typing approach based on known markers is time-
- 7 consuming and less reproducible. We survey deep learning models of automated cell
- 8 type identification.

9

- 4.5.1. DigitalDLSorter
- DigitalDLSorter [51] was proposed to identify and quantify the immune cells infiltrated in
- tumors captured in bulk RNA-seq, utilizing single-cell RNA-seq data.
- 14 *Model.* DigitalDLSorter is a 4-layer DNN (**Fig. 4A**) (2 hidden layers of 200 neurons each
- and an output of 10 cell types). The DigitalDLSorter is trained with two single-cell
- datasets: breast cancers [94] and colorectal cancers [95]. For each cell, it is determined
- to be tumor cell or non-tumor cell using RNA-seq based CNV method [94], followed by
- using xCell algorithm [96] to determine immune cell types for non-tumor cells. Different
- 19 pseudo bulk (from 100 cells) RNA-seq datasets were prepared with known mixture
- 20 proportions to train the DNN. The output of DigitalDLSorter is the predicted proportions
- of cell types in the input bulk sample.
- 22 Result. DigitalDLSorter was first tested on simulated bulk RNA-seq samples.
- 23 DigitalDLSorter achieved excellent agreement (linear correlation of 0.99 for colorectal
- cancer, and good agreement in quadratic relationship for breast cancer) at predicting cell
- 25 types proportion. The proportion of immune and nonimmune cell subtypes of test bulk

TCGA samples was predicted by DigitalDLSorter and the results showed a very good correlation to other deconvolution tools including TIMER [94], ESTIMATE [97], EPIC [98] and MCPCounter [99]. Using DigitalDLSorter predicted CD8+ (good prognosis for overall and disease-free survival) and Monocytes-Macrophages (MM, indicator for protumoral activity) proportions, it is found that patients with higher CD8+/MM ratio had better survival for both cancer types than those with lower CD8+/MM ratio. Both EPIC and MCPCounter yielded non-significant survival associations using their cell proportion estimate.

4.5.2. scCapsNet

- scCapsNet [52] is an interpretable capsule network designed for cell type prediction. The paper showed that the trained network could be interpreted to inform marker genes and regulatory modules of cell types.
- *Model.* scCapsNet takes log-transformed, normalized expressions as input and follows the general CapsNet model described in Section 3.4. Capsule v_k represents the probability of a single cell x_n belonging to cell type k, which will be used for cell-type classification. Once trained, the interpretation of marker genes and regulatory modules can be achieved by determining first the important primary capsules for each cell type and then the most significant genes for each important primary capsule (identified based on c_{kl} directly). To determine the genes that are important for an important primary capsule l, genes are ranked base on the scores of the first principal component computed from the columns of $W_{P,l}$ in Eq. (15) and then the markers are obtained by a greedy search along with the ranked list for the best classification performance.
- Results. scCapsNet's performance was evaluated on human PBMCs [100] and mouse retinal bipolar cells [101] datasets and shown to have comparable accuracies (99% and

97% respectively) with DNN and other popular ML algorithms (SVM, random forest, LDA and nearest neighbor). However, the interpretability of scCapsNet was demonstrated extensively. First, examining the coupling coefficients for each cell type showed that only a few primary capsules have high values and thus are effective. Subsequently, a set of core genes were identified for each effective capsule using the greedy search on the PC-score ranked gene list. GO enrichment analysis showed that these core genes were enriched in cell-type-related biological functions. Mapping the expression data into space spanned by PCs of the columns of $W_{P,l}$ corresponding to all core genes uncovered regulatory modules that would be missed by the T-SNE of gene expressions, which demonstrated the effectiveness of the embeddings learned by scCapsNet in capturing the functionally important features.

12

13

1

2

3

4

5

6

7

8

9

10

11

4.5.3. netAE: network-enhanced autoencoder

- netAE [102] is a VAE-based semi-supervised cell type prediction model (**Figs. 2A, 4C**)
- that deals with scenarios of having a small number of labeled cells.
- 16 <u>Model.</u> netAE works with UMI counts and assumes a ZINB distribution for x_{gn} as in Eq.
- 17 (25) in scVI. However, netAE adopts the general VAE loss as in Eq. (6) with two function-
- 18 specific loss as

$$L(\mathbf{\Theta}) = -\mathcal{L}(\mathbf{\Theta}) + \lambda_1 \sum_{n \in S} Q(\mathbf{z}_n) + \lambda_2 \sum_{n \in S_L} log f(y_n | \mathbf{z}_n), \tag{36}$$

- where S is a set of indices for all cells and S_L is a subset of S for only cells with cell type
- 20 labels, Q is modified Newman and Girvan modularity [103] that quantifies cluster strength
- using \mathbf{z}_n , f is the softmax function, and y_n is the cell type label. The second loss in Eq.
- 22 (36) functions as a clustering constraint and the last term is the cross-entropy loss that
- 23 constrains the cell type classification.

1 Results: netAE was compared with popular dimension reduction methods including scVI, ZIFA, PCA and AE as well as a semi-supervised method scANVI [104]. For different 2 dimensionality reduction methods, cell type classification from latent features of cells was 3 4 carried out using KNN and logistic regression. The effect of different labeled samples 5 sizes on classification performance was also investigated, where the sample size varied 6 from as few as 10 cells to 70% of all cells. Among 3 test datasets (mouse brain cortex, 7 human embryo development, and mouse hematopoietic stem and progenitor cells), netAE outperformed most of the baseline methods. Latent features were visualized using 8 9 t-SNE and cell clusters by netAE were tighter than those by other embedding spaces. 10 netAE also showed consistency of better cell-type classification with improved cell type 11 clustering. This suggested that the latent spaces learned with added modularity constraint 12 in the loss helped identify clusters of similar cells. Ablation study by removing each of the 13 three loss terms in Eq. (36) showed a drop of cell-type classification accuracy, suggesting 14 all three were necessary for the optimal performance.

15

16

17

18

19

20

21

22

23

24

4.5.4. scDGN - supervised adversarial alignment of single-cell RNA-seq data

scDGN [53], or Single Cell Domain Generalization Network (**Fig. 4G**), is a domain adversarial network that aims to accurately assign cell types of single cells while performing batch removal (domain adaptation) at the same time. It benefits from the superior ability of domain adversarial learning to learn representations that are invariant to technical confounders.

<u>Model.</u> scDGN takes the log-transformed, normalized expression as the input and has

three main modules: i) an encoder $(E_{\phi}(x_n))$ for dimension reduction of scRNA-seq data, ii) cell-type classifier, $C_{\phi_C}\Big(E_{\phi}(x_n)\Big)$ with parameters ϕ_C , and iii) domain (batch)

discriminator, $D_{\phi_{\mathbb{D}}}\left(E_{\phi}(x_n)\right)$. The model has a Siamese design and the training takes a pair of cells (x_1, x_2) , each from the same or different batches. The encoder network contains two hidden layers with 1146 and 100 neurons. $C_{\phi_{\mathcal{C}}}$ classifies the cell type and $D_{\phi_{\mathbb{D}}}$ predicts whether x_1 and x_2 are from the same batch or not. The overall loss is denoted by

$$L(\boldsymbol{\phi}, \boldsymbol{\phi}_{C}, \boldsymbol{\phi}_{D}) = L_{C}\left(C_{\boldsymbol{\phi}_{C}}\left(E_{\boldsymbol{\phi}}(\boldsymbol{x}_{1})\right)\right) - \lambda L_{D}\left(D_{\boldsymbol{\phi}_{D}}\left(E_{\boldsymbol{\phi}}(\boldsymbol{x}_{1})\right), D_{\boldsymbol{\phi}_{D}}\left(E_{\boldsymbol{\phi}}(\boldsymbol{x}_{2})\right)\right), \quad (37)$$

where L_C is the cross-entropy loss, L_D is a contrastive loss as described in [105]. Notice that (37) has an adversarial formulation and minimizing this loss maximizes the misclassification of cells from different batches, thus making them indistinguishable. Similar to GAN training, scDGN is trained to iteratively solve: $\hat{\boldsymbol{\phi}}_D = \operatorname{argmin}_{\boldsymbol{\phi}_D} L(\hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\phi}}_C, \boldsymbol{\phi}_D)$ and $(\hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\phi}}_C) = \operatorname{argmin}_{\boldsymbol{\phi}, \boldsymbol{\phi}_C} L(\boldsymbol{\phi}, \boldsymbol{\phi}_C, \hat{\boldsymbol{\phi}}_D)$. Results. scDGN was tested for classifying cell types and aligning batches ranging in size

from 10 to 39 cell types and from 4 to 155 batches. The performance was compared to a series of deep learning and traditional ML methods, including Lin et al. DNN [72], CaSTLe [106], MNN [32], scVI [17], and Seurat [10]. scDGN outperformed all other methods in the classification accuracy on a subset of scQuery datasets (0.29), PBMC (0.87), and 4 of the six Seurat pancreatic datasets (0.86 - 0.95). PCA visualization of the learned data representations demonstrated that scDGN overcame the batch differences and clearly separated cell clusters based on cell types, while other methods were vulnerable to batch effects. In summary, scDGN is a supervised adversarial alignment method to eliminate the batch effect of scRNA-seq data and create cleaner representations of cell types.

4.6. Biological function prediction

12

13

14

15

16

17

18

19

20

21

- 1 Predicting biological functions and responses to treatment at single cell level or cell types
- 2 is critical to understand cellular system functioning and potent responses to stimulations.
- 3 DL models are capable of capture gene-gene relationship and their property in the latent
- 4 space. Several surveyed models demonstrate exciting results to learn complex biological
- 5 functions and outcomes.

9

10

11

12

13

14

15

16

18

19

20

21

22

23

24

7 4.6.1. CNNC: convolutional neural network for coexpression

8 CNNC [54] is proposed to infer causal interactions between genes from scRNA-seq data.

Model. CNNC is a Convolutional Neural Network (CNN) (**Fig. 4F**), one of the most popular DL models. CNNC takes expression levels of two genes from many cells and transforms them into a 32 x 32 image-like normalized empirical probability function (NEPDF), which measures the probabilities of observing different coexpression levels between the two genes. CNNC includes 6 convolutional layers, 3 max-pooling layers, 1 flatten layer, and

one output layer. All convolution layers have 32 kernels of size 3x3. Depending on the

application, the output layer can be designed to predict the state of interaction (Yes/No)

between the genes or the causal interaction between the input genes (no interaction,

gene A regulates gene B, or gene B regulates gene A).

Result. CNNC was trained to predict transcription factor (TF)-Gene interactions using the mESC data from scQuery [107], where the ground truth interactions were obtained from the ChIP-seq dataset from the GTRD database [108]. The performance was compared with DNN, count statistics [109], and mutual information-based approach [110]. CNNC was shown to have more than 20% higher AUPRC than other methods and reported almost no false-negative for the top 5% predictions. CNNC was also trained to predict the pathway regulator-target gene pairs. The positive regulator-gene pairs were obtained

from KEGG [111], Reactome [112], and negative samples were genes pairs that appeared in pathways but do not interacted. CNNC was shown to have better performance of predicting regulator-gene pairs on both KEGG and Reactome pathways than other methods including Pearson correlation, count statistics, GENIE3 [113], Mutual information, Bayesian directed network (BDN), and DREMI [110]. CNNC was also applied for causality prediction between two genes, that is if two genes regulate each other and if they do, which gene is the regulator. The ground truth causal relationships were also obtained from the KEGG and Reactome datasets. Again, CNNC reported better performance than BDN, the common method developed to learn casual relationships from gene expression data. CNNC was finally trained to assign 3 essential cell functions (cell cycle, circadian rhythm, and immune system) to genes. This is achieved by training CNNC to predict pairs of genes from the same function (e.g. Cell Cycle defined by mSigDB from gene set enrichment analysis (GSEA) [114]) as 1 and all other pairs as 0. The performance was compared with "guilt by association" and DNN, and CNNC was shown to have more than 4% higher AUROC and reported all positives for the top 10% predictions.

17 18

19

20

21

22

23

24

16

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

4.6.2. scGen, a generative model to predict perturbation response of single cells across cell types

scGen [115] is designed to learn cellular responses to certain perturbations such as drug treatment and gene knockout from single-cell expression data, and then predict cellular responses to the same perturbation for a new sample or a new cell type. The novelty of scGen is that it learns the cellular response in the latent space instead of the expression data space.

1 Model. ScGen follows the general VAE (Figs. 2A, 4C) for scRNA-seg data but uses the "latent space arithmetics" to learn perturbations' response. Given scRNA-seg samples of 2 3 perturbed (denoted as p) and unperturbed cells (denoted as unp), a VAE model is trained. Then, the latent space representation \mathbf{z}_p and \mathbf{z}_{unp} are obtained for the perturbed and 4 5 unperturbed cells. Following the notion that VAE could map nonlinear operations (e.g., 6 perturbation) in the data space to linear operations in the latent space, ScGen estimates the response in the latent space as $\delta = \bar{z}_p - \bar{z}_{unp}$, where $\bar{z}_{.}$ is the average representation 7 8 of samples from the same cell type or different cell types. Then, given the latent representation $\mathbf{z'}_{unp}$ of an unperturbed cell for a new sample from the same or different 9 cell type, the latent representation of the corresponding perturbed cell can be predicted 10 as $\mathbf{z'}_{p} = \mathbf{z'}_{unp} + \boldsymbol{\delta}$. The expression of the perturbed cell can also be estimated by feeding 11 $\mathbf{z'}_p$ into the VAE decoder. The scGen can also be expanded to samples and treatment 12 across two species (using orthologues between species). When scGen is trained for 13 14 species 1 (s_1) with both perturbed and unperturbed cells and species 2 (s_2) with only 15 unperturbed cells, the latent code for the perturbed cells from s_2 can be predicted as $\mathbf{z}_{s_2,p} = \frac{1}{2} \left(z_{s_1,p} + z_{s_2,unp} + \delta_s + \delta_p \right)$ where $\boldsymbol{\delta}_p = z_{s_1,unp} - z_{s_1,p}$ captures the response of 16 perturbation and $\delta_s = z_{s_1} - z_{s_2}$ represents the difference between species. 17 18 Result. scGen was applied to predict perturbation of out-of-samples response in human PBMCs data, and scGen showed a higher average correlation (R²= 0.948) between 19 20 predicted and real data for six cell types [116]. Compared with other methods including CVAE [117], style transfer GAN [118], linear approaches based on vector arithmetics (VA) 21 22 similar in [119] and PCA+VA, scGen predicted full distribution of ISG15 gene (strongest 23 regulated gene by IFN-β) response to IFN-β [116], while others might predict mean

(CAVE and style transfer GAN) but failed to produce the full distribution. scGen was also tested on predicting the intestinal epithelial cells' response to infections [120]. For early transit-amplifying cells, scGen showed good prediction (R²=0.98) for both *Heligmosomoides polygyrus* and *Salmonella* infections. Finally, scGen was evaluated for perturbation across species using scRNA-seq data set by Hagai et al. [121], which comprises bone marrow-derived mononuclear phagocytes from mice, rats, rabbits, and pigs perturbed with lipopolysaccharide (LPS). scGen's predictions of LPS perturbation responses were shown to be highly correlated (R² = 0.91) with the real responses.

5. Conclusions

We systematically survey 25 DL models according to the challenges they address. We categorize major DL model statements into VAE, AE, and GAN with a unified mathematic formulation in Section 3 (graphic model representation in Fig. 2), which will guide readers to focus on the DL model selection, training strategies, and loss functions in each algorithm. Specifically, the differences in loss functions are highlighted in each DL model's applications to meet specific objectives. DL/ML models that 25 surveyed models are evaluated against are presented in Fig. 3, providing a straightforward way for readers to pick up the most suitable DL model at a specific step for their own scRNA-seq data analysis. All evaluation methods are listed in Table 3, as we foresee Table 3 to be an easy recipe book for researchers to establish their scRNA-seq pipeline. In addition, a summary of all the 25 DL models concerning their DL models, evaluation metrics, implementation environment, downloadable source codes, features, and application notes is presented in Table 1a and 1b. Taken together, this survey provides a rich

resource to select a DL model for proper research applications, and we expect to inspire new DL model developments for scRNA-seq analysis.

One advantage of DL for scRNA-seq repeatedly demonstrated in many of these surveyed papers is DL's ability to scale to a large number of cells, thanks to the stochastic gradient descent algorithm. For imputation, DCA shows linear scalability with the number of cells, and sclGAN and Deeplinpute are demonstrated to scale to 100K cells while non-DL algorithms including SAVER and SCRABBLE fail due to excessive memory usage and runtime [19]. A similar favorable scalability result has been echoed for batch normalization by DESC and iMAP, clustering by scDeepCluster, and multi-functional analysis by scVI, LDVAE, SAUCIE, and scScope. Overall, the advantage of DL in scalability becomes more apparent over non-DL approaches after the number of cells exceeds thousands. However, many of these comparisons exclude the time for determining DL models' hyperparameters. Although iMAP shows that the model is robust against model hyperparameters, determinination of optimal hyperparameters in DL models has not been comprehensively studied for these scRNA-seq tasks.

This review focuses on surveying common DL models, such as AE, VAE, and GAN, and their model variations or combinations to address single-cell data analysis challenges. With the advancement of multi-omics single-cell technologies, new single-cell data types and DL models will be introduced to the single-cell analysis pipeline, such as cyTOF using SAUCIE [15], spatial transcriptome using DNN [122], and CITE-seq that simultaneously generates read counts for surface protein expression along with gene expression [123, 124]. Other than 3 most common unsupervised DL models using AE, VAE, and GAN, this review also discusses supervised network frameworks including

CapsNet (e.g. scCapsNet [52]), CNN (e.g. CNNC [54]), and domain adaption learning (e.g. scDGN [53]). It is expected that more DL models and learning paradigms will be developed and implemented for the most challenging steps for scRNA-seq data, including but not limited to, multi-omics data integration and data interpretation. For example, integrating protein-protein interaction graphs into DL models has been shown for its advantages of biological knowledge and nonlinear interactions embedded in the graphs [125-127]. Indeed, a recently published scRNA-seq analysis pipeline, scGNN [128], incorporates 3 iterative autoencoders (including one graph autoencoder) and successfully demonstrated Alzheimer's disease-related neural development and differentiation mechanisms. We expect that our careful organization of this review will provide a basic understanding of DL models for scRNA-seq and inspire innovative applications of DL models for single cell analysis.

Funding

This article's publication costs were supported partially by the National Institutes of Health (CTSA 1UL1RR025767-01 to YC, R01GM113245 to YH, NCI Cancer Center Shared Resources P30CA54174 to YC, and K99CA248944 to YCC); National Science Foundation (2051113 to YFJ); Cancer Prevention and Research Institute of Texas (RP160732 to YC, RP190346 to YC and YH); and the Fund for Innovation in Cancer Informatics (ICI Fund to YCC and YC). The funding sources had no role in the design of the study; collection, analysis, and interpretation of data; or in writing the manuscript.

Authors' contributions

- 1 YH, YC, MF, and YFJ conceived the study. MF, ZL, TZ, MMH, YCC, ZY, KP, SJ, JZ, SJG,
- 2 YFJ, YC and YH summarized resources, wrote, and approved the final version of the
- 3 paper.

References

- 7 1. Lahnemann D, Koster J, Szczurek E, McCarthy DJ, Hicks SC, Robinson MD, Vallejos CA,
- 8 Campbell KR, Beerenwinkel N, Mahfouz A et al: Eleven grand challenges in single-cell
- 9 data science. Genome Biol 2020, **21**(1):31.
- 10 2. Vitak SA, Torkenczy KA, Rosenkrantz JL, Fields AJ, Christiansen L, Wong MH, Carbone L,
- 11 Steemers FJ, Adey A: **Sequencing thousands of single-cell genomes with combinatorial**
- 12 **indexing**. *Nat Methods* 2017, **14**(3):302-308.
- 13 3. Wu H, Wang C, Wu S: Single-Cell Sequencing for Drug Discovery and Drug
- 14 **Development**. *Curr Top Med Chem* 2017, **17**(15):1769-1777.
- 4. Kinker GS, Greenwald AC, Tal R, Orlova Z, Cuoco MS, McFarland JM, Warren A, Rodman
- 16 C, Roth JA, Bender SA et al: Pan-cancer single-cell RNA-seq identifies recurring
- programs of cellular heterogeneity. *Nat Genet* 2020, **52**(11):1208-1218.
- 18 5. Navin NE: The first five years of single-cell cancer genomics and beyond. Genome Res
- 19 2015, **25**(10):1499-1507.
- 20 6. Suva ML, Tirosh I: Single-Cell RNA Sequencing in Cancer: Lessons Learned and
- 21 **Emerging Challenges**. *Mol Cell* 2019, **75**(1):7-12.
- 22 7. Mannarapu M, Dariya B, Bandapalli OR: Application of single-cell sequencing
- technologies in pancreatic cancer. Mol Cell Biochem 2021, 476(6):2429-2437.
- 24 8. Wauters E, Van Mol P, Garg AD, Jansen S, Van Herck Y, Vanderbeke L, Bassez A, Boeckx
- B, Malengier-Devlies B, Timmerman A et al: Discriminating mild from critical COVID-19
- by innate and adaptive immune single-cell profiling of bronchoalveolar lavages. Cell
- 27 Res 2021, **31**(3):272-290.
- 28 9. Bost P, Giladi A, Liu Y, Bendjelal Y, Xu G, David E, Blecher-Gonen R, Cohen M, Medaglia
- 29 C, Li H et al: Host-Viral Infection Maps Reveal Signatures of Severe COVID-19 Patients.
- 30 *Cell* 2020, **181**(7):1475-1488 e1412.

- 1 10. Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, 3rd, Hao Y, Stoeckius
- M, Smibert P, Satija R: Comprehensive Integration of Single-Cell Data. Cell 2019,
- 3 **177**(7):1888-1902 e1821.
- 4 11. Wolf FA, Angerer P, Theis FJ: SCANPY: large-scale single-cell gene expression data
- 5 **analysis**. *Genome Biol* 2018, **19**(1):15.
- 6 12. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Ł, Polosukhin I:
- 7 **Attention is all you need**. In: Advances in neural information processing systems: 2017.
- 8 5998-6008.
- 9 13. Karpathy A, Toderici G, Shetty S, Leung T, Sukthankar R, Fei-Fei L: Large-scale video
- 10 classification with convolutional neural networks. In: Proceedings of the IEEE
- 11 conference on Computer Vision and Pattern Recognition: 2014. 1725-1732.
- 12 14. Deng L, Liu Y: **Deep learning in natural language processing**: Springer; 2018.
- 13 15. Amodio M, van Dijk D, Srinivasan K, Chen WS, Mohsen H, Moon KR, Campbell A, Zhao Y,
- Wang X, Venkataswamy M et al: Exploring single-cell data with deep multitasking neural
- 15 **networks**. *Nat Methods* 2019, **16**(11):1139-1145.
- 16. Srinivasan S, Leshchyk A, Johnson NT, Korkin D: **A hybrid deep clustering approach for**
- 17 robust cell type profiling using single-cell RNA-seq data. RNA 2020, 26(10):1303-1319.
- 18 17. Lopez R, Regier J, Cole MB, Jordan MI, Yosef N: Deep generative modeling for single-
- 19 **cell transcriptomics**. *Nat Methods* 2018, **15**(12):1053-1058.
- 20 18. Eraslan G, Simon LM, Mircea M, Mueller NS, Theis FJ: Single-cell RNA-seq denoising
- using a deep count autoencoder. *Nat Commun* 2019, **10**(1):390.
- 22 19. Xu Y, Zhang Z, You L, Liu J, Fan Z, Zhou X: sclGANs: single-cell RNA-seg imputation
- using generative adversarial networks. *Nucleic Acids Res* 2020, **48**(15):e85.
- 24 20. Arisdakessian C, Poirion O, Yunits B, Zhu X, Garmire LX: DeepImpute: an accurate, fast,
- 25 and scalable deep neural network method to impute single-cell RNA-seq data. Genome
- 26 Biol 2019, **20**(1):211.
- 27 21. Tran HTN, Ang KS, Chevrier M, Zhang X, Lee NYS, Goh M, Chen J: A benchmark of batch-
- 28 effect correction methods for single-cell RNA sequencing data. Genome Biol 2020,
- 29 **21**(1):12.
- 30 22. Petegrosso R, Li Z, Kuang R: Machine learning and statistical methods for clustering
- 31 single-cell RNA-sequencing data. Brief Bioinform 2020, 21(4):1209-1223.

- 1 23. Abdelaal T, Michielsen L, Cats D, Hoogduin D, Mei H, Reinders MJT, Mahfouz A: A
- 2 comparison of automatic cell identification methods for single-cell RNA sequencing
- 3 data. Genome Biol 2019, 20(1):194.
- 4 24. Wang J, Zou Q, Lin C: A comparison of deep learning-based pre-processing and
- 5 clustering approaches for single-cell RNA sequencing data. Briefings in Bioinformatics
- 6 2021.
- 7 25. Picelli S, Bjorklund AK, Faridani OR, Sagasser S, Winberg G, Sandberg R: Smart-seq2 for
- 8 sensitive full-length transcriptome profiling in single cells. Nat Methods 2013,
- 9 **10**(11):1096-1098.
- 10 26. Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, Tirosh I, Bialas AR,
- 11 Kamitaki N, Martersteck EM et al: Highly Parallel Genome-wide Expression Profiling of
- 12 Individual Cells Using Nanoliter Droplets. Cell 2015, 161(5):1202-1214.
- 13 27. Eisenstein M: Single-cell RNA-seq analysis software providers scramble to offer
- solutions. *Nat Biotechnol* 2020, **38**(3):254-257.
- 15 28. Chen G, Ning B, Shi T: Single-Cell RNA-Seq Technologies and Related Computational
- 16 **Data Analysis**. *Front Genet* 2019, **10**:317.
- 17 29. Vallejos CA, Marioni JC, Richardson S: BASiCS: Bayesian Analysis of Single-Cell
- 18 **Sequencing Data**. *PLoS Comput Biol* 2015, **11**(6):e1004333.
- 19 30. Lun AT, Bach K, Marioni JC: Pooling across cells to normalize single-cell RNA
- sequencing data with many zero counts. Genome Biol 2016, 17:75.
- 21 31. Hafemeister C, Satija R: Normalization and variance stabilization of single-cell RNA-seq
- data using regularized negative binomial regression. *Genome Biol* 2019, **20**(1):296.
- 23 32. Haghverdi L, Lun ATL, Morgan MD, Marioni JC: Batch effects in single-cell RNA-
- sequencing data are corrected by matching mutual nearest neighbors. Nat Biotechnol
- 25 2018, **36**(5):421-427.
- 33. Butler A, Hoffman P, Smibert P, Papalexi E, Satija R: Integrating single-cell transcriptomic
- data across different conditions, technologies, and species. Nat Biotechnol 2018,
- **36**(5):411-420.
- 34. Korsunsky I, Millard N, Fan J, Slowikowski K, Zhang F, Wei K, Baglaenko Y, Brenner M, Loh
- 30 PR, Raychaudhuri S: Fast, sensitive and accurate integration of single-cell data with
- 31 **Harmony**. *Nat Methods* 2019, **16**(12):1289-1296.

- 1 35. Peng T, Zhu Q, Yin P, Tan K: **SCRABBLE: single-cell RNA-seq imputation constrained**
- 2 by bulk RNA-seq data. Genome Biol 2019, **20**(1):88.
- 3 36. Huang M, Wang J, Torre E, Dueck H, Shaffer S, Bonasio R, Murray JI, Raj A, Li M, Zhang
- 4 NR: **SAVER**: gene expression recovery for single-cell RNA sequencing. Nat Methods
- 5 2018, **15**(7):539-542.
- 6 37. Li WV, Li JJ: An accurate and robust imputation method scimpute for single-cell RNA-
- 7 seq data. Nat Commun 2018, 9(1):997.
- 8 38. Roweis ST, Saul LK: Nonlinear dimensionality reduction by locally linear embedding.
- 9 Science 2000, **290**(5500):2323-2326.
- 10 39. Welch JD, Hartemink AJ, Prins JF: SLICER: inferring branched, nonlinear cellular
- trajectories from single cell RNA-seq data. *Genome Biol* 2016, **17**(1):106.
- 12 40. Linderman GC, Rachh M, Hoskins JG, Steinerberger S, Kluger Y: Fast interpolation-based
- 13 t-SNE for improved visualization of single-cell RNA-seq data. Nat Methods 2019,
- 14 **16**(3):243-245.
- 15 41. Becht E, McInnes L, Healy J, Dutertre CA, Kwok IWH, Ng LG, Ginhoux F, Newell EW:
- Dimensionality reduction for visualizing single-cell data using UMAP. Nat Biotechnol
- 17 2018.
- 18 42. Subelj L, Bajec M: Unfolding communities in large complex networks: combining
- defensive and offensive label propagation for core extraction. Phys Rev E Stat Nonlin
- 20 Soft Matter Phys 2011, **83**(3 Pt 2):036103.
- 21 43. Traag VA, Waltman L, van Eck NJ: From Louvain to Leiden: guaranteeing well-
- 22 **connected communities**. *Sci Rep* 2019, **9**(1):5233.
- 23 44. Wang B, Zhu J, Pierson E, Ramazzotti D, Batzoglou S: Visualization and analysis of
- single-cell RNA-seq data by kernel-based similarity learning. Nat Methods 2017,
- **14**(4):414-416.
- 26 45. Finak G, McDavid A, Yajima M, Deng J, Gersuk V, Shalek AK, Slichter CK, Miller HW,
- 27 McElrath MJ, Prlic M et al: MAST: a flexible statistical framework for assessing
- 28 transcriptional changes and characterizing heterogeneity in single-cell RNA
- 29 **sequencing data**. *Genome Biol* 2015, **16**:278.
- 30 46. Kharchenko PV, Silberstein L, Scadden DT: Bayesian approach to single-cell differential
- 31 **expression analysis**. *Nat Methods* 2014, **11**(7):740-742.

- 1 47. Miao Z, Deng K, Wang X, Zhang X: **DEsingle for detecting three types of differential**
- 2 expression in single-cell RNA-seq data. *Bioinformatics* 2018, **34**(18):3223-3224.
- 3 48. Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, Bengio
- 4 Y: Generative adversarial networks. *Communications of the ACM* 2020, **63**(11):139-144.
- 5 49. Gulrajani I, Ahmed F, Arjovsky M, Dumoulin V, Courville A: Improved training of
- 6 wasserstein gans. arXiv preprint arXiv:170400028 2017.
- 7 50. Arjovsky M, Chintala S, Bottou L: Wasserstein gan. arXiv 2017. arXiv preprint
- 8 arXiv:170107875 2017, **30**.
- 9 51. Torroja C, Sanchez-Cabo F: DigitaldIsorter: Deep-Learning on scRNA-Seq to
- Deconvolute Gene Expression Data. Front Genet 2019, **10**:978.
- 11 52. Wang L, Nie, R., Yu, Z. et al.: An interpretable deep-learning architecture of capsule
- 12 networks for identifying cell-type gene expression programs from single-cell RNA-
- 13 **sequencing data**. *Nat Mach Intell* 2020, **2**:693-703.
- 14 53. Ge S, Wang H, Alavi A, Xing E, Bar-Joseph Z: Supervised Adversarial Alignment of
- 15 **Single-Cell RNA-seq Data**. *J Comput Biol* 2021, **28**(5):501-513.
- 16 54. Yuan Y, Bar-Joseph Z: Deep learning for inferring gene relationships from single-cell
- 17 **expression data**. *Proc Natl Acad Sci U S A* 2019.
- 18 55. Eraslan G, Avsec Z, Gagneur J, Theis FJ: Deep learning: new computational modelling
- 19 **techniques for genomics**. *Nat Rev Genet* 2019, **20**(7):389-403.
- 20 56. Patel ND, Nguang SK, Coghill GG: **Neural network implementation using bit streams**.
- 21 *IEEE Trans Neural Netw* 2007, **18**(5):1488-1504.
- 22 57. van Dijk D, Sharma R, Nainys J, Yim K, Kathail P, Carr AJ, Burdziak C, Moon KR, Chaffer
- 23 CL, Pattabiraman D et al: Recovering Gene Interactions from Single-Cell Data Using
- 24 **Data Diffusion**. *Cell* 2018, **174**(3):716-729 e727.
- 25 58. Wang J, Agarwal D, Huang M, Hu G, Zhou Z, Ye C, Zhang NR: **Data denoising with transfer**
- learning in single-cell transcriptomics. *Nat Methods* 2019, **16**(9):875-878.
- 27 59. Badsha MB, Li R, Liu B, Li YI, Xian M, Banovich NE, Fu AQ: Imputation of single-cell gene
- expression with an autoencoder neural network. Quant Biol 2020, 8(1):78-94.
- 29 60. Yu B, Chen C, Qi R, Zheng R, Skillman-Lawrence PJ, Wang X, Ma A, Gu H: scGMAI: a
- Gaussian mixture model for clustering single-cell RNA-Seq data based on deep
- 31 **autoencoder**. *Brief Bioinform* 2020.

- 1 61. Lin P, Troup M, Ho JW: CIDR: Ultrafast and accurate clustering through imputation for single-cell RNA-seq data. *Genome Biol* 2017, **18**(1):59.
- 3 62. Berthelot D, Schumm, T. and Metz, L.: **BEGAN: Boundary Equilibrium Generative**4 **Adversarial Networks**. *arXiv* 2017.
- 5 63. Berthelot D, Schumm T, Metz L: **Began: Boundary equilibrium generative adversarial**6 **networks**. *arXiv preprint arXiv:170310717* 2017.
- 7 64. Wang T, Johnson TS, Shao W, Lu Z, Helm BR, Zhang J, Huang K: **BERMUDA: a novel**
- 8 deep transfer learning method for single-cell RNA sequencing batch correction
- 9 reveals hidden high-resolution cellular subtypes. *Genome Biol* 2019, **20**(1):165.
- 10 65. Borgwardt KM, Gretton A, Rasch MJ, Kriegel HP, Scholkopf B, Smola AJ: Integrating
- structured biological data by Kernel Maximum Mean Discrepancy. *Bioinformatics* 2006,
- 12 **22**(14):e49-57.
- 13 66. Crow M, Paul A, Ballouz S, Huang ZJ, Gillis J: Characterizing the replicability of cell types
- defined by single cell RNA-sequencing data using MetaNeighbor. Nat Commun 2018,
- 15 **9**(1):884.
- 16 67. Polanski K, Young MD, Miao Z, Meyer KB, Teichmann SA, Park JE: **BBKNN: fast batch**
- alignment of single cell transcriptomes. *Bioinformatics* 2020, **36**(3):964-965.
- 18 68. Li X, Wang K, Lyu Y, Pan H, Zhang J, Stambolian D, Susztak K, Reilly MP, Hu G, Li M: **Deep**
- 19 learning enables accurate clustering with batch effect removal in single-cell RNA-seq
- 20 **analysis**. *Nat Commun* 2020, **11**(1):2338.
- 21 69. Guo X, Gao, L., Liu, X., and Yin, J.: Improved deep embedded clustering with local
- structure preservation. Proc 26th International Joint Conference on Artificial Integlligence
- 23 2017:1753-1759.
- 24 70. Hie B, Bryson B, Berger B: Efficient integration of heterogeneous single-cell
- transcriptomes using Scanorama. *Nat Biotechnol* 2019, **37**(6):685-691.
- 26 71. Wang D, Hou S, Zhang L, Wang X, Liu B, Zhang Z: iMAP: integration of multiple single-
- cell datasets by adversarial paired transfer networks. *Genome Biol* 2021, **22**(1):63.
- 28 72. Lin C, Jain S, Kim H, Bar-Joseph Z: **Using neural networks for reducing the dimensions**
- of single-cell RNA-Seq data. Nucleic Acids Res 2017, 45(17):e156.
- 30 73. Rashid S, Shah S, Bar-Joseph Z, Pandya R: Dhaka: Variational Autoencoder for
- 31 Unmasking Tumor Heterogeneity from Single Cell Genomic Data. *Bioinformatics* 2019.

- 1 74. Tirosh I, Izar B, Prakadan SM, Wadsworth MH, 2nd, Treacy D, Trombetta JJ, Rotem A.
- 2 Rodman C, Lian C, Murphy G et al: Dissecting the multicellular ecosystem of metastatic
- 3 melanoma by single-cell RNA-seq. *Science* 2016, **352**(6282):189-196.
- 4 75. Zahn H, Steif A, Laks E, Eirew P, VanInsberghe M, Shah SP, Aparicio S, Hansen CL:
- 5 Scalable whole-genome single-cell library preparation without preamplification.
- 6 Nature Methods 2017, **14**(2):167-173.
- 7 76. Ding J, Condon A, Shah SP: Interpretable dimensionality reduction of single cell
- 8 transcriptome data with deep generative models. *Nat Commun* 2018, **9**(1):2002.
- 9 77. Gronbech CH, Vording MF, Timshel PN, Sonderby CK, Pers TH, Winther O: scVAE:
- variational auto-encoders for single-cell gene expression data. Bioinformatics 2020,
- **36**(16):4415-4422.
- 12 78. Wang D, Gu J: VASC: Dimension Reduction and Visualization of Single-cell RNA-seq
- Data by Deep Variational Autoencoder. Genomics Proteomics Bioinformatics 2018,
- 14 **16**(5):320-331.
- 15 79. Jang E. GSaPB: Categorical reparameterization with gumbel-softmax. arXiv 2016.
- 16 80. Tian T, Wan, J., Song, Q. et al.: Clustering single-cell RNA-seq data with a model-based
- deep learning approach. Nat Mach Intell 2019, 1.
- 18 81. Regev A, Teichmann SA, Lander ES, Amit I, Benoist C, Birney E, Bodenmiller B, Campbell
- 19 P, Carninci P, Clatworthy M et al: The Human Cell Atlas. Elife 2017, 6.
- 20 82. Xie J, Girshick R, Farhadi A: Unsupervised deep embedding for clustering analysis. In:
- 21 International conference on machine learning: 2016. PMLR: 478-487.
- 22 83. Marouf M, Machart P, Bansal V, Kilian C, Magruder DS, Krebs CF, Bonn S: Realistic in
- 23 silico generation and augmentation of single-cell RNA-seq data using generative
- adversarial networks. *Nat Commun* 2020, **11**(1):166.
- 25 84. Miyato TaK, M: **cGANs with projection discriminator**. *Preprint* 2018.
- 26 85. Zappia L, Phipson B, Oshlack A: **Splatter: simulation of single-cell RNA sequencing data**.
- 27 Genome Biol 2017, **18**(1):174.
- 28 86. Lindenbaum O, Stanley, J. S., Wolf, G. and Krishnaswamy, S.: Geometry-based data
- 29 **generation**. Advances in Neural Information Processing Systems 2018.
- 30 87. Svensson V, Gayoso A, Yosef N, Pachter L: Interpretable factor models of single-cell
- 31 RNA-seq via variational autoencoders. *Bioinformatics* 2020, **36**(11):3418-3421.

- 1 88. Levine JH, Simonds EF, Bendall SC, Davis KL, Amir el AD, Tadmor MD, Litvin O, Fienberg
- 2 HG, Jager A, Zunder ER et al: Data-Driven Phenotypic Dissection of AML Reveals
- 3 **Progenitor-like Cells that Correlate with Prognosis**. *Cell* 2015, **162**(1):184-197.
- 4 89. Qiu X, Mao Q, Tang Y, Wang L, Chawla R, Pliner HA, Trapnell C: Reversed graph
- 5 **embedding resolves complex single-cell trajectories**. *Nat Methods* 2017, **14**(10):979-982.
- 6 90. McInnes L, Healy, J. & Melville, J.: Umap: uniform manifold approximation and
- 7 projection for dimension reduction. *ArXiv* 2018.
- 8 91. van der Maaten LH, G.: Visualizing data using t-SNE. J Mach Learn 2008, 9:2579-2605.
- 9 92. Moon KRea: PHATE: a dimensionality reduction method for visualizing trajectory
- structures in high-dimensional biological data. bioRxiv 2017.
- 11 93. Deng Y, Bao F, Dai Q, Wu LF, Altschuler SJ: **Scalable analysis of cell-type composition**
- from single-cell transcriptomics using deep recurrent learning. Nat Methods 2019,
- **16**(4):311-314.
- 14 94. Chung W, Eum HH, Lee HO, Lee KM, Lee HB, Kim KT, Ryu HS, Kim S, Lee JE, Park YH et
- 15 al: Single-cell RNA-seq enables comprehensive tumour and immune cell profiling in
- primary breast cancer. *Nat Commun* 2017, **8**:15081.
- 17 95. Li H, Courtois ET, Sengupta D, Tan Y, Chen KH, Goh JJL, Kong SL, Chua C, Hon LK, Tan
- 18 WS et al: Reference component analysis of single-cell transcriptomes elucidates
- cellular heterogeneity in human colorectal tumors. *Nat Genet* 2017, **49**(5):708-718.
- 20 96. Aran D, Hu Z, Butte AJ: xCell: digitally portraying the tissue cellular heterogeneity
- 21 **landscape**. *Genome Biol* 2017, **18**(1):220.
- 22 97. Yoshihara K, Shahmoradgoli M, Martinez E, Vegesna R, Kim H, Torres-Garcia W, Trevino
- V, Shen H, Laird PW, Levine DA et al: Inferring tumour purity and stromal and immune
- 24 cell admixture from expression data. *Nat Commun* 2013, **4**:2612.
- 25 98. Racle J, de Jonge K, Baumgaertner P, Speiser DE, Gfeller D: Simultaneous enumeration
- of cancer and immune cell types from bulk tumor gene expression data. *Elife* 2017, 6.
- 99. Becht E, Giraldo NA, Lacroix L, Buttard B, Elarouci N, Petitprez F, Selves J, Laurent-Puig P,
- Sautes-Fridman C, Fridman WH et al: Estimating the population abundance of tissue-
- infiltrating immune and stromal cell populations using gene expression. Genome Biol
- 30 2016, **17**(1):218.

- 1 100. Zheng GX, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, Ziraldo SB, Wheeler TD,
- 2 McDermott GP, Zhu J et al: Massively parallel digital transcriptional profiling of single
- 3 **cells**. *Nat Commun* 2017, **8**:14049.
- 4 101. Shekhar K, Lapan SW, Whitney IE, Tran NM, Macosko EZ, Kowalczyk M, Adiconis X, Levin
- 5 JZ, Nemesh J, Goldman M et al: Comprehensive Classification of Retinal Bipolar
- 6 Neurons by Single-Cell Transcriptomics. *Cell* 2016, **166**(5):1308-1323 e1330.
- 7 102. Dong Z, Alterovitz G: netAE: semi-supervised dimensionality reduction of single-cell
- 8 RNA sequencing to facilitate cell labeling. *Bioinformatics* 2021, **37**(1):43-49.
- 9 103. Newman ME: Modularity and community structure in networks. Proc Natl Acad Sci U S
- 10 *A* 2006, **103**(23):8577-8582.
- 11 104.Xu C, Lopez R, Mehlman E, Regier J, Jordan MI, Yosef N: **Probabilistic harmonization**
- and annotation of single-cell transcriptomics data with deep generative models. *Mol*
- 13 Syst Biol 2021, **17**(1):e9620.
- 14 105. Hadsell R, Chopra S, LeCun Y: **Dimensionality reduction by learning an invariant**
- mapping. In: 2006 IEEE Computer Society Conference on Computer Vision and Pattern
- 16 Recognition (CVPR'06): 2006. IEEE: 1735-1742.
- 17 106.Lieberman Y, Rokach L, Shay T: CaSTLe-classification of single cells by transfer
- learning: harnessing the power of publicly available single cell RNA sequencing
- experiments to annotate new experiments. *PloS one* 2018, **13**(10):e0205499.
- 20 107. Alavi A, Ruffalo M, Parvangada A, Huang Z, Bar-Joseph Z: A web server for comparative
- analysis of single-cell RNA-seq data. *Nat Commun* 2018, **9**(1):4768.
- 22 108. Yevshin I, Sharipov R, Valeev T, Kel A, Kolpakov F: GTRD: a database of transcription
- factor binding sites identified by ChIP-seg experiments. Nucleic Acids Research 2016,
- 24 **45**(D1):D61-D67.
- 25 109. Wang YXR, Waterman MS, Huang HY: Gene coexpression measures in large
- 26 heterogeneous samples using count statistics. P Natl Acad Sci USA 2014,
- 27 **111**(46):16371-16376.
- 28 110. Krishnaswamy S, Spitzer MH, Mingueneau M, Bendall SC, Litvin O, Stone E, Pe'er D, Nolan
- 29 GP: Systems biology. Conditional density-based analysis of T cell signaling in single-
- 30 **cell data**. *Science* 2014, **346**(6213):1250689.
- 31 111. Kanehisa M, Furumichi M, Tanabe M, Sato Y, Morishima K: **KEGG: new perspectives on**
- 32 genomes, pathways, diseases and drugs. *Nucleic Acids Res* 2017, **45**(D1):D353-D361.

- 1 112. Fabregat A, Jupe S, Matthews L, Sidiropoulos K, Gillespie M, Garapati P, Haw R, Jassal B,
- 2 Korninger F, May B et al: The Reactome Pathway Knowledgebase. Nucleic Acids Res
- 3 2018, **46**(D1):D649-D655.
- 4 113. Huynh-Thu VA, Irrthum A, Wehenkel L, Geurts P: Inferring regulatory networks from
- 5 **expression data using tree-based methods**. *PloS one* 2010, **5**(9):e12776.
- 6 114. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A,
- 7 Pomeroy SL, Golub TR, Lander ES et al: Gene set enrichment analysis: a knowledge-
- 8 based approach for interpreting genome-wide expression profiles. Proc Natl Acad Sci
- 9 *USA* 2005, **102**(43):15545-15550.
- 10 115. Lotfollahi M, Wolf FA, Theis FJ: scGen predicts single-cell perturbation responses. Nat
- 11 *Methods* 2019, **16**(8):715-721.
- 12 116. Kang HM, Subramaniam M, Targ S, Nguyen M, Maliskova L, McCarthy E, Wan E, Wong S,
- Byrnes L, Lanata CM et al: Multiplexed droplet single-cell RNA-sequencing using
- natural genetic variation. *Nat Biotechnol* 2018, **36**(1):89-94.
- 15 117. Duvenaud D, Maclaurin D, Iparraguirre J, Bombarell R, Hirzel T, Aspuru-Guzik A, Adams R:
- Advances in Neural Information Processing Systems 28. Cortes C, Lawrence ND, Lee
- 17 DD, Sugiyama M, Garnett R, Eds 2015:2224-2232.
- 18 118. Amodio M, Krishnaswamy S: **MAGAN: Aligning biological manifolds**. In: *International*
- 19 Conference on Machine Learning: 2018. PMLR: 215-223.
- 20 119. Ghahramani A, Watt FM, Luscombe NM: Generative adversarial networks simulate gene
- expression and predict perturbations in single cells. *bioRxiv* 2018:262501.
- 22 120. Haber AL, Biton M, Rogel N, Herbst RH, Shekhar K, Smillie C, Burgin G, Delorey TM, Howitt
- MR, Katz Y et al: A single-cell survey of the small intestinal epithelium. Nature 2017,
- **551**(7680):333-339.
- 25 121. Hagai T, Chen X, Miragaia RJ, Rostom R, Gomes T, Kunowska N, Henriksson J, Park JE,
- 26 Proserpio V, Donati G et al: Gene expression variability across cells and species shapes
- 27 **innate immunity**. *Nature* 2018, **563**(7730):197-202.
- 28 122. Maseda F, Cang Z, Nie Q: DEEPsc: A Deep Learning-Based Map Connecting Single-
- 29 **Cell Transcriptomics and Spatial Imaging Data**. *Front Genet* 2021, **12**:636743.
- 30 123. Musu Y, Liang C, Deng M: Clustering single cell CITE-seq data with a canonical
- 31 correlation based deep learning method. *bioRxiv* 2021.

- 1 124.Zhou Z, Ye C, Wang J, Zhang NR: Surface protein imputation from single cell transcriptomes by deep neural networks. *Nat Commun* 2020, **11**(1):651.
- 3 125. Ramirez R, Chiu Y-C, Hererra A, Mostavi M, Ramirez J, Chen Y, Huang Y, Jin Y-F:
- 4 Classification of Cancer Types Using Graph Convolutional Neural Networks. Frontiers
- 5 in Physics 2020, **8**(203).
- 6 126. Ramirez R, Chiu YC, Zhang S, Ramirez J, Chen Y, Huang Y, Jin YF: Prediction and
- 7 interpretation of cancer survival using graph convolution neural networks. Methods
- 8 2021, **192**:120-130.
- 9 127. Battaglia PW, Hamrick JB, Bapst V, Sanchez-Gonzalez A, Zambaldi V, Malinowski M,
- Tacchetti A, Raposo D, Santoro A, Faulkner R: **Relational inductive biases, deep learning,**
- and graph networks. arXiv preprint arXiv:180601261 2018.
- 12 128. Wang J, Ma A, Chang Y, Gong J, Jiang Y, Qi R, Wang C, Fu H, Ma Q, Xu D: scGNN is a
- 13 novel graph neural network framework for single-cell RNA-Seq analyses. *Nat Commun*
- 14 2021, **12**(1):1882.
- 15 129.Peng Y, Baulier E, Ke Y, Young A, Ahmedli NB, Schwartz SD, Farber DB: Human
- embryonic stem cells extracellular vesicles and their effects on immortalized human
- 17 **retinal Muller cells**. *PLoS One* 2018, **13**(3):e0194004.
- 130. Stoeckius M, Hafemeister C, Stephenson W, Houck-Loomis B, Chattopadhyay PK, Swerdlow
- 19 H, Satija R, Smibert P: Simultaneous epitope and transcriptome measurement in single
- 20 **cells**. *Nat Methods* 2017, **14**(9):865-868.
- 21 131.La Manno G, Gyllborg D, Codeluppi S, Nishimura K, Salto C, Zeisel A, Borm LE, Stott SRW,
- Toledo EM, Villaescusa JC et al: Molecular Diversity of Midbrain Development in Mouse,
- 23 Human, and Stem Cells. Cell 2016, 167(2):566-580 e519.
- 24 132. Azizi E, Carr AJ, Plitas G, Cornish AE, Konopacki C, Prabhakaran S, Nainys J, Wu K,
- 25 Kiseliovas V, Setty M et al: Single-Cell Map of Diverse Immune Phenotypes in the Breast
- 26 **Tumor Microenvironment**. *Cell* 2018, **174**(5):1293-1308 e1236.
- 27 133. Chu LF, Leng N, Zhang J, Hou Z, Mamott D, Vereide DT, Choi J, Kendziorski C, Stewart R,
- Thomson JA: Single-cell RNA-seq reveals novel regulators of human embryonic stem
- 29 cell differentiation to definitive endoderm. Genome Biol 2016, 17(1):173.
- 30 134. Baron M, Veres A, Wolock SL, Faust AL, Gaujoux R, Vetere A, Ryu JH, Wagner BK, Shen-
- 31 Orr SS, Klein AM et al: A Single-Cell Transcriptomic Map of the Human and Mouse

- Pancreas Reveals Inter- and Intra-cell Population Structure. Cell Syst 2016, 3(4):346-
- 2 360 e344.
- 3 135. Camp JG, Sekine K, Gerber T, Loeffler-Wirth H, Binder H, Gac M, Kanton S, Kageyama J,
- 4 Damm G, Seehofer D et al: Multilineage communication regulates human liver bud
- 5 **development from pluripotency**. *Nature* 2017, **546**(7659):533-538.
- 6 136. Muraro MJ, Dharmadhikari G, Grun D, Groen N, Dielen T, Jansen E, van Gurp L, Engelse
- 7 MA, Carlotti F, de Koning EJ et al: A Single-Cell Transcriptome Atlas of the Human
- 8 **Pancreas**. *Cell Syst* 2016, **3**(4):385-394 e383.
- 9 137. Darmanis S, Sloan SA, Zhang Y, Enge M, Caneda C, Shuer LM, Hayden Gephart MG, Barres
- BA, Quake SR: A survey of human brain transcriptome diversity at the single cell level.
- 11 Proc Natl Acad Sci U S A 2015, **112**(23):7285-7290.
- 12 138. Tirosh I, Venteicher AS, Hebert C, Escalante LE, Patel AP, Yizhak K, Fisher JM, Rodman C,
- Mount C, Filbin MG et al: Single-cell RNA-seq supports a developmental hierarchy in
- 14 **human oligodendroglioma**. *Nature* 2016, **539**(7628):309-313.
- 15 139. Patel AP, Tirosh I, Trombetta JJ, Shalek AK, Gillespie SM, Wakimoto H, Cahill DP, Nahed
- BV, Curry WT, Martuza RL et al: Single-cell RNA-seq highlights intratumoral
- heterogeneity in primary glioblastoma. Science 2014, **344**(6190):1396-1401.
- 18 140. Zahn H, Steif A, Laks E, Eirew P, VanInsberghe M, Shah SP, Aparicio S, Hansen CL:
- 19 Scalable whole-genome single-cell library preparation without preamplification. *Nat*
- 20 *Methods* 2017, **14**(2):167-173.
- 21 141. Pollen AA, Nowakowski TJ, Shuga J, Wang X, Leyrat AA, Lui JH, Li N, Szpankowski L,
- Fowler B, Chen P et al: Low-coverage single-cell mRNA sequencing reveals cellular
- 23 heterogeneity and activated signaling pathways in developing cerebral cortex. Nat
- 24 Biotechnol 2014, **32**(10):1053-1058.
- 25 142.Xin Y, Kim J, Okamoto H, Ni M, Wei Y, Adler C, Murphy AJ, Yancopoulos GD, Lin C,
- Gromada J: RNA Sequencing of Single Human Islet Cells Reveals Type 2 Diabetes
- 27 **Genes**. Cell Metab 2016, **24**(4):608-615.
- 28 143. Yan L, Yang M, Guo H, Yang L, Wu J, Li R, Liu P, Lian Y, Zheng X, Yan J et al: Single-cell
- 29 RNA-Seq profiling of human preimplantation embryos and embryonic stem cells. *Nat*
- 30 Struct Mol Biol 2013, **20**(9):1131-1139.

- 1 144. Chevrier S, Levine JH, Zanotelli VRT, Silina K, Schulz D, Bacac M, Ries CH, Ailles L, Jewett
- 2 MAS, Moch H et al: An Immune Atlas of Clear Cell Renal Cell Carcinoma. Cell 2017,
- 3 **169**(4):736-749 e718.
- 4 145. Han X, Wang R, Zhou Y, Fei L, Sun H, Lai S, Saadatpour A, Zhou Z, Chen H, Ye F et al:
- 5 **Mapping the Mouse Cell Atlas by Microwell-Seq.** *Cell* 2018, **172**(5):1091-1107 e1017.
- 6 146. Hrvatin S, Hochbaum DR, Nagy MA, Cicconet M, Robertson K, Cheadle L, Zilionis R, Ratner
- 7 A, Borges-Monroy R, Klein AM et al: Single-cell analysis of experience-dependent
- 8 transcriptomic states in the mouse visual cortex. *Nat Neurosci* 2018, **21**(1):120-129.
- 9 147. Joost S, Zeisel A, Jacob T, Sun X, La Manno G, Lonnerberg P, Linnarsson S, Kasper M:
- 10 Single-Cell Transcriptomics Reveals that Differentiation and Spatial Signatures Shape
- Epidermal and Hair Follicle Heterogeneity. Cell Syst 2016, 3(3):221-237 e229.
- 12 148. Paul F, Arkin Y, Giladi A, Jaitin DA, Kenigsberg E, Keren-Shaul H, Winter D, Lara-Astiaso D,
- Gury M, Weiner A et al: Transcriptional Heterogeneity and Lineage Commitment in
- 14 **Myeloid Progenitors**. *Cell* 2015, **163**(7):1663-1677.
- 15 149. Buettner F, Natarajan KN, Casale FP, Proserpio V, Scialdone A, Theis FJ, Teichmann SA,
- Marioni JC, Stegle O: Computational analysis of cell-to-cell heterogeneity in single-cell
- 17 RNA-sequencing data reveals hidden subpopulations of cells. Nat Biotechnol 2015,
- 18 **33**(2):155-160.
- 19 150. Biase FH, Cao X, Zhong S: **Cell fate inclination within 2-cell and 4-cell mouse embryos**
- revealed by single-cell RNA sequencing. Genome Res 2014, 24(11):1787-1796.
- 21 151. Deng Q, Ramskold D, Reinius B, Sandberg R: Single-cell RNA-seg reveals dynamic,
- random monoallelic gene expression in mammalian cells. Science 2014, 343(6167):193-
- 23 196.
- 24 152. Klein AM, Mazutis L, Akartuna I, Tallapragada N, Veres A, Li V, Peshkin L, Weitz DA,
- 25 Kirschner MW: Droplet barcoding for single-cell transcriptomics applied to embryonic
- 26 **stem cells**. *Cell* 2015, **161**(5):1187-1201.
- 27 153. Goolam M, Scialdone A, Graham SJL, Macaulay IC, Jedrusik A, Hupalowska A, Voet T,
- Marioni JC, Zernicka-Goetz M: Heterogeneity in Oct4 and Sox2 Targets Biases Cell Fate
- 29 in **4-Cell Mouse Embryos**. *Cell* 2016, **165**(1):61-74.
- 30 154.Kim JK, Kolodziejczyk AA, Ilicic T, Teichmann SA, Marioni JC: Characterizing noise
- 31 structure in single-cell RNA-seq distinguishes genuine from technical stochastic
- allelic expression. *Nat Commun* 2015, **6**:8687.

- 1 155. Usoskin D, Furlan A, Islam S, Abdo H, Lonnerberg P, Lou D, Hjerling-Leffler J, Haeggstrom
- J, Kharchenko O, Kharchenko PV et al: Unbiased classification of sensory neuron types
- 3 by large-scale single-cell RNA sequencing. *Nat Neurosci* 2015, **18**(1):145-153.
- 4 156. Zeisel A, Munoz-Manchado AB, Codeluppi S, Lonnerberg P, La Manno G, Jureus A,
- 5 Marques S, Munguba H, He L, Betsholtz C et al: Brain structure. Cell types in the mouse
- 6 cortex and hippocampus revealed by single-cell RNA-seq. Science 2015,
- 7 **347**(6226):1138-1142.
- 8 157. Yu Z, Liao J, Chen Y, Zou C, Zhang H, Cheng J, Liu D, Li T, Zhang Q, Li J et al: Single-Cell
- 9 Transcriptomic Map of the Human and Mouse Bladders. J Am Soc Nephrol 2019,
- 10 **30**(11):2159-2176.
- 11 158. Tusi BK, Wolock SL, Weinreb C, Hwang Y, Hidalgo D, Zilionis R, Waisman A, Huh JR, Klein
- 12 AM, Socolovsky M: Population snapshots predict early haematopoietic and erythroid
- 13 **hierarchies**. *Nature* 2018, **555**(7694):54-60.
- 14 159. Pijuan-Sala B, Griffiths JA, Guibentif C, Hiscock TW, Jawaid W, Calero-Nieto FJ, Mulas C,
- 15 Ibarra-Soria X, Tyser RCV, Ho DLL et al: A single-cell molecular map of mouse
- gastrulation and early organogenesis. *Nature* 2019, **566**(7745):490-495.
- 17 160. Cao J, Spielmann M, Qiu X, Huang X, Ibrahim DM, Hill AJ, Zhang F, Mundlos S, Christiansen
- 18 L, Steemers FJ et al: The single-cell transcriptional landscape of mammalian
- organogenesis. *Nature* 2019, **566**(7745):496-502.
- 20 161. Setty M, Tadmor MD, Reich-Zeliger S, Angel O, Salame TM, Kathail P, Choi K, Bendall S,
- Friedman N, Pe'er D: Wishbone identifies bifurcating developmental trajectories from
- 22 **single-cell data**. *Nat Biotechnol* 2016, **34**(6):637-645.
- 23 162. Nestorowa S, Hamey FK, Pijuan Sala B, Diamanti E, Shepherd M, Laurenti E, Wilson NK,
- 24 Kent DG, Gottgens B: A single-cell resolution map of mouse hematopoietic stem and
- progenitor cell differentiation. *Blood* 2016, **128**(8):e20-31.
- 163. Cao J, Packer JS, Ramani V, Cusanovich DA, Huynh C, Daza R, Qiu X, Lee C, Furlan SN,
- 27 Steemers FJ et al: Comprehensive single-cell transcriptional profiling of a multicellular
- 28 **organism**. *Science* 2017, **357**(6352):661-667.
- 29 164. Strehland AaG, J.: Cluster ensembles---a knowledge reuse framework for combining
- 30 multiple partitions. J Mach Learn Res 2002, **3**:583-617.
- 31 165. McDaid AF, Greene D, Hurley N: Normalized mutual information to evaluate overlapping
- 32 community finding algorithms. arXiv preprint arXiv:11102515 2011.

- 1 166. MacKay DJ, Mac Kay DJ: Information theory, inference and learning algorithms:
- 2 Cambridge university press; 2003.
- 3 167. Hubert LaA, P.: **Comparing Partitions**. *Journal of Classification* 1985, **2**:193-218.
- 4 168. Buttner M, Miao Z, Wolf FA, Teichmann SA, Theis FJ: **A test metric for assessing single-**5 **cell RNA-seq batch correction**. *Nat Methods* 2019, **16**(1):43-49.
- 6 169. Rosenberg A, Hirschberg J: V-measure: A conditional entropy-based external cluster
- 7 **evaluation measure**. In: Proceedings of the 2007 joint conference on empirical methods in
- 8 natural language processing and computational natural language learning (EMNLP-CoNLL):
- 9 2007. 410-420.

1 Figure Captions

- 2 Figure 1. Single cell data analysis steps for both conventional ML methods
- 3 (bottom) and DL methods (top). Depending on the input data and analysis objectives,
- 4 major scRNA-seq analysis steps are illustrated in the center flow chart (color boxes) with
- 5 conventional ML approaches along with optional analysis modules below each analysis
- 6 step. Deep learning approaches are categorized as Deep Neural Network, Generative
- 7 Adversarial Network, Variational Autoencoder, and Autoencoder. For each DL approach,
- 8 optional algorithms are listed on top of each step.

9

- 10 Figure 2. Graphical models of the major surveyed DL models including A)
- 11 Variational Autoencoder **B)** Autoencoder; and **C)** Generative Adversarial Network

12

- 13 **Figure 3**. **Algorithm comparison grid.** DL methods surveyed in the paper are listed on
- the left-hand side, and some in the column. Algorithms selected to compare in each DL
- method are marked by "■" at each cross-point.

16

- 17 **Figure 4. DL model network illustration. A)** Deep neural network, **B)** Autoencoder, **C)**
- 18 Variational autoencoder, **D**) Autoencoder with recursive imputer, **E**) Generative
- adversarial network, **F**) Convolutional neural network, and **G**) Domain adversarial neural
- 20 network.

21

Tables

 Table 1a.
 Deep Learning algorithms reviewed in the paper

Арр	Algorithm	Models	Evaluation	Environment	Codes	Refs
Imputation						
•	DCA	AE	DREMI	Keras, Tensorflow, scanpy	https://github.com/theislab/d	[18]
	SAVER-X	AE+TL	t-SNE, ARI	R/sctransfer	https://github.com/jingshuw/ SAVERX	[58]
	DeepImpute	DNN	MSE, Pearson's correlation	Keras/Tensorf low	https://github.com/lanagarm ire/DeepImpute	[20]
	LATE	AE	MSE	Tensorflow	https://github.com/audreyqy fu/LATE	[59]
	scGAMI	AE	NMI, ARI, HS and CS	Tensorflow	https://github.com/QUST- AIBBDRC/scGMAI/	[60]
	sclGANs	GAN	ARI, ACC, AUC, and F-score	PyTorch	https://github.com/xuyungan g/scIGANs	[19]
Batch corr	ection					
	BERMUDA	AE+TL	kBET, the entropy of Mixing, SI	PyTorch	https://github.com/txWang/B ERMUDA	[64]
	DESC	AE	ARI, KL	Tensorflow	https://github.com/eleozzr/desc	[68]
	iMAP	AE+GAN	kBET, LISI	PyTorch	https://github.com/Svvord/i MAP	[71]
Clustering,	, latent representati	on, dimension i	reduction, and data augmentation			
	Dhaka	VAE	ARI, Spearman Correlation	Keras/Tensorf low	https://github.com/Microsoft Genomics/Dhaka	[73]
	scvis	VAE	KNN preservation, log-likelihood	Tensorflow	https://bitbucket.org/jerry00/ scvis-dev/src/master/	[76]
	scVAE	VAE	ARI	Tensorflow	https://github.com/scvae/sc vae	[77]
	VASC	VAE	NMI, ARI, HS, and CS	H5py, keras	https://github.com/wang- research/VASC	[78]
	scDeepCluster	AE	ARI, NMI, clustering accuracy	Keras, Scanpy	https://github.com/ttgump/sc DeepCluster	[80]

	cscGAN	GAN	t-SNE, marker genes, MMD, AUC	Scipy, Tensorflow	https://github.com/imsb- uke/scGAN	[83]	
Multi-functional models (IM: imputation, BC: batch correction, CL: clustering)							
	scVI	VAE	IM: L₁ distance; CL: ARI, NMI, SI; BC: Entropy of Mixing	PyTorch, Anndata	https://github.com/YosefLab /scvi-tools	[17]	
	LDVAE	VAE	Reconstruction errors	Part of scVI	https://github.com/YosefLab /scvi-tools	[87]	
	SAUCIE	AE	IM: R² statistics; CL: SI;BC: modified kBET; Visualization: Precision/Recall	Tensorflow	https://github.com/Krishnas wamyLab/SAUCIE/	[15]	
	scScope	AE	<pre>IM:Reconstruction errors; BC: Entropy of mixing; CL: ARI</pre>	Tensorflow, Scikit-learn	https://github.com/Altschule rWu-Lab/scScope	[93]	
Cell type Identification							
[DigitalDLSorter	DNN	Pearson correlation	R/Python/Ke ras	https://github.com/cartof/digit alDLSorter	[51]	
	scCapsNet	CapsNet	Cell-type Prediction accuracy	Keras, Tensorflow	https://github.com/wanglf19/ scCaps	[52]	
	netAE	VAE	Cell-type Prediction accuracy, t- SNE for visualization	pyTorch	https://github.com/LeoZDong/netAE	[102]	
	scDGN	DANN	Prediciton accuracy	pyTorch	https://github.com/SongweiG e/scDGN	[53]	
Function analysis							
	CNNC	CNN	AUROC, AUPRC, and accuracy	Keras, Tensorflow	https://github.com/xiaoyeye/ CNNC	[54]	
	scGen	VAE	Correlation, visualization	Tensorflow	https://github.com/theislab/s cgen	[115]	

DL Model keywords: AE: autoencoder, AE+TL: autoencoder with transfer learning, AE: variational autoencoder, GAN: Generative adversarial network, CNN: convolutional neural network, DNN: deep neural network, DANN: domain adversarial neural network, CapsNet: capsule neural network

Table 1b. Comparison of Deep Learning algorithms reviewed in the paper

Арр	Algorithm	Feature	Application notes			
Imputation	n					
	DCA	 Strongest recovery of the top 500 genes Choices of noise models, including NB, and ZINB Outperform other existing methods in capturing cell population structure 	 AE integrated into the Scanpy framework High scalability of AE, up to millions of cells This method was compared to SAVER, scImpute, and MAGIC 			
	SAVER-X	 Pretraining from existing data sets (transfer learning) Decomposes the variation into three components 	 SAVER-X pretraining on PBMCs outperformed other denoising methods, including DCA, scVI, scImpute, and MAGIC SAVER-X was also applied for cross-species analysis 			
	DeepImpute	 Divide-and-conquer approach, using a bank of DNN models Reduced complexity by learning smaller sub-network Minimized overfitting by removing target genes from input 	 DeepImpute had the highest overall accuracy and offered shorter computation time than other methods like MAGIC, DrImpute, ScImpute, SAVER, VIPER, and DCA DeepImpute showed benefits in improving clustering results and identifying significantly differentially expressed genes Scalable and faster training time 			
	LATE	 Takes the log-transformed expression as input No explicit distribution assumption on input data 	 LATE outperforms other existing methods like MAGIC, SAVER, DCA, scVI, particularly when the ground truth contains only a few or no zeros Better scalability than DCA and scVI up to 1.3 million cells with 10K genes 			
	scGAMI	 A model designed for clustering but it includes an AE Uses fast independent component analysis algorithm: FastICA 	 Significantly improved the clustering performance in eight of seventeen selected scRNA-seq datasets scGMI's scalability needs to be improved 			
	scIGANs	 Trains a GAN model to generate samples with imputed expressions 	 This framework forces the model to reconstruct the real samples and discriminate between real and generated samples. Best reported performance in clustering compared to DCA, DeepImpute, SAVER, scImpute, MAGIC Scalable over 100K cells, also robust to small datasets 			
Batch cor	rection					
	BERMUDA	Preserves batch-specific biological signals through transfer-learning Preserves batch-specific cell populations	 Outperform other methods like mnnCorrect, BBKNN, Seurat, and scVI Removes batch effects even when the cell population compositions across different batches are vastly different Scalable by using mini-batch gradient descent algorithm during training 			

DESC	 Removes batch effect through clustering with the hypothesis that batch differences in expressions are smaller than true biological variations Does not require explicit batch information for batch removal 	 DESC is effective in removing the batch effect, whereas CCA, MNN, Seurat 3.0, scVI, BERMUDA, and scanorama were sensitive to batch definitions DESC is biologically interpretable and can reveal both discrete and pseudo-temporal structures of cells Small memory footprint and GPU enabled
iMAP	 iMAP combines AE and GAN for batch effect removal It consists of two processing stages, each including a separate DL model 	 iMAP was shown to separate the batch-specific cell types but mix batch shared cell types and outperformed other existing batch correction methods including Harmony, scVI, fastMNN, Seurat Capable handling datasets from Smart-seq2 and 10X Genomics platforms Demonstrated the stability over hyperparameters, and scalability for thousands of cells.
	ntation, dimension reduction, and data a	
Dhaka	 It was proposed to reduce the dimension of scRNA-seq data for efficient stratification of tumor subpopulations 	 Dhaka was shown to have an ARI higher than most other comparing methods including t-SNE, PCA, SIMLR, NMF, an autoencoder, MAGIC, and scVI Dhaka can represent an evolutionary trajectory
scvis	 VAE network that learns low-dimensional representations Capture both local and global neighboring structures 	 scvis was tested on the simulated data and outperformed t-SNE scvis is much more scalable than BH t-SNE (t-SNE takes O(M log(M)) time and O(M log(M)) space) for very large dataset (>1 million cells)
scVAE	 scVAE includes multiple VAE models for denoising gene expression levels and learning the low-dimensional latent representation Gaussian Mixture VAE (GMVAE) with negative binomial distribution achieved the highest lower bound and ARI 	 GMVAE was also compared with Seurat and shown to perform better, however, scVAE performed no better than scVI or scvis Algorithm applicable to large datasets with million cells
VASC	 Another VAE for dimension reduction and latent representation Explicitly model dropout with a Zero- inflated layer 	 VASC was compared with PCA, t-SNE, ZIFA, and SIMLR on 20 datasets In the study of embryonic development from zygote to blast cells, VASC shithe owed better performance to model embryo developmental progression VASC is reported to handle a large number of cells or cell types Demonstrated model stability
scDeepCluster	 AE network that simultaneously learns feature representation and performs clustering via explicit modeling of cell clusters 	 It was tested on the simulation data with different dropout rates and compared with DCA, MPSSC and SIMLR CIDR, PCA + k-mean, scvis and DEC significantly outperforming them Running time of scDeepCluster scales linearly with the number of cells, suitable for large scRNA-seq datasets

cscGAN	 It was designed to augment the existing scRNA-seq samples by generating expression profiles of specific cell types or subpopulations The cscGAN learns the expression patterns of a specific subpopulation (few cells), and simultaneously learns the cells from all populations (a large number of cells). 	 cscGAN was shown to generate high-quality scRAN-seq data for specific cell types. The augmentation in this method improved the identification of rare cell types and the ability to capture transitional cell states from trajectory analysis Better scalability than SUGAR (Synthesis Using Geometrically Aligned Random-walks) Capable re-establish developmental trajectories through pseudo-time analysis via cscGAN data augmentation
Multi-functional models		
scVI	 Designed to address a range of fundamental analysis tasks, including batch correction, visualization, clustering, and differential expression Integrated a normalization procedure and batch correction 	 ScVI was shown to be faster than most non-DL algorithms and scalable to handle twice as many cells as non-DL algorithms with a fixed memory For imputation, ScVI, together with other ZINB-based models, performed better than methods using alternative distributions Similar scalability as DCA
LDVAE	 Adaption of scVI to improve the model interpretability 	 For LDVAE the variations along the different axes of the latent variable establish direct linear relationships with input genes.
SAUCIE	It is applied to the normalized data instead of count data	 Results showed that SAUCIE had a better or comparable performance with other approaches SAUCIE has better scalability and faster runtimes than any of the other models Applications with single-cell CyTOF datasets
scScope	AE with recurrent steps designed for imputation and batch correction	 It was compared with PCA, MAGIC, ZINB-WaVE, SIMLR, AE, scVI, and DCA Efficiently identify cell subpopulations from complex datasets with high dropout rates, large numbers of subpopulations and rare cell types For scalability and training speed, scScope was shown to offer scalability (for >100K cells) with high efficiency (faster than most of the approaches)
Cell type Identification		
DigitalDLSorter	 A deconvolution model with 4-layer DNN Designed to identify and quantify the immune cells infiltrated in tumors captured in bulk RNA-seq, utilizing single-cell RNA-seq data 	 DigitalDLSorter achieved excellent agreement (linear correlation of 0.99 for colorectal cancer, and good agreement in quadratic relationship for breast cancer) at predicting cell type proportion.
scCapsNet	 It takes log-transformed, normalized expressions as input and follows the general CapsNet model 	 Interpretable capsule network designed for cell type prediction scCapsNet makes the deep-learning black box transparent through the direct interpretation of internal parameters

netAE	 VAE-based semi-supervised cell type prediction model Aiming at learning a low dimensional space from which the original space can be accurately reconstructed 	 Deals with scenarios of having a small number of labeled cells. netAE outperformed most of the baseline methods including scVI, ZIFA, PCA and AE as well as a semi-supervised method scANVI
scDGN	 scDGN takes the log-transformed, normalized expression as the input Supervised domain adversarial network 	 scDGN was tested for classifying cell types and aligning batches scDGN outperformed many deep learning and traditional machine learning methods in classification accuracy, including DNN, CaSTLe, MNN, scVI, and Seurat
Function analysis		
CNNC	 CNNC takes expression levels of two genes from many cells and transforms them into a 32 x 32 image-like normalized empirical probability function Inferring causal interactions between genes from scRNA-seq 	 CNNC outperforms prior methods for inferring TF-gene and protein-protein interactions, causality inference, and functional assignments Was shown to have more than 20% higher AUPRC than other methods and reported almost no false-negative for the top 5% predictions
scGen	 ScGen follows the general VAE for scRNA-seq data but uses the "latent space arithmetics" to learn perturbations' response Designed to learn cell response to certain perturbation (drug treatment, gene knockout, etc) 	 Compared with other methods including CVAE, style transfer GAN, linear approaches based on vector arithmetics(VA) and PCA+VA, scGen predicted full distribution of ISG15 gene (strongest regulated gene by IFN-b) response to IFN-b scGen can be used to translate the effect of a stimulation trained in study A to how stimulated cells would look in study B, given a control sample set

Abbreviation: NB: negative binomial distribution; ZINB: zero-inflated negative binomial distribution; TF: Transcription factor;

Table 2a: Simulated single-cell data/algorithms

1

2 3

Title	Algorithm	# Cells	Simulation methods	Reference
Splatter	DCA, DeepImpute, PERMUDA, scDeepCluster, scVI, scScope, solo	~2000	Splatter/R	[85]
CIDR	sclGAN	50	CIDR simulation	[61]
NB+dropout	Dhaka	500	Hierarchical model of NB/Gamma + random dropout	
Bulk RNA- seq	SAUCIE	1076	1076 CCLE bulk RNAseq + dropout conditional on the expression level	
SIMLR	scScope	1 million	SIMLR, high-dimensional data generated from latent vector	[44]
SUGAR	cscGAN	3000	Generating high dimensional data that follows a low dimensional manifold	[86]

Table 2b: Human single-cell data sources used by different DL algorithms

Title	Algorithm	Cell origin	# Cells	Data Sources	Reference
68k PBMCs	DCA SAVER-X LATE, scVAE, scDeepCluster, scCapsNet, scDGN	Blood	68,579	10X Single Cell Gene Expression Datasets	
Human pluripotent	DCA	hESCs	1,876	GSE102176	[129]
CITE-seq	SAVER-X	Cord blood mononuclear cells	8,005	GSE100866	[130]
Midbrain and Dopaminergic Neuron Development	SAVER-X	Brain/ embryo ventral midbrain cells	1,977	GSE76381	[131]
НСА	SAVER-X	Immune cell, Human Cell Atlas	500,000	HCA data portal	
Breast tumor	SAVER-X	Immune cell in tumor micro- environment	45,000	GSE114725	[132]
293T cells	DeepImpute, iMAP	Embryonic kidney	13,480	10X Single Cell Gene Expression Datasets	
Jurkat	DeepImpute, iMAP	Blood/ lymphocyte	3,200	10X Single Cell Gene Expression Datasets	
ESC, Time- course	scGAN	ESC	350, 758	GSE75748	[133]

Baron-Hum-1	scGMAI, VASC	Pancreatic islets	1,937	GSM2230757	[134]
Baron-Hum-2	scGMAI, VASC	Pancreatic islets	1,724	GSM2230758	[134]
Camp	scGMAI, VASC	Liver cells	303	GSE96981	[135]
CEL-seq2	PERMUDA, DESC	Pancreas/Islet s of Langerhans		GSE85241	[136]
Darmanis	scGMAI, sclGAN, VASC	Brain/cortex	466	GSE67835	[137]
Tirosh-brain	Dhaka, scvis	Oligodendrogli oma	>4800	GSE70630	[138]
Patel	Dhaka	Primary glioblastoma cells	875	GSE57872	[139]
Li	scGMAI, VASC	Blood	561	GSE146974	[68]
Tirosh-skin	scvis	melanoma	4645	GSE72056	[74]
xenograft 3, and 4	Dhaka	Breast tumor	~250	EGAS00001002170	[140]
Petropoulos	VASC/netAE	Human embryos	1,529	E-MTAB-3929	
Pollen	scGMAI, VASC		348	SRP041736	[141]
Xin	scGMAI, VASC	Pancreatic cells (α-, β-, δ-)	1,600	GSE81608	[142]
Yan	scGMAI, VASC	embryonic stem cells	124	GSE36552	[143]
PBMC3k	VASC, scVI	Blood	2,700	SRP073767	[100]
CyTOF, Dengue	SAUCIE	Dengue infection	11 M, ~42 antibodies	Cytobank, 82023	[15]
CyTOF, ccRCC	SAUCIE	Immunue profile of 73 ccRCC patients	3.5 M, ~40 antibodies	Cytobank: 875	[144]
CyTOF, breast	SAUCIE	3 patients		Flow Repository: FR- FCM-ZYJP	[132]
Chung, BC	DigitalDLSorter	Breast tumor	515	GSE75688	[94]
Li, CRC	DigitalDLSorter	Colorectal cancer	2,591	GSE81861	[95]
Pancreatic datasets	scDGN	Pancreas	14693	SeuratData	
Kang, PBMC	scGen	PBMC stimulated by INF-β	~15,000	GSE96583	[116]

Table 2c: Mouse single-cell data sources used by different DL algorithms

Title	Algorithm	Cell origin	# Cells	Data Sources	Reference
Brain cells from E18 mice	DCA, SAUCIE	Brain Cortex	1,306,127	10x: Single Cell Gene Expression Datasets	
Midbrain and Dopaminergic Neuron Development	SAVER-X	Ventral Midbrain	1907	GSE76381	[131]
Mouse cell atlas	SAVER-X		405,796	GSE108097	[145]
neuron9k	DeepImpute	Cortex	9128	10x: Single Cell Gene Expression Datasets	
Mouse Visual Cortex	DeepImpute	Brain cortex	114601	GSE102827	[146]
murine epidermis	DeepImpute	Epidermis	1422	GSE67602	[147]
myeloid progenitors	LATE DESC, SAUCIE	Bone marrow	2,730	GSE72857	[148]
Cell-cycle	sclGAN	mESC	288	E-MTAB-2805	[149]
A single-cell survey		Intestine	7721	GSE92332	[120]
Tabula Muris	iMAP	Mouse cells	>100K		
Baron-Mou-1	VASC	Pancreas	822	GSM2230761	[134]
Biase	scGMAI, VASC	Embryos/SMA RTer	56	GSE57249	[150]
Biase	scGMAI, VASC	Embryos/Fluidi gm	90	GSE59892	[150]
Deng	scGMAI, VASC	Liver	317	GSE45719	[151]
Klein	VASC scDeepCluster sclGAN	Stem Cells	2,717	GSE65525	[152]
Goolam	VASC	Mouse Embryo	124	E-METAB-3321	[153]
Kolodziejczyk	VASC	mESC	704	E-MTAB-2600	[154]
Usoskin	VASC	Lumbar	864	GSE59739	[155]
Zeisel	VASC, scVI, SAUCIE, netAE	Cortex, hippocampus	3,005	GSE60361	[156]
Bladder cells	scDeepCluster	Bladder	12,884	GSE129845	[157]
HEMATO	scVI	Blood cell	>10,000	GSE89754	[158]
retinal bipolar cells	scVI, scCapsNet	retinal	~25,000	GSE81905	[101]

	SAUCIE				
Embryo at 9 time points	LDAVE	embryos from E6.5 to E8.5	116,312	GSE87038	[159]
Embryo at 9 time points	LDAVE	embryos from E9.5 to E13.5	~2 millions	GSE119945	[160]
CyTOF,	SAUCIE	Mouse thymus	200K, ~38 antibodies	Cytobank: 52942	[161]
Nestorowa	netAE	hematopoietic stem and progenitor cells	1,920	GSE81682	[162]
small intestinal epithelium	scGen	Infected with Salmonella and worm H. polygyrus	1,957	GSE92332	[120]

Table 2d: Single-cell data derived from other species

Title	Algorithm	Species	Tissue	# Cells	SRA/GEO	Reference
Worm neuron cells ¹	scDeepCluster	C. elegans	Neuron	4,186	GSE98561	[163]
Cross species, stimulation with LPS and dsRNA	scGen	Mouse, rat, rabbit, and pig	bone marrow- derived phagocyte	5,000 to 10,000 /species	13 accessions in ArrayExpress	[121]

^{1.} Processed data is available at https://github.com/ttgump/scDeepCluster/tree/master/scRNA-seq%20data

Table 2e: Large single-cell data source used by various algorithms

Title	Sources	Notes
10X Single-cell gene expression dataset	https://support.10xgenomics.com/single- cell-gene-expression/datasets	Contains large collection of scRNA- seq dataset generated using 10X system
Tabula Muris	https://tabula-muris.ds.czbiohub.org/	Compendium of scRNA-seq data from mouse
HCA	https://data.humancellatlas.org/	Human single-cell atlas
MCA	https://figshare.com/s/865e694ad06d585 7db4b, or GSE108097	Mouse single-cell atlas
scQuery	https://scquery.cs.cmu.edu/	A web server cell type matching and key gene visualization. It is also a source for scRNA-seq collection (processed with common pipeline)
SeuratData	https://github.com/satijalab/seurat-data	List of datasets, including PBMC and human pancreatic islet cells
cytoBank	https://cytobank.org/	Community of big data cytometry

Table 3. Evaluation metrics used in surveyed DL algorithms

Evaluation Method	s used in surveyed DL algorithms Equations	Explanation
Pseudobulk RNA-seq		Average of normalized (log2-transformed) scRNA-seq counts across cells is calculated and then correlation coefficient between the pseudobulk and the actual bulk RNA-seq profile of the same cell type is evaluated.
Mean squared error (MSE)	$MSE = \frac{1}{n} \sum_{i=1}^{n} (x_i - \hat{x}_i)^2$	MSE assesses the quality of a predictor, or an estimator, from a collection of observed data x , with \hat{x} being the predicted values.
Pearson correlation	$\rho_{X,Y} = \frac{cov(X,Y)}{\sigma_X \sigma_Y}$	where cov() is the covariance, σ_X and σ_Y are the standard deviation of X and Y , respectively.
Spearman correlation	$\rho_s = \rho_{r_X, r_Y} = \frac{cov(r_X, r_Y)}{\sigma_{r_X} \sigma_{r_Y}}$	The Spearman correlation coefficient is defined as the Pearson correlation coefficient between the rank variables, where r_X is the rank of X.
Entropy of accuracy, H _{acc} [21]	$H_{acc} = -\frac{1}{M} \sum_{i=1}^{M} \sum_{j=1}^{N_i} p_i(x_j) \log(p_i(x_j))$	Measures the diversity of the ground-truth labels within each predicted cluster group. $p_i(x_j)$ (or $q_i(x_j)$) are the proportions of cells in the j^{th} ground-truth cluster (or predicted cluster) relative to the total number of cells in the j^{th} predicted cluster (or ground-truth clusters), respectively.
Entropy of purity, <i>H</i> _{pur} [21]	$H_{pur} = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{M_i} q_i(x_j) \log \left(q_i(x_j) \right)$	Measures the diversity of the predicted cluster labels within each ground-truth group
Entropy of mixing [32]	$E = \sum_{i=1}^{C} p_i \log(p_i)$	This metric evaluates the mixing of cells from different batches in the neighborhood of each cell. C is the number of batches, and p_i is the proportion of cells from batch i among N nearest cells.
Mutual Information (MI) [164]	$MI(U,V) = \sum_{i=1}^{ U } \sum_{j=1}^{ V } P_{UV}(i,j) \log \left(\frac{P_{UV}(i,j)}{P_{U}(i)P_{V}(j)} \right)$	where $P_U(i) = \frac{ U_i }{N}$ and $P_V(j) = \frac{ V_j }{N}$. Also, define the joint distribution probability is $P_{UV}(i,j) = \frac{ U_i \cap V_j }{N}$. The MI is a measure of mutual dependency between two cluster assignments U and V .
Normalized Mutual Information (NMI) [165]	$NMI(U,V) = \frac{2 \times MI(U,V)}{[H(U) + H(V)]}$	where $H(U) = \sum P_U(i) \log(P_U(i))$, $H(V) = \sum P_V(i) \log(P_V(i))$. The NMI is a normalization of the MI score between 0 and 1.

Kullback–Leibler (KL) divergence [166]	$D_{KL}(P Q) = \sum_{x \in \chi} P(x) \log \left(\frac{P(x)}{Q(x)} \right)$
Jaccard Index	$J(U,V) = \frac{\lfloor U \cap V \rfloor}{\lfloor U \cup V \rfloor}$
Fowlkes-Mallows Index for two clustering algorithms (FM)	$FM = \sqrt{\frac{TP}{TP + FP}} \times \frac{TP}{TP + FN}$
Rand index (RI)	$RI = (a+b)/\binom{n}{2}$
Adjusted Rand index (ARI) [167]	$ARI = \frac{RI - E[RI]}{\max(RI) - E[RI]}$
Silhouette index	$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$
Maximum Mean Discrepancy (MMD) [65]	$MMD(F, p, q) = \sup_{f \in F} \ \mu_p - \mu_q\ _f$
k-Nearest neighbor batch-effect test (kBET) [168]	$a_n^k = \sum_{l=1}^L \frac{(N_{nl}^k - k \cdot f_l)^2}{k \cdot f_l} \sim X_{L-1}^2$

Local Inverse Simpson's

Index (LISI) [34]

 $\frac{1}{\lambda(n)} = \frac{1}{\sum_{l=1}^{L} (p(l))^2}$

where discrete probability distributions P and Q are defined on the same probability space χ . This relative entropy is the measure for directed divergence between two distributions.

 $0 \le J(U,V) \le 1$. J = 1 if clusters *U* and *V* are the same. If *U* are *V* are empty, J is defined as 1.

TP as the number of pairs of points that are present in the same cluster in both U and V; FP as the number of pairs of points that are present in the same cluster in U but not in V; FN as the number of pairs of points that are present in the same cluster in V but not in U; and TN as the number of pairs of points that are in different clusters in both U and V.

Measure of constancy between two clustering outcomes, where a (or b) is the count of pairs of cells in one cluster (or different clusters) from one clustering algorithm but also fall in the same cluster (or different clusters) from the other clustering algorithm.

ARI is a corrected-for-chance version of RI, where E[RI] is the expected Rand Index.

where a(i) is the average dissimilarity of i^{th} cell to all other cells in the same cluster, and b(i) is the average dissimilarity of i^{th} cell to all cells in the closest cluster. The range of s(i) is [-1,1], with 1 to be well-clustered and -1 to be completely misclassified.

MMD is a non-parametric distance between distributions based on the reproducing kernel Hilbert space, or, a distance-based measure between two distribution p and q based on the mean embeddings μ_p and μ_q in a reproducing kernel Hilbert space F.

Given a dataset of N cells from L batches with N_l denoting the number of cells in batch l, N_{nl}^k is the number of cells from batch l in the k-nearest neighbors of cell n, f_l is the global fraction of cells in batch l, or $f_l = \frac{N_l}{N}$, and X_{L-1}^2 denotes the X^2 distribution with L-1 degrees of freedom. It uses a X^2 -based test for random neighborhoods of fixed size to determine the significance ("well-mixed").

This is the inverse Simpson's Index in the k-nearest neighbors of cell n for all batches, where p(l) denotes the proportion of batch l in the k-nearest neighbors. The score reports the

Homogeneity	$HS = 1 - \frac{H(P(U V))}{H(P(U))}$
	H(P(U))

Completeness
$$CS = 1 - \frac{H(P(V|U))}{H(P(V))}$$

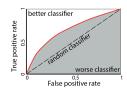
V-Measure [169]
$$V_{\beta} = \frac{(1+\beta)HS \times CS}{\beta HC + CS}$$

Precision, recall
$$Precision = \frac{TP}{TP + FP}$$
, $recall = \frac{TP}{TP + FN}$

Accuracy
$$Accuracy = \frac{TP + TN}{N}$$

F₁-score
$$F_1 = \frac{2 \operatorname{Precision} \cdot \operatorname{Recall}}{\operatorname{Precision} + \operatorname{Recall}}$$

AUC, RUROC



effective number of batches in the k-nearest neighbors of cell n.

where H() is the entropy, and U is the ground-truth assignment and V is the predicted assignment. The HS range from 0 to 1, where 1 indicates perfectly homogeneous labeling.

Its values range from 0 to 1, where 1 indicates all members from a ground-truth label are assigned to a single cluster.

where β indicates the weight of HS. V-Measure is symmetric, *i.e.* switching the true and predicted cluster labels does not change V-Measure.

TP: true positive, FP: false positive, FN, false negative.

N: all samples tested, TN: true negative

A harmonic mean of precision and recall. It can be extended to F_{β} where β is a weight between precision and recall (similar to V-measure).

Area Under Curve (grey area). Receiver operating characteristic (ROC) curve (red line). A similar measure can be performed on the Precision-Recall curve (PRC), or AUPRC. Precision-Recall curves summarize the trade-off between the true positive rate and the positive predictive value for a predictive model (mostly for an imbalanced dataset).

