



OPEN

A topology-preserving dimensionality reduction method for single-cell RNA-seq data using graph autoencoder

Zixiang Luo^{1,4}, Chenyu Xu^{2,4}, Zhen Zhang³✉ & Wenfei Jin¹✉

Dimensionality reduction is crucial for the visualization and interpretation of the high-dimensional single-cell RNA sequencing (scRNA-seq) data. However, preserving topological structure among cells to low dimensional space remains a challenge. Here, we present the single-cell graph autoencoder (scGAE), a dimensionality reduction method that preserves topological structure in scRNA-seq data. scGAE builds a cell graph and uses a multitask-oriented graph autoencoder to preserve topological structure information and feature information in scRNA-seq data simultaneously. We further extended scGAE for scRNA-seq data visualization, clustering, and trajectory inference. Analyses of simulated data showed that scGAE accurately reconstructs developmental trajectory and separates discrete cell clusters under different scenarios, outperforming recently developed deep learning methods. Furthermore, implementation of scGAE on empirical data showed scGAE provided novel insights into cell developmental lineages and preserved inter-cluster distances.

Single-cell RNA sequencing (scRNA-seq) is an ideal approach for investigating cell-cell variation. Conventional dimensionality reduction techniques such as principal component analysis (PCA) and t-Distributed Stochastic Neighbor Embedding (t-SNE)¹ were implemented on scRNA-seq data for visualization and downstream analyses, significantly increasing our understanding of cellular heterogeneity and development progress. The recent emergence of massively parallel scRNA-seq such as droplet platforms enabled interrogation of millions of cells in complex biological systems^{2–5}, which provide a fantastic potential for dissection of tissue and cellular micro-environment, identification of rare/new cell types, inference of developmental lineages, and elucidation of the mechanism of cellular response to stimulations⁶. However, the data generated by massively parallel scRNA-seq are of high dropout and high noise with complex structure, which posed a series of challenges on dimensionality reduction. Particularly, it is a big challenge to preserve the complex topological structure among cells.

Many dimensionality reduction methods have been developed or introduced for scRNA-seq data analyses in the past several years. Recently developed competitive methods include DCA⁷, scVI⁸, scDeepCluster⁹, PHATE¹⁰, SAUCIE¹¹, scGNN¹², ZINB-WaVE¹³ and Ivis¹⁴. Among them, deep learning showed the greatest potentials. For instance, DCA, scDeepCluster, Ivis, and SAUCIE adapted the autoencoder to denoise, visualize and cluster the scRNA-seq data. However, these deep learning-based models only embedded the distinct cell features while ignoring the cell–cell relationships, which limited their ability to reveal the complex topological structure among cells and made them difficult to elucidate the developmental trajectory. The recently proposed graph autoencoder¹⁵ is very promising as it preserves the long-distance relationships among data in a latent space. In this study, we developed the single-cell graph autoencoder (scGAE). It improved the graph autoencoder to preserving global topological structure among cells. We further extended the scGAE for visualization, trajectory inference, and clustering. Analyses of simulated data and empirical data showed that scGAE outperformed the other competitive methods.

¹Shenzhen Key Laboratory of Gene Regulation and Systems Biology, School of Life Sciences, Southern University of Science and Technology, Shenzhen 518055, China. ²Department of Electric Engineering, Iowa State University, Ames, IA 50011, USA. ³Department of Mathematics, International Center for Mathematics, National Center for Applied Mathematics (Shenzhen), Guangdong Provincial Key Laboratory of Computational Science and Material Design, Southern University of Science and Technology, Shenzhen 518055, China. ⁴These authors contributed equally: Zixiang Luo and Chenyu Xu. ✉email: zhangz@sustech.edu.cn; jinwf@sustech.edu.cn

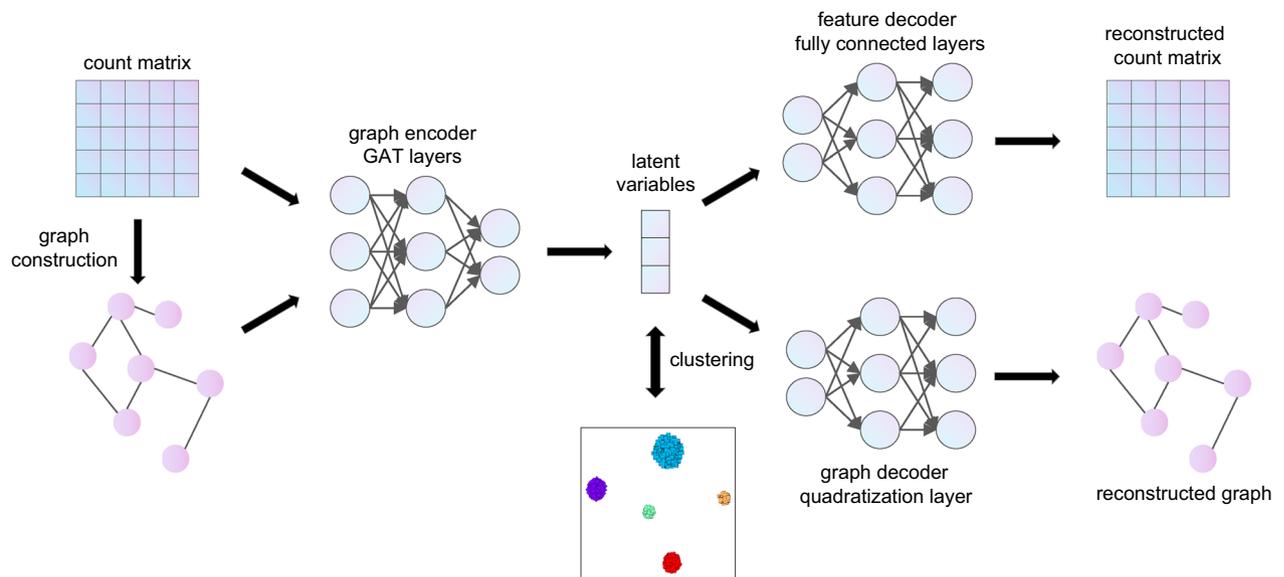


Figure 1. The model architecture of scGAE. The normalized count matrix represents the gene expression level in each cell. The adjacency matrix is constructed by connecting each cell to its K nearest neighbors. The encoder takes the count matrix and the adjacency matrix as inputs and generates low-dimensional latent variables. The feature decoder reconstructs the count matrix. The graph decoder reconstructs the adjacency matrix. Clustering is performed on the latent variables.

Results

The model architecture of scGAE. scGAE combines the advantage of the deep autoencoder and graphical model to embed the topological structure of high-dimensional scRNA-seq data to a low-dimensional space (Fig. 1). After getting the normalized count matrix, scGAE builds the adjacency matrix among cells by K-nearest-neighbor algorithm. The encoder maps the count matrix to a low-dimensional latent space by graph attentional layers¹⁶. scGAE decodes the embedded data with a feature decoder and a graph decoder. The feature decoder reconstructs the count matrix to preserve the feature information; The graph decoder recovers the adjacency matrix and preserves the topological structure information. It decodes the embedded data to the spaces with the same dimension as original data by minimizing the distance between the input data and the reconstructed data (see “Methods”). We use deep clustering to learn the data embedding and do cluster assignment simultaneously¹⁷, generating a clustering-friendly latent representation (Supplementary Fig. S1). The implementation and usage of scGAE can be found on Github: <https://github.com/ZixiangLuo1161/scGAE>.

Visualization of scGAE embedded data and comparison to other methods. To systematically evaluate the performance of scGAE, we summarized four representative scenarios (scenario1: cells in continuous differentiation lineages; scenario2: cells in differentiation lineages where cells concentrate at the center of each branch; scenario3: distinct cell populations with apparent differences; and scenario4: distinct cell populations with small population differences) (Fig. 2 left). We used Splatter¹⁸ and PROSSTT¹⁹ to simulate scRNA-seq data in four scenarios. For scGAE, the data was visualized by tSNE after projected to a latent space. Compared with other methods, scGAE better captured the complex structures in the data (Fig. 2). In scenario1 and scenario2, scGAE almost entirely reproduced the differentiation lineages (Fig. 2a,b), while other methods only revealed some local structures and failed to exhibit the overall structure of simulated data. The results of tSNE and SAUCIE exhibited distinct clusters but lost lineage relationship in scenario2. In scenario3 and 4, scGAE almost perfectly preserved the compact cell clusters and inter-cluster distances in the simulated data, while the clusters inferred by other methods are dispersed, and the topological structure among these clusters was not preserved (Fig. 2c,d). Only scGAE separated all the clusters while the other methods mixed different types of cells when the differences between clusters are small (Fig. 2d). Based on these observations, scGAE perfectly reproduced the differentiation lineages and distinct clusters in the simulated data, indicating scGAE outperforms other competitive methods in restoring the relationship between cells.

Trajectory inference and cell clustering based on scGAE embedded data. We further quantitatively evaluated the performance of scGAE for trajectory inference tasks. The scGAE and other competitive methods were used to perform dimensionality reduction on the developmental lineage data simulated by PROSSTT (scenario1 and 2). We conducted trajectory inference on these embedded data using DPT²⁰. The Kendall correlation coefficient²¹ between the inferred trajectories and the ground truth was calculated to measure their similarity. Because scDeepCluster is a clustering method, we didn't include it for trajectory inference tasks. The results showed that scGAE, scGNN, and scVI better recovered the original trajectory than the other competitive methods on both scenario1 and 2 (Fig. 3a,b). Compared with scenario1, the data is not uniformly

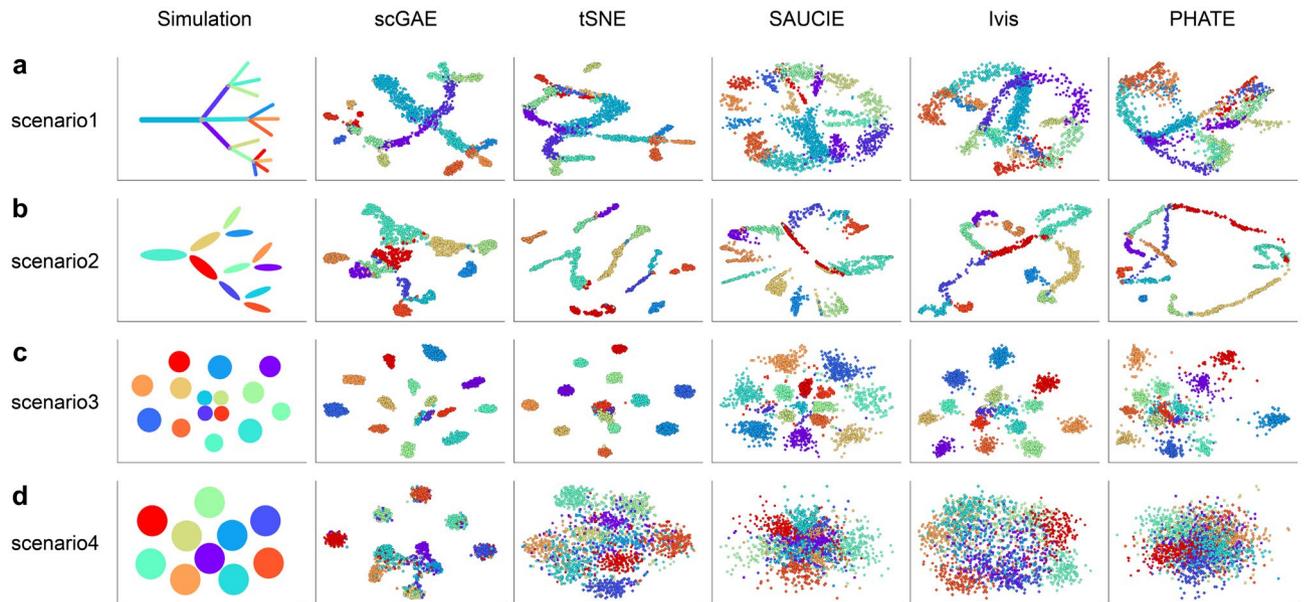


Figure 2. Visualization of the four simulated datasets by scGAE, tSNE, SAUCIE, Ivis, and PHATE. Each color represents a cell subpopulation in the simulated dataset. **(a)** scenario1: cells in continuous differentiation lineages. **(b)** scenario2: cells in differentiation lineages where cells concentrate at the center of each branch. **(c)** scenario3: distinct cell populations with apparent population differences. **(d)** scenario4: distinct cell populations with small population differences.

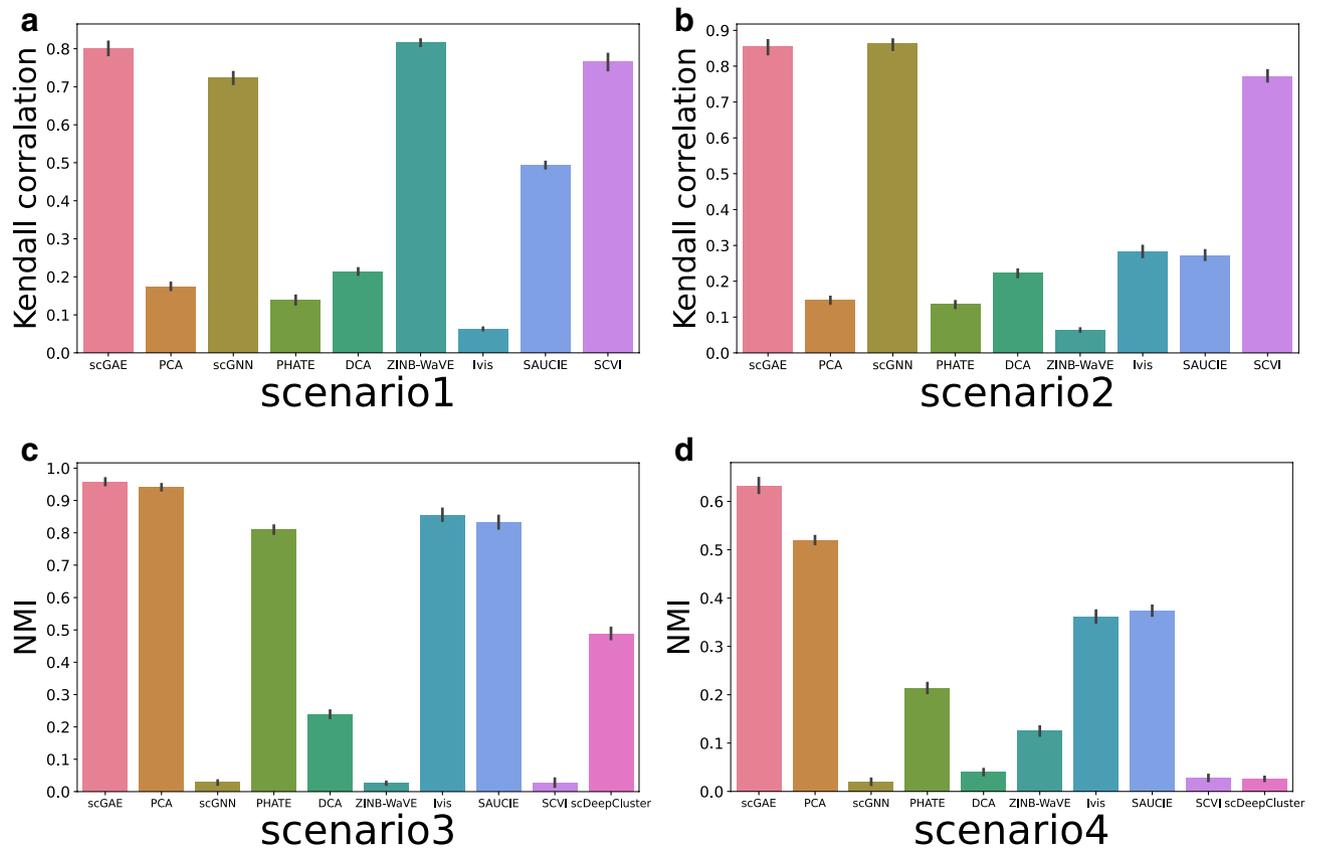


Figure 3. Quantitative evaluation of scGAE and several other competitive methods on clustering and trajectory inference tasks. In scenario1 **(a)** and scenario2 **(b)**, the Kendall correlation between the ground truth and inferred trajectory was calculated. In scenario3 **(c)** and scenario4 **(d)**, the normalized mutual information (NMI) measures the difference between the ground truth and the inferred clusters.

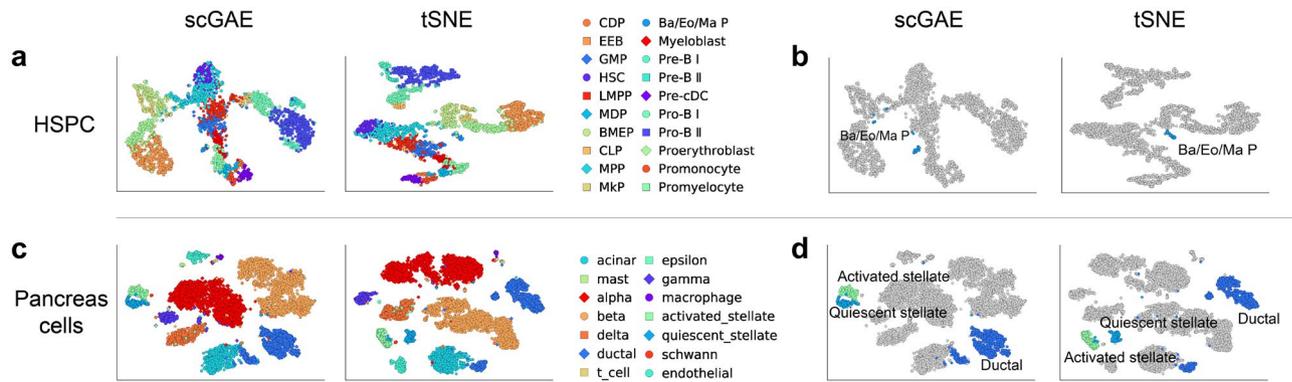


Figure 4. Analyses of two real datasets. (a) Visualization of HSPC cells by scGAE and tSNE (b) scGAE identified the multiple subpopulations in previous reported Ba/Eo/MaP. (c) Visualization of pancreas cells by scGAE and tSNE. (d) The close distance between two stellate states and the short distance between ductal subtypes recovered by scGAE.

distributed along the developmental trajectory in scenario2. Most methods have a lower Kendall correlation, but two graph neural network based methods and scVI still have good performances. It shows that the graph-based structure can well preserve the relationship among data. Next, we evaluated the performance of scGAE and other competitive methods on cell clustering tasks with data simulated by Splatter (scenario3 and 4). We performed Louvain clustering on these embedded data. Normalized mutual information (NMI) was used to measure the difference between inferred clusters and ground truth. The results showed that scGAE was the best among these methods (Fig. 3c,d, Supplementary Fig. S2). Although scVI, ZINB-Wave, and scGNN performed well for trajectory inference (Fig. 3a,b), they got a low score in the cell clustering task (Fig. 3c,d). The inconsistency between data structure imposed in existing methods and simulated data structure might contribute the differences of performance. Some methods such as scGAE assume no prior hypothesis on the data, which may facilitate their performances in all cases. Also, different data preprocessing approaches might affect the results. For the methods that takes normalized data as input, we normalized data using the Seurat R package. While the three method that dropped most only accept raw data as input. Moreover, when there are noises, scGAE can do better than these three methods in the low-dimensional cell clustering. This may be because scGAE optimize clustering and latent representation simultaneously in one shot.

To test the effect of zero-inflation, we varied the parameters in scenario 4 for $dropout.shape = -1$, $dropout.mid$ range in $(-0.5, 0, 0.5, 1)$. The corresponding dropout rates are $12 \pm 0.3\%$, $17 \pm 0.4\%$, $23 \pm 0.5\%$, and $30 \pm 0.6\%$. The corresponding normalized mutual information (NMI) is 0.62, 0.62, 0.65, and 0.61. The result shows that scGAE is robust again zero-inflation. Overall, scGAE performed well for both trajectory inference and cell clustering in four scenarios.

scGAE identified novel subpopulations that shaped hematopoietic lineage relationship. Single cell analysis of hematopoietic stem and progenitor cells (HSPCs) have significantly increased our understanding of the early cell subpopulations and developmental trajectory during hematopoiesis^{5,22–27}. We further used scGAE to analyze HSPCs scRNA-seq data from our previous study⁵ (Fig. 4a). We found the previous identified Basophil/Eosinophil/Mast progenitors (Ba/Eo/MaP) has been classified into multiple subpopulations (Fig. 4b). It indicates that the cells in Ba/Eo/MaP may have different differentiation potentials at early phase. While the other competitive methods did not identify the subpopulations in Ba/Eo/MaP (Supplemental Figs. S3a, S4a), supporting scGAE has the highest statistical power to identify the substructure in the scRNA-seq data.

scGAE preserved topological structure among human pancreatic cells populations. The function of the pancreas hinges on complex interactions among distinct cell types and cell populations. We re-analyzed the scRNA-seq data of human pancreatic cells from Baron et al.²⁸. Although the pancreatic cell subpopulations identified by scGAE are the same as the original study, we found the distances and topological structures among cell types inferred by scGAE better fit our knowledge (Fig. 4c). For instance, the activated stellate and quiescent stellate showed similar expression profiles and phenotypes²⁹. scGAE revealed the close relationship between two cell populations better than the other methods (Fig. 4d and Supplemental Figs. S3b, S4b). scGAE also preserved the short distance between two ductal subtypes, while some methods including tSNE project them into a longer distance. Moreover, scGAE clearly separated other cell populations while SAUCIE, Ivis, and PHATE mixed some of the clusters. Overall, scGAE preserved the topological structure among different cell populations, which greatly benefit our understanding of the cellular relationships.

Discussion

Because of the high noises of scRNA-seq data and complicated cellular relationships, preserving the topological structure of scRNA-seq data in low-dimensional space is still a challenge. We proposed scGAE which is a promising topology-preserving dimensionality reduction method. It generates a low-dimensional representation that better preserves both the global structure and local structure of the high-dimensional scRNA-seq data. The

key innovation of scGAE is to embed the structure information and feature information simultaneously using a multitask graph autoencoder. It is suitable for analyzing the data both in lineages and clusters. The learned latent representation benefits various downstream analyses, including clustering, trajectory inference, and visualization. The analyses on both simulated data and empirical data suggested scGAE accurately preserved the topological structures of data.

scGNN¹² is another tool that utilizes graph autoencoder for single cell RNA-seq data dimensionality reduction. scGAE is designed to perform dimensionality reduction while being friendly for further clustering and trajectory inference. scGNN is designed to do multi-tasks for modeling heterogeneous cell–cell relationships and their underlying complex gene expression patterns. It consists of four types of autoencoders with appropriate regularizations and iterations among these autoencoders. From the performance perspective, scGAE and scGNN have similar performance on the trajectory inference while scGAE has better performance on clustering. From the computational perspective, the running time of scGAE is much shorter than scGNN and memory cost is slightly lower than scGNN. This is due to the iterative process in scGNN, which is more time-consuming and requires more computational resources.

As an early study adapting graph autoencoder for dimensionality reduction of scRNA-seq data, this approach is likely to be significantly improved in the future. Firstly, because the complex data structure is hard to be directly embedded into two-dimensional space by graph autoencoder, we embedded the scRNA-seq data into an intermediate dimension and used tSNE to visualize the embedded data into a two-dimensional space. However, tSNE focuses more on local information, and it sometimes fails to correctly recover the global structure, which may distort the topological structure in the data. A better visualization method is needed to preserve the topological structure of scRNA-seq data. Secondly, the graph in scGAE is constructed by the K-nearest neighbor (KNN) algorithm that relies on a predefined parameter K. However, the optimal K varies among different datasets and different parts of a dataset. Constructing an optimal graph is challenging due to the difficulty in determining a suitable K, which could be our potential future endeavors. Thirdly, scGAE has a moderate time cost but a relatively high memory cost compared with other statistics model and deep learning methods without graph-based layers (Supplementary Figs. S5–S7). This is caused by the recursive neighborhood expansion across layers in graph neural network³⁰. In the future, we will investigate more efficient architectures such as GNN with graph sampling³⁰ to reduce the time and memory cost.

Methods

Joint graph autoencoder. The graph autoencoder is a type of artificial neural network for unsupervised representation learning on graph-structured data¹⁵. The graph autoencoder often has a low-dimensional bottleneck layer so that it can be used as a model for dimensionality reduction. Let the inputs be single-cell graphs of node matrices X and adjacency matrices A . In our joint graph autoencoders³¹, there is one encoder E for the whole graph and two decoders D_X and D_A for nodes and edges respectively. In practice, we first encode the input graph into a latent variable $h = E(X, A)$, and then we decode h into the reconstructed node matrix $X_r = D_X(h)$ and the reconstructed adjacency matrix $A_r = D_A(h)$. The objective of learning process is to minimize the reconstruction loss

$$L_r = \lambda \|X - X_r\|_2^2 + (1 - \lambda) \|A - A_r\|_2^2,$$

where the weight λ is a hyper-parameter. In our experiments, λ is set to be 0.6.

We used the Python package Spektral³² to implement our model. There are many types of graph neural networks that can be used as the encoder or decoder. Hereby, to extract the features of a node with the aid of its neighbors, we apply graph attention layers as default in the encoder. Other graph neural networks such as GCN³³, GraphSAGE³⁴ and TAGCN³⁵ can also be implemented as the encoder in scGAE. The feature decoder D_X is a four-layer fully connected neural network with 64, 256, 512 nodes in hidden layers.

The edge decoder consists of a fully connected layer followed by the composition of quadratization and activation:

$$A_r = D_A(h) = \sigma(ZZ^T),$$

where $Z = \sigma(W_h)$ arises as an output of a fully connected layer with the weight matrix W , and $\sigma(x) = \max(0, x)$ is the rectified linear unit.

Deep-clustering embedding. Motivated by Yang et al.³⁶, we use a two-stage method. The first stage is to pre-train scGAE by minimizing L_r . The resulting neural network parameters are set as the initialization of the second stage, which we call alter-training. The loss function in the alter-training stage comprises both reconstruction error L_r and clustering cost $L_c = L_c(h, \mu)$:

$$L = L_r + \gamma L_c,$$

where μ is a collection of clustering centroids, and γ is a hyper-parameter set as 2.5 in our experiments.

The alter-training consists of doing the following two steps alternately:

1. Given a collection of clustering centroids μ , update network parameters by minimizing L ;
2. Compute the embedded data h using the updated network, and do clustering in the embedded space to obtain new centroids μ ;

In experiments, we use the pre-trained network to generate the initial embedded data which are clustered to obtain the initial centroids by Louvain³⁷. There are various choices for the loss L_c and the clustering algorithm in the second step¹⁷. In practice, we compute the new centroids μ by minimizing L_c using the stochastic gradient descent. A good choice of L_c is the soft assignment loss, which is the KL divergence of empirical clustering assignment distribution Q from a target distribution P . This is motivated by t-SNE¹ which uses a proper distribution Q in low dimensional space in order to inherit the clustering property from the high dimensional space. Given an embedded point h_i and a centroid μ_j , Q is defined as Student's t -distribution $q_{ij} = \frac{(1 + \|h_i - \mu_j\|^2)^{-1}}{\sum_{j'} (1 + \|h_i - \mu_{j'}\|^2)^{-1}}$. An

ideal target distribution should have the following properties: (1) improve cluster purity, (2) put more emphasis on data points assigned with high confidence, and (3) prevent large clusters from distorting the hidden feature space. In experiments, we follow DEC³⁸ choose P as $p_{ij} = \frac{q_{ij}^2 / \sum_i q_{ij}}{\sum_{j'} q_{ij}^2 / \sum_i q_{ij}}$.

Evaluation metric. Clustering results are measured by Normalized Mutual Information (NMI)³⁹. Given the knowledge of the ground truth class assignments U and our clustering algorithm assignment V on n data points, NMI measures the agreement of the two assignment, ignoring permutations. NMI is defined as

$$\text{NMI}(U, V) = \frac{1}{\text{mean}(H(U), H(V))} \sum_{i=1}^{|U|} \sum_{j=1}^{|V|} \frac{|U_i \cap V_j|}{N} \log \left(\frac{n|U_i \cap V_j|}{|U_i||V_j|} \right),$$

where $H(U) = -\sum_{i=1}^{|U|} \frac{|U_i|}{n} \log \left(\frac{|U_i|}{n} \right)$ is the entropy.

Trajectory inference results are measured by Kendall correlation coefficient. We define an order among the set of observations $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$: any pair of observations (x_i, y_i) and (x_j, y_j) , where $i < j$ are said to be concordant if either both $x_i > x_j$ and $y_i > y_j$ hold or both $x_i < x_j$ and $y_i < y_j$ hold; otherwise they are said to be discordant. Denote the number of concordant pairs as N_{conco} and the number of discordant pairs as N_{discon} . Kendall correlation coefficient is defined as

$$\tau = \frac{2(N_{conco} - N_{discon})}{n(n-1)}.$$

Data simulation. We simulated five scRNA-seq datasets using Splatter R package (data1, data3, and data4) and PROSSTT Python package (data2 and data5). The cells in data1 and data5 are in the linear distribution along the developmental trajectory. The cells in data2 have a skewed distribution where cells concentrate at the center of each branch. The cells in data3 and data4 are in distinct clusters with moderate and small cluster differences, respectively. All datasets have 2000 cells and 5000 genes. Data1, data2, data3, and data4 were simulated for scenario1 to scenario4 for data visualization. Data5, data2, data3, and data4 are used for the evaluation of scGAE on trajectory inference and cell clustering tasks.

Data preprocessing. The scRNA-seq data preprocessing was conducted using scTransform⁴⁰ in The Seurat package⁴¹. The pre-processed count matrix was used to construct the single-cell graph, where the nodes represent cells, and the edges represent the relationships between cells. The cell graph is built by the K-nearest neighbor (KNN) algorithm⁴² in the Scikit-learn Python package⁴³. The default K is pre-defined as 35 in this study and adjusted according to the datasets in our experiments. The generated adjacency matrix is a 0–1 matrix, where 1 represents being connected, and 0 represents no connection.

Empirical scRNA-seq data. We analyzed two different scRNA-seq datasets, namely HSPCs data and pancreatic cells data. HSPCs data and pancreatic cells data represent cells showing lineages relationship and cells showing distinct clusters, respectively. The HSPCs data are single-cell transcriptome data of FACS sorted CD34+ cells from human bone marrow mononuclear cells, accessible in the national genomics data center (HRA000084) and described in our previous study⁹. The pancreases cells data contains 10,000 single-cell transcriptomes with 14 distinct cell clusters, download from GEO (GSE84133)²⁸.

Competitive methods. Nine competitive methods, namely scDeepCluster, DCA, scVI, PCA, Ivis, SAUCIE, scGNN, ZINB-Wave, and PHATE, were compared with scGAE. Among these methods, scDeepCluster, DCA, scVI, Ivis, scGNN, and SAUCIE are deep learning based and showed the greatest potential. These methods usually generate hidden variables for downstream analysis, including visualization, clustering, and trajectory inference. The raw count matrix was used as input for DCA, scVI, scGNN, ZINB-Wave and scDeepCluster. For methods that take normalized data as input (scGAE, SAUCIE, PCA, Ivis, and PHATE), scTransform was used for data preprocessing. Each software was run following its manual and with default parameters. For SAUCIE, Ivis, and DCA, we first performed PCA to reduce the dimension to 100, 50, and 32 PCs, respectively. Ivis, SAUCIE, and PHATE directly generate the 2-dimensional embeddings. The cell clustering and trajectory inference were performed on the two-dimensional embeddings. scGNN and ZINB-Wave generated 128 and 10 dimensional embeddings. Both scGAE and PCA embedded simulated data to ten dimensions and embedded empirical data to 20 dimensions due to the complex structure of the empirical data. We performed tSNE to visualize data for these methods.

Data availability

The hematopoietic stem and progenitor cells (HSPCs) data is available in the Genome Sequence Archive in BIG Data Center, under accession numbers HRA000084. The data of human pancreatic cells is available through NCBI GEO with the accession number GSE84133.

Code availability

Accession codes The code and software of scGAE are available on GitHub (<https://github.com/ZixiangLuo1161/scGAE>).

Received: 23 May 2021; Accepted: 17 September 2021

Published online: 08 October 2021

References

- van der Maaten, L. & Hinton, G. Visualizing data using t-SNE. *J. Mach. Learn. Res.* **9**, 2579–2605 (2008).
- Jaitin, D. A. *et al.* Massively parallel single-cell RNA-seq for marker-free decomposition of tissues into cell types. *Science* **343**, 776–779 (2014).
- Klein, A. M. *et al.* Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell* **161**, 1187–1201 (2015).
- Macosko, E. Z. *et al.* Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. *Cell* **161**, 1202–1214 (2015).
- Qin, P. *et al.* Integrated decoding hematopoiesis and leukemogenesis using single-cell sequencing and its medical implication. *Cell Discov.* **7**, 1–17 (2021).
- Shapiro, E., Biezuner, T. & Linnarsson, S. Single-cell sequencing-based technologies will revolutionize whole-organism science. *Nat. Rev. Genet.* **14**, 618–630 (2013).
- Eraslan, G., Simon, L. M., Mircea, M., Mueller, N. S. & Theis, F. J. Single-cell RNA-seq denoising using a deep count autoencoder. *Nat. Commun.* **10**, 1–14 (2019).
- Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for single-cell transcriptomics. *Nat. Methods* **15**, 1053–1058 (2018).
- Tian, T., Wan, J., Song, Q. & Wei, Z. Clustering single-cell RNA-seq data with a model-based deep learning approach. *Nat. Mach. Intell.* **1**, 191–198 (2019).
- Moon, K. R. *et al.* Visualizing structure and transitions in high-dimensional biological data. *Nat. Biotechnol.* **37**, 1482–1492 (2019).
- Amodio, M. *et al.* Exploring single-cell data with deep multitasking neural networks. *Nat. Methods* **16**, 1139–1145 (2019).
- Wang, J. *et al.* scgcn is a novel graph neural network framework for single-cell RNA-seq analyses. *Nat. Commun.* **12**, 1–11 (2021).
- Risso, D., Perraudeau, F., Gribkova, S., Dudoit, S. & Vert, J.-P. A general and flexible method for signal extraction from single-cell RNA-seq data. *Nat. Commun.* **9**, 1–17 (2018).
- Szibert, B., Cole, J. E., Monaco, C. & Drozdov, I. Structure-preserving visualisation of high dimensional single-cell datasets. *Sci. Rep.* **9**, 1–10 (2019).
- Kipf, T. N. & Welling, M. Variational graph auto-encoders. *stat* **1050**, 21 (2016).
- Velickovic, P. *et al.* Graph attention networks. *stat* **1050**, 4 (2018).
- Min, E. *et al.* A survey of clustering with deep learning: From the perspective of network architecture. *IEEE Access* **6**, 39501–39514 (2018).
- Zappia, L., Phipson, B. & Oshlack, A. Splatter: Simulation of single-cell RNA sequencing data. *Genome Biol.* **18**, 1–15 (2017).
- Papadopoulos, N., Gonzalo, P. R. & Söding, J. Probst: Probabilistic simulation of single-cell RNA-seq data for complex differentiation processes. *Bioinformatics* **35**, 3517–3519 (2019).
- Haghverdi, L., Büttner, M., Wolf, F. A., Büttner, F. & Theis, F. J. Diffusion pseudotime robustly reconstructs lineage branching. *Nat. Methods* **13**, 845 (2016).
- Kendall, M. G. A new measure of rank correlation. *Biometrika* **30**, 81–93 (1938).
- Velten, L. *et al.* Human haematopoietic stem cell lineage commitment is a continuous process. *Nat. Cell Biol.* **19**, 271–281 (2017).
- Buenrostro, J. D. *et al.* Integrated single-cell analysis maps the continuous regulatory landscape of human hematopoietic differentiation. *Cell* **173**, 1535–1548 (2018).
- Hay, S. B., Ferchen, K., Chetal, K., Grimes, H. L. & Salomonis, N. The human cell atlas bone marrow single-cell interactive web portal. *Exp. Hematol.* **68**, 51–61 (2018).
- Karamitros, D. *et al.* Single-cell analysis reveals the continuum of human lympho-myeloid progenitor cells. *Nat. Immunol.* **19**, 85–97 (2018).
- Tusi, B. K. *et al.* Population snapshots predict early haematopoietic and erythroid hierarchies. *Nature* **555**, 54–60 (2018).
- Zheng, S., Papalexi, E., Butler, A., Stephenson, W. & Satija, R. Molecular transitions in early progenitors during human cord blood hematopoiesis. *Mol. Syst. Biol.* **14**, e8041 (2018).
- Baron, M. *et al.* A single-cell transcriptomic map of the human and mouse pancreas reveals inter- and intra-cell population structure. *Cell Syst.* **3**, 346–360 (2016).
- Bachem, M. G., Zhou, S., Buck, K., Schneiderhan, W. & Siech, M. Pancreatic stellate cells—role in pancreas cancer. *Langenbeck's Arch. Surg.* **393**, 891–900 (2008).
- Chen, J., Ma, T. & Xiao, C. Fastgcn. Fast learning with graph convolutional networks via importance sampling. In *International Conference on Learning Representations* (2018).
- Lerique, S., Abitbol, J. L. & Karsai, M. Joint embedding of structure and features via graph convolutional networks. *Appl. Netw. Sci.* **5**, 1–24 (2020).
- Grattarola, D. & Alippi, C. Graph neural networks in tensorflow and keras with spektral [application notes]. *IEEE Comput. Intell. Mag.* **16**, 99–106 (2021).
- Kipf, T. N. & Welling, M. Semi-supervised classification with graph convolutional networks. In *Proceedings of the 5th International Conference on Learning Representations, ICLR '17* (2017).
- Hamilton, W. L., Ying, R. & Leskovec, J. Inductive representation learning on large graphs. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS'17*, 1025–1035 (Curran Associates Inc., 2017).
- Du, J., Zhang, S., Wu, G., Moura, J. M. & Kar, S. Topology adaptive graph convolutional networks. arXiv preprint [arXiv:1710.10370](https://arxiv.org/abs/1710.10370) (2017).
- Yang, B., Fu, X., Sidiropoulos, N. D. & Hong, M. Towards k-means-friendly spaces: Simultaneous deep learning and clustering. In *International conference on machine learning*, 3861–3870 (PMLR, 2017).
- Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in large networks. *J. Stat. Mech. Theory Exp.* **2008**, P10008 (2008).
- Xie, J., Girshick, R. & Farhadi, A. Unsupervised deep embedding for clustering analysis. In *International conference on machine learning*, 478–487 (PMLR, 2016).

39. Shannon, C. E. A mathematical theory of communication. *Bell Syst. Tech. J.* **27**, 379–423 (1948).
40. Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biol.* **20**, 1–15 (2019).
41. Stuart, T. *et al.* Comprehensive integration of single-cell data. *Cell* **177**, 1888–1902 (2019).
42. Peterson, L. E. K-nearest neighbor. *Scholarpedia* **4**, 1883 (2009).
43. Pedregosa, F. *et al.* Scikit-learn: Machine learning in python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011).

Acknowledgements

The work of Jin was supported by National Key R&D Program of China (2018YFC1004500), National Natural Science Foundation of China (81872330, 31741077), the Shenzhen Innovation Committee of Science and Technology (JCYJ20170817111841427, ZDSYS20200811144002008), the Shenzhen Science and Technology Program (KQTD20180411143432337), and Center for Computational Science and Engineering, Southern University of Science and Technology. The work of Zhang was partially supported by the NSFC Grant (Nos. 11731006, 12071207), the Guangdong Basic and Applied Basic Research Foundation (2021A1515010359) and the Guangdong Provincial Key Laboratory of Computational Science and Material Design (No. 2019B030301001).

Author contributions

W.J. and Z.Z. conceived and designed the project. Z.L. and C.X. developed the algorithm, coded the program and performed the data analysis. W.J. and Z.L. wrote the manuscript with inputs from all authors.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-021-99003-7>.

Correspondence and requests for materials should be addressed to Z.Z. or W.J.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021