

Tim Zhang¹, Amirali Amirsoleimani², Mostafa Rahimi Azghadi³, Jason K Eshraghian⁴, Roman Genov⁵, and Yu Xia¹

¹Department of Bioengineering, McGill University

²Department of Electrical Engineering and Computer Science, York University

³College of Science and Engineering, James Cook University

⁴Department of Electrical and Computer Engineering, UC Santa Cruz

⁵Department of Electrical and Computer Engineering, University of Toronto

February 3, 2023

SSCAE: A Neuromorphic SNN Autoencoder for sc-RNA-seq Dimensionality Reduction

Tim Zhang¹, Amirali Amirsoleimani², Jason K. Eshraghian³, Mostafa Rahimi Azghadi⁴, Roman Genov⁵, and Yu Xia¹

¹Department of Bioengineering, McGill University, Montreal H3A 0E9, Canada

²Department of Electrical Engineering and Computer Science, York University, Toronto ON M3J 1P3, Canada

³Department of Electrical and Computer Engineering, UC Santa Cruz, CA 95064, United States

⁴College of Science and Engineering, James Cook University, QLD 4811, Australia

⁵Department of Electrical and Computer Engineering, University of Toronto, Toronto M5S, Canada

Abstract—Single-cell RNA sequencing is an emerging technique in the field of biology that departs radically from the previous assumption of gene-expression homogeneity within a tissue. The large quantity of data generated by this technology enables discoveries of cellular biology and disease mechanics that were previously not possible, and calls for accurate, scalable, and efficient processing pipelines. In this work, we propose SSCAE (spiking single-cell autoencoder), a novel SNN-based autoencoder for sc-RNA-seq dimensionality reduction. We apply this architecture to a variety of datasets, and the results show that it can match and surpass the performance of current state-of-the-art techniques. Moreover, the potential of this technique lies in its ability to be scaled up and to take advantage of neuromorphic hardware, circumventing the memory bottleneck that currently limits the size of sequencing datasets that can be processed.

Index Terms—Single Cell RNA, Next-gen sequencing, Spiking Neural Network, Deep Learning

I. INTRODUCTION

RNA sequencing has been a fundamental technique in the field of biology and genomics, giving researchers the ability to quantitatively analyze the mRNA molecules within samples, thus enabling downstream studies in cellular processes. Sequencing technologies have undergone multiple iterations of improvements to throughput, speed, and accuracy. The recent popularization of next-gen sequencing (NGS) has fueled the emergence of the single-cell RNA sequencing (sc-RNA-seq) field, pioneered in 2009 [1]. Departing from the conventional bulk RNA-seq technologies only capable of reflecting gene expression profiles on the population level, single-cell RNA sequencing (sc-RNA-seq) allows for the reading of transcripts on the individual cell level.

Recently, it was shown that gene expression is heterogeneous within a population [2], [3], departing from the previous assumption that cell populations within the tissue are homogeneous. It is deduced that expression heterogeneities can lead to cell differentiation [4] and cancer progression [5], [6].

The general data processing pipeline is shown in Fig.1. The raw counts are first sequenced with technologies such as NGS followed by alignment and deduplication. Preprocessing is performed to remove low-quality cells and normalize data.

Dimensionality reduction is then applied to the processed reads-matrices before potential clustering and classification. Visualizing the cells in lower dimensions is essential for downstream applications in order to identify cell groups and cell lineages.

The main challenge facing sc-RNA-seq data analysis is the amount of noise present in the data, which has proven to be much more significant than noise in bulk RNA sequencing [7]. Common dimensionality-reduction methods, including PCA and t-SNE, are widely applied to sc-RNA-seq datasets. However, such conventional techniques suffer from such noisy data which leads to significant performance degradation. Additionally, such techniques also struggle with preserving large-scale information such as intercluster relationships [8].

Recent advances in machine learning and ANNs (artificial neural networks) have enabled new sc-RNA-seq data processing techniques based on autoencoders [9]–[15], and tends to better capture the underlying data representations compared to rule-based heuristics. A hurdle facing deep autoencoder architectures and large-scale neural networks for sc-RNA-seq processing is the immense amount of data that is needed, and potentially limits scalability on modern hardware accelerators.

Spiking neural networks (SNNs) are a new breed of ANN that more closely resemble a biological neural network, where information is communicated between neurons in sparse binary spike trains instead of high precision activations, as in conventional deep neural networks. It is considered to be the third generation of neural networks [16]. Such a network can be much faster to execute on neuromorphic hardware [17], [18] and also boasts high noise tolerance [19], [20], both of which are characteristics that offer practical benefits to sc-RNA-seq.

In this work, we propose SSCAE (spiking single-cell autoencoder), a novel SNN-based autoencoder architecture for single-cell RNA sequencing (sc-RNA-seq) data dimensionality reduction. We apply this architecture on a variety of datasets, and the results show that it is able to match and surpass the performance of current state-of-art techniques. The objective of this work is to demonstrate the efficacy of using neuromorphic computing in the field of bioinformatics and we demonstrate a series of advantages that it may offer in applications. To

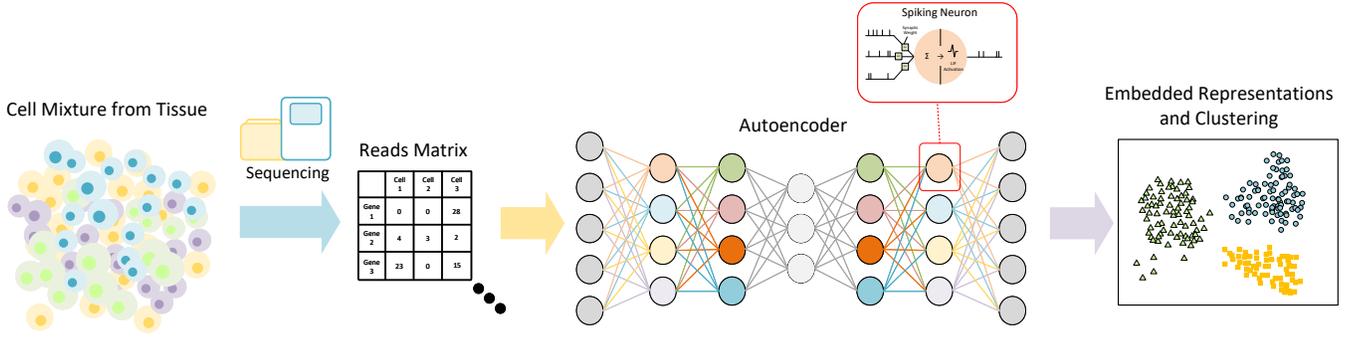


Fig. 1: SSCAE pipeline overview: Cell mixtures collected from tissues are first sequenced, then the reads matrix is passed as input to the spiking autoencoder. The cells embedded in the lower-dimensional space are visualized and clustered.

the best of our knowledge, this is the first work to utilize neuromorphic algorithms in the field of bioinformatics, as applied to single-cell RNA sequencing data.

II. METHODS

A. Network Architecture

Our work builds the autoencoder using the `snnTorch` framework from [21]. A leaky integrate-and-fire (LIF) neuron model used is modelled mathematically as in Eq. 1, where U is the membrane potential, β is the decay rate, W is the weight matrix, S_{out} is the output spike generated by the neuron, and θ is the firing threshold of the neuron. Direct-input-encoding scheme is used as in [22], where input data is treated as a current X injected into a neuron. Backpropagation through time (BPTT) is used to train the network (Eq.2), where \mathcal{L} is the loss, $W[s]$ is the weight at time $[s]$, noting that W is the same across all time s . In the forward pass, a Heaviside step function is used to model neuronal spiking. As it is a non-differentiable function, a surrogate gradient function is required to approximate the Heaviside function for gradient calculation. The arctan surrogate gradient is chosen to mitigate the “dead neuron” and “vanishing gradient” problem [21], [23], as it provided the best results when compared with other surrogate gradient functions during testing.

$$U[t] = \underbrace{\beta U[t-1]}_{\text{decay}} + \underbrace{WX[t]}_{\text{input}} - \underbrace{S_{\text{out}}[t-1]\theta}_{\text{reset}} \quad (1)$$

$$\frac{\partial \mathcal{L}}{\partial W} = \sum_t \frac{\partial \mathcal{L}[t]}{\partial W} = \sum_t \sum_{s \leq t} \frac{\partial \mathcal{L}[t]}{\partial W[s]} \frac{\partial W[s]}{\partial W} \quad (2)$$

The architecture consists of a 3-layer encoder network and a 3-layer decoder network, with a latent dimension of 2. A dropout layer is added at the input to counter the effect of dropout noise from sequencing. Inspired by [9], [24], a zero-inflation (ZI) layer is added after the decoder, which models the dropout event with a probability distribution of $e^{-\tilde{y}^2}$, where \tilde{y} is the reconstructed value of a gene. Since backpropagation cannot function on probability distributions, a Gumbel-softmax distribution is used for reparameterization [9], [25]. The firing threshold is set to a large value $\theta \rightarrow \infty$

for the final layer and its membrane potential is used as the reconstructed value.

B. Datasets

Multiple datasets were used for validating SSCAE, all of which are acquired from the collection at Sanger Institute. Results from the four most demonstrative datasets are presented to illustrate the behavior of SSCAE under a variety of conditions, cell-types, sample sizes, and difficulties, comparing SSCAE against the three other most commonly used methods: PCA, t-SNE, and UMAP. Out of these, two datasets are obtained from [26] encompassing 14 cell-types from the transcriptomic map of human pancreases, one is obtained from [27] of human embryonic cells, and one is from [28] sequencing human embryonic cells.

Each dataset is preprocessed using the Bioconductor package in R [29]. Raw reads are first quality-controlled to remove low-quality sequences, followed by normalization and log-transformation.

C. Performance Metrics

To quantitatively and qualitatively benchmark our unsupervised learning algorithm, we employ the following metrics. For NMI and ARI, K-means is first applied to the dimensionality-reduced data.

1) *Silhouette Coefficient*: The silhouette coefficient is a common technique for evaluating supervised learning algorithms defined by Eq. 3, where C_I is the number of points in class I , $d(i, j)$ is the distance between data points i and j in class I . The coefficient measures the quality of the clusters by calculating the ratio between the intra-cluster distances (compactness) with the inter-cluster distances (separation), ranging from -1 to 1, where a score closer to 1 indicates more distinct clusters.

$$\begin{aligned} a(i) &= \frac{1}{|C_I| - 1} \sum_{j \in C_I, i \neq j} d(i, j) \\ b(i) &= \min_{J \neq I} \frac{1}{|C_J|} \sum_{j \in C_J} d(i, j) \\ s(i) &= \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}, \text{ if } |C_I| > 1 \end{aligned} \quad (3)$$

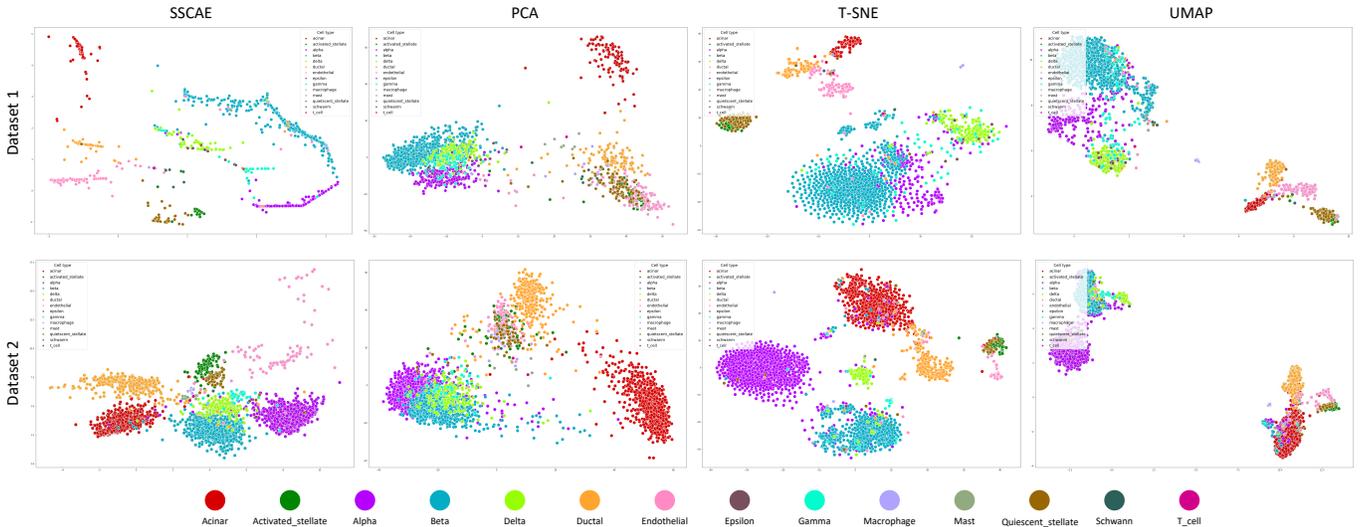


Fig. 2: The lower-dimensional representation of the cells, color-coded by the ground-truth cell-type labels.

2) *Normalized Mutual Information (NMI)*: NMI is defined as in Eq. 4, where X is the predicted label, Y is the ground truth label, $I(X, Y)$ is the mutual information between X and Y , $H(X)$ and $H(Y)$ are the entropies of X and Y .

$$\text{NMI}(\mathbf{X}, \mathbf{Y}) = \frac{I(\mathbf{X}, \mathbf{Y})}{\sqrt{H(\mathbf{X})H(\mathbf{Y})}} \quad (4)$$

3) *Adjusted Rand Index*: ARI is another common metric of cluster validity, often used in conjunction with the NMI and is used to measure the similarity between the predicted clusters and ground truth clusters, calculated from the contingency table.

4) *Preservation of Pairwise Distances*: Many of the downstream applications for sc-RNA-seq data not only requires reduced dimensionality for the distinct clusters, but they also demand the embedded cell populations to preserve the local and global structure of the original data [8]. Hence, we plot the pairwise distances between points in the latent-space versus the pairwise distances in the original higher-dimensional space, and calculate the Pearson correlation coefficient. Moreover, we qualitatively compare the distribution of distances in the original higher-dimensional space against the distribution of distances in the latent space.

III. EXPERIMENTAL RESULTS

We ran SSCAE against the three other most popular techniques on two of the most challenging datasets from Hemberg Group’s collection. The dimensionality-reduced representations are shown in Fig. 2(a), each datapoint is color-coded according to the cell-type. Visually, all four techniques show a certain degree of class separation, with non-linear techniques, SSCAE, t-SNE, and UMAP generating more distinct clusters than PCA (the linear technique) across the two datasets. Closer inspection reveals that SSCAE generally provides better cluster separation than t-SNE and UMAP, as evidenced by the

inability of the latter two to resolve Schwann cells (dark green) and quiescent-stellate (brown) cells. Additionally, the SSCAE shows significantly less overlap between alpha (purple), beta (cyan), gamma (light blue), and delta (light green) compared to the other methods.

This observation is supported by quantitative measurements as seen in Fig.3. Across datasets, SSCAE achieved the highest score in all three metrics. t-SNE and UMAP exhibited similar behaviour and shared similar NMI scores, with UMAP generally providing better results than t-SNE. Silhouette scores show a greater advantage of SSCAE, validating the more compact clusters and better separations.

	SSCAE	PCA	t-SNE	UMAP
NMI	0.75	0.57	0.58	0.65
ARI	0.63	0.35	0.29	0.33
Silhouette	0.41	0.14	0.20	0.25

	SSCAE	PCA	t-SNE	UMAP
NMI	0.73	0.57	0.63	0.64
ARI	0.65	0.37	0.40	0.39
Silhouette	0.49	0.23	0.20	0.16

Fig. 3: A heatmap showing the clustering metrics from each of the four techniques across two datasets.

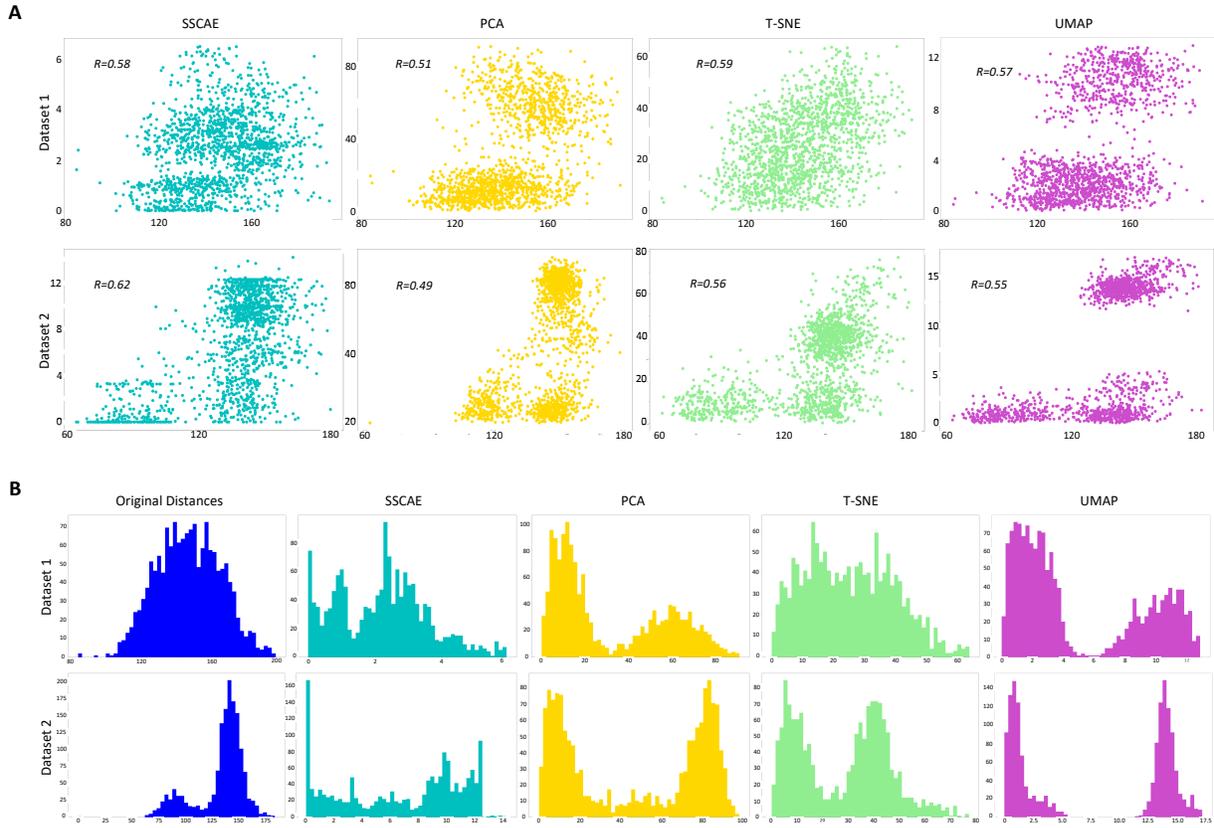


Fig. 4: A: Distances in original higher-dimensional space vs. distances in embedded space, with Pearson coefficients. B: Distributions of distances in original higher-dimensional space compared against distributions in embedded space.

Subsequently, the distance preservation properties are examined to investigate SSSCAE’s ability to capture the original dataset’s local and global structure. As shown in Fig.4A, all four techniques show a positive correlation between distances in the original higher-dimensional space versus the distances in the embedding. Their respective Pearson correlation coefficients are relatively similar. It is worth noting that, despite similar scores, PCA and UMAP both display discontinuities in the plots suggesting inconsistencies in the data structure preservation. This is also evident in Fig.4B, where the distance distributions are examined in original data space and compared with the embedded space. In dataset 1, t-SNE and SSSCAE’s distributions best resemble the original data’s normal distribution. In dataset 2, the four techniques displayed similar behaviour, while PCA, t-SNE and UMAP display a bimodal behaviour, which can be visually identified in Fig.2 as the data forms certain large clusters very far from each other, but each large cluster contains many overlapping unresolved smaller clusters.

IV. DISCUSSIONS AND FUTURE OUTLOOK

Through this work, SSSCAE has demonstrated its efficacy in sc-RNA-seq dimensionality reduction, resolving data clusters while preserving the data structure. The potential of this technique lies in its ability to be scaled up and to take

advantage of the benefits of neuromorphic hardware such as IBM TrueNorth [30], Intel Loihi [31], and emerging in-memory computing hardware [32], which may alleviate the memory bottleneck currently limiting the scalability of sc-RNA-seq processing for larger datasets. For future work, the exact power consumption and computational time savings will be evaluated on neuromorphic hardware and be compared against existing techniques. Additionally, the noise tolerance properties should be further investigated on larger datasets.

V. CONCLUSION

In this work, we proposed the novel SSSCAE spiking autoencoder architecture designed for sc-RNA-seq data processing. We implemented this pipeline on 2 datasets and examined the results through several aspects, showing that our technique is able to generate more compact clusters while maintaining better separation between clusters. Furthermore, the local and global data structure is more accurately preserved with SSSCAE which is essential for downstream applications. Additionally, this work demonstrates the efficacy of neuromorphic computing in the field of bioinformatics, bringing a myriad of potential benefits including improved noise tolerance, processing speed and computational efficiency. The hardware-algorithm co-optimization should be further investigated and the potential benefits will be quantified.

REFERENCES

- [1] F. Tang, C. Barbacioru, Y. Wang, E. Nordman, C. Lee, N. Xu, X. Wang, J. Bodeau, B. B. Tuch, A. Siddiqui *et al.*, “mrna-seq whole-transcriptome analysis of a single cell,” *Nature methods*, vol. 6, no. 5, pp. 377–382, 2009.
- [2] A. K. Shalek, R. Satija, J. Shuga, J. J. Trombetta, D. Gennert, D. Lu, P. Chen, R. S. Gertner, J. T. Gaubblomme, N. Yosef *et al.*, “Single-cell rna-seq reveals dynamic paracrine control of cellular variation,” *Nature*, vol. 510, no. 7505, pp. 363–369, 2014.
- [3] S. Huang, “Non-genetic heterogeneity of cells in development: more than just noise,” *Development*, vol. 136, no. 23, pp. 3853–3862, 2009.
- [4] A. S. Cuomo, D. D. Seaton, D. J. McCarthy, I. Martinez, M. J. Bonder, J. Garcia-Bernardo, S. Amaty, P. Madrigal, A. Isaacson, F. Buettner *et al.*, “Single-cell rna-sequencing of differentiating ips cells reveals dynamic genetic effects on gene expression,” *Nature communications*, vol. 11, no. 1, pp. 1–14, 2020.
- [5] P. Bischoff, A. Trinks, B. Obermayer, J. P. Pett, J. Wiederspahn, F. Uhlitz, X. Liang, A. Lehmann, P. Jurmeister, A. Elsner *et al.*, “Single-cell rna sequencing reveals distinct tumor microenvironmental patterns in lung adenocarcinoma,” *Oncogene*, vol. 40, no. 50, pp. 6748–6758, 2021.
- [6] Y. Zhang, D. Wang, M. Peng, L. Tang, J. Ouyang, F. Xiong, C. Guo, Y. Tang, Y. Zhou, Q. Liao *et al.*, “Single-cell rna sequencing in cancer research,” *Journal of Experimental & Clinical Cancer Research*, vol. 40, no. 1, pp. 1–17, 2021.
- [7] B. Hwang, J. H. Lee, and D. Bang, “Single-cell rna sequencing technologies and bioinformatics pipelines,” *Experimental & molecular medicine*, vol. 50, no. 8, pp. 1–14, 2018.
- [8] E. Becht, L. McInnes, J. Healy, C.-A. Dutertre, I. W. Kwok, L. G. Ng, F. Ginhoux, and E. W. Newell, “Dimensionality reduction for visualizing single-cell data using umap,” *Nature biotechnology*, vol. 37, no. 1, pp. 38–44, 2019.
- [9] D. Wang and J. Gu, “Vasc: dimension reduction and visualization of single-cell rna-seq data by deep variational autoencoder,” *Genomics, proteomics & bioinformatics*, vol. 16, no. 5, pp. 320–331, 2018.
- [10] D. Tran, H. Nguyen, B. Tran, C. La Vecchia, H. N. Luu, and T. Nguyen, “Fast and precise single-cell data analysis using a hierarchical autoencoder,” *Nature communications*, vol. 12, no. 1, pp. 1–10, 2021.
- [11] T. A. Geddes, T. Kim, L. Nan, J. G. Burchfield, J. Y. Yang, D. Tao, and P. Yang, “Autoencoder-based cluster ensembles for single-cell rna-seq data analysis,” *BMC bioinformatics*, vol. 20, no. 19, pp. 1–11, 2019.
- [12] D. Talwar, A. Mongia, D. Sengupta, and A. Majumdar, “Autoimpute: Autoencoder based imputation of single-cell rna-seq data,” *Scientific reports*, vol. 8, no. 1, pp. 1–11, 2018.
- [13] R. Lopez, J. Regier, M. B. Cole, M. I. Jordan, and N. Yosef, “Deep generative modeling for single-cell transcriptomics,” *Nature methods*, vol. 15, no. 12, pp. 1053–1058, 2018.
- [14] J. Ding, A. Condon, and S. P. Shah, “Interpretable dimensionality reduction of single cell transcriptome data with deep generative models,” *Nature communications*, vol. 9, no. 1, pp. 1–13, 2018.
- [15] E. Lin, S. Mukherjee, and S. Kannan, “A deep adversarial variational autoencoder model for dimensionality reduction in single-cell rna sequencing analysis,” *BMC bioinformatics*, vol. 21, no. 1, pp. 1–11, 2020.
- [16] W. Maass, “Networks of spiking neurons: the third generation of neural network models,” *Neural networks*, vol. 10, no. 9, pp. 1659–1671, 1997.
- [17] M. Pfeiffer and T. Pfeil, “Deep learning with spiking neurons: opportunities and challenges,” *Frontiers in neuroscience*, p. 774, 2018.
- [18] M. R. Azghadi, C. Lammie, J. K. Eshraghian, M. Payvand, E. Donati, B. Linares-Barranco, and G. Indiveri, “Hardware implementation of deep network accelerators towards healthcare and biomedical applications,” *IEEE Transactions on Biomedical Circuits and Systems*, vol. 14, no. 6, pp. 1138–1159, 2020.
- [19] S. Park, D. Lee, and S. Yoon, “Noise-robust deep spiking neural networks with temporal information,” in *2021 58th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2021, pp. 373–378.
- [20] J. K. Eshraghian, X. Wang, and W. D. Lu, “Memristor-based binarized spiking neural networks: Challenges and applications,” *IEEE Nanotechnology Magazine*, vol. 16, no. 2, pp. 14–23, 2022.
- [21] J. K. Eshraghian, M. Ward, E. Neftci, X. Wang, G. Lenz, G. Dwivedi, M. Bennamoun, D. S. Jeong, and W. D. Lu, “Training spiking neural networks using lessons from deep learning,” *arXiv preprint arXiv:2109.12894*, 2021.
- [22] B. Rueckauer, I.-A. Lungu, Y. Hu, M. Pfeiffer, and S.-C. Liu, “Conversion of continuous-valued deep networks to efficient event-driven networks for image classification,” *Frontiers in neuroscience*, vol. 11, p. 682, 2017.
- [23] W. Fang, Z. Yu, Y. Chen, T. Huang, T. Masquelier, and Y. Tian, “Deep residual learning in spiking neural networks,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 21 056–21 069, 2021.
- [24] E. Pierson and C. Yau, “Zifa: Dimensionality reduction for zero-inflated single-cell gene expression analysis,” *Genome biology*, vol. 16, no. 1, pp. 1–10, 2015.
- [25] E. Jang, S. Gu, and B. Poole, “Categorical reparameterization with gumbel-softmax,” *arXiv preprint arXiv:1611.01144*, 2016.
- [26] M. Baron, A. Veres, S. L. Wolock, A. L. Faust, R. Gaujoux, A. Vetere, J. H. Ryu, B. K. Wagner, S. S. Shen-Orr, A. M. Klein *et al.*, “A single-cell transcriptomic map of the human and mouse pancreas reveals inter- and intra-cell population structure,” *Cell systems*, vol. 3, no. 4, pp. 346–360, 2016.
- [27] M. Goolam, A. Scialdone, S. J. Graham, I. C. Macaulay, A. Jedrusik, A. Hupalowska, T. Voet, J. C. Marioni, and M. Zernicka-Goetz, “Heterogeneity in oct4 and sox2 targets biases cell fate in 4-cell mouse embryos,” *Cell*, vol. 165, no. 1, pp. 61–74, 2016.
- [28] L. Yan, M. Yang, H. Guo, L. Yang, J. Wu, R. Li, P. Liu, Y. Lian, X. Zheng, J. Yan *et al.*, “Single-cell rna-seq profiling of human preimplantation embryos and embryonic stem cells,” *Nature structural & molecular biology*, vol. 20, no. 9, pp. 1131–1139, 2013.
- [29] R. C. Gentleman, V. J. Carey, D. M. Bates, B. Bolstad, M. Dettling, S. Dudoit, B. Ellis, L. Gautier, Y. Ge, J. Gentry *et al.*, “Bioconductor: open software development for computational biology and bioinformatics,” *Genome biology*, vol. 5, no. 10, pp. 1–16, 2004.
- [30] F. Akopyan, J. Sawada, A. Cassidy, R. Alvarez-Icaza, J. Arthur, P. Merolla, N. Imam, Y. Nakamura, P. Datta, G.-J. Nam *et al.*, “Truenorth: Design and tool flow of a 65 mw 1 million neuron programmable neurosynaptic chip,” *IEEE transactions on computer-aided design of integrated circuits and systems*, vol. 34, no. 10, pp. 1537–1557, 2015.
- [31] M. Davies, N. Srinivasa, T.-H. Lin, G. Chinya, Y. Cao, S. H. Choday, G. Dimou, P. Joshi, N. Imam, S. Jain *et al.*, “Loihi: A neuromorphic manycore processor with on-chip learning,” *Ieee Micro*, vol. 38, no. 1, pp. 82–99, 2018.
- [32] M. Rahimi Azghadi, Y.-C. Chen, J. K. Eshraghian, J. Chen, C.-Y. Lin, A. Amirsoleimani, A. Mehonic, A. J. Kenyon, B. Fowler, J. C. Lee *et al.*, “Complementary metal-oxide semiconductor and memristive hardware for neuromorphic computing,” *Advanced Intelligent Systems*, vol. 2, no. 5, p. 1900189, 2020.