

Gene Regulatory Network Inference Based on Convolutional GRU

Siying Du*

School of Information, Yunnan Normal University, Kunming, China

1960089397@qq.com

**Corresponding author*

Keywords: Convolutional GRU, Gene Regulatory Network, Deep Learning

Abstract: Time-course single-cell RNA sequencing (scRNA-seq) data reflect gene expression changes over time, offering a valuable resource for exploring dynamic gene interactions and building dynamic gene regulatory networks (GRNs). However, most existing methods are typically designed for bulk RNA sequencing (bulk RNA-seq) data and cannot be directly applied to time-course scRNA-seq data. Addressing this issue, we present CGGRN, an approach based on convolutional gated recurrent unit (GRU) for inferring GRNs. CGGRN transforms time-course data into images, including raw pairwise gene images and neighborhood images, and aggregates them with time point information into a four-dimensional tensor. The tensor is then fed into the convolutional GRU to capture features for each gene pair and reconstruct the GRN. We conducted trials on four time-course scRNA-seq datasets using CGGRN, and the outcomes show that CGGRN surpasses existing models in constructing GRN.

1. Introduction

Gene regulatory networks are among the most important and common biological networks in biology. Their core function is to capture and describe the various complex processes that influence gene expression, thereby affecting the physiological state and behavior patterns of cells. A major objective of systems biology is the accurate depiction of these intracellular regulatory associations, effectively incorporating genes, transcription factors, and their interplays into a dynamic network architecture. By reconstructing gene regulatory networks, scientists can gain deeper insights into the fundamental mechanisms of gene function, further advancing the understanding of cellular functions and the study of complex disease mechanisms [1, 2].

High-throughput sequencing technologies and efficient cell separation techniques have laid a solid foundation for modern single-cell sequencing platforms. RNA sequencing (RNA-seq) technology has made in-depth studies of the entire transcriptome possible, driving many important biological discoveries and has become one of the widely used technologies in biomedical research. Common sequencing data are mainly categorized into bulk RNA-seq data and scRNA-seq data. Bulk RNA-seq technology performs high-throughput sequencing of RNA from millions of cells, providing the average expression level of each gene [3, 4], making it suitable for revealing the overall trends in gene expression. However, the limitation of bulk sequencing lies in its inability to accurately quantify

the RNA content of low-abundance cells and potential bias in results when heterogeneous cell populations are present in the sample.

In contrast, scRNA-seq technology allows for the isolation of individual cells, transcript capture, and library construction for sequencing, enabling unprecedented resolution for analyzing the basic biological characteristics of cell populations and biological systems. It provides more precise, cell-level information. The application of this technology has enabled researchers to unravel the heterogeneity and dynamic changes in gene expression at the single-cell level, significantly improving the precision and detail of GRN construction. By combining advanced computational methods, such as deep learning, the inference and analysis of GRNs have become more efficient and accurate, providing powerful support for uncovering the complex mysteries of life.

In recent years, scRNA-seq technology has made considerable advancements, allowing for unbiased, reproducible, high-resolution, and high-throughput transcriptomic analysis of individual cells [5-7]. Unlike conventional bulk RNA sequencing, single-cell RNA sequencing captures gene activity at the resolution of individual cells, enabling the identification of transcriptomic features specific to various cell types in biological tissues. This technology provides a comprehensive view of gene expression heterogeneity between cells and facilitates the reconstruction of cell type-specific GRNs [8-11].

With the development of deep learning technology, more and more GRN reconstruction methods have emerged, particularly for single-cell sequencing data. However, many existing methods are mainly designed for static data, which makes them limited when handling time-course data and difficult to capture time features. As a result, some researchers have attempted to use pseudo-time methods to infer GRNs from scRNA-seq data, but these methods are not suitable for true time-course data and fail to fully consider the dependencies and dynamic changes between time points. Therefore, how to design methods capable of handling time-course data and effectively capturing the dynamic changes of gene regulatory networks remains a challenging problem.

Researchers can use deep learning techniques to interpret data from complex distributions or interaction patterns, enabling the effective reconstruction of GRNs from single-cell time-course data. For example, the dynGENIE3 [12] method utilizes Ordinary Differential Equations (ODE) to describe dynamic changes in gene expression. The TDL [13] method proposed by Yuan et al. converts the data into image form and aggregates it with time points into a three-dimensional tensor as the model input. However, this method does not consider neighborhood context information, which may lead to false positives. dynDeepDRIM [14], while converting gene pairs into image form, also generates neighborhood images of gene pairs and inputs them along with time points into CNNs to predict interactions between gene pairs. However, this method requires considerable computational resources.

Inspired by the above works, we propose a gene regulatory network method based on convolutional GRU. Initially, the data is transformed into images, including raw gene pair images and neighborhood images, which are then combined with time points into a four-dimensional tensor. This tensor is subsequently fed into a convolutional GRU to capture the time and spatial characteristics of the data and infer the gene regulatory network. The key contributions of this method are:

- (1) The data is converted into images, and a convolutional GRU model is employed to learn the time and spatial features, followed by the inference of the GRN.
- (2) CGGRN is compared with other existing methods across four distinct time-course datasets, and experimental results demonstrate that CGGRN outperforms current models in gene regulatory network inference.

2. Proposed Framework

We propose a gene regulatory network inference method based on convolutional GRU, named

CGGRN. The method first constructs the main images and neighborhood images of gene pairs, which are then aggregated with time point information into a four-dimensional tensor as the model input. Specifically, for a gene pair (m, n) , its image at time point t is denoted as $I_{m,n}^t$, referred to as the main image. Additionally, neighborhood information of the gene pair is considered, mainly consisting of the gene's own image and the neighborhood images of the gene pair, represented as $I_{m,m}^t, I_{n,n}^t$, and $\{I_{m,u_1}^t, I_{m,u_2}^t, \dots, I_{m,u_r}^t, I_{n,v_1}^t, I_{n,v_2}^t, \dots, I_{n,v_r}^t\}$, where $\{u_1, u_2, \dots, u_r\}$ and $\{v_1, v_2, \dots, v_r\}$ refer to the top r genes that show significant covariance with the two genes in the pair. After generating the images, for each time point in the data, the main images and neighborhood images of the gene pairs are stacked to form a three-dimensional tensor. Then, the three-dimensional tensors of each time point are aggregated, ultimately forming a four-dimensional tensor, which serves as the model input. Next, the convolutional GRU method is employed, leveraging the strengths of both convolutional and recurrent neural networks to effectively capture the time and spatial features of the data. Ultimately, the gene regulatory network is inferred.

3. Experiments and Analysis

In this section, CGGRN is compared with other existing techniques using four time-course datasets from mice and humans. To ensure the randomness of data splitting, the entire dataset is randomly divided into training and testing sets with an 8:2 ratio. Within the training set, 20% of the data is further randomly selected as a validation set to monitor the model's performance in real-time, and the best model is saved based on the performance on the validation set. AUROC is chosen as the evaluation metric for model performance. The results demonstrate that CGGRN performs better than the existing methods.

3.1 Datasets

We selected four time-course scRNA-seq datasets from mice and humans, which are detailed in Table 1. These include the mouse embryonic stem cell datasets (mESC1 and mESC2) [15, 16], and the human embryonic stem cell datasets (hESC1 and hESC2) [17, 18]. The human embryonic stem cell and mESC1 datasets contain 36 transcription factors (TFs), while mESC2 contains 38 TFs. We classify TF-gene pairs as positive if significant peak signals are detected in the promoter region of the target gene, and negative if no significant peaks are observed.

Table 1: Detailed information of time-course scRNA-seq datasets

Datasets	Points	Genes	Cells
hESC1	5	26178	1529
hESC2	6	19189	758
mESC1	9	23481	3456
mESC2	4	24175	2717

3.2 Evaluation Metrics

In this study, we use the area under the receiver operating characteristic curve (AUROC) to evaluate the model's performance. AUROC measures the overall performance of the classifier by calculating the area under the receiver operating characteristic curve, reflecting the model's classification ability at different classification thresholds. Specifically, the closer the AUROC value is to 1, the better the model's performance and the more accurate the classification.

3.3 Comparative Experiments and Analysis

We assessed the effectiveness of CGGRN in inferring GRN on four time-course datasets and compared it with several existing methods, including dynGENIE3, TDL-3DCNN, TDL-LSTM, and dynDeepDRIM. dynGENIE3 describes the dynamic changes in gene expression by using ODE. TDL-3D CNN and TDL-LSTM are two methods within TDL, which use 3D CNN and long short-term memory networks (LSTM), respectively, to predict the complex regulatory relationships between genes. The dynDeepDRIM method converts each gene pair into a main image and neighborhood image, and combines the time point information to integrate these into a four-dimensional tensor, which is then used as the input to a CNN for GRN prediction.

The results of the comparison are shown in Figure 1, with different methods represented by different colors, the x-axis indicating the datasets, and the y-axis showing AUROC values. It is clear from the figure that CGGRN consistently delivers the highest AUROC values across the four time-course datasets. Additionally, on the hESC2 dataset, CGGRN's AUROC value is approximately 3% higher than that of dynDeepDRIM. The experimental results demonstrate that CGGRN exhibits high accuracy and effectiveness in gene regulatory network inference, outperforming existing methods.

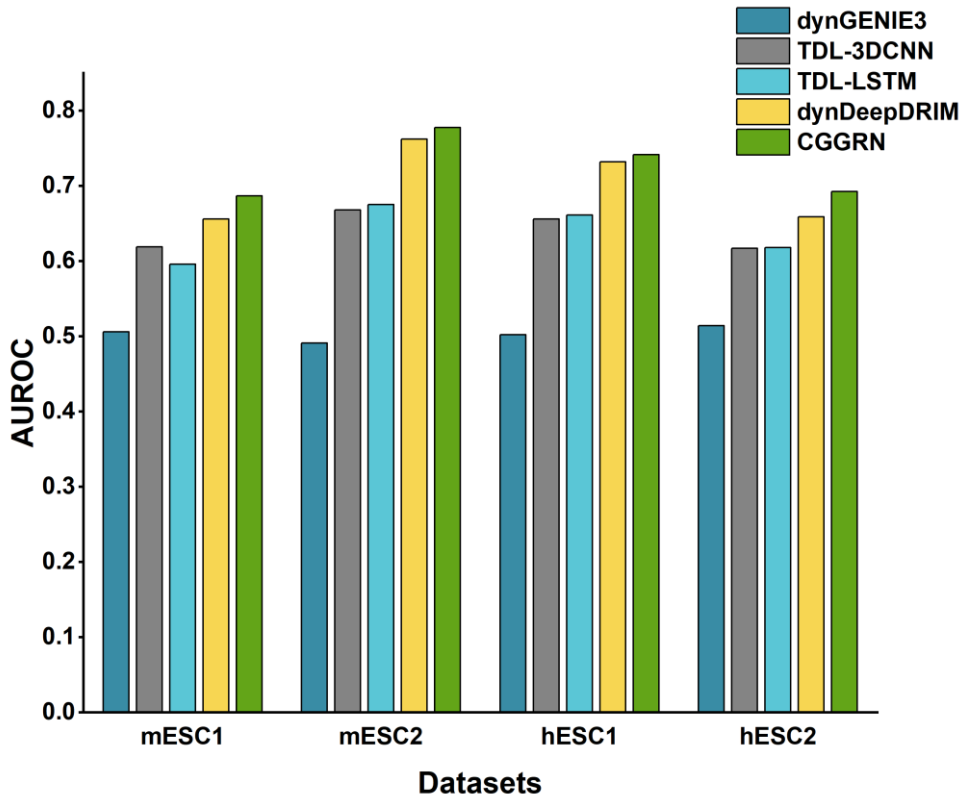


Figure 1: Comparison of CGGRN with other approaches on four time-course scRNA-seq datasets.

4. Conclusion

Gene regulation in organisms is often dynamic and accompanied by time-dependent information. Nonetheless, most current gene regulatory network inference methods are mainly built for static data, making them ill-suited for time-course data. To overcome this limitation, we propose a GRN inference approach based on convolutional GRU. Specifically, we first generate the main images of gene pairs and neighborhood images, which are then aggregated with time point information into a four-dimensional tensor, providing the model with spatial features and time features. These four-

dimensional tensors are then passed as input to the convolutional GRU network, where the convolutional layers effectively capture the spatial features of the data, and the GRU layers learn the time features. Finally, we use the trained model to infer the GRN. To validate the effectiveness of the proposed method, we evaluated the performance of CGGRN on four single-cell time-course datasets, using AUROC as the evaluation metric. The experimental results indicate that CGGRN outperforms existing methods in GRN inference, demonstrating its advantages in handling time-course data.

There is a causal relationship between transcription factors and their target genes, and this relationship is not static; it dynamically evolves with changes in the developmental stages of the organism. Therefore, how to effectively incorporate and model this causal relationship, and how to use it for precise gene regulatory network inference, is a direction worth exploring in our future research.

References

- [1] L. F. Iglesias-Martinez, W. Kolch, and T. Santra, "BGRMI: A method for inferring gene regulatory networks from time-course gene expression data and its application in breast cancer research," *Scientific Reports*, vol. 6, no. 1, p. 37140, 2016.
- [2] X. Zhang, J. Zhao, J.-K. Hao, X.-M. Zhao, and L. Chen, "Conditional mutual inclusive information enables accurate quantification of associations in gene regulatory networks," *Nucleic acids research*, vol. 43, no. 5, pp. e31-e31, 2015.
- [3] E. Shapiro, T. Biezuner, and S. Linnarsson, "Single-cell sequencing-based technologies will revolutionize whole-organism science," *Nature Reviews Genetics*, vol. 14, no. 9, pp. 618-630, 2013.
- [4] O. Stegle, S. A. Teichmann, and J. C. Marioni, "Computational and analytical challenges in single-cell transcriptomics," *Nature Reviews Genetics*, vol. 16, no. 3, pp. 133-145, 2015.
- [5] A. Zeisel et al., "Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq," *Science*, vol. 347, no. 6226, pp. 1138-1142, 2015.
- [6] B. Treutlein et al., "Reconstructing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq," *Nature*, vol. 509, no. 7500, pp. 371-375, 2014.
- [7] S. Islam et al., "Quantitative single-cell RNA-seq with unique molecular identifiers," *Nature methods*, vol. 11, no. 2, pp. 163-166, 2014.
- [8] A. Seb  Pedr  s et al., "Cnidarian cell type diversity and regulation revealed by whole-organism single-cell RNA-Seq," *Cell*, vol. 173, no. 6, pp. 1520-1534. e20, 2018.
- [9] P. W. Hook et al., "Single-cell RNA-seq of mouse dopaminergic neurons informs candidate gene selection for sporadic Parkinson disease," *The American Journal of Human Genetics*, vol. 102, no. 3, pp. 427-446, 2018.
- [10] S. Aibar et al., "SCENIC: single-cell regulatory network inference and clustering," *Nature methods*, vol. 14, no. 11, pp. 1083-1086, 2017.
- [11] A. Karbalayghareh, U. Braga-Neto, and E. R. Dougherty, "Intrinsically Bayesian robust classifier for single-cell gene expression time series in gene regulatory networks," in *Proceedings of the 8th ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*, 2017, pp. 766-767.
- [12] V. A. Huynh-Thu and P. Geurts, "dynGENIE3: dynamical GENIE3 for the inference of gene networks from time series expression data," *Scientific reports*, vol. 8, no. 1, p. 3384, 2018.
- [13] Y. Yuan and Z. Bar-Joseph, "Deep learning of gene relationships from single cell time-course expression data," *Briefings in bioinformatics*, vol. 22, no. 5, p. bbab142, 2021.
- [14] Y. Xu, J. Chen, A. Lyu, W. K. Cheung, and L. Zhang, "dynDeepDRIM: a dynamic deep learning model to infer direct regulatory interactions using time-course single-cell gene expression data," *Briefings in Bioinformatics*, vol. 23, no. 6, p. bbac424, 2022.
- [15] S. Semrau, J. E. Goldmann, M. Soumillon, T. S. Mikkelsen, R. Jaenisch, and A. Van Oudenaarden, "Dynamics of lineage commitment revealed by single-cell transcriptomics of differentiating embryonic stem cells," *Nature communications*, vol. 8, no. 1, p. 1096, 2017.
- [16] A. M. Klein et al., "Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells," *Cell*, vol. 161, no. 5, pp. 1187-1201, 2015.
- [17] L.-F. Chu et al., "Single-cell RNA-seq reveals novel regulators of human embryonic stem cell differentiation to definitive endoderm," *Genome biology*, vol. 17, pp. 1-20, 2016.
- [18] S. Petropoulos et al., "Single-cell RNA-seq reveals lineage and X chromosome dynamics in human preimplantation embryos," *Cell*, vol. 165, no. 4, pp. 1012-1026, 2016.