

# HALink: A VAE-Based and Hybrid GCN–GAT Architecture for Inferring Single-Cell Gene Regulatory Networks

Chaowang Lan<sup>1\*</sup>, Yulong Yuan<sup>1</sup>, Jingxin Wu<sup>1</sup>, Huiwu Zhang<sup>1</sup>, Xingyu Ji<sup>1</sup>, and Caihua Liu<sup>2</sup>

<sup>1</sup> School of Artificial Intelligence, Guilin University of Electronic Technology, No. 1 Jinji Road, 541004, Guangxi, China  
[chaowanglan@guet.edu.cn](mailto:chaowanglan@guet.edu.cn)

<sup>2</sup> Department of Computer Science, Cornell University, New York, United States of America

**Abstract.** Gene regulatory networks (GRNs) describe the complex interactions that control gene expression and are key indicators of cellular function and disease progression. Although numerous methods have been developed for inferring GRNs from single-cell transcriptomic data, current methods have two limitations: sensitivity to data sparsity and technical noise in single-cell RNA sequencing, and inadequate integration of both global structural information and local dependencies within GRNs. To address these limitations, we propose HALink. This method integrating a variational autoencoder-based model and Top-K sparsification strategy to address the data sparsity, reduce technical noise, and insufficient prior knowledge. Furthermore, it utilizes a hybrid Graph Convolutional Network–Graph Attention Network (GCN–GAT) architecture that simultaneously captures both global structural information and local dependencies within GRNs, enabling more comprehensive and accurate GRN construction. Experiments results demonstrate that both the VAE-based model with sparsification strategy and the hybrid GCN–GAT architecture are able to improve the prediction performance, while their integrated implementation yields synergistic improvements. Compared with the state-of-art methods, HALink achieves superior performance on most experimental cases across four benchmark datasets. Our method facilitates the identification of key gene regulatory networks underlying life processes and enables the discovery of biologically meaningful insights into disease pathogenesis.

**Keywords:** Gene regulatory network · Variational Autoencoder · sparsification strategy · hybrid architecture.

## 1 Introduction

Gene regulatory networks (GRNs) represent molecular interactions between transcription factors (TFs) and their target genes, capturing the dynamics of transcriptional regulation [1, 2]. These networks play vital roles in revealing cellular

functions, developmental processes, and deciphering genomic regulatory mechanisms [3–5]. Therefore, accurately inferring GRNs is of great importance for a wide range of biological applications [6], particularly in understanding complex human disease mechanisms, advancing disease prevention, diagnosis, and treatment, and aiding in the discovery of potential drug targets [7–9]. Many researches exemplify progress in this area. For example, Wang et al. developed the NetID method to construct regulatory modules in hematopoietic differentiation and uncovering critical regulators of cell fate decisions [10]. Similarly, Li et al. introduced STREAM to reveal key regulatory interactions associated with Alzheimer’s disease and lymphoma.[11]

The rapid advancement of high-throughput single-cell RNA sequencing technology [12], the inference of GRNs has emerged as a major focus in analyzing single-cell RNA data, leading to the development of diverse computational methods [13]. Existing GRN inference methods can generally be classified into two categories: correlation-based methods and graph neural network-based methods. The correlation-based methods employ statistical metrics, such as Pearson coefficient [14] and mutual information (MI) [15], to quantify pairwise gene co-expression patterns. These methods establish regulatory edges when the statistical metrics surpass predetermined thresholds. However, such methods often fail to capture nonlinear regulatory relationships and are highly sensitive to noise and sparsity data. With recent progress in graph representation learning, GNN-based methods have been increasingly applied to infer GRNs from single-cell RNA-seq data [16, 17]. For example, the GNNLink model employs graph neural networks to model regulatory relationships between genes [18]. The scMGATGRN model, which is based on graph attention networks, incorporates multi-view and view-level attention mechanisms to optimize gene representations [19]. GCLink combines graph contrastive learning with graph attention networks to infer GRNs [20]. DeepRIG [21] applies the graph-based deep learning model to infer GRNs. The GENELink model leverages known regulatory information to infer the gene pair interactions and predict potential regulatory edges [22, 23], achieving highly accuracy and stability of GRNs. Although there are many methods have been developed to infer GRNs, they have two principal limitations. Firstly, the single-cell RNA sequencing data is sparsity [24] and noise [25], which have significance negative effect on inferring GRNs. Secondly, the state-of-art methods typically emphasize either global structural information or local dependencies of graph, failing to effectively integrate both perspectives for inferring GRNs.

To overcome these limitations, we develop HALink, a novel method that synergistically integrates the complementary strengths of Graph Convolutional Networks (GCN) and Graph Attention Networks (GAT), to infer GRNs. This method comprises two key stages. Firstly, a variational autoencoder-based (VAE-based) model is employed to generate a sparse adjacency matrix and a Top-K sparsification strategy is utilized to eliminate the impact of data noise and sparsity. This sparse adjacency matrix enables to alleviate data sparsity and insufficient prior information, thereby significantly enhancing input data quality for subsequent graph learning. Secondly, a hybrid GCN-GAT architecture is utilized

to capture global structural information and local dependencies within GRNs. The architecture incorporates learnable weighting factors to dynamically balance the contributions from both GCN and GAT. The innovations of HALink has two aspects: (1) the combination of VAE-based model and the Top-K sparsification strategy effectively reduces noise and redundancy of the data, improving the quality of input data; (2) the hybrid GCN-GAT architecture capturing both global and local information within GRNs, enhancing the performance of inferring GRN. The experimental results demonstrate that HALink achieves significant performance improvement across multiple single-cell datasets.

## 2 Methods

The framework of HALink is presented in Fig.1. This method consists of two main steps: (1) Adjacency matrix generation and sparsification. and (2) Hybrid GCN-GAT architecture construction. The first step aims to mitigate the effects of data sparsity and noise, while the second step is designed to integrate both global structural and local dependency within GRNs. Given a gene expression matrix  $X \in \mathbb{R}^{N \times M}$ , where  $N$  and  $M$  correspond to the number of gene and sample, respectively.  $A^P \in \mathbb{R}^{N \times M}$  represents the gene-gene interaction prior adjacency matrix. The prior adjacency matrix is derived from [22].

### 2.1 Adjacency matrix generation and sparsification

The single-cell data often is sparsity and contains large number of noise, which have significance negative influence to infer GRNs. To address these problems, we apply the VAE-based model, which enable to fill in the missing parts of the prior information of the data and captures potential regulatory relationships, to generate an adjacency matrix. Then the Top-K sparsification is employed to further reduce noise and unnecessary computations by retaining the strongest gene interactions.

**VAE-based model for generating adjacency matrix** The VAE is a generative model that learns latent representations to recreate data. Inspired by VAE, we develop a VAE-based model to generate adjacency matrix. This model maps the gene expression matrix into latent variables, then recreating the gene-gene adjacency matrix. The framework of the VAE-based model is presented in Fig. 1A. For a gene expression matrix  $X$ , it is mapped to a Gaussian distribution with mean  $\mu$  and log  $\sigma^2$ . Then using the reparameterization trick, latent variables  $z$  are sampled as the equation 2.1:

$$z = \mu + \epsilon \cdot \sigma, \quad \epsilon \sim \mathcal{N}(0, I) \quad (1)$$

The loss function consists of a reconstruction loss and the KL divergence. The reconstruction loss calculates the difference between the original gene similarity

matrix and the reconstructed gene-gene adjacency matrix. The KL divergence measures the difference between the Gaussian distribution generated by the encoder and the standard normal Gaussian distribution. The loss function  $\mathcal{L}_{\text{total}}$  is calculated by equation 2.1.

$$\mathcal{L}_{\text{total}} = \|\hat{A} - X \cdot X\|^2 - 0.5 \cdot (1 + \log \sigma^2 - \mu^2 - \sigma^2) \quad (2)$$

Where  $\hat{A}$  is the output of the VAE-based model (also the reconstructed gene-gene adjacency matrix). This objective guides the training of VAE to learn the latent regulatory relationships among genes.

**The Top-K sparsification strategy** The reconstructed gene-gene adjacency matrix reflects the regulatory relationships between genes. For two genes with a low regulatory relationship, it can be indicated that there is no significant interaction between them. In order to identify the strong regulatory relationship gene pairs, we apply the Top-K sparsification strategy. This strategy aims to convert the gene-gene adjacency matrix  $\hat{A}$  into sparse adjacency matrix  $A^S$  by retaining the top-K strongest regulatory edges for each gene and truncating all weaker edges. given tow genes  $i$  and  $j$ , the  $i$ th line and  $j$ th column of spare adjacency matrix  $A_{ij}^S$  is set 1 if the distance between gene  $i$  and gene  $j$  is one of the  $N$  smallest for gene  $i$ , and to 0 otherwise.

## 2.2 Hybrid GCN-GAT architecture construction

The GCN enables each node to aggregate information from its multiple neighbors through iterative message passing. Therefore, it integrates information across the entire graph structure, capturing the global structural information of the graph. The GAT employs a self-attention mechanism to assign differential weights to the neighbors of each node through, thereby capturing the most informative local dependencies within the graph. In order to capture both global structural information and local dependencies, we develop a hybrid GCN-GAT architecture, which combines the GCN and GAT network, to infer GRNs.

Given a spare adjacency matrix  $A^S$  and a prior adjacency matrix  $A^P$ . These two matrix combines into a fusion adjacency matrix  $A^F$ . The equation of calculating fusion adjacency matrix is in equation 2.2:

$$A^F = A^S + A^P \quad (3)$$

i This fusion adjacency matrix is the origin input of the GCN and GAT branch. The symmetric normalization and random-walk normalization are applied to normalize the fusion adjacency matrix in GCN and GAT branch, respectively.

The GCN performs convolution operations on the fusion adjacency matrix  $A^F$  to iteratively learn node embeddings layer by layer. The computation of each GCN layer is defined by equation 4:

$$H_{GCN}^{(l+1)} = \sigma(A^F H_{GCN}^{(l)} W_{GCN}^{(l)}) \quad (4)$$

Where  $H_{GCN}^{(l+1)}$  represents the node embeddings at the  $l$ -th layer,  $W_{GCN}^{(l)}$  is the weight matrix at the  $l$ -th layer, and  $\sigma$  is the activation function. The output of the GCN, which is defined by  $\mathbf{h}_{GCN}$ , is the node embeddings that contains the information across the graph structure.

The GAT assigns different neighbor weights to each node through the self-attention mechanism, allowing for the dynamic adjustment of the influence of each neighboring node. The formula of each GAT layer is shown in in equation 2.2:

$$H_{GAT}^{(l+1)} = \sigma(\alpha^l H_{GAT}^{(l)} W_{GAT}^{(l)}) \quad (5)$$

Where  $\alpha^l$  is the attention coefficient matrix that calculated dynamically by the network and  $W_{GCN}^{(l)}$  is the weight matrix at the  $l$ -th layer. We define the output of GAT is  $\mathbf{h}_{GAT}$ .

To integrate the strengths of both GCN and GAT, we design a feature fusion layer that combines the outputs of GCN and GAT through a learnable weighting factor  $\eta$ . The final node representations are obtained by a weighted combination of the outputs of both branches:

$$\mathbf{h}_{\text{final}} = \eta \cdot \mathbf{h}_{GCN} + (1 - \eta) \cdot \mathbf{h}_{GAT} \quad (6)$$

In the feature fusion layer, The weighting coefficient  $\eta$  is implemented as a learnable scalar parameter and is initialized to 0.5, such that both branches contribute equally at the beginning of training. This mechanism allows the model to automatically adjust the contributions of the two branches based on the specific requirements of the task and data. To train this parameter, we use gradient descent and update it through backpropagation.

**Low-dimensional encoding and prediction** After the feature fusion layers, we input gene pairs into two identical channels (Channel 1 and Channel 2) to process different types of information, which further helps in learning the low-dimensional representations of the nodes. The model fuses different information through these channels to enhance the representation capability of gene nodes. Then, the dot product is used as a scoring function to evaluate the similarity between the pair of genes based on the learned representations.

## 3 Results

### 3.1 Dataset and preprocessing

Seven single-cell RNA sequencing (scRNA-seq) datasets, which is collected by BEELINE et al. [26], are utilized to evaluate the performance of our method. These datasets include two human cell types and five mouse cell types. The human cell types are (i) human embryonic stem cells (hESC) [27] and (ii) human hepatocytes (hHEP) [28]. The mouse cell types are (i) mouse embryonic stem cells (mESC) [29], (ii) erythroid-lineage mouse hematopoietic stem cells (mHSC-E), (iii) granulocyte-monocyte-lineage mouse hematopoietic stem cells (mHSC-GM), (iv) lymphoid-lineage mouse hematopoietic stem cells (mHSC-L), and (v)

mouse dendritic cells (mDC) [30]. Except for one available ground-truth network based on loss-of-function/gain-of-function (LOF/GOF) experiments [31], which is limited to mESC, all of the above datasets are associated with three types of ground-truth regulatory networks: STRING database [32], non-specific ChIP-seq [33–35], and cell-type-specific ChIP-seq [31, 36, 37]. All datasets are publicly available from the Gene Expression Omnibus (GEO), with the following accession numbers: GSE75748 (hESC), GSE81252 (hHEP), GSE48968 (mDC), GSE98664 (mESC), and GSE81682 (mHSC) [38]. Detailed information about each dataset and its corresponding ground-truth networks is provided in Supplementary Table S5.

We adopted the preprocessing method proposed by Pratapa et al. [26] for each scRNA-seq dataset and focused on inferring gene interactions regulated by transcription factors (TFs). The specific strategies for splitting training, validation, and testing sets are provided in the Supplementary Material.

### 3.2 Evaluation metrics

To comprehensively evaluate the performance of the proposed hybrid GCN-GAT architecture for inferring GRNs, we employ AUROC (Area Under the Receiver Operating Characteristic Curve) and AUPRC (Area Under the Precision-Recall Curve) as evaluation metrics. AUROC quantifies a model’s overall capability in binary classification by measuring its ability to distinguish between positive and negative samples. The higher the AUROC score, the better the classification performance of the model. AUPRC is employed to evaluate models on imbalanced datasets, as it assesses performance by precision and recall. AUPRC provides more information than AUROC in measuring on minority categories.

**Selecting the best top-K strongest regulatory genes** In subsection The Top-K Sparsification Strategy, we proposed a Top-K strategy to identify the top  $K$  strongest regulatory genes. In order to select the best  $K$  value, we set the Top-K threshold for sparsification in the range from 20 to 60. The performances of our method on different dataset under different  $K$  values are shown in Supplementary Figure S2 and S3. According to these figures, although our method has varying performance at different  $k$  values, it achieves the best performance in most datasets when  $k$  is 45. Therefore, we recommend that setting the optimal  $K$  value to 45 by default. This selection is based on the following considerations: First, from a biological prior perspective, in known real-world GRNs, regulatory edges typically account for only 1% to 5% of all possible connections. We empirically verified that the edges generated with  $Top - K \in [20, 60]$  fall within this proportion, aligning with the expected sparsity of such networks. To prevent the selected top  $k$  gene pairs from overlapping with the positive samples in the subsequent test set, we first check if these gene pairs are present in the test set. If they are, we remove them to avoid prematurely knowing the answers from the test set.

### 3.3 The VAE-based model and hybrid GCN-GAT architecture enhance the performance of inferring GRN

To evaluate the contribution of each component, we conduct ablation studies on the VAE-based module with Top-K sparsification and the hybrid GCN-GAT architecture. Detailed AUROC and AUPRC results across all datasets are reported in Supplementary Tables S1–S4 and S6–S9. The results demonstrate that both components contribute to consistent performance improvements over the baseline GENELink model, and the full model achieves the best performance in the majority of experimental cases.

**Comparing with other methods** In this subsection, a comprehensive benchmarking study is conducted to evaluate the performance of our method against seven infer GRN methods. These comparison methods included two classical unsupervised approaches—Pearson correlation coefficient (PCC) [14] and mutual information (MI) [15]—along with five state-of-the-art graph neural network-based methods: DeepRIG [21], GNNLink [18], scMGATGRN [19], GCLink [20], and GENELink [22]. The overall AUROC score of these methods is presented in Fig. 2.

As shown in Fig. 2, our method achieves superior performance over all methods on both the Specific and LOF/GOF datasets. Notably, our method attains an AUROC above 0.92 across all cell types, whereas other methods reach this threshold only in a limited experimental case. On the Non-Specific dataset, our method outperforms other methods in three cell types—mHSC-E, mHSC-GM, and mHSC-L—while exhibiting comparable performance to the best method in hESC cell types. On the STRING dataset, HALink yields the highest predictive accuracy in mHSC-GM and mHSC-L cell types. Across all 44 experimental cases, our method exceeded an AUROC of 0.9 in over half of the experimental cases (24/44), in contrast to other methods, which surpassed this level in a few experimental cases. These results demonstrate that our method enables accurate and robust reconstruction of GRNs under diverse benchmarking scenarios.

The AUPRC score of all these methods are provided in Supplementary Figure S1, which further corroborate the consistent outperformance of our approach over existing methods.

## 4 Discussion and Conclusion

In this study, we propose a hybrid GCN-GAT architecture combined with a VAE-based model and a Top-K sparsification strategy for inferring GRN. On one hand, the VAE-based model reconstructs the adjacency matrix and, together with the sparsification strategy, effectively mitigates the sparsity and noise issues inherent in single-cell data, thereby enhancing the biological plausibility and stability of the input structure. On the other hand, the hybrid GCN-GAT architecture with learnable weighting factors enables dynamic fusion of global structural and local

neighborhood information, significantly improving the model’s ability to capture complex regulatory relationships and enhancing inference accuracy.

Experimental results on multiple single-cell RNA sequencing datasets demonstrate that the proposed method outperforms traditional approaches in terms of accuracy, robustness, and generalization, validating the synergistic effect of the VAE-based sparsification module and the hybrid GCN–GAT architecture. However, the method still has certain limitations, such as its dependence on the quality of training data and relatively high computational cost on large-scale datasets. Future work may focus on integrating multi-omics data and optimizing the model architecture to improve generalization and computational efficiency, as well as exploring its application to dynamic gene regulatory network modeling.

## 5 Supplementary Material

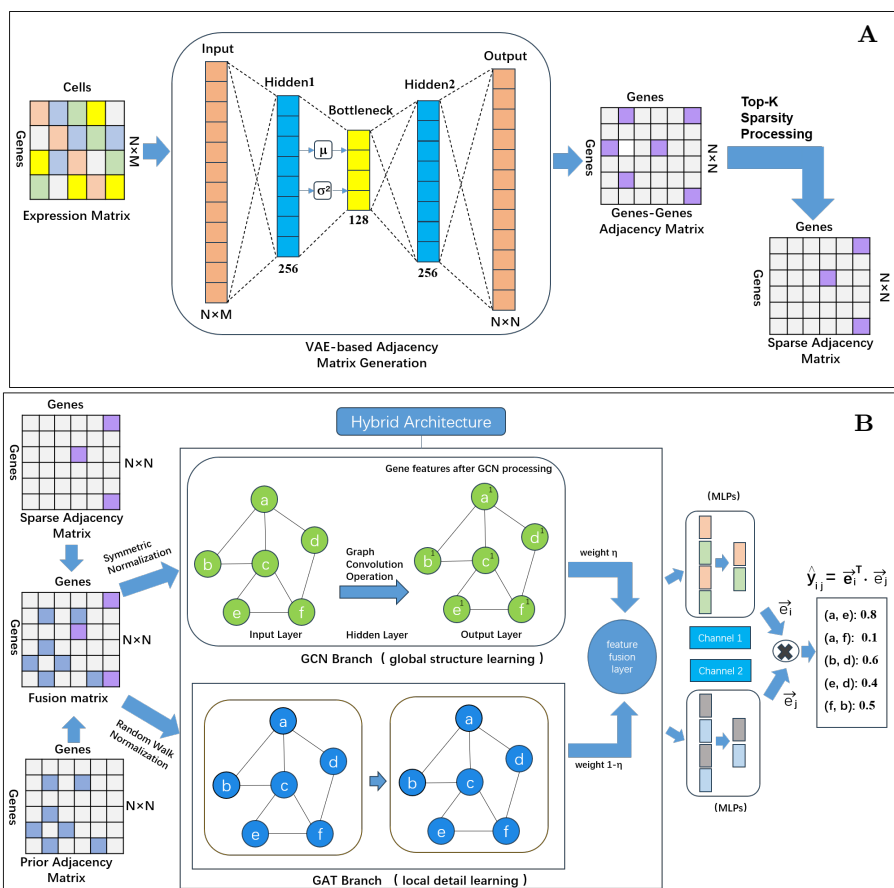
The code will be uploaded upon the acceptance of this manuscript. Below is the link to the supplementary material <https://doi.org/10.6084/m9.figshare.31732408>

## References

1. Alberto de la Fuente. What are gene regulatory networks?. In *Handbook of research on computational methodologies in gene regulatory networks*, pages 1–27. IGI Global, 2010.
2. Daniele Mercatelli, Laura Scalambra, Luca Triboli, Forest Ray, and Federico M. Giorgi. Gene regulatory network inference resources: A practical overview. *Biochimica et Biophysica Acta (BBA)-Gene Regulatory Mechanisms*, 1863(6):194430, 2020. Elsevier.
3. Eric H. Davidson and Douglas H. Erwin. Gene regulatory networks and the evolution of animal body plans. *Science*, 311(5762):796–800, 2006.
4. Douglas H. Erwin and Eric H. Davidson. The evolution of hierarchical gene regulatory networks. *Nature Reviews Genetics*, 10(2):141–148, 2009. Nature Publishing Group UK London.
5. Junlin Xu, Changcheng Lu, Shuting Jin, Yajie Meng, Xiangzheng Fu, Xiangxiang Zeng, Ruth Nussinov, and Feixiong Cheng. Deep learning-based cell-specific gene regulatory networks inferred from single-cell multiome data. *Nucleic Acids Research*, 53(5):gkaf138, 2025. Oxford University Press.
6. Frank Emmert-Streib, Matthias Dehmer, and Benjamin Haibe-Kains. Gene regulatory networks and their applications: understanding biological and medical problems in terms of networks. *Frontiers in Cell and Developmental Biology*, 2:38, 2014. Frontiers Media SA.
7. Benzhe Su, Weiwei Wang, Xiaohui Lin, Shenglan Liu, and Xin Huang. Identifying the potential miRNA biomarkers based on multi-view networks and reinforcement learning for diseases. *Briefings in Bioinformatics*, 25(1):bbad427, 2024.
8. Mengyuan Zhao, Wenying He, Jijun Tang, Quan Zou, and Fei Guo. A comprehensive overview and critical evaluation of gene regulatory network inference technologies. *Briefings in Bioinformatics*, 22(5):bbab009, 2021.
9. Guo Mao, Ruigeng Zeng, Jintao Peng, Ke Zuo, Zhengbin Pang, and Jie Liu. Reconstructing gene regulatory networks of biological function using differential equations of multilayer perceptrons. *BMC Bioinformatics*, 23(1):503, 2022.

10. Weixu Wang, Yichen Wang, Ruiqi Lyu, and Dominic Grün. Scalable identification of lineage-specific gene regulatory networks from metacells with NetID. *Genome Biology*, 25(1):275, 2024.
11. Yang Li, Anjun Ma, Yizhong Wang, Qi Guo, Cankun Wang, Hongjun Fu, Bingqiang Liu, and Qin Ma. Enhancer-driven gene regulatory networks inference from single-cell RNA-seq and ATAC-seq data. *Briefings in Bioinformatics*, 25(5):bbae369, 2024.
12. Eva Hedlund and Qiaolin Deng. Single-cell RNA sequencing: technical advancements and biological applications. *Molecular Aspects of Medicine*, 59:36–46, 2018.
13. Riet De Smet and Kathleen Marchal. Advantages and limitations of current network inference methods. *Nature Reviews Microbiology*, 8(10):717–729, 2010.
14. Alicia T. Specht and Jun Li. LEAP: constructing gene co-expression networks for single-cell RNA-sequencing data using pseudotime ordering. *Bioinformatics*, 33(5):764–766, 2017.
15. Thalia E. Chan, Michael P. H. Stumpf, and Ann C. Babbie. Gene regulatory network inference from single-cell data using multivariate information measures. *Cell Systems*, 5(3):251–267, 2017.
16. Jiayi Dong, Jiahao Li, and Fei Wang. Deep learning in gene regulatory network inference: a survey. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 2024.
17. Mengyuan Zhao, Wenying He, Jijun Tang, Quan Zou, and Fei Guo. A hybrid deep learning framework for gene regulatory network inference from single-cell transcriptomic data. *Briefings in Bioinformatics*, 23(2):bbab568, 2022.
18. Guo Mao, Zhengbin Pang, Ke Zuo, Qinglin Wang, Xiangdong Pei, Xinhai Chen, and Jie Liu. Predicting gene regulatory links from single-cell RNA-seq data using graph neural networks. *Briefings in Bioinformatics*, 24(6):bbad414, 2023.
19. Lin Yuan, Ling Zhao, Yufeng Jiang, Zhen Shen, Qinhu Zhang, Ming Zhang, Chun-Hou Zheng, and De-Shuang Huang. scMGATGRN: a multiview graph attention network-based method for inferring gene regulatory networks from single-cell transcriptomic data. *Briefings in Bioinformatics*, 25(6):bbae526, 2024.
20. Weiming Yu, Zerun Lin, Miaofang Lan, and Le Ou-Yang. GCLink: a graph contrastive link prediction framework for gene regulatory network inference. *Bioinformatics*, 41(3):btaf074, 2025.
21. Jiacheng Wang, Yaojia Chen, and Quan Zou. Inferring gene regulatory network from single-cell transcriptomes with graph autoencoder model. *PLOS Genetics*, 19(9):e1010942, 2023.
22. Guangyi Chen and Zhi-Ping Liu. Graph attention network for link prediction of gene regulations from single-cell RNA-sequencing data. *Bioinformatics*, 38(19):4522–4529, 2022.
23. Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. *Advances in Neural Information Processing Systems*, 31, 2018.
24. Steffen Albrecht, Tommaso Andreani, Miguel A Andrade-Navarro, and Jean Fred Fontaine. Single-cell specific and interpretable machine learning models for sparse scChIP-seq data imputation. *Plos One*, 17(7):e0270043, 2022.
25. HaiYun Wang, JianPing Zhao, ChunHou Zheng, and YanSen Su. scDSSC: deep sparse subspace clustering for scRNA-seq data. *PLOS Computational Biology*, 18(12):e1010772, 2022.
26. Aditya Pratapa, Amogh P. Jalihal, Jeffrey N. Law, Aditya Bharadwaj, and T. M. Murali. Benchmarking algorithms for gene regulatory network inference from single-cell transcriptomic data. *Nature Methods*, 17(2):147–154, 2020.

27. Kishan Kc, Rui Li, Feng Cui, Qi Yu, and Anne R. Haake. GNE: a deep learning framework for gene network inference by aggregating biological information. *BMC Systems Biology*, 13:1–14, 2019.
28. J. Gray Camp, Keisuke Sekine, Tobias Gerber, Henry Loeffler-Wirth, Hans Binder, Malgorzata Gac, Sabina Kanton, Jorge Kageyama, Georg Damm, Daniel Seehofer, et al. Multilineage communication regulates human liver bud development from pluripotency. *Nature*, 546(7659):533–538, 2017.
29. Tetsutaro Hayashi, Haruka Ozaki, Yohei Sasagawa, Mana Umeda, Hiroki Danno, and Itoshi Nikaido. Single-cell full-length total RNA sequencing uncovers dynamics of recursive splicing and enhancer RNAs. *Nature Communications*, 9(1):619, 2018.
30. Alex K. Shalek, Rahul Satija, Joe Shuga, John J. Trombetta, Dave Gennert, Diana Lu, Peilin Chen, Rona S. Gertner, Jellert T. Gaublomme, Nir Yosef, et al. Single-cell RNA-seq reveals dynamic paracrine control of cellular variation. *Nature*, 510(7505):363–369, 2014.
31. Huilei Xu, Caroline Baroukh, Ruth Dannenfels, Edward Y. Chen, Christopher M. Tan, Yan Kou, Yujin E. Kim, Ihor R. Lemischka, and Avi Ma’ayan. ESCAPE: database for integrating high-content published data collected from human and mouse embryonic stem cells. *Database*, 2013:bat045, 2013.
32. Damian Szklarczyk, Annika L. Gable, David Lyon, Alexander Junge, Stefan Wyder, Jaime Huerta-Cepas, Milan Simonovic, Nadezhda T. Doncheva, John H. Morris, Peer Bork, et al. STRING v11: protein–protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Research*, 47(D1):D607–D613, 2019.
33. Luz Garcia-Alonso, Christian H. Holland, Mahmoud M. Ibrahim, Denes Turei, and Julio Saez-Rodriguez. Benchmark and integration of resources for the estimation of human transcription factor activities. *Genome Research*, 29(8):1363–1375, 2019.
34. Zhi-Ping Liu, Canglin Wu, Hongyu Miao, and Hulin Wu. RegNetwork: an integrated database of transcriptional and post-transcriptional regulatory networks in human and mouse. *Database*, 2015:bav095, 2015.
35. Heonjong Han, Jae-Won Cho, Sangyoung Lee, Ayoung Yun, Hyojin Kim, Dasom Bae, Sunmo Yang, Chan Yeong Kim, Muyoung Lee, Eunbeen Kim, et al. TRRUST v2: an expanded reference database of human and mouse transcriptional regulatory interactions. *Nucleic Acids Research*, 46(D1):D380–D386, 2018.
36. Jill E. Moore, Michael J. Purcaro, Henry E. Pratt, Charles B. Epstein, Noam Shores, Jessika Adrian, Trupti Kawli, Carrie A. Davis, Alexander Dobin, et al. Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature*, 583(7818):699–710, 2020.
37. Shinya Oki, Tazro Ohta, Go Shioi, Hideki Hatanaka, Osamu Ogasawara, Yoshihiro Okuda, Hideya Kawaji, Ryo Nakaki, Jun Sese, and Chikara Meno. ChIP-Atlas: a data-mining suite powered by full integration of public ChIP-seq data. *EMBO Reports*, 19(12):e46255, 2018.
38. Sonia Nestorowa, Fiona K. Hamey, Blanca Pijuan Sala, Evangelia Diamanti, Mairi Shepherd, Elisa Laurenti, Nicola K. Wilson, David G. Kent, and Berthold Göttgens. A single-cell resolution map of mouse hematopoietic stem and progenitor cell differentiation. *Blood, The Journal of the American Society of Hematology*, 128(8):e20–e31, 2016.



**Fig. 1.** The overview of the framework of HALink. (A) **Adjacency matrix generation and sparsification.** A VAE-based model is applied to generate a latent adjacency matrix representing regulatory relationships, which fills in the missing parts of prior information. Meanwhile, Top-K sparsification retains the strongest gene-gene interactions, further reducing noise and unnecessary computations. (B) **Hybrid GCN-GAT architecture construction.** We propose a hybrid GCN-GAT architecture that combines the strengths of graph convolutional networks (GCN) and graph attention networks (GAT). GCN excels at capturing the global structure of the graph, while GAT focuses on capturing fine-grained relationships between nodes. By combining both, we can leverage the complementary advantages of these networks to enhance the learning of potential regulatory relationships between genes.

		TFs+500								TFs+1000							
		HALink	GENELink	DeepRIG	GNNLink	sdMGATGRN	GCLink	MI	PCC	HALink	GENELink	DeepRIG	GNNLink	sdMGATGRN	GCLink	MI	PCC
Specific	hESC	0.96	0.81	0.75	0.83	0.87	0.84	0.5	0.52	0.96	0.82	0.74	0.85	0.87	0.84	0.51	0.51
	hHEP	0.94	0.82	0.77	0.84	0.89	0.88	0.48	0.48	0.95	0.86	0.76	0.82	0.89	0.87	0.45	0.48
	mDC	0.97	0.72	0.48	0.65	0.76	0.78	0.52	0.52	0.97	0.73	0.54	0.7	0.77	0.78	0.58	0.51
	mESC	0.92	0.89	0.75	0.83	0.91	0.9	0.5	0.5	0.93	0.9	0.74	0.82	0.92	0.9	0.52	0.49
	mHSC-E	0.93	0.87	0.71	0.82	0.92	0.89	0.51	0.48	0.95	0.9	0.74	0.79	0.9	0.9	0.41	0.47
	mHSC-GM	0.95	0.89	0.68	0.89	0.89	0.9	0.34	0.51	0.96	0.9	0.76	0.85	0.93	0.91	0.33	0.5
	mHSC-L	0.94	0.84	0.73	0.83	0.86	0.85	0.45	0.53	0.94	0.84	0.78	0.84	0.86	0.86	0.54	0.55
Non-Specific	hESC	0.74	0.58	0.75	0.74	0.66	0.72	0.52	0.56	0.77	0.64	0.78	0.76	0.69	0.73	0.52	0.56
	hHEP	0.72	0.72	0.86	0.75	0.71	0.72	0.59	0.57	0.74	0.72	0.84	0.82	0.71	0.74	0.59	0.59
	mDC	0.77	0.7	0.75	0.8	0.8	0.8	0.6	0.62	0.81	0.73	0.71	0.87	0.79	0.8	0.57	0.56
	mESC	0.6	0.58	0.71	0.77	0.79	0.81	0.54	0.57	0.63	0.59	0.71	0.87	0.83	0.82	0.53	0.56
	mHSC-E	0.86	0.74	0.72	0.8	0.76	0.8	0.58	0.61	0.89	0.79	0.74	0.81	0.73	0.84	0.66	0.6
	mHSC-GM	0.91	0.78	0.75	0.75	0.75	0.85	0.66	0.66	0.95	0.84	0.7	0.78	0.76	0.86	0.73	0.61
	mHSC-L	0.95	0.76	0.61	0.61	0.64	0.76	0.6	0.6	0.96	0.75	0.63	0.64	0.58	0.78	0.6	0.63
Lofgof	mESC	0.96	0.81	0.73	0.84	0.72	0.85	0.6	0.53	0.95	0.83	0.77	0.83	0.71	0.85	0.6	0.52
STRING	hESC	0.74	0.82	0.89	0.92	0.88	0.85	0.57	0.66	0.74	0.83	0.88	0.92	0.89	0.87	0.62	0.67
	hHEP	0.69	0.81	0.91	0.92	0.91	0.84	0.56	0.64	0.7	0.83	0.89	0.91	0.92	0.86	0.6	0.65
	mDC	0.75	0.83	0.91	0.92	0.89	0.88	0.62	0.65	0.81	0.86	0.77	0.91	0.9	0.88	0.58	0.63
	mESC	0.6	0.83	0.91	0.92	0.91	0.86	0.54	0.59	0.65	0.84	0.9	0.92	0.91	0.88	0.59	0.61
	mHSC-E	0.73	0.84	0.82	0.92	0.9	0.87	0.57	0.72	0.88	0.86	0.69	0.88	0.89	0.87	0.65	0.71
	mHSC-GM	0.91	0.83	0.7	0.9	0.89	0.87	0.64	0.75	0.94	0.87	0.81	0.9	0.87	0.88	0.67	0.75
	mHSC-L	0.95	0.87	0.52	0.92	0.72	0.86	0.74	0.71	0.96	0.81	0.64	0.79	0.83	0.85	0.75	0.64

**Fig. 2.** The AUROC scores of HALink and seven other methods were compared across four ground truth networks: cell-type-specific ChIP-seq, non-specific ChIP-seq, LOF/GOF, and STRING.