

Exploring Graph-Based Neural Networks for Modeling Long-Range Epigenetic Gene Regulation

Brown University Computational Biology Undergraduate Senior Honors
Thesis

Author: Xavier Loinaz

Advisor: Ritambhara Singh

2021

Abstract

Long-range spatial interactions among genomic regions are critical for regulating gene expression, and their disruption has been associated with a host of diseases. However, when modeling the effects of regulatory factors, most deep learning models either neglect long-range interactions or fail to capture the inherent 3D structure of the underlying genomic organization. To address these limitations, in this thesis we present two graph-based neural network architectures: GC-MERGE, a **Graph Convolutional Model for Epigenetic Regulation of Gene Expression**, as well as XL-MERGE, **Xavier Loinaz’s Model for Epigenetic Regulation of Gene Expression**. Both integrate measurements of both the spatial genomic organization and local regulatory factors, specifically histone modifications, for predicting gene expression. This formulation enables these models to incorporate crucial information about long-range interactions via a natural encoding of spatial interactions into a graph representation. We apply GC-MERGE and XL-MERGE to datasets for the GM12878 (lymphoblastoid), K562 (myelogenous leukemia), and HUVEC (human umbilical vein endothelial) cell lines and demonstrate predictive performance comparable to state-of-the-art methods for GC-MERGE, and superior performance to state-of-the-art methods for XL-MERGE. In addition, we give evidence that GC-MERGE is interpretable in terms of the observed biological regulatory factors, highlighting both the histone modifications and the interacting genomic regions contributing to a gene’s predicted expression. We provide model explanations for several exemplar genes and validate them with evidence from the literature. These models not only present a novel setups for predicting gene expression by integrating multimodal datasets in a graph convolutional framework, but also enable interpretation of the biological mechanisms driving the model’s predictions. The code for GC-MERGE is available at: <https://github.com/rsinghlab/GC-MERGE>.

Acknowledgements

I would like to thank my advisor for this project, Prof. Ritambhara Singh, for her mentorship, enthusiasm, and understanding throughout my work with her. Thank you for giving me such a positive undergraduate research experience, and helping inspire me to further continue research in the field of computational biology in my career.

I would also like to thank members of the Singh Lab who I worked with who helped me along the way. Thank you to Jeremy Bigness for onboarding me onto his project, his accessibility in answering all of my questions, and his understanding when I did not always get things right at first. Thank you to Shalin Patel for his major contributions in accelerating our project pipelines and always being so quick to answer me on Slack about questions I had.

Thank you to my research mentors I have had throughout my past as well. I want to thank Prof. David Borton, Prof. Vicki Colvin, Xiaoting Guo, Jake Villanova, Dr. Benny Coyac, Prof. Jill Helms, and Dr. Eric Sabelman for helping nurture my development as a researcher, having patience with me, and giving me enthusiasm to pursue research further.

I am also in gratitude to my friends, who allowed me to be myself when I would get stressed about school and/or research, and my parents, who have financially supported me throughout college and given me access to these opportunities.

.....
Author's signature

Contents

Introduction	6
1.1 Background	6
1.1.1 Longe-Range Gene Regulation	6
1.1.2 Related Past Work	6
1.1.3 Challenges	7
1.2 Novel Graph-Based Methods to Model Long-Range Epigenetic Gene Regulation	7
1.2.1 GC-MERGE	7
1.2.2 XL-MERGE	8
1.2.3 Model Interpretation	8
GC-MERGE	10
2.1 Methods	10
2.1.1 Graph Convolutional Networks (GCNs)	10
2.1.2 Interpretation of GC-MERGE	11
2.2 Experimental Setup	12
2.2.1 Overview of Datasets	12
2.2.2 Graph Construction and Data Integration	13
2.2.3 Baseline Models	13
2.2.4 Evaluation Metrics	14
2.3 Results	15
2.3.1 Gene Expression Prediction Results	15
2.3.2 Interpretation Results	15
2.4 Discussion	17
XL-MERGE	19
3.1 Methods	19
3.1.1 Model Architecture	19
3.2 Experimental Setup	20
3.2.1 Overview of Datasets and Data Integration	20
3.2.2 Baseline Models	20
3.2.3 Model Ablation Analyses	21
3.2.4 Evaluation Metrics	22
3.3 Results	22
3.3.1 Gene Expression Prediction Results	22
3.3.2 Model Ablation Results	22
3.4 Discussion	23
3.4.1 Adding Interpretability for XL-MERGE	23
3.4.2 Regression Task Results for XL-MERGE	23
3.4.3 Normalizing Genic Region Input Data	23

3.4.4 Different Training Rates for Different Parts of the Model	24
Bibliography	25

Appendix

GC-MERGE Model Details	30
GC-MERGE Additional Dataset Details and Results	32
XL-MERGE Model Details and Results	40

List of Figures

- 2.1 **Model comparison and evaluation.** GC-MERGE gives state-of-the-art performance for the regression task on all three studied cell lines. (a) The Pearson correlation coefficients (PCC) obtained by running GC-MERGE on each cell line are displayed. (b) GC-MERGE outperforms all the baseline models for each of the three cell lines. Scores are calculated as the average of 10 runs and standard deviations are denoted by error bars. . . . 14
- 2.2 **Model explanations for exemplar genes. Top:** For (a) SIDT1, designated as node 60561 (yellow circle), the subgraph of neighbor nodes is displayed. The size of each neighbor node correlates with its predictive importance as determined by GNNExplainer. Nodes in red denote regions corresponding to known enhancer regions regulating SIDT1 [13] (note that multiple interacting fragments can be assigned to each node, see Supplementary Tables B.1 and B.2). All other nodes are displayed in gray. Nodes with importance scores corresponding to outliers have been removed for clarity. **Bottom:** The scaled feature importance scores for each of the five core histone marks used in this study are shown in the bar graph. Results also presented for (b) AKR1B1, (c) LAPTM5, and (d) TOP2B. 16

1. Introduction

1.1 Background

1.1.1 Long-Range Gene Regulation

Gene regulation determines the fate of every cell, and its disruption leads to diverse diseases ranging from cancer to neurodegeneration [17, 26]. Although specialized cell types – from neurons to cardiac cells – exhibit different gene expression patterns, the information encoded by the linear DNA sequence remains virtually the same in all non-reproductive cells of the body. Therefore, the observed differences in cell type must be encoded by elements extrinsic to sequence, commonly referred to as epigenetic factors. Epigenetic factors found in the local neighborhood of a gene typically include histone marks (also known as histone modifications). These marks are naturally occurring chemical additions to histone proteins that control how tightly the DNA strands are wound around the proteins and the recruitment or occlusion of transcription factors. However, the focus of attention in genomics has shifted increasingly to the study of long-range epigenetic regulatory interactions that result from the three-dimensional organization of the genome [25]. For example, one early study demonstrated that chromosomal rearrangements, some located as far as 125 kilo-basepairs (kbp) away, disrupted the region downstream of the PAX6 transcription unit causing Aniridia (absence of the iris) and related eye anomalies [16]. Thus, chromosomal rearrangement can not only directly affect the expression of proximal genes but can also indirectly affect a gene located far away by perturbing its regulatory (e.g., enhancer-promoter) interactions. This observation indicates that while local regulation of genes is informative, studying long-range gene regulation is critical to understanding cell development and disease. However, experimentally testing for all possible combinations of long-range and short-range regulatory factors for $\sim 20,000$ genes is infeasible given the vast size of the search space. Therefore, computational and data-driven approaches are necessary to efficiently search this space and reduce the number of testable hypotheses due to the sheer scope of the problem.

1.1.2 Related Past Work

Recently, deep learning frameworks have been applied to predict gene expression from histone modifications, and their empirical performance has often exceeded the previous machine learning methods [4, 6, 14]. Among their many advantages, deep neural networks perform automatic feature extraction by efficiently exploring feature space and then finding nonlinear transformations of the weighted averages of those features. This formulation is especially relevant to complex biological systems since they are inherently nonlinear. For instance, Singh *et al.* [30] introduced DeepChrome, which used a convolutional neural network (CNN) to aggregate five types of histone mark ChIP-seq signals in a 10,000 bp region around the transcription start site (TSS) of each gene. Using a

similar setup, they next introduced attention layers to their model [29], yielding a comparable performance but with the added ability to visualize feature importance within the local neighborhood of a gene. These methods framed the gene expression problem as a binary classification task in which the gene was either active or inactive. Agarwal *et al.* [1] introduced Xpresso, a CNN framework that operated on the promoter sequences of each gene and 8 other annotated features associated with mRNA decay to predict steady-state mRNA levels. This model focused primarily on the regression task, such that each prediction corresponded to the logarithm of a gene’s expression. While all the studies listed above accounted for combinatorial interactions among features at the local level, they did not incorporate long-range regulatory interactions known to play a critical role in differentiation and disease [17, 26].

1.1.3 Challenges

Modeling these long-range interactions is a challenging task due to two significant reasons. First, we cannot confidently pick an input size for the genomic regions as regulatory elements can control gene expression from various distances. Second, inputting a large region will introduce sparsity and noise into the data, making the learning task difficult. A potential solution to this problem is to incorporate information from long-range interaction networks captured from experiments like Hi-ChIP [19] and Hi-C [34]. These assays use high-throughput sequencing to measure 3D genomic structure, in which each read pair corresponds to an observed 3D contact between two genomic loci. While Hi-ChIP focuses only on spatial interactions mediated by a specific protein, Hi-C captures the global interactions of all genomic regions. Recently, Zeng *et al.* [39] combined a CNN, encoding promoter sequences, with a fully connected network using Hi-ChIP datasets to predict gene expression values. The authors then evaluated the relative contributions of the promoter sequence and promoter-enhancer submodules to the model’s overall performance. While this method incorporated long-range interaction information, its use of HiChIP experiments narrowed this information to spatial interactions mediated by H3K27ac and YY1. Furthermore, CNN models only capture the local topological patterns instead of modeling the underlying spatial structure of the data. Thus, the interpretation of their model was limited to local sequence features.

1.2 Novel Graph-Based Methods to Model Long-Range Epigenetic Gene Regulation

1.2.1 GC-MERGE

To address the previously mentioned issues, we developed a **Graph Convolutional Model of Epigenetic Regulation of Gene Expression** (GC-MERGE), a graph-based deep learning framework that integrates 3D genomic data with existing histone mark signals to predict gene expression. Figure ?? provides a schematic of our overall approach. Unlike previous methods, our model incorporates genome-wide interaction information by using the Hi-C data. To accomplish this, we use a graph convolutional network (GCN) to capture the underlying spatial structure. GCNs are particularly well-suited to representing spatial relationships, as a Hi-C map can be represented as an adjacency matrix of an undirected graph $G \in \{V, E\}$. Here, V nodes represent the genomic regions and E edges repre-

sent their interactions. Our prediction task formulation captures the local and spatial relationships between the histone marks and gene expression. While some methods use many other types of features, such as promoter sequences [1, 39], we focus our efforts solely on histone modifications and extract their relationship to the genes. Even with this simplified set of features, we show that our model’s performance exceeds that of certain baseline methods for the gene expression prediction task.

Limitations of GC-MERGE

Although GC-MERGE introduces an effective way to integrate potential enhancer and repressor epigenetic interactions to better predict gene expression for certain genic regions, there were a couple significant problems with its task formulation. First of all, in order to have standardized identification covering all regions of the genome, the genomic regions in the dataset for GC-MERGE are regularly divided at 10kb intervals. Thus, placement of an actual gene sequence within this interval could be variable — the gene sequence could be in the middle of this 10kb interval or it could be at the boundary. If the genic region is at the boundary and represented by the 10kb interval, significant regulatory information could be lost coming from epigenetic marks from one side of the genic region not covered in the 10kb interval. Secondly, GC-MERGE does not optimally extract epigenetic information from long-range interactions through our graph formulation of a genic region and its long-range interacting regions. Neighboring nodes’ information is passed through the network as an average of histone mark levels over their entire 10kb regions, which may not give us the most information in determining various 10kb regions that may contribute as enhancers/repressors for regulating gene expression.

1.2.2 XL-MERGE

We develop **X**avier **L**oinaz’s **M**odel of **E**pigenetic **R**egulation of **G**ene **E**xpression (XL-MERGE) as a way of addressing the two mentioned limitations of GC-MERGE. To account for the issue of gene sequences’ variable location within genomic intervals used to predict gene expression, we introduce a dataset from Singh *et al.* [30] that counts epigenetic marks for certain genes within 5kb of their respective transcription start sites. By ensuring that the transcription start site is centered within the interval we use to predict gene expression this reduces variability. Additionally, to better extract neighboring nodes’ epigenetic information to predict potential long-range enhancer and repressor interactions, we introduce a convolutional operation with maxpooling in order to extract important histone mark patterns relative to one another and introduce some positional invariance for detecting these patterns. This alone leads to significantly improved gene expression prediction based off neighboring nodes as will be later discussed. The performance metrics of this model are also shown to be better than GC-MERGE by a statistically significant margin.

1.2.3 Model Interpretation

Another significant contribution of this work is to enable biologists to determine at the genic level which regulatory interactions – local or distal – most affect the gene’s expression and which histone marks modulate these interactions. By making the model’s predictive drivers more transparent, this information can suggest promising hypotheses

and guide new research directions. To that effect, we perform an interpretation of the GC-MERGE’s predictions that quantifies the relative importance of the underlying biological regulatory factors driving each gene’s output. We integrate the GNNExplainer method [37] within our modeling framework to highlight not only the important node features (histone modifications) but also the important edges (long-range interactions) that contribute to determining a particular gene’s predicted expression. In this thesis, we apply our method to the three cell lines from Rao *et al.* [22] – GM12878 (lymphoblastoid), K562 (myelogenous leukemia), and HUVEC (human umbilical vein endothelial cells). While solving the gene expression prediction as a regression problem is more valuable for the community, our interpretation framework required us to formulate it as a classification task. Therefore, we perform both regression and classification tasks and demonstrate state-of-the-art performance. Furthermore, we show that our framework allows biologists to tease apart the cumulative effects of different regulatory mechanisms at the genic level. Table ?? places the proposed framework among state-of-the-art deep learning models and lists each model’s properties.

2. GC-MERGE

This chapter outlines the methods, experimental setup, and results for GC-MERGE, and it also contains a brief discussion.

2.1 Methods

2.1.1 Graph Convolutional Networks (GCNs)

Graph convolutional networks (GCNs) are a generalization of convolutional neural networks (CNNs) to graph-based relational data that is not natively structured in Euclidean space [18]. Due to the expressive power of graphs, GCNs have been applied across a wide variety of domains, including recommender systems [12], and social networks [21]. The prevalence of graph-based datasets in biology has made these models a popular choice for tasks like modeling protein-protein interactions [36], stem cell differentiation [2], and chemical reactivity for drug discovery [32].

We use the GraphSAGE formulation [11] as our GCN for its relative simplicity and its capacity to learn generalizable, inductive representations not limited to a specific graph. The input to the model is represented as a graph $G \in \{V, E\}$, with nodes V and edges E , and a corresponding adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ [18], where N is the number of nodes. For each node v , there is also an associated feature vector \mathbf{x}_v . The goal of the network is to learn a state embedding $\mathbf{h}_v^K \in \mathbb{R}^d$ for v , which is obtained by aggregating information over v 's neighborhood K times. Here, d is the dimension of the embedding vector. This state embedding is then fed through a fully-connected network to produce an output \hat{y}_v , which can then be applied to downstream classification or regression tasks.

Within this framework, the first step is to initialize each node with its input features. In our case, the feature vector $\mathbf{x}_v \in \mathbb{R}^m$ is obtained from the ChIP-seq signals corresponding to the five ($m = 5$) core histone marks (H3K4me1, H3K4me3, H3K9me3, H3K36me3, and H3K27me3) in our dataset:

$$\mathbf{h}_v^0 = \mathbf{x}_v \quad (2.1)$$

Next, to transition from the $(k - 1)^{th}$ layer to the k^{th} hidden layer in the network for node v , we apply an aggregation function to the neighborhood of each node. This aggregation function is analogous to a convolution operation over regularly structured Euclidean data such as images. A standard convolution function operates over a grid and represents a pixel as a weighted aggregation of its neighboring pixels. Similarly, a graph convolution performs this operation over the neighbors of a node in a graph. In our case, the aggregation function calculates the mean of the neighboring node features:

$$\mathbf{h}_{\mathcal{N}(v)}^k = \sum_{u \in \mathcal{N}(v)} \frac{\mathbf{h}_u^{k-1}}{|\mathcal{N}(v)|} \quad (2.2)$$

Here, $\mathcal{N}(v)$ represents the adjacency set of node v . We update the node’s embedding by concatenating the aggregation with the previous layer’s representation to retain information from the original embedding. Next, just as done in a standard convolution operation, we take the matrix product of this concatenated representation with a learnable weight matrix to complete the weighted aggregation step. Finally, we apply a non-linear activation function, such as ReLU, to capture the higher-order non-linear interactions among the features:

$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \left[\mathbf{h}_{\mathcal{N}(v)}^k \parallel \mathbf{h}_v^{k-1} \right] \right), \forall k \in \{1, \dots, K\} \quad (2.3)$$

Here, \parallel represents concatenation, σ is a non-linear activation function, and \mathbf{W}_k is a learnable weight parameter. After this step, each node is assigned a new embedding. After K iterations, the node embedding encodes information from the neighbors that are K -hops away from that node:

$$\mathbf{z}_v = \mathbf{h}_v^K \quad (2.4)$$

Here, \mathbf{z}_v is the final node embedding after K iterations. For regression, we feed \mathbf{z}_v into a fully connected network and output a prediction $\hat{y}_v \in \mathbb{R}$, representing a real-valued expression level. We use the mean squared error (MSE) as the loss function. The model architecture is summarized in Supplementary Figure A.1.

2.1.2 Interpretation of GC-MERGE

Although a model’s architecture is integral to its performance, just as important is understanding how the model arrives at its predictions. Neural networks in particular have sometimes been criticized for being “black box” models, such that no insight is provided into how the model operates. Most graph-based interpretability approaches either approximate models with simpler models whose decisions can be used for explanations [23] or use an attention mechanism to identify relevant features in the input that guide a particular prediction [35]. In general, these methods, along with gradient-based approaches [28, 33] or DeepLift [27], focus on the explanation of important node features and do not incorporate the structural information of the graph. However, a recent method called *Graph Neural Net Explainer* (or GNNExplainer) [37], given a trained GCN, can identify a small subgraph as well as a small subset of features that are crucial for a particular prediction. The authors demonstrate its interpretation capabilities on simulated and real-world graphs.

In order to apply this method, the problem must be constructed as a classification task. Therefore, we feed the learned embedding \mathbf{z}_v in Equation 2.4 into a fully connected network and output a prediction \hat{y}_v for each target node using a *Softmax* layer to compute probabilities for each class c . Here, class $c \in \{0, 1\}$ corresponds to whether the gene is either off/inactive ($c = 0$) or on/active ($c = 1$). We use the true binarized gene expression value $y_v \in \{0, 1\}$ by thresholding the expression level relative to the median as the target predictions (as done previously [4, 30, 29, 39]), using a negative log-likelihood (NLL) loss to train the model.

Next, we integrate the GNNExplainer module into our classifier framework. GNNExplainer maximizes the mutual information between the probability distribution of the model’s class predictions over all nodes and the probability distribution of the class predictions for a particular node conditioned on some fractional masked subgraph of neigh-

boring nodes and features. Subject to regularization constraints, it jointly optimizes the fractional node and feature masks, determining the extent to which each element informs the prediction for a particular node.

Specifically, given a node v , the goal is to learn a subgraph $G_s \subseteq G$ and a feature mask $X_s = \{x_j \mid v_j \in G_s\}$ that contribute the most to driving the full model’s prediction of \hat{y}_v . To achieve this objective, the algorithm learns a mask that maximizes the mutual information (MI) between the original model and the masked model. Mathematically, this objective function is as follows:

$$\max_{G_s} MI(Y, (G_s, X_s)) = H(Y) - H(Y \mid G_s, X_s) \quad (2.5)$$

where H is the entropy of a distribution. Since this is computationally intractable with an exponential number of graph masks, GNNExplainer optimizes the following quantity using gradient descent:

$$\min_{M, N} - \sum_{c=1}^C \mathbb{1}_{\{y=c\}} \log(P_\phi(Y = y \mid G = A_c \odot \sigma(M), X = X_c \odot \sigma(N))) \quad (2.6)$$

where c represents the class, A_c represents the adjacency matrix of the computation graph, M represents the subgraph mask, and N represents the feature mask. The importance scores of the nodes and features are obtained by applying the sigmoid function to the subgraph and feature masks, respectively. Finally, the element-wise entropies of the masks are calculated and inserted as regularization terms into the loss function. Therefore, in the context of our model, GNNExplainer learns which genomic regions (via the subgraph mask) and which features (via the feature mask) are most important in driving the model’s predictions.

2.2 Experimental Setup

2.2.1 Overview of Datasets

GC-MERGE requires the following information: (1) Interactions between the genomic regions (Hi-C contact maps); (2) Histone mark signals representing the regulatory signals (ChIP-seq measurements); (3) Expression levels for each gene (RNA-seq measurements). For each gene in a particular region, the first two datasets are the inputs into our proposed model, whereas gene expression is the predicted target. We formulate the problem as both regression and classification tasks. We take the base-10 logarithm of the gene expression values for the regression task, adding a pseudo-count of 1. For the classification task, we binarize the gene expression values as either 0 (off) or 1 (on) using the median as the threshold, consistent with previous studies [4, 30, 29, 39]. Constructing a binary classifier enables us to integrate the GNNExplainer interpretive mechanism with our framework.

We focused on three human cell lines from Rao *et al.* [22]: (1) GM12878, a lymphoblastoid cell line with a normal karyotype, (2) K562, a myelogenous leukemia cell line, and (3) HUVEC, a human umbilical vein endothelial cell line. For each of these cell lines, we accessed RNA-seq expression and ChIP-Seq signal datasets for five uniformly profiled histone marks from the REMC repository [24]. These histone marks include (1) H3K4me1, associated with enhancer regions; (2) H3K4me3, associated with promoter regions; (3) H3K9me3, associated with heterochromatin; (4) H3K36me3, associated with

actively transcribed regions; and (5) H3K27me3, associated with polycomb repression. We chose these marks because of the wide availability of the relevant data as well as for ease of comparison with previous studies [30, 29, 39].

2.2.2 Graph Construction and Data Integration

Our main innovation is formulating the graph-based prediction task to integrate two very different data modalities (histone mark signals and Hi-C interaction frequencies). We represented each genomic region with a node and connected edges between it and the nodes corresponding to its neighbors (bins with non-zero entries in the adjacency matrix) to construct the graph. Due to the large size of the Hi-C graph, we subsampled neighbors to form a subgraph for each node we fed into the model. While there are methods to perform subsampling on large graphs using a random node selection approach (e.g., [38]), we used a simple strategy of selecting the top j neighbors with the highest Hi-C interaction frequency values. We empirically selected the value $j = 10$ for the number of neighbors. A smaller number of neighbors (i.e., $j = 5$) resulted in decreased performance while selecting more neighbors proved prohibitive due to memory constraints.

To integrate the Hi-C datasets (preprocessing details in Supplementary Section B) with the RNA-seq and ChIP-seq datasets, we obtained the average ChIP-seq signal for each of the five core histone marks over the chromosomal region corresponding to each node. In this way, a feature vector of length five was associated with each node. For the RNA-seq data, we took each gene’s transcriptional start site (TSS) and assigned it to the node corresponding to the chromosomal region in which the TSS is located. We applied a mask during the training phase so that the model made predictions only on nodes corresponding to chromosomal regions with genes. If multiple genes were assigned to the same node, we took the median of the expression levels. We assigned 70% of the nodes to the training set, 15% to the validation set, and 15% to the testing set. We provide the details of the hyperparameter tuning in Supplementary Section A.2.

2.2.3 Baseline Models

We compared GC-MERGE with the following deep learning baselines for gene expression prediction formulated as both regression and classification tasks:

- **Multi-layer perceptron (MLP)**: A simple MLP comprised of three fully-connected layers. In this framework, the model predictions for each node do not incorporate feature information from the node’s neighbors.
- **Shuffled neighbor model**: GC-MERGE applied to shuffled Hi-C matrices, such that the neighbors of each node are randomized. The shuffled neighbor and MLP baselines can be viewed as proxies for the importance of including information from long-range regulatory interactions for similarly processed inputs.
- **Convolutional neural network (CNN)**: A convolutional neural network based on DeepChrome [30]. This model takes 10 kb regions corresponding to the genomic regions demarcated in the Hi-C data and subdivides each region into 100 bins. Each bin is associated with five channels, which correspond to the ChIP-seq signals of the same five core histone marks in the present study. A standard convolution is applied to the channels, followed by a fully-connected network.

For the regression task, the range of the outputs is the set of continuous real numbers. For the classification task, a *Softmax* function is applied to the models’ output to yield

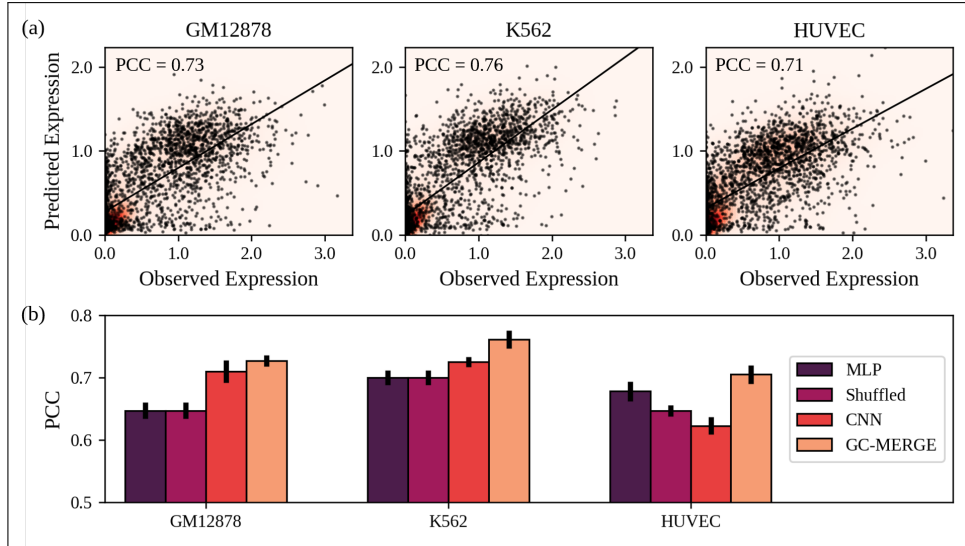


Figure 2.1: **Model comparison and evaluation.** GC-MERGE gives state-of-the-art performance for the regression task on all three studied cell lines. (a) The Pearson correlation coefficients (PCC) obtained by running GC-MERGE on each cell line are displayed. (b) GC-MERGE outperforms all the baseline models for each of the three cell lines. Scores are calculated as the average of 10 runs and standard deviations are denoted by error bars.

a binary prediction. For the CNN baseline, genomic regions are subdivided into smaller 100-bp bins, consistent with Singh *et al.* [30]. However, GC-MERGE and the baselines other than the CNN average the histone modification signals over the entire 10 kb region.

We also implemented GC-MERGE on higher resolution ChIP-seq datasets (1000-bp bins), which we fed through a linear embedding module to form features for the Hi-C nodes. We did not observe an improvement in the performance for the high-resolution input (Supplementary Figure B.1). Additionally, we compared our results to the published results of two other recent deep learning methods, Xpresso by Agarwal *et al.* [1] and DeepExpression by Zeng *et al.* [39], when such comparisons were possible, since in some cases the experimental data sets were unavailable or the code provided did not run

2.2.4 Evaluation Metrics

We measured the regression task performance of all the models by calculating the Pearson correlation coefficient (PCC), which quantifies the correlation between the true and predicted gene expression values in the test set. For interpretation, we adapted our model to perform classification. Therefore, we also evaluated the classification performance using two metrics: the area under the receiver operating characteristic curve (AUROC) and the area under the precision-recall curve (AUPR).

2.3 Results

2.3.1 Gene Expression Prediction Results

We evaluate GC-MERGE and the baseline models on the regression task for the GM12878, K562, and HUVEC cell lines. Figure 2.1(a) shows the predicted versus true gene expression values for GC-MERGE and Figure 2.1(b) compares our model’s performance with the baselines. Note that we determine the Pearson correlation coefficient (PCC) by taking the average of ten runs and denote the standard deviation by the error bars on the graph. For GM12878, the Pearson correlation coefficient of GC-MERGE predictions (PCC = 0.73) exceeds that of the other baselines. Furthermore, we note that our model performance also compares favorably to Xpresso (PCC \approx 0.65) [1], a CNN model that uses promoter sequence and 8 features associated with mRNA decay to predict gene expression. For K562, GC-MERGE again outperforms all alternative baseline models (PCC = 0.76). In addition, GC-MERGE performance also exceeds that of Xpresso (PCC \approx 0.71) [1] as well as DeepExpression (PCC = 0.65) [39], a CNN model that uses promoter sequence data as well as spatial information from H3K27ac and YY1 Hi-ChIP experiments. Our model gives better performance (PCC = 0.71) for HUVEC as well. Neither Xpresso nor DeepExpression studied this cell line.

These results strongly suggest that including spatial information can improve gene expression predictive performance over methods solely using local histone mark features as input. We emphasize that this prediction task allows us to model the relationships between the histone marks, 3D structure of the DNA, and gene expression. Therefore, a good performance indicates that the model can leverage the existing data to learn these connections. One of our main goals is to extract these relationships from the model and present GC-MERGE as a hypothesis driving tool for understanding epigenetic regulation.

2.3.2 Interpretation Results

To determine the underlying biological factors driving the model’s predictions, we integrate the GNNExplainer method, designed for classification tasks, into our modeling framework. Adapting our GC-MERGE model to the classification task also resulted in state-of-the-art performance (Supplementary Figure B.2) achieving 0.87, 0.88, and 0.85 AUPR scores for GM12878, K562, and HUVEC, respectively. Once trained, we show that our classification model can determine which spatial interactions are most critical to a gene’s expression and the histone marks that are most important. For GM12878, a lymphoblastoid cell line, we selected four genes: SIDT1, AKR1B1, LAPTM5, and TOP2B as exemplar genes. These genes are among the most highly expressed genes in our data set, and they have also been experimentally shown to be controlled by several long-range promoter-enhancer interactions [13]. To illustrate the validity of our approach, we perform analyses for each of these genes and corroborate our results using previous studies from the literature. Supplementary Table B.1 lists the chromosomal coordinates and corresponding node identifiers for each gene.

- **SIDT1** encodes a transmembrane dsRNA-gated channel protein and is part of a larger family of proteins necessary for systemic RNA interference [7, 20]. This gene has also been implicated in chemoresistance to the drug gemcitabine in adenocarcinoma cells [7] and is regulated by at least three chromosomal regions [13, 20]. In Figure 2.2(a), we show that for SIDT1, the model makes use of all three genomic re-

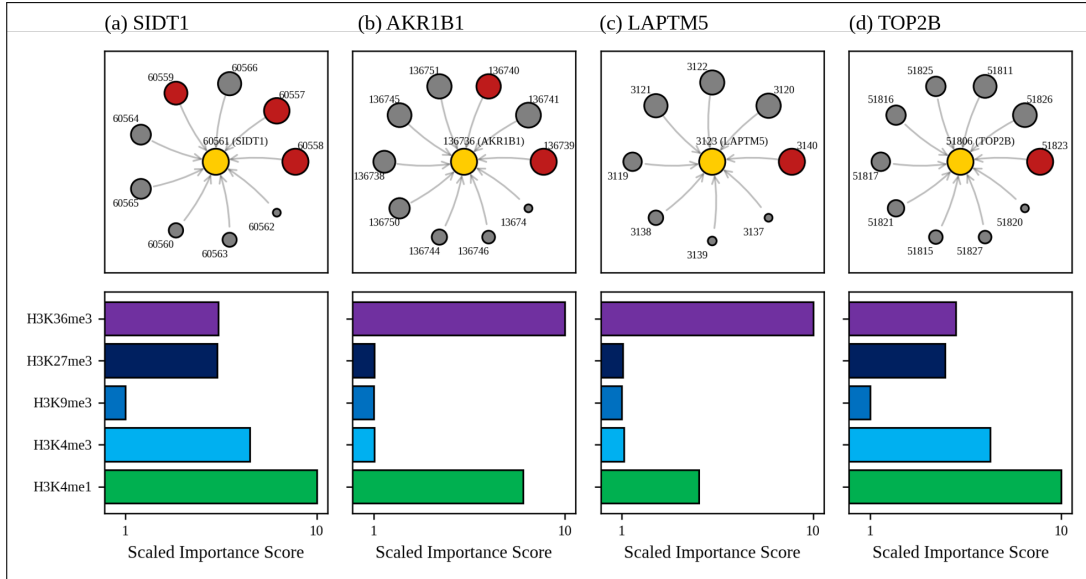


Figure 2.2: **Model explanations for exemplar genes.** **Top:** For (a) SIDT1, designated as node 60561 (yellow circle), the subgraph of neighbor nodes is displayed. The size of each neighbor node correlates with its predictive importance as determined by GNNExplainer. Nodes in red denote regions corresponding to known enhancer regions regulating SIDT1 [13] (note that multiple interacting fragments can be assigned to each node, see Supplementary Tables B.1 and B.2). All other nodes are displayed in gray. Nodes with importance scores corresponding to outliers have been removed for clarity. **Bottom:** The scaled feature importance scores for each of the five core histone marks used in this study are shown in the bar graph. Results also presented for (b) AKR1B1, (c) LAPTM5, and (d) TOP2B.

gions known to have regulatory effects by assigning high importance scores to those nodes (indicated by the size of the node). In addition, we plot the importance scores assigned to the histone marks (node features) that are most important in driving the model’s predictions. From the bar graph, it is apparent that H3K4me1 and H3K4me3 are the two most important features in determining the model’s prediction. This histone mark profile has been associated with regions flanking transcription start sites (TSS) in highly expressed genes [24, 8].

- **AKR1B1** encodes an enzyme that belongs to the aldo-keto reductase family. It has also been identified as a key player in complications associated with diabetes [5, 20] and is regulated by at least two chromosomal regions [13]. As seen in Figure 2.2(b), the model strongly bases its predictions for AKR1B1 on both of the regions known to have regulatory effects (location information in Supplementary Table B.2). We also show that H3K36me3 and H3K4me1 are the two histone marks with the highest scaled importance scores. This chromatin state signature is correlated with genic enhancers of highly expressed genes [24].
- **LAPTM5** encodes a receptor protein that spans the lysosomal membrane [20]. It is highly expressed in immune cells and plays a role in the downregulation of T and B cell receptors and the upregulation of macrophage cytokine production [10] as well as interacts with at least one regulatory sequence [13]. In Figure 2.2(c), the genomic region corresponding to the node with the highest scaled importance score has been experimentally shown to interact with LAPTM5 (see Supplementary

Table B.2) [13] and its histone mark profile is characteristic of genic enhancers [8].

- **TOP2B** encodes DNA topoisomerase II beta, a protein that controls the topological state of DNA during transcription and replication [20]. It transiently breaks and then reforms duplex DNA, relieving torsional stress. Mutations in this enzyme can lead to B cell immunodeficiency [3] and it has been shown to interact with at least two regulatory regions [13]. Figure 2.2(d) shows that the most important neighbor node has been corroborated by experiments to have a regulatory role for that gene [13] and its histone mark profile is indicative of regions flanking TSS [8].

To confirm that the node importance scores obtained from GNNExplainer do not merely reflect the relative magnitudes of the Hi-C counts or the distances between genomic regions, we investigated the relationships among the Hi-C counts, genomic distances, and scaled importance scores for all four exemplar genes (Supplementary Figures B.3 – B.6). We observe that the scaled importance scores do not correspond to the Hi-C counts or the pairwise genomic distances. For example, for *SIDT1*, the three experimentally validated interacting nodes achieve the highest importance scores (10, 9.55, and 7.73). However, they do not correspond to the regions with the highest Hi-C counts (154.78, 412.53, and 170.55 for each of the three known regulatory regions while the highest count is 602.84). In addition, although they are close to the *SIDT1* gene region (40, 20, and 30 kbp away), there are other nodes at the same or closer distances that do not have promoter-enhancer interactions. Therefore, we show that by modeling the histone modifications and the spatial configuration of the genome, GC-MERGE infers connections that could serve as important hypothesis-driving observations for gene regulatory experiments.

2.4 Discussion

We present GC-MERGE, a graph-based deep learning model, which integrates both local and long-range epigenetic data using a graph convolutional network framework to predict gene expression and explain its drivers. We demonstrate its state-of-the-art performance for the gene expression prediction task, outperforming the baselines on the GM12878, K562, and HUVEC cell lines. We also determine the relative contributions of histone modifications and long-range interactions for four highly expressed genes, showing that our model recapitulates known experimental results in a biologically interpretable manner.

With respect to potential future work for GC-MERGE, our framework can be applied on additional cell lines as high-quality Hi-C data sets become available. Incorporating other features, such as promoter sequence, would also be natural extensions. One avenue of particular importance would be to develop more robust methods for interpreting GCNs. For example, while the GNNExplainer model is a theoretically sound framework and yields an unbiased estimator for the importance scores of the subgraph nodes and features, there is variation in the interpretation scores generated over multiple runs. Furthermore, with larger GCNs, the optimization function utilized in GNNExplainer is challenging to minimize in practice. For some iterations, the importance scores converge with little differentiation and the method fails to arrive at a compact representation. This may be due to the relatively small penalties the method applies with respect to constraining the optimal size of the mask and the entropy of the distribution. We plan to address this issue in the future by implementing more robust forms of regularization.

In summary, GC-MERGE demonstrates proof-of-principle for using GCNs to predict

gene expression using both local epigenetic features and long-range spatial interactions. Interpretation of this model allows us to propose plausible biological explanations of the key regulatory factors driving gene expression as well as provide guidance regarding promising hypotheses and new research directions.

3. XL-MERGE

This chapter outlines the methods, experimental setup, and results for XL-MERGE, and it also contains a brief discussion.

3.1 Methods

3.1.1 Model Architecture

XL-MERGE has a similar architecture to GC-MERGE in that it is also a GraphSAGE-formulated graph convolutional network, the mathematical formulation for which is described in the GC-MERGE chapter, but its main differences have to do with pre-embeddings that are created to give better representations of the nodes within our graph formulation, as well as pre-embeddings to better represent genic regions. A schematic of XL-MERGE is shown in Figure C.1.

Better Representation for Genic Regions

One of the drawbacks of GC-MERGE was that in its formulation the 10kb genic regions it was using to predict gene expression were not necessarily centered around the gene’s respective TSS (transcription start site). For XL-MERGE, we attempt to avoid this problem through integrating a new dataset where histone mark data is centered around each TSS. This can be obtained through the dataset used from Singh *et al.* [30]. Additionally, rather than just taking the average of each histone mark across the entire 10kb region as in GC-MERGE, we pass more granular data from the 10kb region (100 bins of 100-base pair regions within the entire 10kb region) into a convolutional layer followed by maxpooling, nonlinear activation, and another linear layer. This serves to capture potential histone mark patterns within the region that could help drive gene expression, and also passes on higher dimensional information for downstream portions of our model than GC-MERGE, perhaps giving greater insight into gene regulation mechanisms.

Better Representation for Neighboring Nodes

Another limitation from GC-MERGE that we seek to improve upon is the capture of relevant histone mark information that could affect long-range interaction between genic regions and potential enhancer/repressor regions. GC-MERGE has a rather simplistic formulation in this respect, since it only takes the average normalized histone mark counts across the entire 10kb region. There are notable issues with this. First of all, we have no way of knowing where along the 10kb long-range regulatory region there is some sort of interaction with the genic region. It could easily be a small portion of this 10kb region, and we would not know where along this region that portion is located. Additionally, assuming that the enhancer/repressor interaction only takes place within a portion of

this 10kb region, by averaging histone mark data across the entire region we introduce significant random noise in our capture, since the histone marks outside of this interacting fragment likely have little to do with regulation of our genic region.

To get around these issues, we can apply a convolutional layer on the 10kb region for each of the neighboring node regions to a particular genic node followed by maxpooling and a nonlinear activation. The convolutional layer seeks to capture certain histone mark patterns within these neighboring node regions that can be relevant for long-range regulation, and the maxpooling gives a degree of positional invariance such that the model is more agnostic to where a certain interacting portion may be located. Therefore, we can gain better capture of potential regulatory interaction motifs for histone marks while also being more discriminatory and filtering out noise from that which we use to find these interactions.

3.2 Experimental Setup

3.2.1 Overview of Datasets and Data Integration

Most of the datasets used for XL-MERGE are covered in section 2.2.1, and their integration is covered in 2.2.2. XL-MERGE uses the same datasets as GC-MERGE and integrates them similarly. One of the additional datasets XL-MERGE uses, however is from Singh *et al.* [30]. This dataset contains histone mark counts in bins of 100 base pairs along each TSS-centered 10kb region via ChIP-Seq experiments for all 5 of the histone marks used in GC-MERGE, along with the corresponding gene catalog ID from the REMC database. Using the gene catalog ID we could map corresponding genic regions to those in GC-MERGE, allowing us to have a TSS-centered representation of genic nodes integrated into our dataset.

Another discrepancy from the GC-MERGE dataset is that feature vectors for each node in the graph were represented differently. Instead of getting the average ChIP-seq signal for each of the 5 core histone marks over the chromosomal region corresponding to each node, the representation was changed to getting the average ChIP-seq signal for each of the 5 core histone marks across each 100 base-pair region along the entire chromosomal region corresponding to the node. This gave 100 bins of the 5 histone marks along each region, making it possible to apply a one-dimensional convolution on each node as referred to earlier.

3.2.2 Baseline Models

In addition to some of the baseline methods run for GC-MERGE (multi-layer perceptron and convolutional neural network), as well as GC-MERGE itself, we were able to run another baseline that was not done for GC-MERGE which provide more evidence for XL-MERGE’s efficacy for predicting gene expression and the potential successful integration of long-range interactions within the genome to drive this more efficacious prediction. This baseline was:

- **AttentiveChrome [29]:** A long short-term memory-based model with attention mechanisms to capture longer-range histone mark dependencies within a genic region and assign appropriate weightings to certain regions within the region. Out of deep learning models we are aware of in the field that are used to predict gene expression from strictly histone marks, AttentiveChrome appears to have the best

predictive power in terms of AUROC (Area Under the Receiver Operating Characteristics). AttentiveChrome even outperforms GC-MERGE by a fair margin. We also ensured that we used the same train-validation-test split for AttentiveChrome as for XL-MERGE when comparing the performance of the two.

Unfortunately, due to time constraints, there have not yet been any regression results run for XL-MERGE. Thus, unlike for GC-MERGE, we are not currently able to directly compare it to Xpresso [1] and DeepExpression [39], since these models only gave results as regression tasks.

3.2.3 Model Ablation Analyses

We also performed a set of model ablation analyses for XL-MERGE, where we zero out the features of an entire portion of our model in order to obtain a more precise analysis of whether or not we capture long-range interactions effectively. For our ablation analysis, we end up zeroing out the TSS-centered local embedding representation for a given gene, such that the predictive power of the model comes purely from neighboring nodes in our graph formulation for that gene. Since in this ablated model the predictive power comes solely from potential long-range interacting nodes, this gives us a less noisy indication as to whether or not these long-range interactions are captured well relative to baselines, rather than having these differences be crowded out by the local genic region driving prediction. The ablation analyses we ended up running were:

- **Standard XL-MERGE ablation:** This was just running our standard XL-MERGE model with the TSS-centered genic node embeddings zeroed out such that genic region information does not contribute to gene expression prediction. Only long-range embedding information should therefore contribute to gene expression prediction.
- **Ablation of XL-MERGE using two-layer perceptrons:** This was an ablation of TSS-centered genic node embeddings except where we slightly modified the mechanism through which we extract information from long-range neighboring nodes. Instead of using a convolution layer followed by maxpooling as in the traditional XL-MERGE model, we use a two-layer perceptron instead. This allows us to compare different mechanisms for which we extract information from long-range interactions.
- **Ablation of XL-MERGE without convolution or multi-layer perceptrons:** This was an ablation of TSS-centered genic node embeddings except there is no specific extraction mechanism for the information of long-range neighboring nodes. That is, we propagate the original normalized histone marks counts through our graph formulation, using that to predict gene expression from neighboring nodes. This provides a null sort of baseline for other ways through which we extract neighboring node gene regulation information.
- **Ablation of XL-MERGE with shuffled neighboring nodes:** Here we run XL-MERGE with the TSS-centered genic node embeddings zeroed out, but in addition to this we shuffle the long-range neighbors corresponding to each genic region. Thus, this can give us insight into whether XL-MERGE utilizes its true neighboring nodes to help drive gene expression prediction with potential regulatory interactions relative to random nodes being used instead.

3.2.4 Evaluation Metrics

Up to this point, XL-MERGE has only been formulated as a binary prediction task, so therefore it is evaluated using the same two classification metrics that were used for GC-MERGE: the area under the receiver operating characteristic curve (AUROC) and the area under the precision-recall curve (AUPR).

3.3 Results

3.3.1 Gene Expression Prediction Results

We evaluate XL-MERGE and the baseline models on the binary classification task for the GM12878, K562, and HUVEC cell lines. Runs for XL-MERGE and each baseline for each cell line were performed 10 times each. Then, we could find the averages across each of these 10 runs. Figure C.2 shows the relative AUROC and AUPR scores computed for XL-MERGE and baselines with a brief analysis, and Figure C.4 shows numerical values for reference.

The results suggest that XL-MERGE is state-of-the-art for binarily predicting gene expression, and that XL-MERGE handily outperforms GC-MERGE in terms of binary gene classification prediction.

3.3.2 Model Ablation Results

We also perform model ablation analyses for XL-MERGE that were described in section 3.2.3. These analyses were each done for the GM12878, K562, and HUVEC cell lines for 10 times for each cell line, and the average AUROC and AUPR scores could be computed across each of these 10 runs. These relative scores are shown in Figure C.3, and Figure C.4 shows actual numerical values for reference.

From these results, it seems that using a convolutional layer along with maxpooling to extract potential long-range interaction information is superior to using a multi-layer perceptron or not using any extraction mechanism at all. This superior performance might suggest that the convolutional neighbor aggregation mechanism is advantageous due its being positionally agnostic to where certain motifs may occur along a potential long-range interaction fragment, and also that the information extracted is less noisy.

There is also strong evidence that XL-MERGE indeed makes use of its potential long-range interacting fragments in its graph construction to drive gene expression prediction, suggesting that it properly incorporates these long-range regulation interactions in its model. The XL-MERGE genic node embedding-ablated model significantly outperforms the genic node embedding-ablated model when neighboring nodes in the graph for a given genic node are shuffled with other random nodes in the graph. Additionally, when neighbors are shuffled in the ablated model, AUROC scores are close to 0.5, suggesting prediction power similar to random guessing when node neighbors are shuffled from the original graph formulation.

3.4 Discussion

Overall, XL-MERGE seems to act as a state-of-the-art model for predicting binary gene expression from histone marks. It provides strong evidence for incorporating long-range regulatory epigenetic interactions to help drive gene expression prediction, thus suggesting the model can potentially identify enhancer/repressor regions. However, this still significant work to further be done to make XL-MERGE more biologically useful as well as comparable to other models, and to improve its robustness.

3.4.1 Adding Interpretability for XL-MERGE

For XL-MERGE, we have yet to add any interpretation angle for analyzing how XL-MERGE makes its predictions, which is critical to make XL-MERGE biologically relevant and useful. For GC-MERGE, we used GNNExplainer, and while it could be good to use GNNExplainer for XL-MERGE, there may be better options. GNNExplainer is computationally useful because it can generate subgraphs most relevant to driving graph neural network predictions in a computationally tractable fashion, but for XL-MERGE only node information one hop away is actually used in predicting gene expression for a particular genic node. Therefore, since we construct our graph such that edges are formed with only the 10 most relevant long-range interacting genomic regions (nodes), there are not a computationally intractable amount of subgraphs or neighboring node combinations that would need to be tested for driving gene expression to generate the most relevant explanation. A simpler, less abstracted method than GNNExplainer, where say we just look at all possible subgraphs for a particular node, might work better because GNNExplainer is commonly known throughout the deep learning community to have issues with generating consistent explanations. This overall presents possible directions for interpretability methods for XL-MERGE that would be advantageous to that of GC-MERGE.

3.4.2 Regression Task Results for XL-MERGE

Another set of results to be added for XL-MERGE would be gene expression prediction regression results. This would make XL-MERGE able to be comparable to other methods, such as Xpresso and DeepExpression, and it would also give XL-MERGE a more information-rich representation for gene expression prediction, rather than just a binary indicator. Additionally, with the regression formulation, downstream analyses can take place for genes that are most highly predicted for expression, or genes that are most highly predicted for no expression.

3.4.3 Normalizing Genic Region Input Data

Another important adjustment that should be made to XL-MERGE involves normalizing the histone mark data for TSS-centered genic regions that is inputted into the model. As of right now, the histone mark data for the genic regions is not normalized in the same way that the histone mark data for long-range interacting regions is. The histone mark data values for the genic regions tend to be higher, thus potentially making XL-MERGE more biased toward the genic region histone mark data. For the future, this issue should

be resolved, and it will perhaps leads to greater influence over gene expression coming from long-range regulatory interactions.

3.4.4 Different Training Rates for Different Parts of the Model

One thing noticed while performing the model ablation analyses of XL-MERGE was that it would take the entire model an apparently shorter time to train than it would during the model ablation analyses. This implies that XL-MERGE takes longer to learn meaningful extractions from long-range interacting regions in the model than it does for the genic regions. In order to best learn from both the genic regions and long-range interacting regions, going forward it may be better to train different parts of the model with different learning rates, or to perhaps ablate one portion of the model during an earlier portion of training as to let the other parts of the model learn their representations first. Thus, all parts of the XL-MERGE model could be used together optimally to better predict gene expression.

Bibliography

- [1] Vikram Agarwal and Jay Shendure. “Predicting mRNA Abundance Directly from Genomic Sequence Using Deep Convolutional Neural Networks”. In: *Cell Reports* 31.7 (May 2020), p. 107663. ISSN: 22111247. DOI: 10.1016/j.celrep.2020.107663. URL: <https://linkinghub.elsevier.com/retrieve/pii/S2211124720306161> (visited on 10/19/2020).
- [2] Ioana Bica et al. “Unsupervised generative and graph representation learning for modelling cell differentiation”. In: *Scientific Reports* 10.1 (Dec. 2020), p. 9790. ISSN: 2045-2322. DOI: 10.1038/s41598-020-66166-8. URL: <http://www.nature.com/articles/s41598-020-66166-8> (visited on 10/23/2020).
- [3] Lori Broderick et al. “Mutations in topoisomerase II result in a B cell immunodeficiency”. In: *Nature Communications* 10.1 (Dec. 2019), p. 3644. ISSN: 2041-1723. DOI: 10.1038/s41467-019-11570-6. URL: <http://www.nature.com/articles/s41467-019-11570-6> (visited on 01/24/2021).
- [4] Chao Cheng et al. “A statistical framework for modeling gene expression using chromatin features and application to modENCODE datasets”. In: *Genome Biology* 12.2 (2011), R15. ISSN: 1465-6906. DOI: 10.1186/gb-2011-12-2-r15. URL: <http://genomebiology.biomedcentral.com/articles/10.1186/gb-2011-12-2-r15> (visited on 10/19/2020).
- [5] K. C. Donaghue et al. “The association of aldose reductase gene (AKR1B1) polymorphisms with diabetic neuropathy in adolescents”. In: *Diabetic Medicine* 22.10 (Oct. 2005), pp. 1315–1320. ISSN: 0742-3071, 1464-5491. DOI: 10.1111/j.1464-5491.2005.01631.x. URL: <http://doi.wiley.com/10.1111/j.1464-5491.2005.01631.x> (visited on 11/05/2020).
- [6] Xianjun Dong et al. “Modeling gene expression using chromatin features in various cellular contexts”. In: *Genome Biology* 13.9 (2012), R53. ISSN: 1465-6906. DOI: 10.1186/gb-2012-13-9-r53. URL: <http://genomebiology.biomedcentral.com/articles/10.1186/gb-2012-13-9-r53> (visited on 10/19/2020).
- [7] Mohamed O. Elhassan, Jennifer Christie, and Mark S. Duxbury. “*Homo sapiens* Systemic RNA Interference-defective-1 Transmembrane Family Member 1 (SIDT1) Protein Mediates Contact-dependent Small RNA Transfer and MicroRNA-21-driven Chemoresistance”. In: *Journal of Biological Chemistry* 287.8 (Feb. 17, 2012), pp. 5267–5277. ISSN: 0021-9258, 1083-351X. DOI: 10.1074/jbc.M111.318865. URL: <http://www.jbc.org/lookup/doi/10.1074/jbc.M111.318865> (visited on 11/05/2020).
- [8] Jason Ernst and Manolis Kellis. “Chromatin-state discovery and genome annotation with ChromHMM”. In: *Nature Protocols* 12.12 (Dec. 2017), pp. 2478–2492. ISSN: 1754-2189, 1750-2799. DOI: 10.1038/nprot.2017.124. URL: <http://www.nature.com/articles/nprot.2017.124> (visited on 10/22/2020).

- [9] Matthias Fey and Jan E. Lenssen. “Fast Graph Representation Learning with PyTorch Geometric”. In: *ICLR Workshop on Representation Learning on Graphs and Manifolds*. 2019.
- [10] Wioletta K. Glowacka et al. “LAPTM5 Protein Is a Positive Regulator of Proinflammatory Signaling Pathways in Macrophages”. In: *Journal of Biological Chemistry* 287.33 (Aug. 2012), pp. 27691–27702. ISSN: 00219258. DOI: 10.1074/jbc.M112.355917. URL: <https://linkinghub.elsevier.com/retrieve/pii/S0021925820477920> (visited on 01/24/2021).
- [11] William L. Hamilton, Rex Ying, and Jure Leskovec. “Inductive Representation Learning on Large Graphs”. In: *Proceedings of the 31st International Conference on Neural Information Processing Systems*. NIPS’17. Long Beach, California, USA: Curran Associates Inc., 2017, pp. 1025–1035. ISBN: 9781510860964.
- [12] Bowen Jin et al. “Multi-behavior Recommendation with Graph Convolutional Networks”. In: *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. SIGIR ’20: The 43rd International ACM SIGIR conference on research and development in Information Retrieval. Virtual Event China: ACM, July 25, 2020, pp. 659–668. ISBN: 978-1-4503-8016-4. DOI: 10.1145/3397271.3401072. URL: <https://dl.acm.org/doi/10.1145/3397271.3401072> (visited on 10/23/2020).
- [13] Inkyung Jung et al. “A compendium of promoter-centered long-range chromatin interactions in the human genome”. In: *Nature Genetics* 51.10 (Oct. 2019), pp. 1442–1449. ISSN: 1061-4036, 1546-1718. DOI: 10.1038/s41588-019-0494-8. URL: <http://www.nature.com/articles/s41588-019-0494-8> (visited on 10/22/2020).
- [14] R. Karlic et al. “Histone modification levels are predictive for gene expression”. In: *Proceedings of the National Academy of Sciences* 107.7 (Feb. 16, 2010), pp. 2926–2931. ISSN: 0027-8424, 1091-6490. DOI: 10.1073/pnas.0909344107. URL: <http://www.pnas.org/cgi/doi/10.1073/pnas.0909344107> (visited on 10/19/2020).
- [15] Diederik P. Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In: *CoRR* abs/1412.6980 (2015).
- [16] Dirk A Kleinjan et al. “Aniridia-associated translocations, DNase hypersensitivity, sequence comparison and transgenic analysis redefine the functional domain of PAX6”. In: *Human molecular genetics* 10.19 (2001), pp. 2049–2059.
- [17] Peter Hugo Lodewijk Krijger and Wouter de Laat. “Regulation of disease-associated gene expression in the 3D genome”. In: *Nature Reviews Molecular Cell Biology* 17.12 (Dec. 2016), pp. 771–782. ISSN: 1471-0072, 1471-0080. DOI: 10.1038/nrm.2016.138. URL: <http://www.nature.com/articles/nrm.2016.138> (visited on 10/18/2020).
- [18] Zhiyuan Liu and Jie Zhou. “Introduction to Graph Neural Networks”. In: *Synthesis Lectures on Artificial Intelligence and Machine Learning* 14.2 (Mar. 19, 2020), pp. 1–127. ISSN: 1939-4608, 1939-4616. DOI: 10.2200/S00980ED1V01Y202001AIM045. URL: <https://www.morganclaypool.com/doi/10.2200/S00980ED1V01Y202001AIM045> (visited on 10/23/2020).
- [19] Maxwell R Mumbach et al. “HiChIP: efficient and sensitive analysis of protein-directed genome architecture”. In: *Nature methods* 13.11 (2016), pp. 919–922.

- [20] National Center for Biotechnology Information National Library of Medicine (US). *Entrez Gene*. <https://www.ncbi.nlm.nih.gov/gene/>. Accessed: 2020-10-22. 1988-.
- [21] Jiezhong Qiu et al. “DeepInf: Social Influence Prediction with Deep Learning”. In: *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. KDD '18: The 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. London United Kingdom: ACM, July 19, 2018, pp. 2110–2119. ISBN: 978-1-4503-5552-0. DOI: 10.1145/3219819.3220077. URL: <https://dl.acm.org/doi/10.1145/3219819.3220077> (visited on 10/23/2020).
- [22] Suhas S.P. Rao, Miriam H. Huntley, Neva C. Durand, et al. “A 3D Map of the Human Genome at Kilobase Resolution Reveals Principles of Chromatin Looping”. In: *Cell* 159.7 (Dec. 2014), pp. 1665–1680. ISSN: 00928674. DOI: 10.1016/j.cell.2014.11.021. URL: <https://linkinghub.elsevier.com/retrieve/pii/S0092867414014974> (visited on 10/25/2020).
- [23] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “Why should i trust you?: Explaining the predictions of any classifier”. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. ACM. 2016, pp. 1135–1144.
- [24] Roadmap Epigenomics Consortium. “Integrative analysis of 111 reference human epigenomes”. In: *Nature* 518.7539 (Feb. 19, 2015), pp. 317–330. ISSN: 0028-0836, 1476-4687. DOI: 10.1038/nature14248. URL: <http://www.nature.com/articles/nature14248> (visited on 10/22/2020).
- [25] M. Jordan Rowley and Victor G. Corces. “Organizational principles of 3D genome architecture”. In: *Nature Reviews Genetics* 19.12 (Dec. 2018), pp. 789–800. ISSN: 1471-0056, 1471-0064. DOI: 10.1038/s41576-018-0060-8. URL: <http://www.nature.com/articles/s41576-018-0060-8> (visited on 10/27/2020).
- [26] Stefan Schoenfelder and Peter Fraser. “Long-range enhancer–promoter contacts in gene expression control”. In: *Nature Reviews Genetics* 20.8 (Aug. 2019), pp. 437–455. ISSN: 1471-0056, 1471-0064. DOI: 10.1038/s41576-019-0128-0. URL: <http://www.nature.com/articles/s41576-019-0128-0> (visited on 10/18/2020).
- [27] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. “Learning important features through propagating activation differences”. In: *arXiv preprint arXiv:1704.02685* (2017).
- [28] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. “Deep inside convolutional networks: Visualising image classification models and saliency maps”. In: *arXiv preprint arXiv:1312.6034* (2013).
- [29] Ritambhara Singh et al. “Attend and Predict: Understanding Gene Regulation by Selective Attention on Chromatin”. In: *Advances in Neural Information Processing Systems* 30 (Dec. 2017), pp. 6785–6795. ISSN: 1049-5258.
- [30] Ritambhara Singh et al. “DeepChrome: deep-learning for predicting gene expression from histone modifications”. In: *Bioinformatics* 32.17 (2016), pp. i639–i648.

- [31] Haitham Sobhy et al. “Highly interacting regions of the human genome are enriched with enhancers and bound by DNA repair proteins”. In: *Scientific Reports* 9.1 (Dec. 2019), p. 4577. ISSN: 2045-2322. DOI: 10.1038/s41598-019-40770-9. URL: <http://www.nature.com/articles/s41598-019-40770-9> (visited on 10/25/2020).
- [32] Mengying Sun et al. “Graph convolutional networks for computational drug development and discovery”. In: *Briefings in Bioinformatics* 21.3 (May 21, 2020), pp. 919–935. ISSN: 1467-5463, 1477-4054. DOI: 10.1093/bib/bbz042. URL: <https://academic.oup.com/bib/article/21/3/919/5498046> (visited on 10/23/2020).
- [33] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. “Axiomatic attribution for deep networks”. In: *arXiv preprint arXiv:1703.01365* (2017).
- [34] Nynke L Van Berkum et al. “Hi-C: a method to study the three-dimensional architecture of genomes.” In: *JoVE (Journal of Visualized Experiments)* 39 (2010), e1869.
- [35] Petar Veličković et al. “Graph attention networks”. In: *arXiv preprint arXiv:1710.10903* (2017).
- [36] Fang Yang et al. “Graph-based prediction of Protein-protein interactions with attributed signed graph embedding”. In: *BMC Bioinformatics* 21.1 (Dec. 2020), p. 323. ISSN: 1471-2105. DOI: 10.1186/s12859-020-03646-8. URL: <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-020-03646-8> (visited on 10/23/2020).
- [37] Rex Ying et al. “GNNExplainer: Generating Explanations for Graph Neural Networks”. In: (2019). arXiv: 1903.03894 [cs.LG].
- [38] Hanqing Zeng et al. “Graphsaint: Graph sampling based inductive learning method”. In: *arXiv preprint arXiv:1907.04931* (2019).
- [39] Wanwen Zeng, Yong Wang, and Rui Jiang. “Integrating distal and proximal information to predict gene expression via a densely connected convolutional neural network”. In: *Bioinformatics* (July 18, 2019). Ed. by Alfonso Valencia, btz562. ISSN: 1367-4803, 1460-2059. DOI: 10.1093/bioinformatics/btz562. URL: <https://academic.oup.com/bioinformatics/advance-article/doi/10.1093/bioinformatics/btz562/5535598> (visited on 10/19/2020).

Appendix

GC-MERGE Model Details	30
GC-MERGE Additional Dataset Details and Results	32
XL-MERGE Model Details and Results	40

A. GC-MERGE Model Details

A.1 Model architecture and training

The model architecture is represented in Figure A.1. Here, the first layer of the model performs a graph convolution on the initial feature embeddings with an output embedding size of 256, followed by application of ReLU, a non-linear activation function. The second layer of the model performs another graph convolution with the same embedding size of 256 on the transformed representations, again followed by application of ReLU. Next, the output is fed into three successive linear layers of sizes 256, 256, and 2, respectively. A regularization step is performed by using a dropout layer with probability 0.5. The model was trained using ADAM, a stochastic gradient descent algorithm [15]. We used the PyTorch Geometric package [9] to implement our code.

A.2 Hyperparameter tuning

Table C.1 details the hyperparameters and the range of values we used to conduct a grid search to determine the optimized model. Specifically, we varied the number of graph convolutional layers, number of linear layers, embedding size for graph convolutional layers, linear layer sizes, and inclusion (or exclusion) of an activation function after the graph convolutional layers. Through earlier iterations of hyperparameter tuning, we also tested the number of neighbors for each node (5 or 10), type of activation functions used for the linear layers of the model (ReLU, LeakyReLU, sigmoid, or tanh), method for accounting for background Hi-C counts, as well as dropout probabilities. Some combinations of hyperparameters were omitted from our grid search because the corresponding model’s memory requirements did not fit on the NVIDIA Titan RTX and Quadro RTX GPUs available to us on Brown University’s Center for Computation and Visualization (CCV) computing cluster. We recorded the loss curves for the training and validation sets over 1000 epochs, by which time the model began to overfit. In addition, the data was split into sets of 70% for training, 15% for validation, and 15% for testing. The optimal hyperparameters for our final model that also proved to be computationally feasible are as follows: 2 graph convolutional layers, 3 linear layers, graph convolutional layer embedding size of 256, linear layer sizes that match that of the graph convolutional layers, and using an activation function (ReLU) after all graph convolutional layers and all linear layers except for the last.

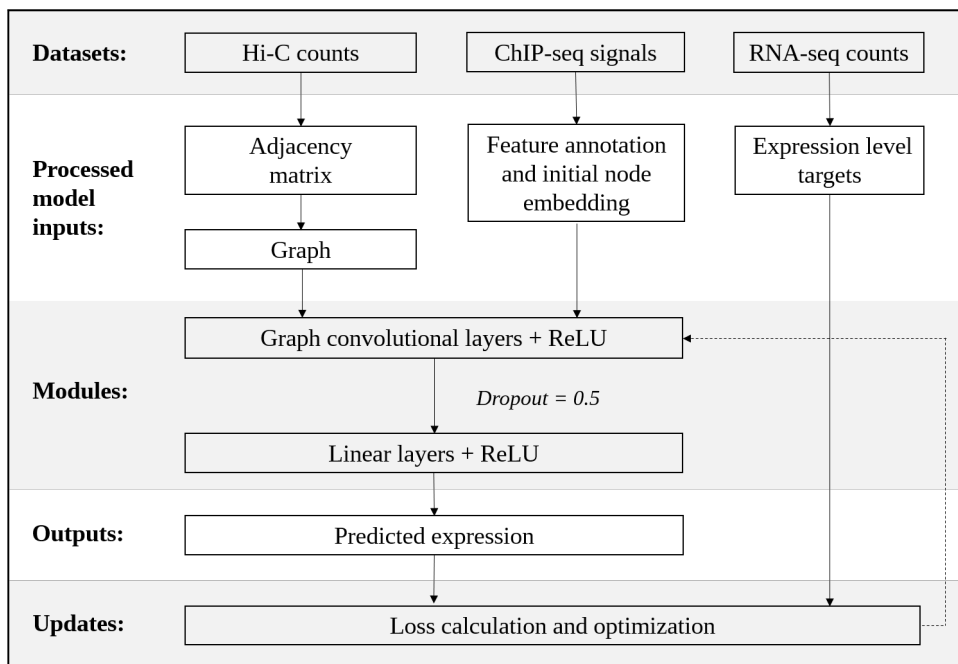


Figure A.1: **Overview of the GCNN model architecture.** The datasets used in our model are Hi-C maps, CHIP-seq signals, and RNA-seq counts. A binarized adjacency matrix is produced from the Hi-C maps by subsampling from the Hi-C matrix, such that only the top 10 neighbors of each node are preserved. The nodes in the graph are annotated with features from the CHIP-seq datasets. Two graph convolutions, each followed by ReLU, are performed. The output is fed into a dropout layer (probability = 0.5), followed by a linear module comprised of three dense layers, in which the first two layers are followed by ReLU. For the regression model, the final output represents the base-10 logarithm of the expression level (with a pseudocount of 1). For the classification model, the output is fed through a *Softmax* layer and then the *argmax* is taken to make the final prediction.

Hyperparameter	Values
Number of graph convolutional layers	1, 2
Number of linear layers	1, 2, 3
Graph convolutional layer embedding sizes	64, 128, 256, 384
Linear layer sizes	Keep sizes of all linear layers constant; alternatively, for each subsequent layer, divide size by 2
Activation function after graph convolutional layers	Include; alternatively, do not include

Table A.1: **Hyperparameter combinations used for tuning in grid search.** A grid search was conducted by varying the following hyperparameters: number of graph convolutional layers, number of linear layers, embedding size for graph convolutional layers, linear layer sizes, and inclusion/exclusion of activation function after the graph convolutional layers.

B. GC-MERGE Additional Dataset Details and Results

For chromosome capture data, we used previously published Hi-C maps at 10 kilobase (kb) resolution for all 22 autosomal chromosomes [22]. We obtained an $N \times N$ symmetric matrix, where each row or column corresponds to a 10 kb chromosomal region. Therefore, each bin coordinate (*row*, *column*) corresponds to the interaction frequency between two respective genomic regions. We applied VC-normalization on the Hi-C maps. In addition, because chromosomal regions located closer together will contact each other more frequently than regions located farther away simply due to chance (rather than due to biologically significant effects), we made an additional adjustment for this background effect. Following Sobhy *et al.* [31], we took the medians of the Hi-C counts for all pairs of interacting regions located the same distance away and used this as a proxy for the background. We subtracted the appropriate median from each Hi-C bin and discarded negative values.

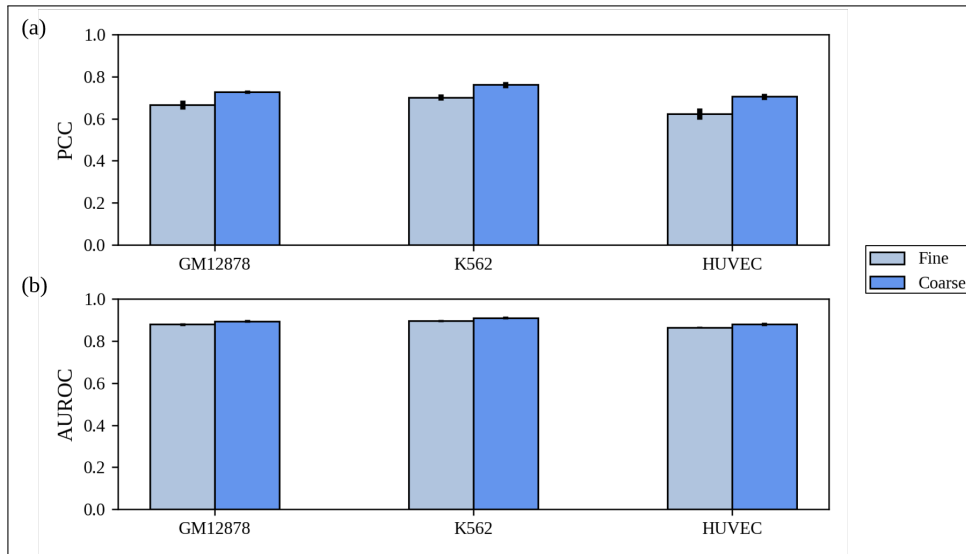


Figure B.1: **Comparison of fine-grained versus coarse-grained ChIP-seq signals for use in GC-MERGE.** For the coarse-grained resolution, ChIP-seq signals were averaged over the entire Hi-C bin (10000 bp resolution). For the fine-grained resolution, ChIP-seq signals were first averaged over 1000 bp bins and then fed into two embedding linear layers followed by ReLU. The output of these embedding layers was then used to feature annotate each node. (a) For the regression task, the fine-grained resolution ChIP-seq data produces performance worse than or comparable to the coarse-grained resolution ChIP-seq data as measured by PCC. (b) For the classification task, the fine-grained resolution ChIP-seq data performs slightly worse than or comparable to that of the coarse-grained resolution ChIP-seq data as measured by AUROC.

Gene	Node Identifier	Node Coordinates	Gene Coordinates
SIDT1	60561	chr3:113249241-113259241	chr3:113532296-113629579
AKR1B1	136736	chr7:134253323-134263323	chr7:134127127-134144036
LAPTM5	3123	chr1:31230000-31240000	chr3:31205316-31230667
TOP2B	51806	chr3:25699241-25709241	chr3:25639475-25706398

Table B.1: **Node coordinates for all exemplar genes: SIDT1, AKR1B1, LAPTM5, and TOP2B.** For each gene, the second and third columns list the corresponding node identifiers and the chromosome coordinates, respectively. The fourth column lists the gene’s actual chromosomal coordinates. Note that the transcription start site was used as the basis for assigning each gene to a node.

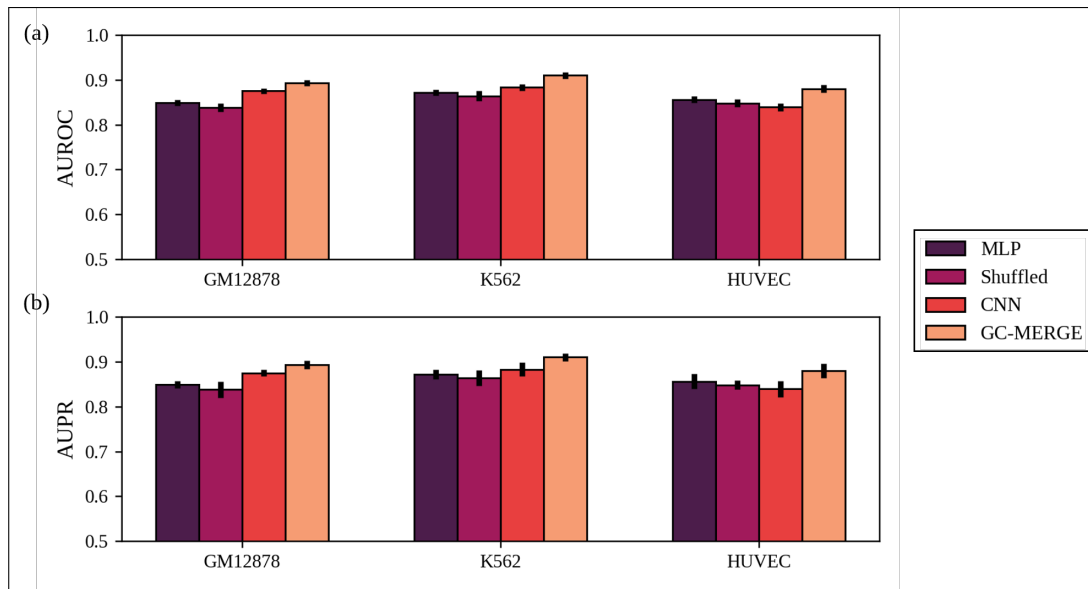


Figure B.2: **Comparison of AUROC and AUPR scores for GC-MERGE and its associated baselines.** GC-MERGE gives state-of-the-art performance for classifying genes as on/active or off/inactive. (a) The AUROC metrics for GM12878, K562, and HUVEC were 0.893, 0.910, and 0.880, respectively. For each of these cell lines, GC-MERGE performance exceeded all other baselines. (b) Using the AUPR metric, GC-MERGE obtains scores of 0.865, 0.884, and 0.848 for GM12878, K562, and HUVEC, respectively. As with the AUROC metric, our model’s performance was the highest among the baselines. Additionally, our AUROC score for K562 (0.91) is comparable to that reported by Zeng *et al.* [39] (0.91). We could not compare scores for the other two cell lines as they do not provide Hi-ChIP data for the cell line to run their model.

Gene	Neighbor Node Identifier	Neighbor Node Coordinates	Interacting Fragment Coordinates
SIDT1	60557	chr3:113209241-113219241	chr3:113251143-113348425
	60558	chr3:113219241-113229241	chr3:113228501-113232053
	60559	chr3:113229241-113239241	chr3:113228501-113232053
	60560	chr3:113239241-113249241	
	60562	chr3:113259241-113269241	
	60563	chr3:113269241-113279241	
	60564	chr3:113279241-113289241	
	60565	chr3:113289241-113299241	
	60566	chr3:113299241-113309241	
AKR1B1	136738	chr7:134273323-134283323	
	136739	chr7:134283323-134293323	chr7:134293046-134298798
	136740	chr7:134293323-134303323	chr7:134293046-134298798
	136741	chr7:134303323-134313323	
	136744	chr7:134333323-134343323	
	136745	chr7:134343323-134353323	
	136746	chr7:134353323-134363323	
	136747	chr7:134363323-134373323	
	136750	chr7:134393323-134403323	
136751	chr7:134403323-134413323		
LAPTM5	3119	chr1:31190000-31200000	
	3120	chr1:31200000-31210000	
	3121	chr1:31210000-31220000	
	3122	chr1:31220000-31230000	
	3137	chr1:31370000-31380000	
	3138	chr1:31380000-31390000	
	3139	chr1:31390000-31400000	
3140	chr1:31400000-31410000	chr1:31401583-31405576	
TOP2B	51811	chr3:25749241-25759241	
	51815	chr3:25789241-25799241	
	51816	chr3:25799241-25809241	
	51817	chr3:25809241-25819241	
	51820	chr3:25839241-25849241	
	51821	chr3:25849241-25859241	
	51823	chr3:25869241-25879241	chr3:25878006-25881223
	51825	chr3:25889241-25899241	
	51826	chr3:25899241-25909241	
51827	chr3:25909241-25919241		

Table B.2: **Neighbor coordinates for SIDT1, AKR1B1, LAPTM5, and TOP2B.** The second column lists the node identifiers for all neighboring nodes of the relevant gene, including neighboring nodes that contain interacting fragments as well as those that do not. The third column third lists the corresponding chromosome coordinates for the node identifier. The fourth column lists the regulatory fragments that interact with each gene as described in Jung *et al.* [13].

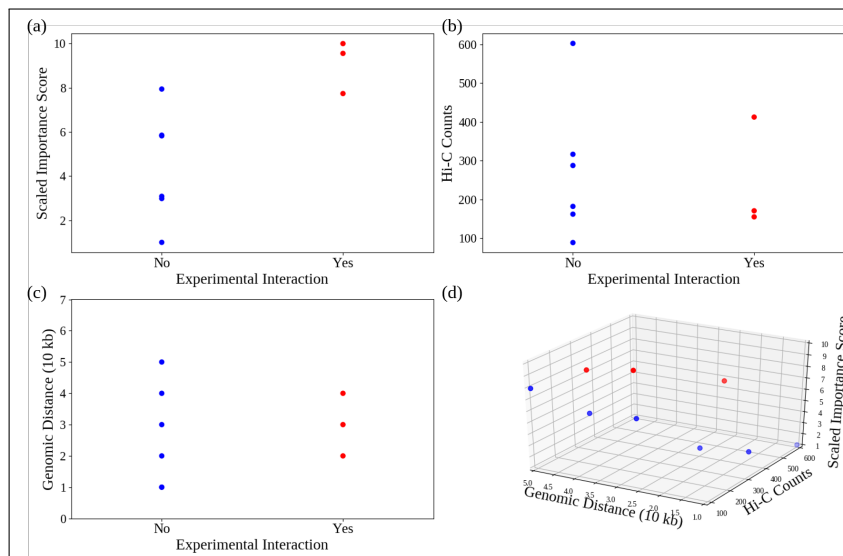


Figure B.3: **Relationships among scaled importance scores, genomic distances, and Hi-C counts for all SIDT1 neighbors.** Nodes corresponding to experimentally validated interacting fragments are denoted in red and all others are denoted in blue. (a) Scaled importance score versus experimental interaction. Experimentally validated interacting fragments are ranked higher on average than non-interacting fragments. (b) Hi-C counts versus experimental interaction. Hi-C counts by themselves are not sufficient to explain the presence of experimentally validated interactions. (c) Genomic distance versus experimental interaction. Genomic distance does not correlate with experimentally validated interactions. (d) 3D plot displaying the relationships among scaled importance scores, genomic distances, and Hi-C counts.

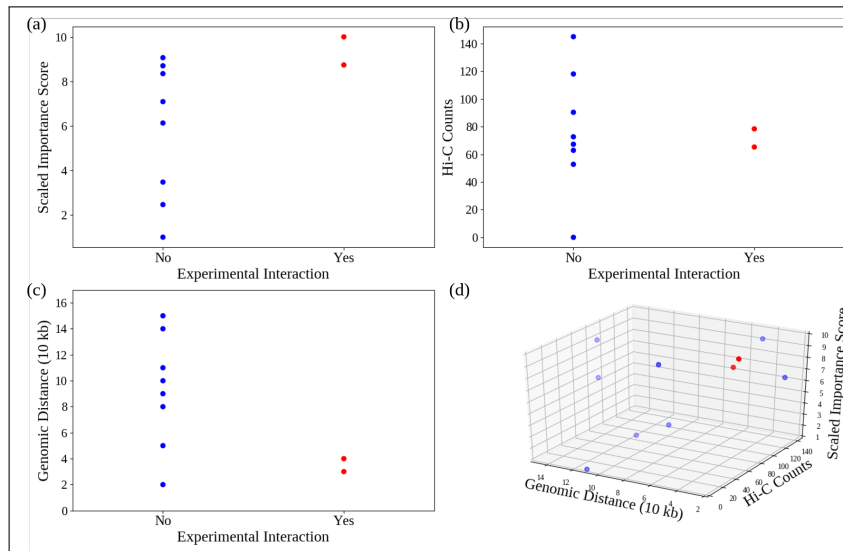


Figure B.4: **Relationships among scaled importance scores, genomic distances, and Hi-C counts for all AKR1B1 neighbors.** Nodes corresponding to experimentally validated interacting fragments are denoted in red and all others are denoted in blue. (a) Scaled importance score versus experimental interaction. Experimentally validated interacting fragments are ranked higher on average than non-interacting fragments. (b) Hi-C counts versus experimental interaction. Hi-C counts by themselves are not sufficient to explain the presence of experimentally validated interactions. (c) Genomic distance versus experimental interaction. Genomic distance does not correlate with experimentally validated interactions. (d) 3D plot displaying the relationships among scaled importance scores, genomic distances, and Hi-C counts.

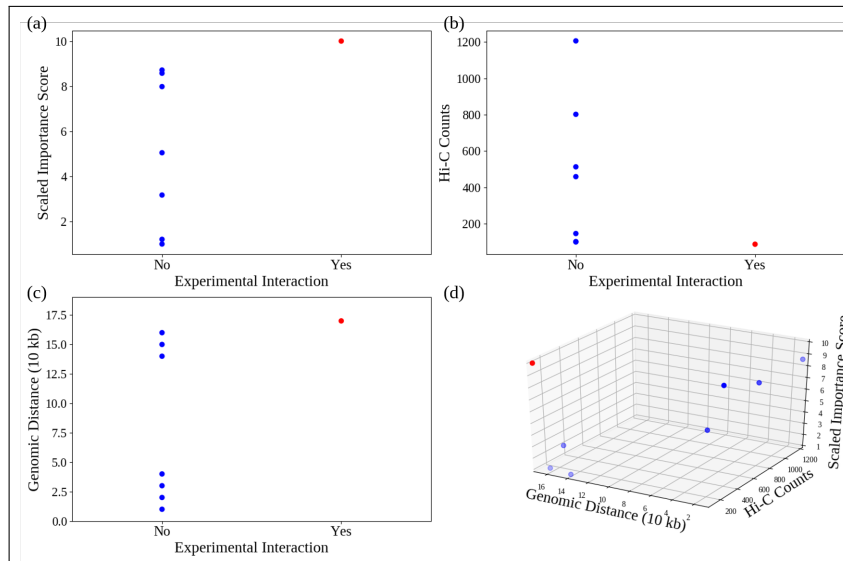


Figure B.5: **Relationships among scaled importance scores, genomic distances, and Hi-C counts for all LAPTM5 neighbors.** Nodes corresponding to experimentally validated interacting fragments are denoted in red and all others are denoted in blue. (a) Scaled importance score versus experimental interaction. Experimentally validated interacting fragments are ranked higher on average than non-interacting fragments. (b) Hi-C counts versus experimental interaction. Hi-C counts by themselves are not sufficient to explain the presence of experimentally validated interactions. (c) Genomic distance versus experimental interaction. Genomic distance does not correlate with experimentally validated interactions. (d) 3D plot displaying the relationships among scaled importance scores, genomic distances, and Hi-C counts.

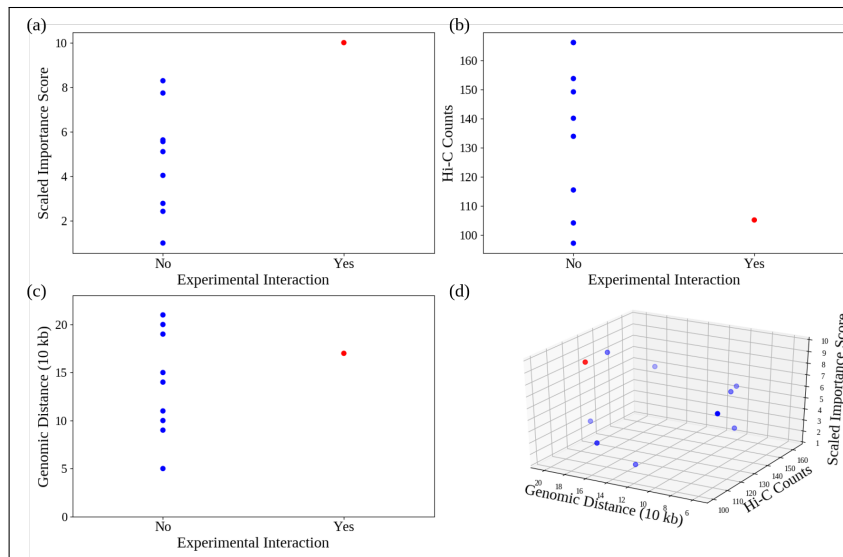


Figure B.6: **Relationships among scaled importance scores, genomic distances, and Hi-C counts for all TOP2B neighbors.** Nodes corresponding to experimentally validated interacting fragments are denoted in red and all others are denoted in blue. (a) Scaled importance score versus experimental interaction. Experimentally validated interacting fragments are ranked higher on average than non-interacting fragments. (b) Hi-C counts versus experimental interaction. Hi-C counts by themselves are not sufficient to explain the presence of experimentally validated interactions. (c) Genomic distance versus experimental interaction. Genomic distance does not correlate with experimentally validated interactions. (d) 3D plot displaying the relationships among scaled importance scores, genomic distances, and Hi-C counts.

C. XL-MERGE Model Details and Results

Hyperparameter	Values
Number of output channels	10, 20, 30
Stride length	1, 2, 5
Kernel size	5, 10, 20
Maxpool size	5, 10

Table C.1: **Hyperparameter combinations used for tuning in grid search for long-range node neighbor convolutional mechanism.** A grid search was conducted by varying the following hyperparameters for the neighboring nodes convolution operation: number of output channels, stride length, kernel size, maxpool size. Some of these combinations were not able to be tested due to memory constraints.

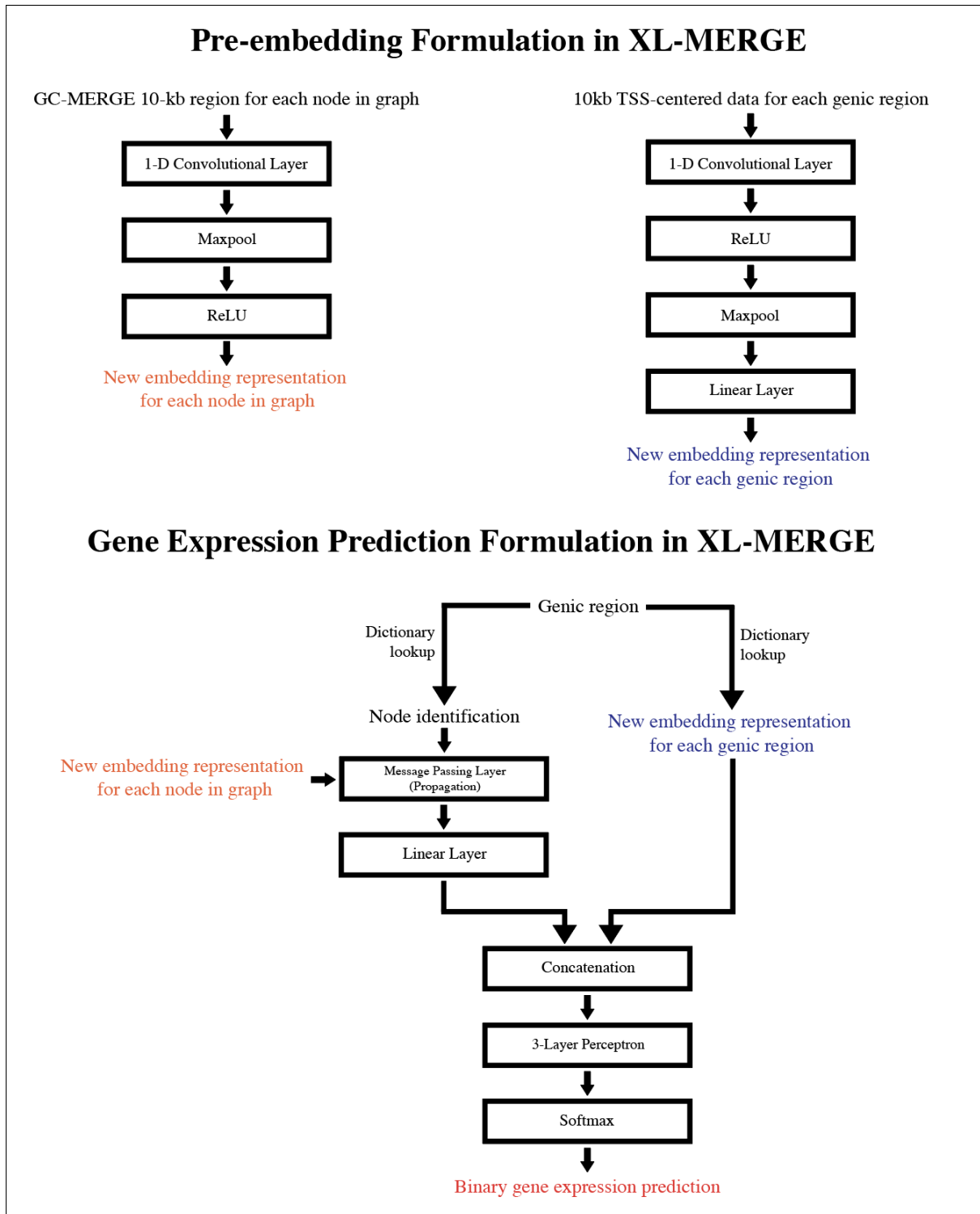


Figure C.1: **Schematic of the architecture of XL-MERGE.** XL-MERGE is similar in architecture to GC-MERGE, but its main differences have to do with the formation of pre-embeddings to better represent genic regions and the long-range interacting nodes of the graph. Convolution followed by maxpooling, nonlinear activation, and a linear layer is applied to the TSS-centered genic region data to create a better representation. For the long-range interacting regions in the graph, convolution followed by nonlinear activation and maxpooling is used to extract more positionally-agnostic, less noisy data from various nodes of the graph.

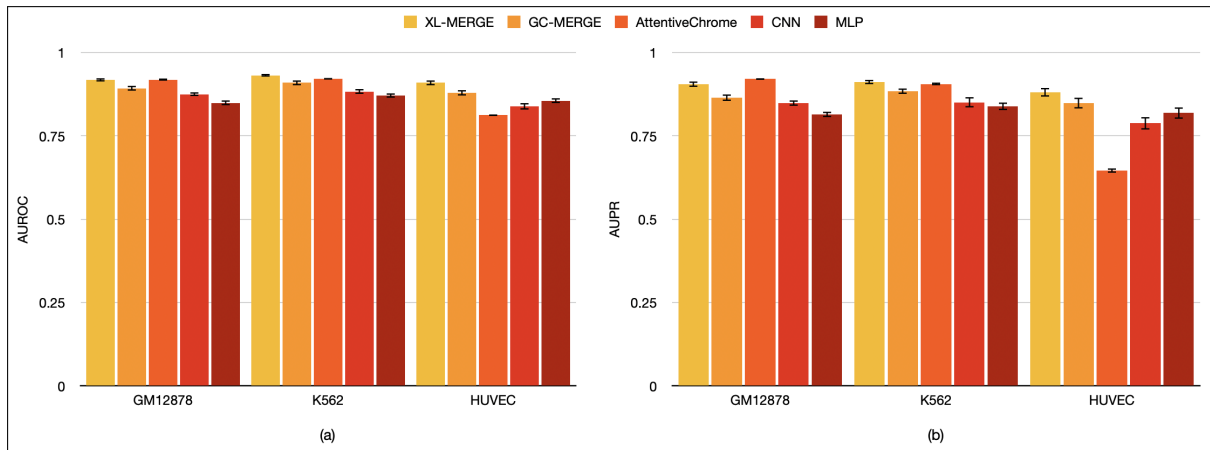


Figure C.2: **Comparison of AUROC and AUPR scores for XL-MERGE and its associated baselines.** XL-MERGE gives evidence that it is a state-of-the-art model for binarily predicting gene expression from histone marks. It also comfortably outperforms GC-MERGE both in terms of AUROC and AUPR. (a) The AUROC metrics for GM12878, K562, and HUVEC were 0.919, 0.932, and 0.910, respectively, which either were tied or were higher than all other baselines run. Additionally, across all three cell lines, XL-MERGE outperforms GC-MERGE by more than 0.02 in terms of AUROC. (b) Using the AUPR metric, XL-MERGE obtains scores of 0.905, 0.912, and 0.881 for GM12878, K562, and HUVEC, respectively. These scores outperform all the other baselines for each cell line, except for AttentiveChrome for GM12878, which scored 0.921. Also, across all three cell lines, XL-MERGE outperforms GC-MERGE by at least 0.028 in terms of AUPR.

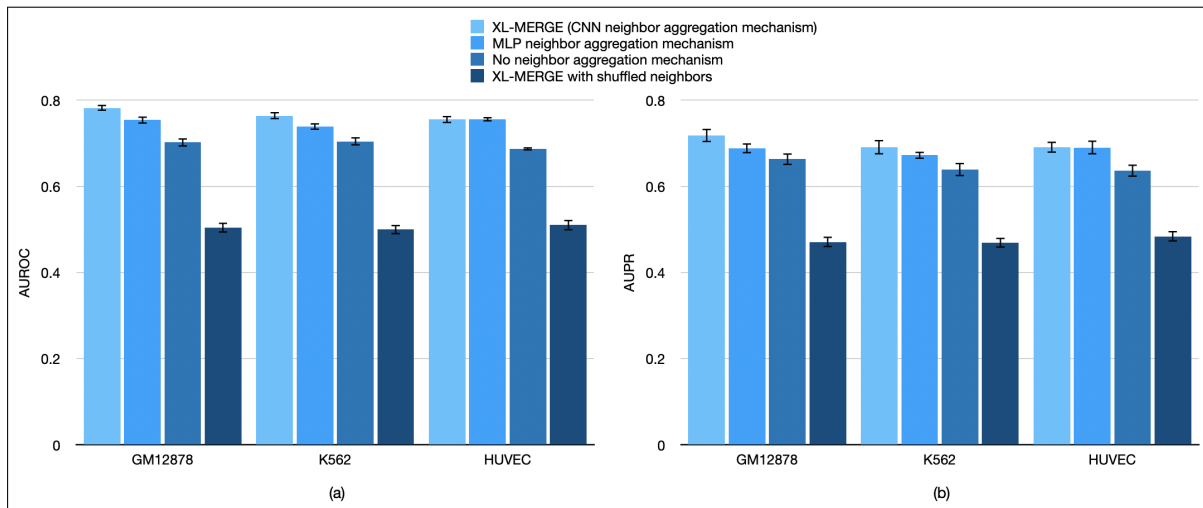


Figure C.3: **Comparison of AUROC and AUPR scores for XL-MERGE with ablated genic node embeddings and associated baselines.** The convolutional mechanism with maxpooling for extracting information from long-range interacting regions seems to be superior than other baselines. Additionally, when long-range interacting neighboring nodes are shuffled, it results in a severe drop-off in importance, strongly suggesting that XL-MERGE successfully integrates long-range interaction information to drive superior gene expression predictions. (a) XL-MERGE with genic node embeddings ablated achieved AUROC scores of 0.782, 0.764, and 0.755 for GM12878, K562, and HUVEC, respectively, and these scores outperform all other respective baselines. Most notably, when the neighbors of genic nodes are shuffled, it results in AUROC scores from 0.5 to 0.51 on average, suggesting prediction little better than random guessing. Thus, our graph construction with significant Hi-C interactions selected for is important for driving gene expression prediction. (b) XL-MERGE with genic node embeddings ablated achieved AUPR scores of 0.718, 0.691, and 0.691 for GM12878, K562, and HUVEC, respectively, and these scores were superior to all other baselines.

Type of Model Being Run	Average Testing Accuracy	Std. Dev. of Accuracy	Average Testing F1	Std. Dev. of F1	Average Testing AUROC	Std. Dev. of AUROC	Average Testing AUPR	Std. Dev. of AUPR
E116 (XL-MERGE)	0.862	0.008	0.852	0.01	0.919	0.005	0.905	0.008
E123 (XL-MERGE)	0.873	0.006	0.864	0.007	0.932	0.004	0.912	0.006
E122 (XL-MERGE)	0.845	0.009	0.839	0.01	0.91	0.007	0.881	0.013
E116 (GC-MERGE)	0.829	0.0096	0.819	0.0109	0.893	0.0068	0.865	0.0092
E123 (GC-MERGE)	0.846	0.0068	0.834	0.0088	0.91	0.0074	0.894	0.0084
E122 (GC-MERGE)	0.81	0.0078	0.802	0.0092	0.88	0.0083	0.848	0.0162
E116 (AttentiveChrome baseline with GC-MERGE processing and data split)	0.875	0.0018	0.873	0.0019	0.919	0.0011	0.921	0.0018
E123 (AttentiveChrome baseline with GC-MERGE processing and data split)	0.88	0.0037	0.876	0.0055	0.922	0.0016	0.906	0.0045
E122 (AttentiveChrome baseline with GC-MERGE processing and data split)	0.731	0.01	0.604	0.0369	0.813	0.0024	0.646	0.006
E116 (CNN baseline run with GC-MERGE processing and data split)	0.828	0.0064	0.816	0.0065	0.875	0.0059	0.849	0.0077
E123 (CNN baseline run with GC-MERGE processing and data split)	0.831	0.0066	0.818	0.0084	0.883	0.0071	0.851	0.0151
E122 (CNN baseline run with GC-MERGE processing and data split)	0.785	0.0093	0.775	0.0117	0.839	0.0091	0.788	0.0186
E116 (MLP baseline with GC-MERGE processing and data split for coarse-grained)	0.784	0.0081	0.772	0.0071	0.849	0.0071	0.815	0.0086
E123 (MLP baseline with GC-MERGE processing and data split for coarse-grained)	0.809	0.0077	0.796	0.0074	0.872	0.0067	0.839	0.0107
E122 (MLP baseline with GC-MERGE processing and data split for coarse-grained)	0.791	0.0043	0.783	0.0063	0.856	0.0072	0.819	0.017
E116 (node-ablated model)	0.714	0.007	0.72	0.012	0.782	0.007	0.718	0.015
E123 (node-ablated model)	0.694	0.006	0.705	0.009	0.764	0.008	0.691	0.017
E122 (node-ablated model)	0.691	0.009	0.701	0.013	0.755	0.009	0.691	0.013
E116 (node-ablated model with MLP pre-aggregation)	0.694	0.0072	0.713	0.009	0.754	0.0083	0.688	0.0121
E123 (node-ablated model with MLP pre-aggregation)	0.676	0.0056	0.696	0.0072	0.739	0.0076	0.672	0.0079
E122 (node-ablated model with MLP pre-aggregation)	0.686	0.0058	0.699	0.007	0.756	0.0054	0.69	0.0158
E116 (node-ablated model with no pre-aggregation mechanism)	0.644	0.0079	0.646	0.012	0.702	0.0094	0.663	0.0133
E123 (node-ablated model with no pre-aggregation mechanism)	0.646	0.0059	0.65	0.0083	0.704	0.0098	0.639	0.016
E122 (node-ablated model with no pre-aggregation mechanism)	0.638	0.005	0.616	0.009	0.687	0.0044	0.636	0.0145
E116 node-ablated with shuffled intrachromosomal neighbors	0.503	0.01	0.479	0.057	0.504	0.012	0.471	0.012
E123 node-ablated with shuffled intrachromosomal neighbors	0.506	0.012	0.43	0.037	0.5	0.011	0.469	0.012
E122 node-ablated with shuffled intrachromosomal neighbors	0.512	0.012	0.42	0.064	0.51	0.012	0.484	0.012

Figure C.4: Comparison of AUROC and AUPR scores for XL-MERGE with ablated genic node embeddings and associated baselines. For reference, here is a table of all the performance metric evaluations for Figures C.1 and C.2. Remember that each row represents a different kind of model modification or baseline, and that for each model modification or baseline 10 runs were performed. The numbers listed are the average metric values obtained over those 10 runs.