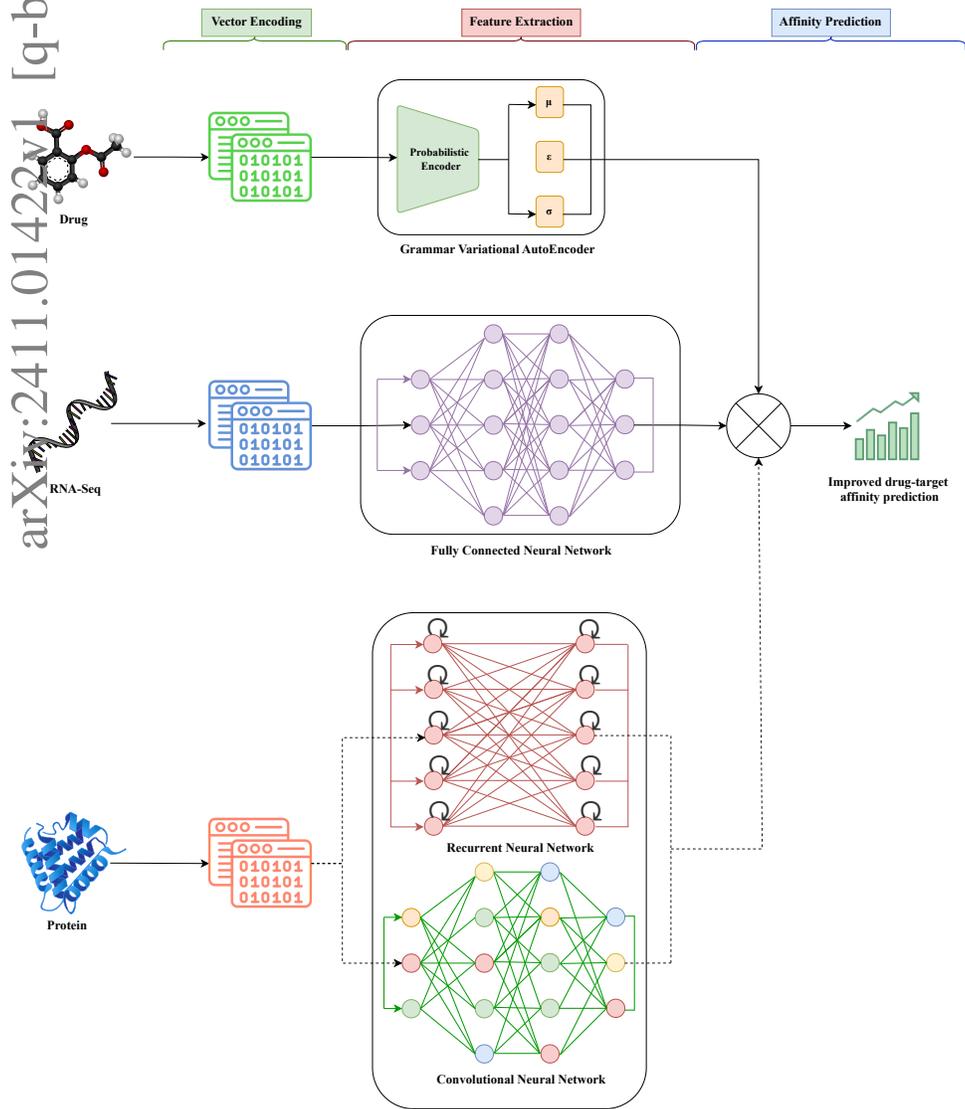


Graphical Abstract

GramSeq-DTA: A grammar-based drug-target affinity prediction approach fusing gene expression information

Kusal Debnath, Pratip Rana, Preetam Ghosh



Highlights

GramSeq-DTA: A grammar-based drug-target affinity prediction approach fusing gene expression information

Kusal Debnath, Pratip Rana, Preetam Ghosh

- GramSeq-DTA combines structural and functional representations of drugs and targets, integrating chemical perturbation data with structural information to enhance drug-target affinity (DTA) prediction accuracy.
- The model uses a Grammar Variational Autoencoder (GVAE) for drug feature extraction besides two different approaches for protein feature extraction, namely Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN).
- Feature extraction from chemical perturbation information is performed using a Fully Connected Neural Network (FCNN).
- When validated on widely used DTA datasets (BindingDB, Davis, and KIBA), the proposed approach outperforms current state-of-the-art models by incorporating both genetic and structural data.

GramSeq-DTA: A grammar-based drug-target affinity prediction approach fusing gene expression information

Kusal Debnath^a, Pratip Rana^b, Preetam Ghosh^{a,*}

^a*Department of Computer Science, Virginia Commonwealth University, , Richmond, 23284, Virginia, USA*

^b*Department of Computer Science, Old Dominion University, , Norfolk, 23529, Virginia, USA*

Abstract

Drug-target affinity (DTA) prediction is a critical aspect of drug discovery. The meaningful representation of drugs and targets is crucial for accurate prediction. Using 1D string-based representations for drugs and targets is a common approach that has demonstrated good results in drug-target affinity prediction. However, these approach lacks information on the relative position of the atoms and bonds. To address this limitation, graph-based representations have been used to some extent. However, solely considering the structural aspect of drugs and targets may be insufficient for accurate DTA prediction. Integrating the functional aspect of these drugs at the genetic level can enhance the prediction capability of the models. To fill this gap, we propose GramSeq-DTA, which integrates chemical perturbation information with the structural information of drugs and targets. We applied a Grammar Variational Autoencoder (GVAE) for drug feature extraction and utilized two different approaches for protein feature extraction: Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN). The chemical perturbation data is obtained from the L1000 project, which provides information on the upregulation and downregulation of genes caused by selected drugs. This chemical perturbation information is processed, and a compact dataset is prepared, serving as the functional feature set of the drugs. By integrating the drug, gene, and target features in the model, our approach outperforms the current state-of-the-art DTA prediction models when validated on widely used DTA datasets (BindingDB, Davis, and KIBA). This work provides a novel and practical approach to DTA prediction by merging the structural and functional aspects of

*Corresponding author

Email address: pghosh@vcu.edu (Preetam Ghosh)

biological entities, and it encourages further research in multi-modal DTA prediction.

Keywords: drug-target affinity, deep learning, grammar-based encoding, chemical perturbation, multi-modal.

1. Introduction

Drug-target affinity (DTA) prediction provides a foundation for modern drug discovery, bringing various benefits to improve efficiency, reduce costs, and increase success rates. The significance of DTA prediction is well-discussed in current literature, emphasizing its role in accelerating the identification of potential drug candidates and minimizing the risk of failure during clinical trials [1, 2]. Recent improvements in computational methods [3, 4, 5, 6] and the availability of relevant data enhance the accuracy and reliability of DTA predictions[7, 8], aiding in the design of efficient therapeutic strategies and effective treatments for many diseases.

To achieve accurate drug-target affinity (DTA) predictions, the way in which both drugs and targets are represented is a critical determinant of the performance of the models. Proper encoding of these molecular entities is essential for capturing the intricate relationships between their structural and functional properties. Early computational approaches to DTA often relied on simplified representations, such as molecular fingerprints for drugs and amino acid sequences for proteins, to feed machine learning models. Although these methods showed some promise, their inability to fully capture the complexity of molecular interactions limited the predictive power of DTA models[9, 10, 11]. As the field evolved, researchers began exploring more sophisticated techniques to better model the structural and chemical properties of drugs and targets, leading to the development of advanced representation methods that significantly improved the predictive accuracy and generalizability of DTA models[12, 13, 14].

Ozturk et al.[15] proposed DeepDTA, where they utilized 1D convolutional neural networks (CNNs) to extract high-level representations of protein sequences and 1D SMILES representations of the compounds. Before this approach, most computational methods treated drug-target affinity prediction as a binary classification problem. DeepDTA redefined the problem as a continuum of binding strength values, providing a broader view of drug-target interactions.

Nguyen et al.[16] advanced the field by representing drugs as graphs instead of linear strings and utilized graph neural networks (GNNs) to predict drug-target

affinity in their proposed deep learning model called GraphDTA. This approach firmly positions the graph-based representation of drugs as a highly effective and reliable method. Building on this trend, Tran et al.[17] proposed the Deep Neural Computation (DeepNC) model, which consists of multiple graph neural network algorithms.

The rise of natural language-based methods in biomedical research has led to further innovations in DTA modeling. Qiu et al.[18] introduced G-K-BertDTA to bridge the gap between the structural and semantic information of molecules. In their approach, drugs were represented as graphs to learn their topological features, and a knowledge-based BERT model was incorporated to obtain the semantic embeddings of the structures, thereby enhancing the feature information.

Nevertheless, there are some limitations to the above-mentioned approaches. Firstly, the information on the relative positions of the constituent atoms and bonds is often missing in the drug encoding approaches adopted in these models. In addition, the functional aspects of those drugs, which can provide relevant insights into their interaction with targets, were also not incorporated.

To address these limitations, we utilized the encoding approach for drugs known as grammar variational autoencoder (GVAE) proposed by Kusner et al.[19]. GVAE discusses the parse tree-based encoding of the drug entities, which allows learning from semantic properties and syntactic rules. This approach can learn a more consistent latent space in which entities with nearby representations decode to discrete similar outputs. In addition, to incorporate the functional aspect of those drugs, we integrated the chemical perturbation information from the L1000 project[20]. In the L1000 project, various chemical entities have been used as perturbagens and tested against multiple human cell lines, primarily linked to several types of cancers, to analyze their gene expression profile. We have utilized these gene expression signatures as the functional feature set for the drugs. Thus, the approach taken in this work utilizes structural and functional representation of drugs, which enhances the drug-target affinity prediction and outperforms the current state-of-the-art methods.

The paper provides a thorough overview of the background research in Section 2, along with a detailed explanation of the methodology, dataset preparation, network architecture, and evaluation metrics in Section 3. Section 4 presents the results of the proposed method on commonly used benchmark datasets and compares its performance to the current state-of-the-art DTA prediction models. Finally, Section 5 addresses the limitations of the study and explores opportunities for future research advancements.

2. Background

2.1. Grammar Variational Autocoder

Gómez-Bombarelli et al. [21] used Gated Recurrent Units (GRUs) and Deep Convolutional Neural Networks (DCNNs) to develop a generative model for molecular entities based on SMILES strings. This model has the potential to encode and decode molecular entities through a continuous latent space, which aids in the exploration of novel molecules with desirable properties in this space. Nevertheless, one major drawback in using string-based representation for molecular entities is their fragility, i.e., minute alteration in the string can lead to complete deviation from the original molecule, even corresponding to the generation of entirely invalid entities. James et al. [22] first proposed the concept of constructing grammar for chemical structures. According to this work, every valid discrete entity can be represented as a parse tree from the given grammar. The advantage of generating parse trees compared to texts is that it ensures the complete validity of the generated entities based on grammar. Thus, Kusner et al. [19] proposed the grammar variational autoencoder (GVAE), which encodes and decodes directly from these parse trees. This approach allows GVAE to learn from syntactic rules as well as to learn semantic properties. Along with its ability to generate valid outputs, this approach can also learn a more coherent latent space in which entities with nearby representations decode to discrete similar outputs.

2.1.1. Context-free Grammar

A context-free grammar (CFG) is conventionally defined as a 4-tuple $G = (V, \Sigma, R, S)$, where V represents a finite set of non-terminal symbols; Σ represents a finite set of terminal symbols, disjoint from V ; R represents a finite set of production rules; S is a unique non-terminal referred to as the start symbol; G represents the grammar that describes a set of trees that can be formed by applying rules in R to leaf nodes until all leaf nodes become terminal symbols in Σ .

The rules R are technically defined as $\alpha \rightarrow \beta$ for $\alpha \in V$ and $\beta \in (V \cup \Sigma)^*$, $*$ denoting the Kleene closure. Practically, these rules are portrayed as a collection of mappings from a solitary non-terminal on the left-hand side in V to a sequence of symbols that can be either terminal or non-terminal by definition. These mappings can be seen as a rule for rewriting.

When a production rule is applied to a non-terminal symbol, it creates a tree structure where the symbols on the right-hand side of the production rule become child nodes for the left-hand side parent. These trees extend from each

non-terminal symbol in V . The language of G is the set of all sequences of terminal symbols that can be generated by traversing the leaf nodes of a tree from left to right. A parse tree is a tree with its root at S and a sequence of terminal symbols as its leaf nodes. The prevalence of context-free languages in computer science is attributed, in part, to the existence of practical parsing algorithms.

2.1.2. Syntactic vs. Semantic Validity

A crucial aspect of grammar-based encoding is that the encoded molecules are syntactically valid, but the semantic validity of these molecules is a matter of discussion. There are some reasons for this phenomenon - a) Some molecules produced by the grammar may be unstable or chemically invalid; for example, a carbon atom cannot make bonds with more than four atoms in a molecule as it has a valency of 4. Nevertheless, the grammar can produce this kind of molecule; b) Assignment of ring-bond digits in SMILES is a non-context-free process. It needs to keep track of the order in which rings are encountered, along with the connectivity between the rings, which can only be determined from the local context of the string. For example, in naphthalene (c1ccc2c(c1)ccc2), the outer ring uses the digit '1', and the inner ring uses '2'. The digits are not nested but rather follow a specific order; c) GVAE can still produce an undetermined sequence if there are existing non-terminal symbols on the stack after processing all logit vectors.

2.2. L1000 Assay

The L1000 project [20] is part of the Library of Integrated Network-Based Cellular Signatures (LINCS) program funded by the National Institutes of Health (NIH). This program aims to catalog and analyze cellular responses to various perturbations to understand how these perturbations modulate cellular functions. This project efficiently manages chemical perturbation data using a structured method, which encircles data generation, processing, and analysis. This project uses various chemical compounds, including FDA-approved drugs, experimental drugs, natural compounds, and other bioactive molecules as perturbagens. Perturbagens selection mainly involves possible connections of these chemicals to numerous biological pathways and disease cross-talks. Various human cell lines, majorly associated with several types of cancers, are chosen to ensure diversity in the biological responses. The L1000 data thus can be used to identify potential new uses for existing drugs or to discover new candidate drugs. Moreover, novel hypotheses can be made that correspond to the possible effects of the new compounds by performing comparative analyses of the gene expression signatures of known drugs. The data also aids in understanding the underlying molecular mechanism

of the diseases by showcasing the effect of different compounds in the alteration of gene expression related to disease pathways.

3. Methodology

3.1. Datasets

3.1.1. Benchmark Datasets

In this study, three datasets are used for the benchmarking purposes: BindingDB [23], Davis [24], and KIBA [25]. The drug-target affinity dataset in the BindingDB database contains experimental binding affinities between small molecules and protein targets and supplementary information on the entities (e.g., ID, Structure, etc.). In the Davis dataset, the targets are kinase proteins, and the drugs are the small molecules (inhibitors) targeting those kinases. Similar to Davis, the KIBA dataset also focuses on kinase proteins and their corresponding inhibitor drugs, but it contains a more significant number of instances than Davis.

In these datasets, the drugs are represented as SMILES strings in these datasets, and the target proteins are represented as amino acid sequences. For BindingDB and Davis datasets, the labels are the K_d (Dissociation constant) values, which indicate the extent of the interactions between each drug-target pair. Meanwhile, a unified KIBA score is used as a label for the KIBA dataset, combining K_d , K_i (Inhibition Constant), and IC_{50} (Half Maximal Inhibitory Coefficient) values for corresponding drug-target pairs. The labels are converted into logarithmic form which helps improve the performance of the model in regression tasks.

3.1.2. L1000 Chemical Perturbation Dataset

The chemical perturbation data available in the L1000 project is documented in raw format. Therefore, the data needs to be processed accordingly for use. The detailed process of preparing the L1000 dataset from the raw data is discussed below:

① The L1000 chemical perturbation data file is loaded where each perturbagen has multiple replicates based on dosage concentration, and each replicate has two lists of associated genes - one for upregulated and the other for downregulated. ② The analysis of the dosage concentration distribution among the replicates shows that samples with a concentration of 10 μM are the most common. Therefore, for standardization purposes, samples with a concentration of 10 μM are selected for further analysis, while the others are excluded. ③ Each unique perturbagen is

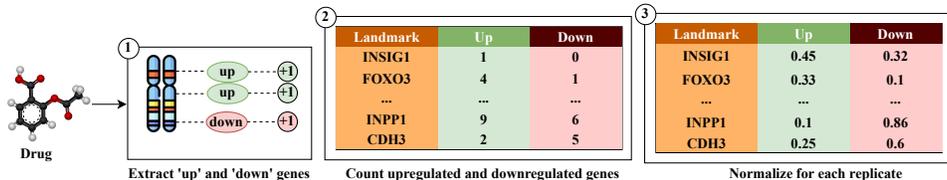


Figure 1: **Preparation of the gene expression dataset.** Gene expression information analyzed on 978 landmark genes for the selected drugs is extracted from the L1000 chemical perturbation data. After considering all the biological replicates of the perturbation analysis, a gene regulation matrix is created for both upregulated and downregulated genes.

mapped into corresponding SMILES representation, which is important for downstream molecular modeling. ④ For each perturbagen p_i , the gene regulatory information is represented as a vector of ‘up’ and ‘down’ regulation values across 978 landmark genes, and the number of times a gene is upregulated or downregulated is counted and normalized by the number of replicates. Let x_{ij}^{up} represent whether gene j is upregulated for perturbagen p_i , and x_{ij}^{down} represent the same for downregulation. The final regulatory vector for each perturbagen is computed as:

$$\mathbf{v}_i = \frac{1}{\text{count}(p_i)} \left(\sum_{j=1}^m x_{ij}^{up}, \sum_{j=1}^m x_{ij}^{down} \right) \quad (1)$$

where $m = 978$ is the number of landmark genes, and $\text{count}(p_i)$ is the number of times perturbagen p_i appears in the dataset. A representative illustration is shown in Figure 1. ⑤ Finally, the vector for each perturbagen is stacked to get the final matrix:

$$\mathbf{V} = \{\mathbf{v}_i\}_{i=1}^k, \quad \mathbf{v}_i \in \mathbb{R}^{978 \times 2} \quad (2)$$

where \mathbf{V} is of shape $k \times 978 \times 2$, and k is the number of unique perturbagens. The processed dataset is an important contribution of this work that can aid future research in RNA-Seq data integration and analysis.

Not all the perturbagens mentioned in the L1000 dataset are entirely present in the benchmark datasets. Therefore, the datasets are processed accordingly, and only those drugs whose corresponding regulatory vector is present in the L1000 dataset are selected. The processing of the datasets resulted in a decrease in the number of total interactions. The summary of all the original and processed benchmark datasets is discussed in Table 1.

Dataset	Compounds	Proteins	Interactions
Original			
BindingDB	22,381	1,860	91,751
Davis	68	379	30,056
KIBA	2,068	229	118,254
Processed			
BindingDB	444	754	18,567
Davis	11	379	4,169
KIBA	12	194	538

Table 1: Dataset statistics - number of compounds, proteins and interactions

3.2. Network Architecture

The complete network consists of 3 main parts - a) Drug encoder, b) RNA-Seq encoder, and c) Protein encoder. Figure 2 shows the schematic diagram of the complete network architecture.

3.2.1. Drug Encoder

For this work, we utilized a pre-trained GVAE model from the study conducted by Zhu et al.[26] that focuses on deep learning-based drug efficacy prediction from transcriptional profiles. One-hot encoded vectors are generated by parsing the SMILES representations of the drugs using a grammar-tree-based approach and then passed into the encoder network. The detailed process of parsing the SMILES and generating one-hot encoded vectors is discussed below:

① SMILES representations are converted into a collection of tokens using a tokenizer. ② The tokenized sequence is then parsed using a grammar adopted from the work of Kusner et al. [19]. This yields a sequence of production rules:

$$G(\tau(S)) = P = \{P_1, P_2, \dots, P_q\} \quad (3)$$

where G is the grammar, $\tau(S)$ is the tokenized sequence and $P = \{P_1, P_2, \dots, P_q\}$ is the sequence of production rules. Each production rule is then mapped to an index I_i in a predefined list of rules. ③ A zero matrix is initialized, denoting the vector to be populated by one-hot encoding:

$$O_{j,I_j} = 1, \quad \forall j \in \{1, 2, \dots, \min(M, q)\}, \quad I_j \in \{1, 2, \dots, N-1\} \quad (4)$$

where O is the encoded one-hot vector of shape $M \times N$, M is the maximum length of sequences, N is the total number of production rules, and q is the number

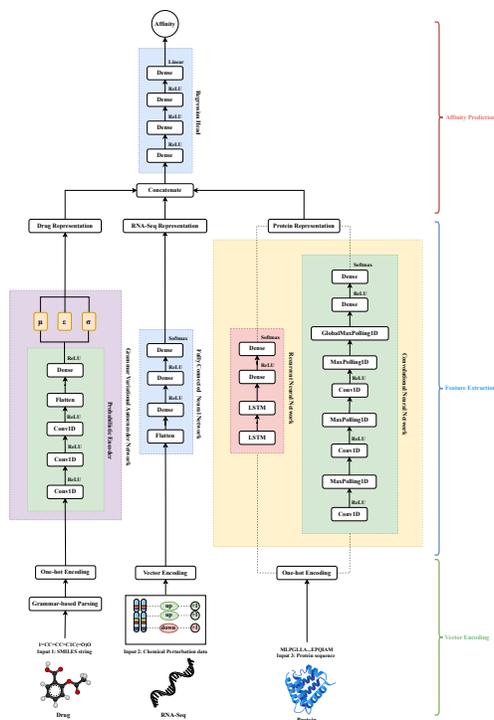


Figure 2: **Network architecture of the proposed model.** The encoded drug information is passed through an GVAE layer, the RNA-Seq information is passed through an FCNN, while the encoded protein information is passed through a series of LSTM layers and 1D CNN layers. Learned representations are concatenated and passed through a FCNN acting as a regression head to predict the affinity.

of productions. If q is smaller than M , the rest of the matrix is padded with an indicator for "end of sequence":

$$O_{j,N-1} = 1, \quad \forall j \in \{q+1, q+2, \dots, M\} \quad (5)$$

In this work, the values of M and N are 277 and 76, respectively. The generated one-hot vectors for each SMILES representations are then passed into the encoder network. The schematic diagram of the overall process is given in Figure 3.

3.2.2. RNA-Seq Encoder

A fully connected neural network (FCNN) is used for the extraction of meaningful features from the high-dimensional RNA-Seq data. As discussed in the

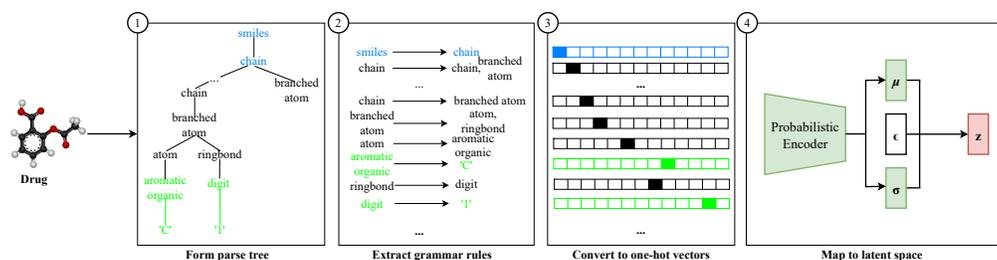


Figure 3: **Encoding of drug SMILES structures.** A parse tree is constructed based on the structural components of SMILES representations. Grammar rules are extracted from the parsed trees. SMILES representations are then converted into one-hot vectors. Finally, the one-hot vectors are transformed into corresponding latent space representations using an encoder network.

L1000 chemical perturbation dataset preparation, the resulting dataset is of shape $k \times 978 \times 2$, where k is the number of unique perturbagens, 978 is the number of landmark genes and 2 indicates the number of columns representing upregulated and downregulated genes. When these vectors are passed through a dense neural network, it learns a condensed and abstract representation of how each perturbagen affects gene expression.

3.2.3. Protein Encoder

Feature extraction from amino acid sequences is achieved using two different types of neural networks: Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN). CNNs are able to identify local motifs and patterns within a sequence by using sliding windows of filters to capture neighboring amino acids. On the other hand, RNNs are capable of extracting long-range dependencies and sequential relationships between amino acids by retaining information from previous positions in the sequence in a step-by-step manner. To encode the sequences, a dictionary of all possible amino acid sequences in the proteins is created. One-hot encoding of a given protein is carried out based on the presence of a particular amino acid in that protein. For standardization, the maximum length of a protein is limited to 1000 sequences. Encoding of all the proteins results in the creation of a vector of $p \times 26 \times 1000$, where p is the number of unique proteins and 26 is the length of the amino acid dictionary.

3.3. Training Settings

The training process is set to run for 500 epochs with an adaptive learning rate that starts at 0.001 using the Adam optimizer. The higher learning rate value is

Parameters	Value
Drug Encoding	
GVAE Encoder Filter Sizes	9, 9, 10
GVAE Encoder Kernel Sizes	9, 9, 11
GVAE Latent Space Dimension	56
RNA-Seq Encoding	
Dense Layers	2
Protein Encoding	
CNN Filter Sizes	32, 64, 96
CNN Kernel Sizes	4, 8, 12
RNN LSTM Layers	2
Regression Head	
Dense Layers	3
Training	
Epochs	500
Learning Rate	0.001
Batch Size	256
Optimizer	Adam

Table 2: Summary of Network Architecture and Training Hyperparameters

chosen to ensure that the training process does not significantly impact the pre-trained weights in the GVAE model. A batch size of 256 has been found to yield the best results, maintaining a balance between memory usage and convergence speed. The summary of the overall network architecture and training hyperparameters are discussed in Table 2.

3.4. Evaluation Metrics

Evaluating deep learning models involves various metrics that capture different aspects of performance. Mean Squared Error (MSE) measures the average squared difference between actual and predicted values, highlighting prediction accuracy. We also used the Concordance Index (C-Index), which is a preferred metric for survival analysis, to evaluate the consistency between predicted risk scores and actual outcomes. Together, these metrics provide a robust framework for model evaluation.

3.4.1. Mean Squared Error (MSE)

Mean Squared Error (MSE) is a common loss function used for regression tasks. It measures the average of the squares of the errors, which are the differences between the predicted and actual values. Mathematically, it is defined as:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (6)$$

where y_i is the actual value, \hat{y}_i is the predicted value, and n is the number of data points.

3.4.2. Concordance Index (C-Index)

The Concordance Index (C-Index) is a metric used primarily in survival analysis to evaluate the predictive accuracy of risk scores. It assesses the degree of concordance between the predicted and actual ordering of event times. The C-Index is calculated as:

$$C = \frac{\text{Number of concordant pairs}}{\text{Number of possible evaluation pairs}} \quad (7)$$

A pair is considered concordant if the predicted and actual orderings of two instances are consistent.

4. Results & Discussions

In this section, we will discuss the performance of the GramSeq-DTA model in detail. The discussion can be divided into the following parts: a) Benchmarking the performance of GramSeq-DTA with integrated RNA-Seq information against the baseline models, b) Advantage of integrating RNA-Seq information to the model, and c) Performance comparison of the proposed model on original and processed datasets

4.1. Benchmarking against baseline models

In order to validate our findings, we conducted a comprehensive performance comparison of GramSeq-DTA, which now includes integrated RNA-Seq information. We compared it against several well-established baseline models: DeepDTA, GraphDTA, DeepNC, and G-K-BertDTA. We used MSE and CI values to assess performance. The model performance is compared across three benchmark datasets - BindingDB, Davis, and KIBA. DeepDTA employs a deep learning framework to capture the complex features of drug-target interactions using

Model	Drug Encoder	RNA-Seq Encoder	Protein Encoder	MSE	CI
DeepDTA	CNN	-	CNN	0.384	0.821
GraphDTA	GINConvNet	-	CNN	0.355	0.819
	GCNNet	-	CNN	0.397	0.786
	GATNet	-	CNN	0.512	0.757
	GAT_GCN	-	CNN	0.384	0.806
DeepNC	GENConv	-	CNN	0.367	0.828
G-K-BertDTA	GINConvNet + Embeddings	-	CNN	0.325	0.832
GramSeq-DTA	GVAE	FCNN	CNN	0.365	0.843
	GVAE	FCNN	RNN	0.355	0.832

Table 3: Performance comparison of GramSeq-DTA with baseline models on the processed BindingDB dataset.

convolutional neural networks. GraphDTA, on the other hand, leverages graph neural networks to represent molecular structures as graphs, enabling it to better capture the topological properties of molecules. DeepNC uses a neural collaborative filtering approach to model interactions, focusing on latent feature extraction. Lastly, G-K-BertDTA integrates graph-based representations with BERT-like architectures to enhance contextual understanding of molecular relationships. Our extensive evaluation, conducted on processed benchmark datasets, showed that GramSeq-DTA consistently outperformed its counterparts in terms of Concordance Index (CI) values, a widely accepted metric for evaluating predictive performance in drug-target affinity modeling. In the BindingDB dataset, GramSeq-DTA outperforms G-K-BertDTA by 1.32% in terms of CI value. Similarly, on the Davis dataset, GramSeq-DTA shows a 0.89% advantage over DeepNC. On the KIBA dataset, GramSeq-DTA exceeds G-K-BertDTA by 2.75%. Importantly, the integration of RNA-Seq data into GramSeq-DTA provided valuable insights into gene expression patterns, contributing to the improved accuracy of the models in predicting drug-target interactions. Detailed results of this comparison can be found in Tables 3, 4, and 5, where the enhanced GramSeq-DTA model demonstrates its robust performance, setting a new standard in the field.

4.2. Advantage of integrating RNA-Seq information

Table 6 indicates that when validated on the processed BindingDB dataset, GramSeq-DTA performs better, with a CI value of 0.843 when integrating RNA-

Model	Drug Encoder	RNA-Seq Encoder	Protein Encoder	MSE	CI
DeepDTA	CNN	-	CNN	0.219	0.779
GraphDTA	GINConvNet	-	CNN	0.187	0.771
	GCNNet	-	CNN	0.214	0.732
	GATNet	-	CNN	0.241	0.713
	GAT_GCN	-	CNN	0.227	0.753
DeepNC	GENConv	-	CNN	0.198	0.789
G-K-BertDTA	GINConvNet + Embeddings	-	CNN	0.169	0.778
GramSeq-DTA	GVAE	FCNN	CNN	0.293	0.796
	GVAE	FCNN	RNN	0.261	0.796

Table 4: Performance comparison of GramSeq-DTA with baseline models on the processed Davis dataset.

Seq information compared to not integrating RNA-Seq information. Similar results are evident in Table 7 and Table 8, where validation is performed on the Davis and KIBA datasets, respectively. CI values of 0.796 and 0.708 are observed for the processed Davis and KIBA datasets, respectively, when RNA-Seq information is integrated. These observations prove that integrating RNA-Seq information with corresponding drug and target structural information can enhance the drug-target affinity prediction ability of the model.

4.3. Performance on original and processed data

Table 9 presents the comparative evaluation of the performance of the proposed model on the original benchmark datasets and the processed benchmark datasets. The results of the model integration with RNA-Seq information are shown for the processed datasets. As shown in Table 1, the number of interactions between the original and the processed datasets differs. Despite losing approximately 80% of data in processing, our model performs better on the processed BindingDB dataset, with a best CI value of 0.843, while for the original BindingDB dataset, the optimum CI value was 0.818. The results on the original and processed Davis dataset are also competitive. During processing the Davis dataset, we lost around 84% of data. Our model indicates a CI value of 0.809 for the original Davis dataset, while for the processed Davis dataset, the CI value is 0.796. Data loss during the processing of the KIBA dataset is 99.5%, which is the highest value among all three datasets. For the KIBA dataset, there is a

Model	Drug Encoder	RNA-Seq Encoder	Protein Encoder	MSE	CI	
DeepDTA	CNN	-	CNN	0.877	0.609	
GraphDTA	GINConvNet	-	CNN	1.061	0.628	
	GCNNNet	-	CNN	0.903	0.631	
	GATNet	-	CNN	0.957	0.609	
	GAT_GCN	-	CNN	0.831	0.671	
DeepNC	GENConv	-	CNN	0.769	0.648	
G-K-BertDTA	GINConvNet + Embeddings		-	CNN	0.693	0.689
GramSeq-DTA	GVAE	FCNN	CNN	0.843	0.708	
	GVAE	FCNN	RNN	1.269	0.688	

Table 5: Performance comparison of GramSeq-DTA with baseline models on the processed KIBA dataset.

Model	Drug Encoder	RNA-Seq Encoder	Protein Encoder	MSE	CI
GramSeq-DTA	GVAE	-	CNN	0.495	0.746
	GVAE	-	RNN	0.495	0.754
	GVAE	FCNN	CNN	0.365	0.843
	GVAE	FCNN	RNN	0.355	0.832

Table 6: Performance comparison of GramSeq-DTA with and without RNA-Seq information integration on the processed BindingDB dataset.

significant difference in CI values (original: 0.823, processed: 0.708) for original and processed datasets. The performance difference for the KIBA dataset can be caused by excessive data loss while processing the original dataset. Based on the performance of the other two process datasets (BindingDB and Davis), with more data available for the processed KIBA dataset, there is a high possibility of getting better performance.

5. Conclusion & Future Directions

In this study, we demonstrated that incorporating chemical perturbation information can enhance drug-target affinity prediction. The core contribution of this research is in its transformation of data from a chemical perturbation assay and

Model	Drug Encoder	RNA-Seq Encoder	Protein Encoder	MSE	CI
GramSeq-DTA	GVAE	-	CNN	0.277	0.705
	GVAE	-	RNN	0.311	0.716
	GVAE	FCNN	CNN	0.293	0.796
	GVAE	FCNN	RNN	0.261	0.796

Table 7: Performance comparison of GramSeq-DTA with and without RNA-Seq information integration on the processed Davis dataset.

Model	Drug Encoder	RNA-Seq Encoder	Protein Encoder	MSE	CI
GramSeq-DTA	GVAE	-	CNN	1.011	0.653
	GVAE	-	RNN	0.876	0.618
	GVAE	FCNN	CNN	0.843	0.708
	GVAE	FCNN	RNN	1.269	0.688

Table 8: Performance comparison of GramSeq-DTA with and without RNA-Seq information integration on the processed KIBA dataset.

utilizing it as an extra modality along with drug and protein structural information. This research could guide future work on a better understanding of affinity prediction among biological entities in the absence of three-dimensional structural information of the entities. Nevertheless, a fundamental limitation of this research is that information from a chemical perturbation assay may not be available for every drug in widely used drug-target affinity benchmark datasets. Therefore, having more data from chemical perturbation assays for additional drugs can further enhance the ability of deep learning models to predict affinity. Future studies should investigate different approaches for transforming the chemical perturbation information. Moreover, introducing advanced feature extraction methods from the biological entities can enhance the prediction. In summary, this work underscores the importance of integrating additional data modalities in drug-target affinity prediction.

CRediT authorship contribution statement

Kusal Debnath: Data Curation, Conceptualization, Methodology, Writing - Original Draft Preparation, Writing - Review & Editing, Visualization. **Pratip Rana:**

Dataset	Model	Drug Encoder	RNA-Seq Encoder	Protein Encoder	MSE	CI
BindingDB						
Original	GramSeq-DTA	GVAE	-	CNN	1.029	0.818
		GVAE	-	RNN	1.061	0.812
Processed	GramSeq-DTA	GVAE	FCNN	CNN	0.365	0.843
		GVAE	FCNN	RNN	0.355	0.832
Davis						
Original	GramSeq-DTA	GVAE	-	CNN	0.446	0.806
		GVAE	-	RNN	0.445	0.809
Processed	GramSeq-DTA	GVAE	FCNN	CNN	0.293	0.796
		GVAE	FCNN	RNN	0.261	0.796
KIBA						
Original	GramSeq-DTA	GVAE	-	CNN	0.272	0.823
		GVAE	-	RNN	0.277	0.823
Processed	GramSeq-DTA	GVAE	FCNN	CNN	0.843	0.708
		GVAE	FCNN	RNN	1.269	0.688

Table 9: Performance GramSeq-DTA without RNA-Seq information integration on the original benchmark datasets.

Conceptualization, Formal Analysis, Investigation, Writing - Original Draft Preparation, Writing - Review & Editing. **Preetam Ghosh:** Conceptualization, Resources, Writing - Review & Editing, Project Administration, Supervision, Funding Acquisition.

Acknowledgment

This work was partially supported by 5R21MH128562-02 (PI: Roberson-Nay), 5R21AA029492-02 (PI: Roberson-Nay), CHRB-2360623 (PI: Das), NSF-2316003 (PI: Cano), VCU Quest (PI: Das) and VCU Breakthroughs (PI: Ghosh) funds awarded to P.G.

Conflict of Interest

The authors declare that they have no conflict of interest regarding the publication of this paper.

Author Biographies



Kusal Debnath is pursuing PhD in Computer Science with a strong focus on AI-driven drug discovery. He obtained his Master of Technology (MTech) degree in Biomedical Engineering from Indian Institute of Technology, Kharagpur, India. His research interests include machine learning applications in bioinformatics and computational biology, and he is passionate about advancing innovations at the intersection of technology and medicine.



Pratip Rana is an assistant professor in the Department of Computer Science at Old Dominion University. He received a Ph.D. degree in Computer Science from Virginia Commonwealth University, USA, in 2020. His research interests include Machine learning, Computational biology, Complex systems, and Modeling & simulation.



Preetam Ghosh is a Professor in the Department of Computer Science and the director of the Biological Networks Lab at Virginia Commonwealth University. He obtained his MS and Ph.D. degrees in Computer Science & Engineering from the University of Texas at Arlington and a BS in Computer Science from Jadavpur University, Kolkata India. His research interests include algorithms, stochastic modeling & simulation, network science and machine learning-related approaches in systems biology and computational epidemiology, and mobile computing-related issues in pervasive grids that have resulted in more than 200 conference and journal articles and several federally funded research projects from NSF, NIH, DoD, and US-VHA. He previously served as the Secretary/Treasurer of ACM SIGBio.

References

- [1] Kit-Kay Mak, Yi-Hang Wong, and Mallikarjuna Rao Pichika. Artificial intelligence in drug discovery and development. *Drug Discovery and Evaluation: Safety and Pharmacokinetic Assays*, pages 1–38, 2023.
- [2] Jintae Kim, Sera Park, Dongbo Min, and Wankyu Kim. Comprehensive survey of recent drug discovery using deep learning. *International Journal of Molecular Sciences*, 22(18):9983, 2021.
- [3] V. Chaitankar, P. Ghosh, E. Perkins, P. Gong, Y. Deng, and C. Zhang. A novel gene network inference algorithm using predictive minimum description length approach. *BMC Systems Biology*, 4(Suppl 1: S7), 2010.
- [4] V. Chaitankar, P. Ghosh, E. Perkins, P. Gong, and C. Zhang. Time lagged information-theoretic approaches to the reverse engineering of gene regulatory networks. *BMC Bioinformatics*, 11(Suppl 6: S19), 2010.
- [5] Joseph J. Nalluri, Debmalya Barh, Vasco Azevedo, and Preetam Ghosh. Mirsig: A consensus-based network inference methodology to identify pan-cancer mirna-mirna interaction signatures. *Scientific Reports*, 7(1), 2017.
- [6] Edian F. Franco, Pratip Rana, Aline Cruz, Víctor V. Calderón, Vasco Azevedo, Rommel T. J. Ramos, and Preetam Ghosh. Performance com-

parison of deep learning autoencoders for cancer subtype detection using multi-omics data. *Cancers*, 13(9), 2021.

- [7] Rocco Meli, Garrett M Morris, and Philip C Biggin. Scoring functions for protein-ligand binding affinity prediction using structure-based deep learning: a review. *Frontiers in bioinformatics*, 2:885983, 2022.
- [8] Seokhyun Moon, Wonho Zhung, and Woo Youn Kim. Toward generalizable structure-based deep learning models for protein–ligand interaction prediction: Challenges and strategies. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 14(1):e1705, 2024.
- [9] Hao Ding, Ichigaku Takigawa, Hiroshi Mamitsuka, and Shanfeng Zhu. Similarity-based machine learning methods for predicting drug–target interactions: a brief review. *Briefings in bioinformatics*, 15(5):734–747, 2014.
- [10] Tapio Pahikkala, Antti Airola, Sami Pietilä, Sushil Shakyawar, Agnieszka Sz wajda, Jing Tang, and Tero Aittokallio. Toward more realistic drug–target interaction predictions. *Briefings in bioinformatics*, 16(2):325–337, 2015.
- [11] Tong He, Marten Heidemeyer, Fuqiang Ban, Artem Cherkasov, and Martin Ester. Simboost: a read-across approach for predicting drug–target binding affinities using gradient boosting machines. *Journal of cheminformatics*, 9:1–14, 2017.
- [12] Xiaoqing Ru, Xiucai Ye, Tetsuya Sakurai, Quan Zou, Lei Xu, and Chen Lin. Current status and future prospects of drug–target interaction prediction. *Briefings in Functional Genomics*, 20(5):312–322, 2021.
- [13] Ming Chen, Yajian Jiang, Xiujuan Lei, Yi Pan, Chunyan Ji, Wei Jiang, and Hongkai Xiong. Drug-target interactions prediction based on signed heterogeneous graph neural networks. *Chinese Journal of Electronics*, 33(1):231–244, 2024.
- [14] Ali K Abdul Raheem and Ban N Dhannoon. Comprehensive review on drug-target interaction prediction-latest developments and overview. *Current Drug Discovery Technologies*, 21(2):56–67, 2024.
- [15] Hakime Öztürk, Arzucan Özgür, and Elif Ozkirimli. Deepdta: deep drug–target binding affinity prediction. *Bioinformatics*, 34(17):i821–i829, 2018.

- [16] Thin Nguyen, Hang Le, Thomas P Quinn, Tri Nguyen, Thuc Duy Le, and Svetha Venkatesh. Graphdta: predicting drug–target binding affinity with graph neural networks. *Bioinformatics*, 37(8):1140–1147, 2021.
- [17] Huu Ngoc Tran Tran, J Joshua Thomas, and Nurul Hashimah Ahamed Has-sain Malim. Deepnc: a framework for drug-target interaction prediction with graph neural networks. *PeerJ*, 10:e13163, 2022.
- [18] Xihe Qiu, Haoyu Wang, Xiaoyu Tan, and Zhijun Fang. Gk bertdta: a graph representation learning and semantic embedding-based framework for drug-target affinity prediction. *Computers in Biology and Medicine*, 173:108376, 2024.
- [19] Matt J. Kusner, Brooks Paige, and José Miguel Hernández-Lobato. Grammar variational autoencoder. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1945–1954. PMLR, 06–11 Aug 2017.
- [20] Aravind Subramanian, Rajiv Narayan, Steven M Corsello, David D Peck, Ted E Natoli, Xiaodong Lu, Joshua Gould, John F Davis, Andrew A Tubelli, Jacob K Asiedu, et al. A next generation connectivity map: L1000 platform and the first 1,000,000 profiles. *Cell*, 171(6):1437–1452, 2017.
- [21] Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato, Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 4(2):268–276, 2018.
- [22] OpenSMILES. The opensmiles format. <http://opensmiles.org/opensmiles.html>, 2015. Accessed: 2024-10-02.
- [23] Tiqing Liu, Yuhmei Lin, Xin Wen, Robert N Jorissen, and Michael K Gilson. Bindingdb: a web-accessible database of experimentally determined protein–ligand binding affinities. *Nucleic acids research*, 35(suppl_1):D198–D201, 2007.
- [24] Mindy I Davis, Jeremy P Hunt, Sanna Herrgard, Pietro Ciceri, Lisa M Wodicka, Gabriel Pallares, Michael Hocker, Daniel K Treiber, and Patrick P

Zarrinkar. Comprehensive analysis of kinase inhibitor selectivity. *Nature biotechnology*, 29(11):1046–1051, 2011.

- [25] Jing Tang, Agnieszka Sz wajda, Sushil Shakyawar, Tao Xu, Petteri Hintsanen, Krister Wennerberg, and Tero Aittokallio. Making sense of large-scale kinase inhibitor bioactivity data sets: a comparative and integrative analysis. *Journal of Chemical Information and Modeling*, 54(3):735–743, 2014.
- [26] Jie Zhu, Jingxiang Wang, Xin Wang, Mingjing Gao, Bingbing Guo, Miaomiao Gao, Jiarui Liu, Yanqiu Yu, Liang Wang, Weikaixin Kong, et al. Prediction of drug efficacy from transcriptional profiles with deep learning. *Nature biotechnology*, 39(11):1444–1452, 2021.