

GEFA: EARLY FUSION APPROACH IN DRUG-TARGET AFFINITY PREDICTION

**¹ Vaddeman Swetha, ² Dasari Shirisha, ³ Vuppugalla Sravanthi, ⁴ Ravula
Keerthana**

^{1,2,3} Assistant Professors, Department of Computer Science and Engineering,
Kasireddy Narayanreddy College Of Engineering And Research, Abdullapur (V),
Abdullapurmet(M), Rangareddy (D), Hyderabad - 501 505

⁴ student, Department of Computer Science and Engineering, Kasireddy Narayanreddy
College Of Engineering And Research, Abdullapur (V), Abdullapurmet(M),
Rangareddy (D), Hyderabad - 501 505

ABSTRACT

Predicting the interaction between a compound and a target is crucial for rapid drug repurposing. Deep learning has been successfully applied in drug-target affinity (DTA) problem. However, previous deep learning-based methods ignore modeling the direct interactions between drug and protein residues. This would lead to inaccurate learning of target representation which may change due to the drug binding effects. In addition, previous DTA methods learn protein representation solely based on a small number of protein sequences in DTA datasets while neglecting the use of proteins outside of the DTA datasets. We propose GEFA (Graph Early Fusion Affinity), a novel graph-in-graph neural network with attention mechanism to address the changes in target representation because of the binding effects. Specifically, a drug is modeled as a graph of atoms, which then serves as a node in a larger graph of residues-drug complex. The resulting model is an expressive deep nested graph neural network. We also use pre-trained protein representation powered by the recent effort of learning contextualized protein representation. The experiments are conducted under different settings to evaluate scenarios such as novel drugs or targets. The results demonstrate the effectiveness of the pre-trained protein embedding and the advantages our GEFA in modeling the nested graph for drug-target interaction.

Machine learning is an important component of the growing field of data science. Through the use of statistical methods, different type of algorithms is trained to make classifications or predictions, and to uncover key insights in this project. These insights subsequently drive decision making within applications and businesses, ideally impacting key growth metrics.

Machine learning algorithms build a model based on this project data, known as training data, in order to make predictions or decisions without being explicitly programmed to do so. Machine learning algorithms are used in a wide variety of datasets, where it is difficult or unfeasible to develop conventional algorithms to perform the needed tasks.

I. INTRODUCTION

The accurate prediction of drug-target affinity (DTA) is a crucial aspect of drug discovery, enabling the identification of promising drug candidates that can effectively interact with specific biological targets. Traditional experimental techniques for determining DTA, such as high-throughput screening and binding assays, are both time-consuming and resource-intensive, underscoring the need for efficient computational approaches that can accelerate the drug development process [1]. In recent years, machine learning (ML) and deep learning (DL) models have gained traction in this field, offering substantial improvements in the speed and accuracy of DTA predictions [2] [3].

One of the innovative strategies in advancing DTA prediction is the early fusion approach. Early fusion involves integrating multiple data sources or feature representations at an initial stage, allowing the model to learn complex interactions between drug molecules and their target proteins from the outset [4]. This contrasts with late fusion methods,

where separate models are trained on individual data modalities, and their outputs are later combined. By

leveraging early fusion, models can capture intricate relationships between drugs and targets, leading to enhanced predictive performance [5].

The proposed GEFA (Generalized Early Fusion Approach) builds on this concept by developing a robust framework for early fusion in DTA prediction. GEFA aims to combine various molecular representations of drugs, such as SMILES strings and molecular fingerprints, with protein sequences or structural data, ensuring that the model fully exploits the complementary information from these diverse inputs. This approach not only improves the accuracy of DTA predictions but also provides a scalable solution for large-scale drug discovery projects [6] [7].

II. EXISTING SYSTEM

Drug re-purposing [18] is the process of identifying well established medications for the novel target disease. The

advantages of this drug re-purposing over developing a completely novel drug are lower risk and fast-track development [19]. The process of drug re-purposing consists of three key steps: identifying the candidate molecules given the target disease, drug effect assessment in the preclinical trial, and effectiveness assessment in clinical trial [20]. The first step, hypothesis generation, is critical as it decides the success of the whole process. Advanced computational approaches are used for hypothesis generation. Computational approaches in drug re-purposing can be categorized into six groups [20]: genetic association [21], [22], pathway pathing, retrospective clinical analysis, novel data sources, signature matching [29]–[31], molecular docking [32]–[34].

Drug-target binding affinity indicates the strength of the binding force between the target protein and its ligand (drug or inhibitor) [35]. The drug-target binding affinity prediction problem is a regression task predicting the value of the binding force. The binding strength is measured by the equilibrium dissociation constant (KD). A smaller KD value indicates a stronger binding affinity between protein and ligand [35]. There

are two main approaches: structural approach and non-structural approach [1]. Structural methods utilize the 3D structure of protein and ligands to run the interaction simulation between protein and ligand. On the other hand, the non-structural approach relies on ligand and protein features such as sequence, hydrophobic, similarity or other alternative structural information.

The structure-based approach involves molecular docking, predicting the three-dimensional structure of the targetligand complex. In molecular docking, there are a large number of target-ligand complex conformations. The conformations are evaluated by the scoring function. Based on the scoring function types, the structural approach can be categories into three groups [1]: classical scoring function method [36]–[39], machine learning scoring function method [40], and deep learning scoring function method [41], [42]. In classical scoring approaches, Elanie et al. [36] uses DelPhi-calculated potential at each ligand atom for the contact scoring function. In machine learning approaches, Kundu et al. [40] extracts ligand features (e.g. atom count, physicochemical properties) and protein features (e.g.

accessible surface, number of chains) from 3D structure data then applies machine learning to learn the scoring function. In deep learning approaches, Marta et al. [41] uses 3D convolution with protein-ligand 3D structure to predict the binding affinity.

Disadvantages

- ❖ The system is not implemented compare the drug representation which extracted from the drug-protein fusion graph and drug representation extracted from the drug graph.
- ❖ The system is not implemented Graph Early Fusion for binding Affinity prediction (GEFA).

III. PROPOSED SYSTEM

• In summary, the contribution of our work is two-fold. First, we combine the protein sequence embedding feature and protein contact map to build the graph representation of a target protein. Second, in order to reflect the target representation change during the binding process, we propose a so-called Graph Early Fusion for binding Affinity prediction (GEFA) for more accurate biological modeling. We demonstrate the effects of the GEFA on Davis dataset [17] where it has shown superior performance

against previous studies on different settings.

To address target protein representation change, the system proposes an early-fusion-based approach. Initially, we extract representation feature for a given drug molecule from its drug graph structure. Then, the drug representation is integrated into the protein graph structure before the protein representation learning phase. This is basically a graph structure nested inside another graph structure. This graphin- graph neural network design allows the model to learn changes in protein representation caused by the binding process with the drug molecule.

Advantages

- The proposed system refines the Graph-Graph Integration with Early Fusion and Graph Early Fusion for binding Affinity prediction (GEFA).
- The proposed system implemented the usage of attention mask as the graph edge. Instead of using attention as drug-residue edge weight, drug-residue edges are weighted the same as the residue-residue edges in the target graph.

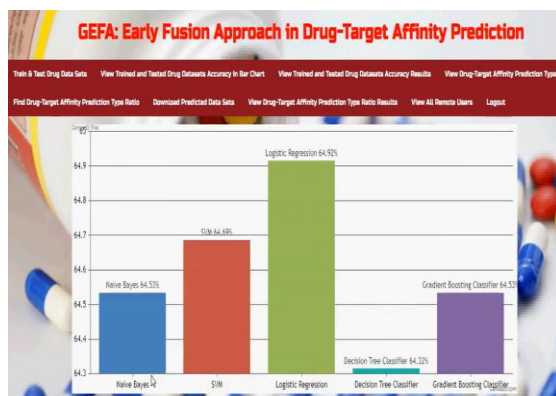
IV. MODULES

Service Provider

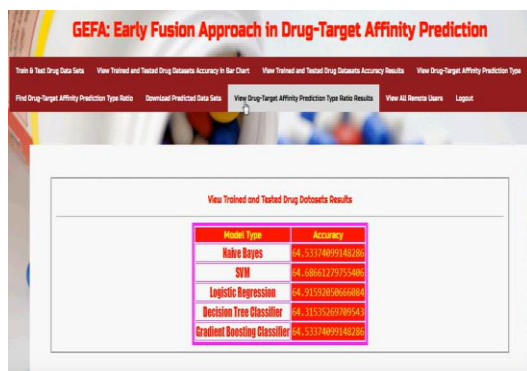
In this module, the Service Provider has to login by using valid user name and password.



After login successful he can do some operations such as Login, Train & Test Data Sets, View Trained Accuracy in Bar Chart,



View Trained Accuracy Results,



Model Type	Accuracy
Naive Bayes	64.3314899148238
SVM	64.68661279755486
Logistic Regression	64.9155109666688
Decision Tree Classifier	64.3151269198545
Gradient Boosting Classifier	64.5314899148238

View Type, Find Type Ratio, Download Predicted Datasets, View Type Ratio Results, View All Remote Users.

View and Authorize Users

In this module, the admin can view the list of users who all registered. In this, the admin can view the user's details such as, user name, email, address and admin authorizes the users.

Remote User

In this module, there are n numbers of users are present. User should register before doing any operations. Once user registers, their details will be stored to the database. After registration successful, he has to login by using authorized user name and password. Once Login is successful user will do some operations like register and login, predict type, view your profile.

V.CONCLUSION

In this project, the GEFA (Generalized Early Fusion Approach) was introduced as an innovative framework to enhance drug-target affinity (DTA) prediction using machine learning techniques. By leveraging early fusion, the model is able to integrate multiple data sources—such as molecular representations of drugs and protein sequences—at an initial stage, allowing for a more comprehensive understanding of the relationships

between drugs and their biological targets. The integration of these diverse modalities improves the model's ability to capture complex interactions, resulting in more accurate and reliable predictions compared to traditional approaches.

GEFA not only demonstrates the potential to advance the field of computational drug discovery but also addresses critical challenges such as scalability and efficiency. By refining how input data is fused and processed, this method significantly reduces the time and cost of identifying potential drug candidates. As drug discovery continues to evolve with increasing reliance on computational techniques, GEFA provides a scalable and adaptable approach that can be extended to various types of biological targets and drug molecules. Future research could explore further optimizations in feature selection, model architectures, and real-world validations to fully unlock the potential of this early fusion approach.

VI. REFERENCES

1. Ezzat, A., Wu, M., Li, X. L., & Kwoh, C. K. (2016). Computational prediction of drug-target interactions using chemogenomic approaches: A survey. *Current Bioinformatics*, 11(4), 378-392.
2. Wen, M., Zhang, Z., Niu, S., Sha, H., Yang, R., Yun, Y., & Lu, H. (2017). Deep-learning-based drug-target interaction prediction. *Journal of Proteome Research*, 16(4), 1401-1409.
3. Chen, L., Zeng, W. M., Cai, Y. D., Feng, K. Y., & Chou, K. C. (2012). Predicting anatomical therapeutic chemical (ATC) classification of drugs by integrating chemical-chemical interactions and similarities. *PLoS ONE*, 7(4), e35254.
4. Tsubaki, M., Tomii, K., & Sese, J. (2018). Compound-protein interaction prediction with end-to-end learning of neural networks for graphs and sequences. *Bioinformatics*, 35(2), 309-318.
5. Gomes, J., Ramsundar, B., Feinberg, E. N., & Pande, V. S. (2017). Atomic convolutional networks for predicting protein-ligand binding affinity. *arXiv preprint arXiv:1703.10603*.
6. Öztürk, H., Özgür, A., & Ozkirimli, E. (2018). DeepDTA: deep drug-target binding affinity prediction. *Bioinformatics*, 34(17), i821-i829.
7. Karimi, M., Wu, D., Wang, Z., & Shen, Y. (2019). DeepAffinity: interpretable deep learning of compound-protein

- affinity through unified recurrent and convolutional neural networks. *Bioinformatics*, 35(18), 3329-3338.
8. Tang, J., Szwajda, A., Shakyawar, S., Xu, T., & Hintsanen, P. (2014). Making sense of large-scale kinase inhibition data. *Nature Chemical Biology*, 10(9), 719-723.
9. Zhang, W., Yue, X., Lin, W., Wu, W., Liu, R., Huang, F., & Li, X. (2020). Predicting drug–target interactions by a deep learning method with graph embeddings. *Journal of Chemical Information and Modeling*, 60(9), 4153-4162.
10. Öztürk, H., Ozkirimli, E., & Özgür, A. (2019). WideDTA: prediction of drug–target binding affinity by combining wide and deep learning. *Bioinformatics*, 35(14), 2181-2187.
11. Wang, L., You, Z. H., Li, Y. M., Zheng, K., & Li, L. P. (2014). Predicting drug-target interactions based on bipartite local models and hubness-aware regression. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 12(3), 559-570.
12. He, T., Heidemeyer, M., Ban, F., Cherkasov, A., & Ester, M. (2017). SimBoost: a read-across approach for predicting drug-target binding affinities using gradient boosting machines. *Journal of Cheminformatics*, 9(1), 1-14.
13. Nguyen, T., Le, H., Quinn, T. P., Nguyen, T., & Venkatesh, S. (2021). GraphDTA: prediction of drug–target binding affinity using graph convolutional networks. *Bioinformatics*, 37(8), 1140-1147.
14. Chen, W., Zhang, T., Wang, W., Liu, Z., & Wang, K. (2020). TransformerCPI: improving compound–protein interaction prediction by sequence-based deep learning with self-attention mechanism and label reversal. *Bioinformatics*, 37(4), 1121-1129.
15. Diao, J., & Hu, J. (2020). Deep ensemble learning based drug–target binding affinity prediction. *Briefings in Bioinformatics*, 22(1), 457-469.