

# Benchmark on Drug Target Interaction Modeling from a Drug Structure Perspective

Anonymous authors

Paper under double-blind review

## Abstract

The prediction modeling of drug-target interactions is crucial to drug discovery and design, which has seen rapid advancements owing to deep learning technologies. Recently developed methods, such as those based on graph neural networks (GNNs) and Transformers, demonstrate exceptional performance across various datasets by effectively extracting structural information. However, the benchmarking of these novel methods often varies significantly in terms of hyperparameter settings and datasets, which limits algorithmic progress. In view of these, we conducted a comprehensive survey and benchmark for drug-target interaction modeling from a structural perspective via integrating tens of explicit (i.e., GNN-based) and implicit (i.e., Transformer-based) structure learning algorithms. We conducted a macroscopical comparison between these two classes of encoding strategies as well as the different featurization techniques that inform molecules’ chemical and physical properties. We then carry out the microscopical comparison between all the integrated models across the six datasets via comprehensively benchmarking their effectiveness and efficiency. To comprehensively assess fairness, we investigate model performance under two experimental scenarios: one with unified hyperparameter settings and the other with individually optimized configurations. Remarkably, the summarized insights from the benchmark studies lead to the design of model combos. We demonstrate that our combos can achieve new state-of-the-art performance on various datasets associated with cost-effective memory and computation.

## 1 Introduction

The prediction modeling of drug-target interactions (DTI) has emerged as an irreplaceable task for efficacious therapeutic interventions. The binding affinity between a drug molecule and its target protein plays a significant role in the design and repurposing of drugs, where a high affinity typically indicates the desired therapeutics, target specificity, long residence, and drug resistance delay (Hughes et al., 2011; Copeland et al., 2006; Swinney, 2004). The precise modeling of DTI can expedite the drug discovery process and circumvent the associated cost (Ashburn & Thor, 2004; Strittmatter, 2014). Deep learning-based frameworks have recently revolutionized this field, enabling more accurate predictions and accelerating the discovery of new compounds by guiding laboratory experiments more efficiently (Wen et al., 2017; Abbasi et al., 2021; Huang et al., 2020a).

Within deep learning frameworks (Öztürk et al., 2018; 2019), drugs are commonly represented using the Simplified Molecular Input Line Entry System (SMILES) (Weininger, 1988a), and proteins are represented as sequences of amino acids. These representations are typically processed using various neural network architectures, such as convolutional neural networks (CNNs) (Krizhevsky et al., 2017; He et al., 2016), recurrent neural networks (RNNs), Transformers, and so on, before being integrated and processed by a multi-layer perceptron (MLP) for DTI prediction. It is notorious that the reliance on sequence-based representations can result in the loss of structural information, which can potentially compromise the DTI predictive capability. From the drug perspective, molecular structure modeling helps identify the specific binding sites (Ma et al., 2011), contributes to predicting pharmacokinetic properties (Ekins et al., 2007), and allows conformational flexibility (Karplus & Kuriyan, 2005).

To address this problem, a number of drug algorithms have been proposed to promote DTI prediction, which can be categorized into explicit and implicit structure learning. First, graph neural networks (GNNs) (Kipf & Welling, 2016; Nguyen et al., 2020) have been widely adopted to learn the molecular structures, owing to their ability to directly operate on graph-based representations of molecules. By explicitly propagating information through the graph, GNNs can learn node and edge features and thereby capture the structural and functional relationships between atoms and bonds. Second, Transformers, originally focused on natural language processing (Vaswani et al., 2017a), have also shown promise in biomedical applications (Huang et al., 2020b; Chen et al., 2020). They rely on self-attention mechanisms to implicitly weight the correlations between different parts of the input SMILES, allowing them to capture long-range dependencies and contextual information.

While these techniques contribute to the learning of drug structures, there is still a key knob under-explored: we lack a systematic study to benchmark their effectiveness and efficiency. Without such a standardized benchmark, it is unachievable to offer fair comparisons and subsequently summarize the design philosophy necessary to inform DTI. There have been several surveys and benchmarks on computational methods for DTI prediction (Öztürk et al., 2018; Huang et al., 2020a; 2021; Xu et al., 2022), which leave out the recent developments of structure learning algorithms and unavoidably fail to focus on drug structure benchmarking. Moreover, although massive efforts (Bal et al., 2024; Zhu et al., 2023; Nguyen et al., 2020) have been made to explore the effectiveness of modeling structural information, they predominantly use their proprietary training hyperparameters, datasets, and evaluation metrics. Due to the various settings, one cannot reach convincing answers as to whether a configuration of structure encoders and/or featurization methods generally performs well. The complex of DTI classification and regression tasks and datasets complicates the benchmark comparison.

In this study, we introduce GTB-DTI, a comprehensive benchmark customized for **G**NN and **T**ransformer-based methodologies for **D**TI prediction. I) We thoroughly examine the implementation details for each category of drug structure learning methods and integrate three widely used datasets for classification and regression tasks, respectively. Then, we harmonize the sensitive hyperparameters across different methods using a greedy search to identify an optimal *sweet spot* configuration. The unified setting lays the foundation for a fair and reproducible benchmark. II) To gain macroscopical insights into the structure encoders and featurization methods, we fix the drug encoder to be either GNN or Transformer-based approaches and benchmark these two strategies in the various settings. We also integrate tens of drug features given their importance to inform molecules’ chemistry and physical properties and evaluate them on the representative datasets. III) To gain macroscopical insights into nuance between 31 concerned models, we conduct the benchmark studies of their effectiveness on the six datasets with the unified setting. Moreover, we assess the efficiency of each method by measuring peak GPU memory usage, running time, and convergences. IV) The comprehensive study finally provides a number of surprising observations: *i*) The CNN encoder accompanied by integer features has close protein embedding performance compared to the Transformer or larger language models, but they are more efficient. *ii*) The explicit and implicit structure encoders for drugs exhibit unequal performances across the different datasets, which suggests their hybrid usage for generalization purposes. *iii*) Inspired by these insights, we conclude with a model combo that leads us to attaining state-of-the-art (SOTA) regression results and performing similarly to SOTA in the DTI classifications. Our combos further deliver cost-effective memory usage and running time as well as faster convergence, which can serve as a new baseline for the following explorations.

## 2 Formulations for Drug-target Interaction Modeling

In this research, we focus on the formulations of recently emerging structure modeling approaches for drug molecules, which could be categorized into explicit methods based on graph neural networks and implicit methods based on Transformers. The target proteins are learned by the sophisticated tools of convolutional/recurrent neural networks (CNNs/RNNs) or Transformers, after which both the molecules’ and proteins’ embeddings are integrated to facilitate interaction prediction. We will also summarize and benchmark the various widely adopted molecule features.

## 2.1 Graph Neural Networks based Methods

A drug molecule is typically represented as a graph  $G = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  and  $\mathcal{E}$  denote the sets of atoms and chemical bonds, respectively. The classical GNN frameworks involve key processes of aggregating and updating node features, collectively referred to as message passing, which can be mathematically represented as follows (Scarselli et al., 2008; Duan et al., 2022):

$$\mathbf{h}_i^{(l+1)} = \text{COMBINE}_{\text{node}}^{(l)} \left( \mathbf{h}_i^{(l)}, \text{AGGREGATE}_{\text{node}}^{(l)} \left( f_{\alpha} \left( \left\{ \mathbf{h}_j^{(l)}, \mathbf{e}_{ij}^{(l)} : j \in \mathcal{N}_i \right\} \right) \right) \right), \quad (1)$$

$$\mathbf{e}_{ij}^{(l+1)} = \text{COMBINE}_{\text{edge}}^{(l)} \left( \mathbf{e}_{ij}^{(l)}, \text{AGGREGATE}_{\text{edge}}^{(l)} \left( g_{\beta} \left( \left\{ \mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} : j \in \mathcal{N}_i \right\} \right) \right) \right), \quad (2)$$

where  $\mathbf{h}_i^{(l)}$  is the feature representation of node  $v_i$  at layer  $l$ ,  $\mathbf{e}_{ij}^{(l)}$  is the feature representation of edge between nodes  $v_i$  and  $v_j$ ,  $\mathcal{N}_i$  refers to the set of neighboring nodes next to node  $v_i$ . Functions  $\text{AGGREGATE}^{(l)}$ ,  $\text{COMBINE}^{(l)}$  aim to aggregate the neighborhood representations and integrate them together with the node features, respectively. Additionally,  $f_{\alpha}$  and  $g_{\beta}$  are feature mapping functions, parameterized by  $\alpha$  and  $\beta$ , respectively. The molecule’s representation can be derived using READOUT a function that processes the set of vertex features  $\mathbf{H}^{(L)}$  at the last layer.

**Graph Convolutional Networks (GCN).** Given a molecule with  $N$  atoms, the adjacency matrix  $\mathbf{A} \in R^{N \times N}$  indicates its connectivity, with  $A_{ij} = 1$  if atom  $v_i$  is adjacent to atom  $v_j$ , and 0 otherwise. Considering the self-connection of atoms, we have  $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ . Let’s  $\mathbf{X} \in R^{N \times C}$  denote the initial atom feature matrix. GCN (Kipf & Welling, 2017) models the message passing as follows:

$$\mathbf{H}^{(l+1)} = \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}^{(l)} \mathbf{W}^{(l)}), \quad (3)$$

where  $\mathbf{H}^{(l)}$  is the node feature matrix at layer  $l$ , starting with  $\mathbf{H}^{(0)} = \mathbf{X}$ . Matrix  $\mathbf{W}^{(l)}$  represents the learnable weights for layer  $l$ ,  $\sigma$  denotes a non-linear activation function, e.g., ReLU, and  $\tilde{\mathbf{D}}$  is a diagonal degree matrix of  $\tilde{\mathbf{A}}$ . A couple of pioneering works have leveraged GCN to facilitate drug-protein interaction prediction (Mukherjee et al., 2022; Tran et al., 2022; Tsubaki et al., 2018; Pan et al., 2023b). For example, DeepGLSTM (Mukherjee et al., 2022) uses mixture-of-depths GCNs to capture drug representations from different scales. CPI (Tsubaki et al., 2018) considers cross-atom distance and introduces the concept of r-radius subgraphs (Costa & Grave, 2010), using r-radius vertices and edges to redefine the structure of graphs.

**Graph Isomorphism Networks (GIN).** GIN excels in learning distinct graph features by approximating the Weisfeiler-Lehman test, enabling it to distinguish a wide range of graph structures (Xu et al., 2018). The message-passing process at the  $(l+1)$ -th layer is of the following form:

$$\mathbf{h}_i^{(l+1)} = \text{MLP}^{(l)}((1 + \epsilon^{(l)})\mathbf{h}_i^{(l)} + \sum_{j \in \mathcal{N}_i} \mathbf{h}_j^{(l)}), \quad (4)$$

where  $\text{MLP}^{(l)}$  is a multi-layer perceptron that parameterizes the update function, and  $\epsilon^{(l)}$  is a learnable parameter. We benchmark several GIN-based drug-target interaction modeling methods. GraphCPI (Quan et al., 2019) and GraphDTA (Nguyen et al., 2020) adopt GIN-based models with batch normalization to obtain the drug representation. SubMDTA (Pan et al., 2023a) uses a subgraph’s generation task and contrastive learning to pretrain a molecular graph encoder with multiple GIN layers for further prediction.

**Graph Attention Networks (GAT).** Unlike fixed-weight aggregation, GAT (Velickovi et al., 2018) employs an attention mechanism to determine neighborhood importance and learn the node embeddings as follows:

$$\mathbf{h}_i^{(l+1)} = \sigma \left( \sum_{j \in i \cup \mathcal{N}_i} \text{softmax}(\text{LeakyReLU}(\mathbf{W}_a^T [\mathbf{W}^{(l)} \mathbf{h}_i^{(l)} || \mathbf{W}^{(l)} \mathbf{h}_j^{(l)}])) \mathbf{W}^{(l)} \mathbf{h}_j^{(l)} \right). \quad (5)$$

$\mathbf{W}_a^T$  denotes attention weights, and  $||$  is a concatenating operation. GraphDTA (Nguyen et al., 2020) and AMMVF (Wang et al., 2023) leverage the multi-head GAT layers to optimize the atom messaging. They integrate GAT with other architectural modules, such as GCN, facilitating a more comprehensive representation of drugs.

**Graph Transformers.** Graph Transformers (Rong et al., 2020; Maziarka et al., 2020) have emerged as powerful alternatives to traditional graph neural networks (GNNs) for molecular representation learning. Unlike conventional GNNs, which rely on message-passing mechanisms to propagate local node information, Graph Transformers leverage self-attention mechanisms to capture both local and global dependencies more effectively. By integrating Message Passing Networks into Transformer-style architectures, these models enhance expressiveness, enabling more comprehensive encoding of molecular structures. This hybrid approach allows Graph Transformers to preserve structural information while benefiting from the flexibility of attention-based learning.

## 2.2 Transformer-based Methods

Besides the graph representation, drugs could also be decorated as SMILES strings (Weininger, 1988b) and encoded similarly to natural language processing. Specifically, after tokenizing SMILES strings, the Transformer model utilizes multi-head attention to model the interactions between different segments of the input and obtain the molecular representations. Positional encodings are also integrated to preserve the sequence order, enhancing the model’s ability to process sequential information effectively. We review and benchmark two typical types of attention mechanisms used for molecular representations.

**Self-Attention**(Huang et al., 2020b; Qian et al., 2023; Yin et al., 2024). Self-attention computes a weighted sum of all input values based on their relevance to each other. Considering an embedding of a SMILES sequence  $\mathbf{H}^{(l)} \in \mathbb{R}^{d \times N}$  at a specific Transformer layer, where  $N$  and  $d$  are token length and dimension, respectively, the attention is calculated by  $\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(\mathbf{Q}\mathbf{K}^T/\sqrt{d_k})\mathbf{V}$ .  $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{d_k \times N}$  and  $\mathbf{V} \in \mathbb{R}^{d_v \times N}$  are projections of the input matrix  $\mathbf{H}^{(l)}$ . Multi-head attention combines these projections across different subspaces for a more detailed analysis. Followed by normalization and feed-forward neural networks, the SMILES embedding is updated,  $\mathbf{H}^{(l+1)}$  and the output from the last layer is treated as molecular representations. Transformer encoders like MolTrans (Huang et al., 2020b) and FOTFCPI (Yin et al., 2024) are adopted to enhance substructure embeddings in proteins and drugs.

**Cross-Attention**(Kurata & Tsukiyama, 2022; Qian et al., 2023). Cross-attention is designed to capture the interaction between the drug and protein sequences, with the query matrix  $\mathbf{Q}$  derived from one sequence and the key and value matrices  $\mathbf{K}, \mathbf{V}$  from another. This mechanism is particularly useful in integrating hybrid representations such as drug graphs and SMILES (Wang et al., 2023), as well as drugs and proteins (Pan et al., 2023b; Kurata & Tsukiyama, 2022).

## 2.3 Feature Processing Methods

Beyond the drugs’ structure or sequence learning with GNNs or Transformers, the extra molecular properties, such as molecular weight, solubility, and lipophilicity, are crucial for building accurate and quantitative drug-target relationship models. We summarize two typical featurization methods.

**Sequence Processing Methods.** Both drugs and proteins are input as strings of ASCII characters, whose features can be extracted using statistical solutions. Integer encoding (Nguyen et al., 2020) simply converts the string to a sequence of integers, which assigns an integer to each character. The N-gram (Dong et al., 2005) captures the statistical dependencies between characters in an input string. Specifically, a 3-gram model breaks down a sequence  $S = \{s_1, s_2, \dots, s_m\}$  into  $\{[s_1, s_2, s_3], [s_2, s_3, s_4], \dots, [s_{m-2}, s_{m-1}, s_m]\}$ , analyzing the relationship between adjacent characters.

**Drug-Unique Featurization Methods.** The additional chemical properties and structural details of SMILES strings are often considered to gain a more comprehensive understanding. Extended-Connectivity Fingerprints (ECFP) (Morgan, 1965; Rogers & Hahn, 2010), involves generating unique identifiers for atoms based on their local chemical environment and iteratively updating these through a hash function to capture a broader molecular context, ultimately producing a set of fingerprints that represent the molecules overall structure. Another approach, RDKit, is used to convert SMILES into molecular graphs (Landrum et al., 2006; Nguyen et al., 2020), where nodes represent the physical and chemical properties of molecules, and bonds are represented by an adjacency matrix. For example, atomic properties such as atom type, degree,

and hydrogen information (like the number of explicit hydrogens) are all crucial for constructing a graph. More detailed properties can be found in Appendix G.

**Embedding Featurization Methods.** Embedding methods are used to translate these discrete sequences into continuous embedding spaces. Notably, SMI2Vec (Quan et al., 2018) and Prot2Vec (Asgari & Mofrad, 2015) convert discrete tokens of drug SMILES and protein sequences into vectors that encapsulate semantic and syntactic similarities, effectively grouping similar tokens together in vector space. Additionally, pre-trained language models (Bal et al., 2024; Lin et al., 2022) are increasingly utilized to leverage large-scale learned patterns, fine-tuned to analyze complex protein data representations effectively.

### 3 A Fair Benchmark Platform Setup

**Benchmark Model and Dataset Selection.** From the perspective of reproducibility, we restrict our analysis to models for which the source code has been publicly released. To enhance the comprehensiveness, credibility, and sophistication of our benchmark, we conduct experiments on more than 30 models, including both GNN-based and Transformer-based methods. These models are derived from papers spanning the years 2018 to 2024. We run these models on 6 frequently evaluated datasets, including both binary interaction classification and continuous affinity regression. For the classification aspect, we utilize datasets including Human (Liu et al., 2015), *Caenorhabditis elegans* (*C. elegans*) (Tsubaki et al., 2018), and DrugBank (Wishart et al., 2008). For regression, we employ the Davis (Davis et al., 2011), KIBA (Tang et al., 2014), and BindingDB datasets (Liu et al., 2007) with dissociation constant (Kd) measures, as processed in Huang et al. (2021). The statistical details of these models and datasets are presented in Appendix B and Table 3, respectively.

**Hyperparameter Configuration.** Given the critical role of hyperparameters in achieving optimal performance, we perform a systematic review of the hyperparameters for all selected models in Appendix E. To ensure that comparisons across models are equitable, we consider comparing each model using its optimal hyperparameters, as reported in the corresponding perspective papers.

**Data Split.** We treat each dataset independently to prevent any information leakage that could arise from training a single model on multiple datasets. For duplicate drugprotein pairs in the regression dataset, only the entry with the maximum affinity score is retained. For duplicates in the classification dataset, all entries are removed if conflicting labels are present; otherwise, a single instance is kept. Following data cleaning, each dataset is split into a training set and a test set. We apply k-fold cross-validation on each training set, where each fold consists of a unique trainingvalidation split, and models trained on different folds are completely independent, thereby eliminating any possibility of cross-fold leakage. The training sets are used to fit the models, the validation sets are used to select the best model during training, and the test sets are for final evaluation.

**Other Training Details.** We train on the training set and see its performance on the validation set at the end of every epoch, and the model that achieves the best validation performance will be saved. After training, we evaluate the saved model on the test set and save their results. We average the final performance metrics across all folds as the final results. Considering the original training epochs, we use 1000 as maximum epochs limitation. To avoid overfitting, we consider an early stop mechanism in training. Given the complexity in the dataset, we use 50 patience for all datasets. We use MSE as an early stopping evaluation metric for regression and F1 for classification. The detailed results are provided in Appendix F.

### 4 A Macroscopic Benchmark on Encoder and Featurization Strategies

**★Encoder Exploration for Drugs and Proteins.** To investigate the influence of different encoding strategies for extracting the structural information of drugs, we employ GIN (Xu et al., 2019) and vanilla Transformer (Vaswani et al., 2017b) as the encoders for drugs. Meanwhile, integer encoding with CNN, n-gram encoding with CNN, and the vanilla Transformer are considered to capture protein’s representations, which are frequently adopted. To leverage the advantages of the pretrained protein information, we include a language model, i.e., Evolutionary Scale Modeling (ESM2) (Lin et al., 2022). The results of various

combinations of drug and protein encoders are shown in Fig. 1. All results are averaged by five-fold cross-validation with an early stop mechanism.

**Obs. 1. GNN and Transformer-based drug encoders exhibit unequal performance depending on DTI tasks.** When the encoder for the protein sequence is fixed, drug features extracted by the GNN structures GIN generally perform better than those by Transformers in regression tasks, but the opposite is true in classification tasks. This disparity may be due to the smaller size of the Human dataset compared to the Davis dataset, which allows for faster convergence in classification tasks than in regression tasks.

**Obs. 2. Transformer models are better but sensitive in extracting features from protein.** Although we only consider the simplest pretrained protein language model of ESM2, it still significantly outperforms other encoders in both tasks. This improvement can likely be attributed to the robust and generalizable representations learned from extensive data by the pretrained model. In addition, the Transformer encoder for the protein achieves the best performance on the classification task but shows unstable performance in the regression task. This is likely due to the smaller size and simpler classification dataset compared to the regression dataset, making the training stop for a fixed early stop threshold.

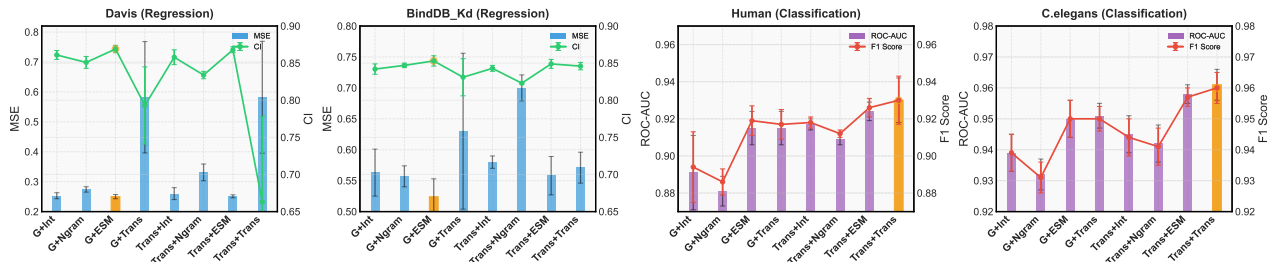


Figure 1: Comparison of different encoding strategies with early stop mechanism for drugs and proteins when the total epoch is 1000, LR is 0.0005, BS is 512, and DR is 0.2. Trans is a Transformer-based model, which is composed of two parts: embedding with the position encoding and the encoder in the Transformer. ESM refers to ESM2.

**Obs 3. Integer encoding appears to be more effective when paired with a CNN as the protein encoder and a fixed drug encoder.** Compared to this specific model configuration, the local context provided by 3-gram encoding does not significantly enhance the model’s predictive performance. This implies that the simple relationships in amino acids’ immediate neighbors, as modeled by Word2Vec, do not capture much useful information compared with simple integer encoding.

**★Featurization Exploration.** Despite the efficacy of GNNs in learning drug structures, the featurization of nodes plays a critical role in capturing both the intrinsic properties of atoms and their contextual relevance. We conduct a detailed analysis of various methods (summarized in Section G of the Appendix) for constructing graph features within the DTI context. The node feature is constructed via various characteristics, such as chemical and physical properties. We categorize each feature into five main classes, e.g., atomic properties (AP), hydrogen information (HI), electron properties (EP), stereochemistry (Ste), and structural information (Str). To better determine which types of features are more effective in capturing the structural information, we conduct an ablation study on the different featurization strategies. We choose GraphDTA (Nguyen et al., 2020) and GraphCPI (Quan et al., 2019) with GIN as our backbone models. The results of feature combinations are reported in Fig. 2.

**Obs. 4. More complex featurization does not necessarily bring a positive effect, and its effectiveness is highly task-dependent.** As shown in Fig. 2, adding features like atomic properties (AP), hydrogen information (HI), and stereochemistry information (Ste) improves performance in the regression task by reducing the MSE loss, suggesting that these features provide valuable information. However, features like electron properties (EP) and structural information (Str) may introduce noise rather than useful information, especially when combined with other features, leading to inconsistent results. Furthermore, in classification tasks, the trend differs, with additional features sometimes negatively impacting performance (blue line) rather than providing benefits. This highlights the importance of careful feature selection, as

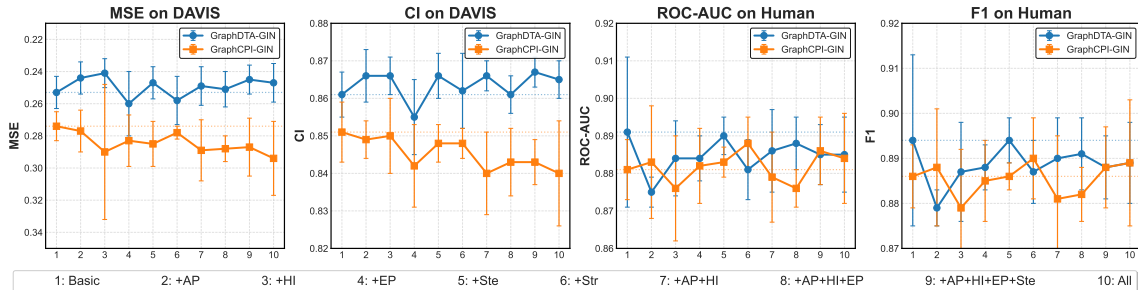


Figure 2: Various performances of GraphDTA-GIN and GraphCPI-GIN versus different features on DAVIS and Human datasets.  $+x$  means that  $x$  is added to the basic featurization. All means using all features.

indiscriminate inclusion of complex features may lead to increased noise, affecting model generalization and robustness.

**Obs. 5. Protein representation plays a crucial role in the effectiveness of different drug featurizations.** As shown in Fig. 2, even when using the same drug featurization strategies, the trends vary depending on the protein representation. This suggests that the way proteins are encoded directly influences how drug features interact with the model. For instance, in GraphDTA, features like atomic properties (AP) and hydrogen information (HI) help improve performance, but in GraphCPI, those same features don’t always provide the same benefits. In some cases, adding more drug features introduces noise rather than useful information. This highlights the fact that drug featurization and protein representation are deeply interconnected, and optimizing one without considering the other may not yield the best results. To build more effective drug-protein interaction models, both components should be considered holistically rather than in isolation.

## 5 A Microscopic Benchmark on DTI Models

**★Benchmark over Effectiveness.** As shown in Table 1 and Table 2, we conduct experiments on models with their optimal hyperparameters across two tasks and three datasets, respectively (see more comprehensive model comparisons in Appendix F). For the unified hyperparameters, we summarized the result in Appendix ???. All results are averaged by five-fold cross-validation with an early stop mechanism.

**Obs. 6. Molecular graphs are better than fingerprints for capturing the graph features of a drug.** In reference to Table 7, it is evident that GNN-based approaches utilizing the molecular graph generally yield superior performance compared with fingerprints (CPI (Tsubaki et al., 2018), BACPI (Li et al., 2022), GANDTI (Wang et al., 2021)). This reinforces the idea that the rich structural and atomic property information inherent to molecular graphs is pivotal for representation extraction, leading to enhanced model performance.

**Obs. 7. Graph structure is a crucial part of extracting a drug’s features.** Different GNNs have distinct performances in both tasks when the protein representation is fixed. Specifically, GIN, with its unique ability to distinguish non-isomorphic graphs, consistently outperforms other models across different protein encoders in regression tasks. Although Transformer-based methods such as MRBDTA are proficient in handling sequential information from SMILES and proteins, the depth of information they capture appears to be marginally less comprehensive than that provided by molecular graph-based approaches. This is substantiated by the superior performance of GNN-based methods, including MGraphDTA, ColdDTA, and SubMDTA, which suggests that GNN captures intricate structural details more effectively.

**★Benchmark over Efficiency** To analyze the training speed and memory usage, we empirically evaluate the peak memory and running time for various methods during the training procedure on one regression dataset and one classification task, respectively. To fairly compare various methods, we set the batch size as 32, as such maximum batch size is adopted by some methods. All results are measured on an RTX 3090 GPU. The memory and running time comparisons are illustrated in Fig. 3.

Table 1: Regression task benchmark on DAVIS, KIBA, and BindingDB\_Kd datasets, respectively. For the GraphDTA and GraphCPI, we only show the one with a specific GNN encoder that has the overall best performance. The best result is highlighted in bold, and the runner-up is underlined. *Avg. Reduction of MSE* is computed by the average (across 3 datasets) of the differences between our model’s MSE and each model’s MSE, divided by the average (across 3 datasets) of each model’s MSE, respectively.

Category	Models	DAVIS			KIBA			BindingDB_Kd			Avg. Reduction of MSE (%)
		MSE	R2	CI	MSE	R2	CI	MSE	R2	CI	
GNN	GraphDTA-GIN	0.253 ± 0.010	0.623 ± 0.015	0.861 ± 0.006	0.255 ± 0.007	-1.840 ± 0.779	0.553 ± 0.019	0.563 ± 0.038	0.693 ± 0.021	0.842 ± 0.007	34.827%
	GraphCPI-GIN	0.274 ± 0.009	0.593 ± 0.013	0.851 ± 0.008	1.681 ± 0.946	-17.724 ± 10.533	0.553 ± 0.094	0.557 ± 0.017	0.696 ± 0.009	0.847 ± 0.003	72.213%
	MGraphDTA	0.232 ± 0.012	0.655 ± 0.018	0.869 ± 0.007	0.032 ± 0.012	0.642 ± 0.133	0.832 ± 0.040	0.529 ± 0.011	0.712 ± 0.006	0.852 ± 0.005	11.980%
	SAGDTA	0.324 ± 0.064	0.518 ± 0.096	0.833 ± 0.027	0.065 ± 0.008	0.279 ± 0.085	0.713 ± 0.032	0.529 ± 0.011	0.712 ± 0.006	0.852 ± 0.005	23.965%
	EmbedDTI	0.280 ± 0.024	0.583 ± 0.036	0.851 ± 0.009	0.289 ± 0.142	-2.217 ± 1.579	0.558 ± 0.038	0.542 ± 0.019	0.705 ± 0.010	0.850 ± 0.004	37.174%
	DeepGLSTM	0.316 ± 0.023	0.529 ± 0.035	0.841 ± 0.007	8.539 ± 7.479	-94.109 ± 83.400	0.514 ± 0.036	0.594 ± 0.061	0.677 ± 0.033	0.840 ± 0.013	92.613%
	CPI	0.402 ± 0.082	0.401 ± 0.122	0.811 ± 0.033	0.052 ± 0.003	0.416 ± 0.036	0.734 ± 0.037	0.762 ± 0.165	0.585 ± 0.090	0.815 ± 0.028	42.599%
	BACPI	0.334 ± 0.015	0.502 ± 0.023	0.827 ± 0.006	0.031 ± 0.004	0.658 ± 0.043	0.831 ± 0.020	0.550 ± 0.010	0.700 ± 0.006	0.845 ± 0.002	23.716%
	DeepNC-HGC	0.309 ± 0.025	0.541 ± 0.037	0.841 ± 0.005	0.080 ± 0.003	0.110 ± 0.036	0.667 ± 0.022	0.572 ± 0.011	0.689 ± 0.006	0.844 ± 0.003	27.367%
	DeepNC-GEN	0.270 ± 0.012	0.597 ± 0.017	0.852 ± 0.009	0.135 ± 0.045	-0.509 ± 0.505	0.608 ± 0.037	0.578 ± 0.020	0.685 ± 0.011	0.840 ± 0.003	28.993%
	DrugBAN	0.242 ± 0.007	0.640 ± 0.010	0.869 ± 0.003	0.029 ± 0.003	0.676 ± 0.032	0.832 ± 0.013	<u>0.465 ± 0.018</u>	<u>0.747 ± 0.010</u>	<u>0.862 ± 0.003</u>	5.163%
	GANDTI	0.318 ± 0.018	0.527 ± 0.027	0.844 ± 0.006	0.030 ± 0.002	0.662 ± 0.026	0.831 ± 0.007	0.621 ± 0.012	0.662 ± 0.006	0.836 ± 0.002	27.967%
	BridgeDPI	1.241 ± 1.432	-0.848 ± 2.133	0.827 ± 0.078	0.325 ± 0.109	0.638 ± 0.121	0.857 ± 0.001	0.514 ± 0.011	0.720 ± 0.006	0.861 ± 0.002	66.442%
	ColdDTA	<u>0.220 ± 0.009</u>	<u>0.672 ± 0.014</u>	<u>0.880 ± 0.004</u>	0.110 ± 0.029	-0.224 ± 0.329	0.673 ± 0.079	0.463 ± 0.008	0.748 ± 0.004	0.866 ± 0.001	11.980%
	SubMDTA	0.289 ± 0.012	0.570 ± 0.018	0.841 ± 0.007	<u>0.029 ± 0.002</u>	<u>0.678 ± 0.025</u>	<u>0.836 ± 0.011</u>	0.532 ± 0.032	0.710 ± 0.017	0.852 ± 0.006	17.882%
	IMAE	0.230 ± 0.009	0.657 ± 0.014	0.874 ± 0.004	0.046 ± 0.018	0.484 ± 0.196	0.781 ± 0.056	0.479 ± 0.012	0.739 ± 0.006	0.863 ± 0.002	7.550%
Transformer	CSDTI	0.331 ± 0.012	0.508 ± 0.017	0.832 ± 0.005	0.088 ± 0.004	0.014 ± 0.041	0.628 ± 0.047	0.768 ± 0.021	0.582 ± 0.012	0.805 ± 0.004	41.196%
	TDGraphDTA	0.222 ± 0.005	0.669 ± 0.008	0.653 ± 0.011	0.091 ± 0.019	-0.009 ± 0.209	0.327 ± 0.125	0.497 ± 0.016	0.729 ± 0.009	0.777 ± 0.005	13.827%
	AMMVF	0.377 ± 0.030	0.439 ± 0.044	0.815 ± 0.005	0.075 ± 0.020	0.161 ± 0.221	0.603 ± 0.141	0.682 ± 0.015	0.628 ± 0.008	0.825 ± 0.002	38.448%
	IIFDTI	0.313 ± 0.018	0.534 ± 0.027	0.836 ± 0.008	0.054 ± 0.013	0.398 ± 0.143	0.691 ± 0.050	0.634 ± 0.024	0.655 ± 0.013	0.832 ± 0.006	30.270%
	ICAN	0.371 ± 0.013	0.448 ± 0.020	0.818 ± 0.006	0.089 ± 0.000	-2.052 ± 0.000	0.500 ± 0.000	0.747 ± 0.031	0.593 ± 0.017	0.813 ± 0.004	42.171%
	MolTrans	0.410 ± 0.136	0.390 ± 0.202	0.812 ± 0.039	4.314 ± 2.290	-47.055 ± 25.515	0.540 ± 0.021	0.695 ± 0.183	0.621 ± 0.100	0.822 ± 0.009	87.119%
	TransformerCPI	0.393 ± 0.022	0.415 ± 0.032	0.802 ± 0.008	0.070 ± 0.003	0.217 ± 0.033	0.800 ± 0.002	0.659 ± 0.040	0.641 ± 0.022	0.829 ± 0.013	37.790%
	MRBDTA	0.241 ± 0.005	0.640 ± 0.008	0.870 ± 0.007	0.050 ± 0.005	0.360 ± 0.058	0.735 ± 0.015	0.507 ± 0.006	0.724 ± 0.003	0.862 ± 0.002	12.531%
	FOTFCPI	0.305 ± 0.012	0.546 ± 0.018	0.839 ± 0.009	0.229 ± 0.180	-1.555 ± 2.003	0.587 ± 0.086	0.567 ± 0.008	0.695 ± 0.004	0.848 ± 0.006	36.603%
	Our combos	<b>0.211 ± 0.007</b>	<b>0.685 ± 0.011</b>	<b>0.886 ± 0.004</b>	<b>0.026 ± 0.004</b>	<b>0.710 ± 0.051</b>	<b>0.849 ± 0.023</b>	<b>0.461 ± 0.006</b>	<b>0.749 ± 0.003</b>	<b>0.869 ± 0.002</b>	0.000%

Table 2: Classification task benchmark on Human, *C.elegans*, and DrugBank datasets, respectively. For the GraphDTA and GraphCPI, we only show the one that has the overall best performance. The best result is highlighted in bold, and the runner-up is underlined. *Avg. Improvement of Accuracy* is computed by the average (across 3 datasets) of the differences between our model’s accuracy and each model’s accuracy, divided by the average (across 3 datasets) of each model’s accuracy, respectively.

Category	Models	Human			C.elegans			Drugbank			Avg. Improvement of Accuracy (%)
		ROC-AUC	Accuracy	F1	ROC-AUC	Accuracy	F1	ROC-AUC	Accuracy	F1	
GNN	GraphDTA-GIN	0.949 ± 0.007	0.885 ± 0.011	0.869 ± 0.011	0.977 ± 0.003	0.929 ± 0.005	0.915 ± 0.005	0.850 ± 0.001	0.783 ± 0.006	0.785 ± 0.005	3.312%
	GraphCPI-GIN	0.941 ± 0.005	0.874 ± 0.007	0.858 ± 0.008	0.971 ± 0.003	0.924 ± 0.008	0.907 ± 0.009	0.838 ± 0.012	0.775 ± 0.010	0.778 ± 0.006	4.275%
	MGraphDTA	0.960 ± 0.004	0.905 ± 0.007	0.893 ± 0.007	0.983 ± 0.002	0.943 ± 0.004	0.931 ± 0.004	0.879 ± 0.004	0.800 ± 0.004	0.806 ± 0.003	1.322%
	SAGDTA	0.957 ± 0.004	0.901 ± 0.005	0.887 ± 0.006	0.966 ± 0.006	0.912 ± 0.014	0.894 ± 0.017	0.819 ± 0.009	0.752 ± 0.010	0.756 ± 0.010	4.600%
	EmbedDTI	0.958 ± 0.003	0.901 ± 0.005	0.888 ± 0.006	0.975 ± 0.003	0.924 ± 0.002	0.908 ± 0.002	0.815 ± 0.007	0.758 ± 0.005	0.765 ± 0.003	3.871%
	DeepGLSTM	0.958 ± 0.004	0.903 ± 0.007	0.890 ± 0.008	0.975 ± 0.004	0.923 ± 0.006	0.906 ± 0.007	0.796 ± 0.014	0.745 ± 0.007	0.752 ± 0.006	4.356%
	CPI	0.951 ± 0.012	0.900 ± 0.010	0.887 ± 0.012	0.955 ± 0.005	0.913 ± 0.007	0.893 ± 0.010	0.739 ± 0.087	0.678 ± 0.072	0.687 ± 0.074	7.708%
	BACPI	0.947 ± 0.003	0.905 ± 0.007	0.893 ± 0.008	0.975 ± 0.003	0.936 ± 0.005	0.921 ± 0.006	0.849 ± 0.004	0.776 ± 0.009	0.782 ± 0.008	2.522%
	DeepNC-HGC	0.932 ± 0.009	0.861 ± 0.015	0.845 ± 0.016	0.970 ± 0.003	0.918 ± 0.004	0.903 ± 0.006	0.809 ± 0.006	0.752 ± 0.006	0.762 ± 0.006	6.006%
	DeepNC-GEN	0.961 ± 0.002	0.907 ± 0.006	0.894 ± 0.006	0.980 ± 0.002	0.932 ± 0.005	0.917 ± 0.008	0.813 ± 0.007	0.736 ± 0.015	0.756 ± 0.007	4.194%
	DrugBAN	0.974 ± 0.002	0.920 ± 0.005	0.910 ± 0.005	0.982 ± 0.002	0.946 ± 0.004	0.935 ± 0.005	0.876 ± 0.004	0.799 ± 0.008	0.801 ± 0.005	0.675%
	GANDTI	0.970 ± 0.002	0.917 ± 0.004	0.906 ± 0.004	0.967 ± 0.003	0.919 ± 0.007	0.901 ± 0.007	0.836 ± 0.014	0.752 ± 0.008	0.763 ± 0.004	3.671%
	BridgeDPI	0.957 ± 0.012	0.887 ± 0.021	0.877 ± 0.020	0.960 ± 0.004	0.882 ± 0.034	0.857 ± 0.040	0.726 ± 0.076	0.644 ± 0.087	0.685 ± 0.047	11.189%
	ColdDTA	0.971 ± 0.002	0.922 ± 0.009	0.912 ± 0.010	0.983 ± 0.003	0.947 ± 0.002	0.936 ± 0.002	<b>0.885 ± 0.004</b>	<b>0.813 ± 0.005</b>	<b>0.816 ± 0.004</b>	0.037%
	SubMDTA	0.971 ± 0.003	0.919 ± 0.006	0.909 ± 0.007	<u>0.985 ± 0.001</u>	0.945 ± 0.007	0.933 ± 0.008	0.861 ± 0.005	0.791 ± 0.005	0.793 ± 0.005	1.055%
	IMAE	0.944 ± 0.004	0.878 ± 0.005	0.863 ± 0.003	0.967 ± 0.004	0.911 ± 0.007	0.892 ± 0.007	0.847 ± 0.004	0.777 ± 0.005	0.780 ± 0.004	4.560%
Transformer	CSDTI	0.905 ± 0.007	0.846 ± 0.007	0.826 ± 0.009	0.910 ± 0.006	0.840 ± 0.011	0.805 ± 0.010	0.774 ± 0.011	0.721 ± 0.006	0.730 ± 0.004	11.467%
	TDGraphDTA	0.977 ± 0.002	0.927 ± 0.005	0.917 ± 0.005	0.984 ± 0.001	0.943 ± 0.007	0.929 ± 0.010	<u>0.880 ± 0.006</u>	<u>0.805 ± 0.006</u>	<u>0.810 ± 0.003</u>	0.299%
	AMMVF	0.962 ± 0.005	0.915 ± 0.007	0.905 ± 0.009	0.984 ± 0.005	0.948 ± 0.006	0.937 ± 0.007	0.692 ± 0.161	0.654 ± 0.088	0.696 ± 0.020	6.595%
	IIFDTI	0.973 ± 0.006	0.920 ± 0.006	0.909 ± 0.008	<b>0.987 ± 0.002</b>	0.948 ± 0.005	0.937 ± 0.005	0.849 ± 0.014	0.777 ± 0.010	0.782 ± 0.011	1.437%
	ICAN	0.971 ± 0.002	0.927 ± 0.005	0.917 ± 0.005	0.977 ± 0.004	0.942 ± 0.003	0.929 ± 0.004	0.839 ± 0.005	0.764 ± 0.005	0.768 ± 0.004	1.899%
	MolTrans	0.979 ± 0.003	0.931 ± 0.002	0.923 ± 0.002	0.980 ± 0.003	0.943 ± 0.004	0.930 ± 0.004	0.868 ± 0.004	0.795 ± 0.005	0.795 ± 0.009	0.525%
	TransformerCPI	0.968 ± 0.003	0.917 ± 0.004	0.906 ± 0.004	0.984 ± 0.001	0.941 ± 0.005	0.929 ± 0.006	0.874 ± 0.007	0.799 ± 0.008	0.803 ± 0.007	0.789%
	MRBDTA	0.971 ± 0.004	0.920 ± 0.007	0.909 ± 0.007	<u>0.985 ± 0.002</u>	<u>0.953 ± 0.002</u>	<u>0.943 ± 0.002</u>	0.866 ± 0.005	0.789 ± 0.006	0.790 ± 0.004	0.789%
	FOTFCPI	<u>0.980 ± 0.003</u>	<b>0.937 ± 0.006</b>	<b>0.929 ± 0.006</b>	<b>0.987 ± 0.001</b>	<u>0.953 ± 0.003</u>	0.942 ± 0.004	0.866 ± 0.002	0.790 ± 0.004	0.793 ± 0.006	0.112%
	Our combos	<b>0.981 ± 0.003</b>	<u>0.936 ± 0.007</u>	<u>0.928 ± 0.008</u>	<b>0.987 ± 0.003</b>	<b>0.954 ± 0.005</b>	<b>0.944 ± 0.006</b>	0.866 ± 0.007	0.793 ± 0.009	0.798 ± 0.005	0.000%

**Obs. 8.** In general, the memory usage of GNN-based methods is smaller than that of Transformer-based methods, which is positively proportional to run time. This difference is primarily due to the self-attention mechanism employed in Transformers, which requires significant memory resources. In contrast, model parameters, such as those in DeepGLSTM, do not exhibit a direct relationship with either runtime or performance.

**★Benchmark over Convergence.** We select the two representative methods from the GNN-based and Transformer-based frameworks, respectively, and evaluate them across six datasets on two tasks. The training



losses are depicted in Fig. 4. To ensure a fair comparison of convergence behavior, we use the previous early stopping setting. Based on the empirical results, we summarize our key observations as follows:

**Obs. 9. GNN-based methods demonstrate quicker convergence compared to Transformer-based methods.** This phenomenon arises from the fact that GNN-based methods have less memory usage and fewer model parameters, leading to larger batch size usage or faster convergence compared with Transformer-based methods.

### 5.1 Our Best Combo of Drug and Protein Encoders

**Deriving Combo from Benchmark Insights.** Based on our benchmark results, we summarize the insights of protein and drug encoder usages and propose a light yet effective architecture, which could be treated as a new strong baseline for future explorations. *Regarding the proteins*, we observe that multi-scale CNNs associated with a mixture of model depths can generally learn the effective protein representations (Yang et al., 2022; Zhu et al., 2023; Fang et al., 2023), which approximate the language model’s accuracy while having lower memory and computation costs. *Regarding the drug molecules*, both GNN- and Transformer-based methods, such as MRBDTA (Zhang et al., 2022), MolTrans (Huang et al., 2020b), and MGraphDTA (Yang et al., 2022) prove promising in DTI tasks. This encourages us to leverage information from hybrid perspectives, i.e., implicit structure (via attention in Transformers) and explicit structure learning (via message passing along edges in GNNs).

Our model design, illustrated in Fig. 5, integrates these components. Specifically, for drug graphs, we adopt a hybrid network that augments the self-attention mechanism with inter-atomic distances and graph adjacency matrices (Maziarka et al., 2020), incorporating both 2D and 3D molecular structural information. Given the projections of molecular input at an attention head, i.e.  $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ , the adjacent matrix  $\mathbf{A} \in \{0, 1\}^{N \times N}$ , and the inter-atomic distances matrix  $\mathbf{D} \in \mathbb{R}^{N \times N}$  obtained using RDkit, the augmented attention is calculated as follows:

$$\text{Multi-Attn} = (\lambda_a \cdot \text{softmax}(\mathbf{QK}^T / \sqrt{d}) + \lambda_d g(\mathbf{D}) + \lambda_g \mathbf{A}) \mathbf{V}, \quad (6)$$

where  $g(\cdot)$  is a row-wise softmax function, and  $\lambda_a, \lambda_d$  and  $\lambda_g$  denote scalars weighting the self-attention, distance, and adjacency matrices, respectively. Besides the implicit and explicit structure learning, we integrate the features from drug SMILES. It is notable that simply utilizing the SMILES representation extracted from a Transformer for downstream tasks does not perform as well as GNN. To align with the protein embedding paradigm, we adopt a simple CNN to unearth potential SMILES information, as suggested in Zhao et al. (2021). Subsequently, due to the fact that cross-attention is more complex and hard to optimize, we implement a straightforward attention mechanism to integrate the representations of the drug graph and SMILES, denoted as  $\mathbf{f}_G$  and  $\mathbf{f}_S$ , respectively, using a weighting parameter  $\lambda$ , as follows:

$$\mathbf{f}_D = \lambda \cdot \mathbf{f}_G + (1 - \lambda) \cdot \mathbf{f}_S, \quad \lambda = \text{MLP}(\text{MLP}(\mathbf{f}_G) + \text{MLP}(\mathbf{f}_S)). \quad (7)$$

Finally, the prediction is obtained by processing the concatenated protein and drug representations through a task-relevant head, as shown in Fig. 5

**Novelty.** As opposed to the previous strategy, which heuristically stacked a large amount of modules of different types, our model design is driven by systematic benchmarking and empirical insights. Through extensive experiments under fair and controlled conditions, we identify key encoder and featurization strategies that consistently outperform others. Notably, we disentangle and quantify the distinct molecular features. For instance, atomic properties and hydrogen information significantly enhance predictive performance, while adding electron properties may introduce noise.

Thus, our combo is not merely an ad hoc combination but a carefully validated design that effectively obtains a superior balance between accuracy and computational efficiency on both classification and regression tasks. The clear empirical guidance to model design offered through this study helps to establish a more principled framework for future work in drug-target interaction modeling and provides a robust, reproducible new baseline for the community.

**Benchmark Comparison to State-of-the-Art Frameworks.** We compare the proposed combos with the SOTA frameworks in Tables 1 and 2, and Figures 3 and 4. It is observed that our model consistently



Figure 3: Model size, memory usage and run time on Davis.

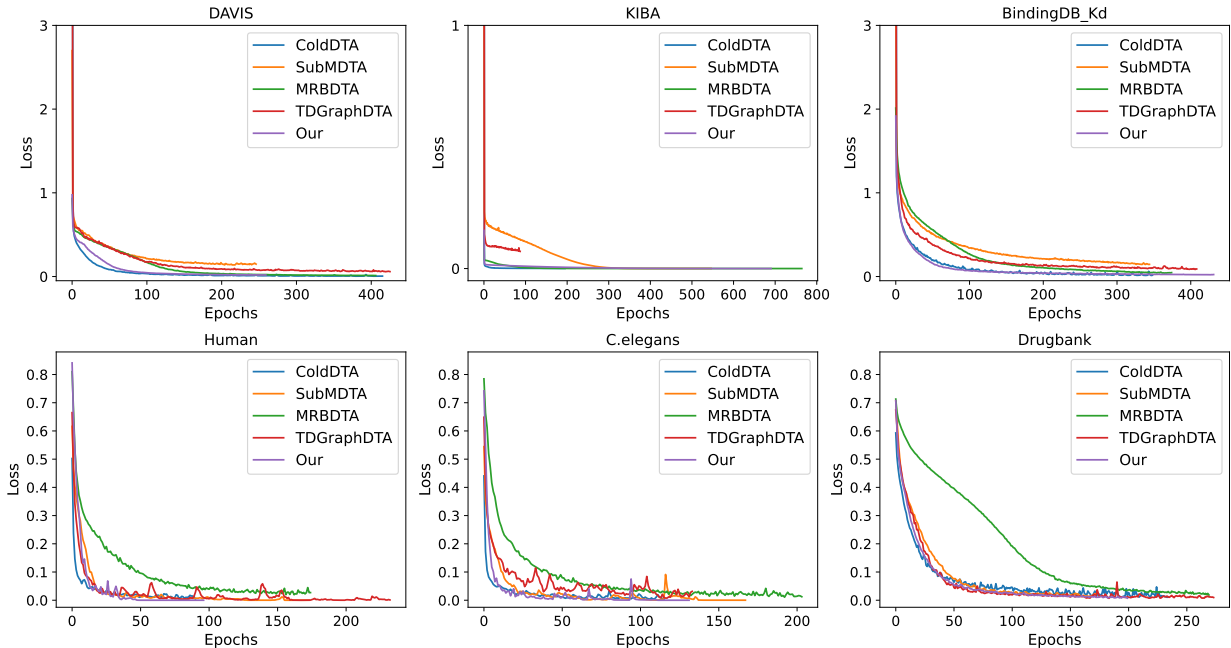


Figure 4: Convergence curves for five selected methods.

achieves the best performance in the regression tasks across three datasets and nearly outperforms most methods in classification tasks. By leveraging the physical conformation information from the molecular graph, our combos converge faster than the other two Transformer-based methods, MRBDTA (Zhang et al., 2022) and TDGraphDTA (Zhu et al., 2023), particularly on the KIBA dataset. Moreover, our model uses

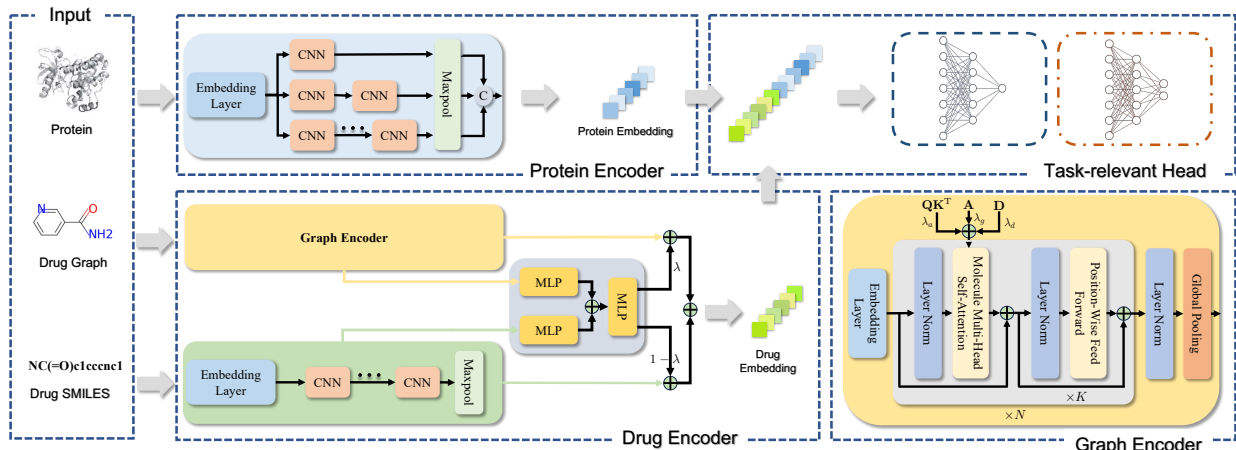


Figure 5: Overview of our proposed model combos.

three times less peak memory and fewer parameters than other Transformer-based methods, enabling faster computation and reduced storage requirements.

## 6 Conclusion

In this work, we establish a benchmark with fair and consistent experimental configurations, aiming to push DTI research, particularly emphasizing the utilization of structural information. Our meticulous approach has entailed thorough exploration of diverse encoder strategies and featurization techniques for both drug molecules and proteins. Moreover, dozens of existing approaches across six representative datasets for both regression and classification tasks are investigated on various metrics, including DTI classification and regression accuracy, peak memory usage, and model convergence. Provided with the comprehensive benchmark results, we propose a novel approach that integrates the strengths of GNN and Transformer-based methods. Our studies on benchmarking and rethinking help lay a solid, practical, and systematic foundation for the DTI community and provide researchers with broader and deeper insights into the intricate dynamics of drug-target interactions.

## References

- Karim Abbasi, Parvin Razzaghi, Antti Poso, Saber Ghanbari-Ara, and Ali Masoudi-Nejad. Deep learning in drug target interaction prediction: current and future perspectives. *Current Medicinal Chemistry*, 28(11):2100–2113, 2021.
- Ehsaneddin Asgari and Mohammad R. K. Mofrad. Continuous distributed representation of biological sequences for deep proteomics and genomics. *PLOS ONE*, 10(11):1–15, 11 2015. doi: 10.1371/journal.pone.0141287. URL <https://doi.org/10.1371/journal.pone.0141287>.
- Ted T. Ashburn and Karl B. Thor. Drug repositioning: identifying and developing new uses for existing drugs. *Nature Reviews Drug Discovery*, 3(8):673–683, Aug 2004. ISSN 1474-1784. doi: 10.1038/nrd1468. URL <https://doi.org/10.1038/nrd1468>.
- Peizhen Bai, Filip Miljković, Bino John, and Haiping Lu. Interpretable bilinear attention network with domain adaptation improves drugtarget prediction. *Nature Machine Intelligence*, 5(2):126–136, 2023.
- Song Bai, Feihu Zhang, and Philip H.S. Torr. Hypergraph convolution and hypergraph attention. *Pattern Recognition*, 110:107637, 2021. ISSN 0031-3203. doi: <https://doi.org/10.1016/j.patcog.2020.107637>. URL <https://www.sciencedirect.com/science/article/pii/S0031320320304404>.

- Rakesh Bal, Yijia Xiao, and Wei Wang. PGraphDTA: Improving Drug Target Interaction Prediction using Protein Language Models and Contact Maps, 2024.
- Lifan Chen, Xiaoqin Tan, Dingyan Wang, Feisheng Zhong, Xiaohong Liu, Tianbiao Yang, Xiaomin Luo, Kaixian Chen, Hualiang Jiang, and Mingyue Zheng. TransformerCPI: improving compound-protein interaction prediction by sequence-based deep learning with self-attention mechanism and label reversal experiments. *Bioinformatics*, 36(16):4406–4414, 05 2020. ISSN 1367-4803. doi: 10.1093/bioinformatics/btaa524. URL <https://doi.org/10.1093/bioinformatics/btaa524>.
- Zhongjian Cheng, Qichang Zhao, Yaohang Li, and Jianxin Wang. IIFDTI: predicting drug-target interactions through interactive and independent features based on attention mechanism. *Bioinformatics*, 38(17):4153–4161, 07 2022. ISSN 1367-4803. doi: 10.1093/bioinformatics/btac485. URL <https://doi.org/10.1093/bioinformatics/btac485>.
- Robert A Copeland, David L Pompliano, and Thomas D Meek. Drug–target residence time and its implications for lead optimization. *Nature reviews Drug discovery*, 5(9):730–739, 2006.
- Fabrizio Costa and Kurt De Grave. Fast neighborhood subgraph pairwise distance kernel. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, ICML’10, pp. 255262, Madison, WI, USA, 2010. Omnipress. ISBN 9781605589077.
- 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.
- Qi-wen Dong, Xiao-long Wang, and Lei Lin. Application of latent semantic analysis to protein remote homology detection. *Bioinformatics*, 22(3):285–290, 11 2005. ISSN 1367-4803. doi: 10.1093/bioinformatics/bti801. URL <https://doi.org/10.1093/bioinformatics/bti801>.
- Keyu Duan, Zirui Liu, Peihao Wang, Wenqing Zheng, Kaixiong Zhou, Tianlong Chen, Xia Hu, and Zhangyang Wang. A comprehensive study on large-scale graph training: Benchmarking and rethinking. *Advances in Neural Information Processing Systems*, 35:5376–5389, 2022.
- S Ekins, J Mestres, and B Testa. In silico pharmacology for drug discovery: applications to targets and beyond. *British journal of pharmacology*, 152(1):21–37, 2007.
- Kejie Fang, Yiming Zhang, Shiyu Du, and Jian He. ColdDTA: Utilizing data augmentation and attention-based feature fusion for drug-target binding affinity prediction. *Computers in Biology and Medicine*, 164: 107372, 2023. ISSN 0010-4825. doi: <https://doi.org/10.1016/j.compbiomed.2023.107372>. URL <https://www.sciencedirect.com/science/article/pii/S0010482523008375>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Jiayue Hu, Wang Yu, Chao Pang, Junru Jin, Nhat Truong Pham, Balachandran Manavalan, and Leyi Wei. DrugormerDTI: Drug Graphormer for drug-target interaction prediction. *Computers in Biology and Medicine*, 161:106946, 2023. ISSN 0010-4825. doi: <https://doi.org/10.1016/j.compbiomed.2023.106946>. URL <https://www.sciencedirect.com/science/article/pii/S0010482523004110>.
- Kexin Huang, Tianfan Fu, Lucas M Glass, Marinka Zitnik, Cao Xiao, and Jimeng Sun. Deeppurpose: a deep learning library for drug–target interaction prediction. *Bioinformatics*, 36(22-23):5545–5547, 2020a.
- Kexin Huang, Cao Xiao, Lucas M Glass, and Jimeng Sun. MolTrans: Molecular Interaction Transformer for drug-target interaction prediction. *Bioinformatics*, 37(6):830–836, 10 2020b. ISSN 1367-4803. doi: 10.1093/bioinformatics/btaa880. URL <https://doi.org/10.1093/bioinformatics/btaa880>.
- Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor W Coley, Cao Xiao, Jimeng Sun, and Marinka Zitnik. Therapeutics data commons: Machine learning datasets and tasks for drug discovery and development. *arXiv preprint arXiv:2102.09548*, 2021.

- James P Hughes, Stephen Rees, S Barrett Kalindjian, and Karen L Philpott. Principles of early drug discovery. *British journal of pharmacology*, 162(6):1239–1249, 2011.
- Yuan Jin, Jiarui Lu, Runhan Shi, and Yang Yang. EmbedDTI: Enhancing the Molecular Representations via Sequence Embedding and Graph Convolutional Network for the Prediction of Drug-Target Interaction. *Biomolecules*, 11(12), 2021. ISSN 2218-273X. doi: 10.3390/biom11121783. URL <https://www.mdpi.com/2218-273X/11/12/1783>.
- Martin Karplus and John Kuriyan. Molecular dynamics and protein function. *Proceedings of the National Academy of Sciences*, 102(19):6679–6685, 2005.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks, 2017.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- Hiroyuki Kurata and Sho Tsukiyama. ICAN: Interpretable cross-attention network for identifying drug and target protein interactions. *PLOS ONE*, 17(10):e0276609, 2022. doi: 10.1371/journal.pone.0276609. URL <https://app.dimensions.ai/details/publication/pub.1152144662>. <https://journals.plos.org/plosone/article/file?id=10.1371/journal.pone.0276609type=printable>.
- Greg Landrum et al. Rdkit: Open-source cheminformatics, 2006. URL <https://www.rdkit.org>.
- Guohao Li, Chenxin Xiong, Ali Thabet, and Bernard Ghanem. Deepergcn: All you need to train deeper gcns, 2020.
- Min Li, Zhangli Lu, Yifan Wu, and YaoHang Li. BACPI: a bi-directional attention neural network for compound-protein interaction and binding affinity prediction. *Bioinformatics*, 38(7):1995–2002, 01 2022. ISSN 1367-4803. doi: 10.1093/bioinformatics/btac035. URL <https://doi.org/10.1093/bioinformatics/btac035>.
- Mufei Li, Jinjing Zhou, Jiajing Hu, Wenxuan Fan, Yangkang Zhang, Yaxin Gu, and George Karypis. Dgl-lifesci: An open-source toolkit for deep learning on graphs in life science. *ACS omega*, 6(41):27233–27238, 2021.
- Xuan Lin, Kaiqi Zhao, Tong Xiao, Zhe Quan, Zhi-Jie Wang, and Philip S. Yu. DeepGS: Deep Representation Learning of Graphs and Sequences for Drug-Target Binding Affinity Prediction. *arXiv preprint arXiv:2003.13902*, abs/2003.13902, 2020. URL <https://api.semanticscholar.org/CorpusID:214728430>.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Sal Candido, et al. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv*, 2022.
- Hui Liu, Jianjiang Sun, Jihong Guan, Jie Zheng, and Shuigeng Zhou. Improving compound-protein interaction prediction by building up highly credible negative samples. *Bioinformatics*, 31(12):i221–i229, 2015.
- 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.
- Dik-Lung Ma, Daniel Shiu-Hin Chan, Paul Lee, Maria Hiu-Tung Kwan, and Chung-Hang Leung. Molecular modeling of drug-dna interactions: virtual screening to structure-based design. *Biochimie*, 93(8):1252–1266, 2011.

- Łukasz Maziarka, Tomasz Danel, Sławomir Mucha, Krzysztof Rataj, Jacek Tabor, and Stanisław Jastrzębski. Molecule attention transformer. *arXiv preprint arXiv:2002.08264*, 2020.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representations of words and phrases and their compositionality. In *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2*, NIPS’13, pp. 3111–3119, Red Hook, NY, USA, 2013. Curran Associates Inc.
- H. L. Morgan. The generation of a unique machine description for chemical structures—a technique developed at chemical abstracts service. *Journal of Chemical Documentation*, 5(2):107–113, 1965. doi: 10.1021/c160017a018. URL <https://doi.org/10.1021/c160017a018>.
- Shrimon Mukherjee, Madhusudan Ghosh, and Partha Basuchowdhuri. DeepGLSTM: Deep Graph Convolutional Network and LSTM based approach for predicting drug-target binding affinity. In *Proceedings of the 2022 SIAM International Conference on Data Mining (SDM)*, pp. 729–737. SIAM, 2022. doi: 10.1137/1.9781611977172.82.
- Thin Nguyen, Hang Le, Thomas P Quinn, Tri Nguyen, Thuc Duy Le, and Svetha Venkatesh. GraphDTA: predicting drugtarget binding affinity with graph neural networks. *Bioinformatics*, 37(8):1140–1147, 10 2020. ISSN 1367-4803. doi: 10.1093/bioinformatics/btaa921. URL <https://doi.org/10.1093/bioinformatics/btaa921>.
- Hakime Öztürk, Arzucan Özgür, and Elif Ozkirimli. Deepdta: deep drug–target binding affinity prediction. *Bioinformatics*, 34(17):i821–i829, 2018.
- Hakime Öztürk, Elif Ozkirimli, and Arzucan Özgür. Widedta: prediction of drug-target binding affinity. *arXiv preprint arXiv:1902.04166*, 2019.
- Shourun Pan, Leiming Xia, Lei Xu, and Zhen Li. SubMDTA: drug target affinity prediction based on substructure extraction and multi-scale features. *BMC Bioinformatics*, 24(1):334, 2023a.
- Yaohua Pan, Yijia Zhang, Jing Zhang, and Mingyu Lu. CSDTI: an interpretable cross-attention network with GNN-based drug molecule aggregation for drug-target interaction prediction. *Applied Intelligence*, 53(22):27177–27190, 2023b.
- Ying Qian, Xinyi Li, Jian Wu, and Qian Zhang. MCL-DTI: using drug multimodal information and bi-directional cross-attention learning method for predicting drugtarget interaction. *BMC bioinformatics*, 24(1):323, 2023.
- Zhe Quan, Xuan Lin, Zhi-Jie Wang, Yan Liu, Fan Wang, and Kenli Li. A system for learning atoms based on long short-term memory recurrent neural networks. In *2018 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pp. 728–733, 2018. doi: 10.1109/BIBM.2018.8621313.
- Zhe Quan, Yan Guo, Xuan Lin, Zhi-Jie Wang, and Xiangxiang Zeng. GraphCPI: Graph Neural Representation Learning for Compound-Protein Interaction. In *2019 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pp. 717–722, Nov 2019. doi: 10.1109/BIBM47256.2019.8983267.
- David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling*, 50(5):742–754, 2010. doi: 10.1021/ci100050t. URL <https://doi.org/10.1021/ci100050t>. PMID: 20426451.
- Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie, Ying WEI, Wenbing Huang, and Junzhou Huang. Self-supervised graph transformer on large-scale molecular data. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 12559–12571. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/94aef38441efa3380a3bed3faf1f9d5d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/94aef38441efa3380a3bed3faf1f9d5d-Paper.pdf).
- Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.

- Stephen M. Strittmatter. Overcoming Drug Development Bottlenecks With Repurposing: Old drugs learn new tricks. *Nature Medicine*, 20(6):590–591, Jun 2014. ISSN 1546-170X. doi: 10.1038/nm.3595. URL <https://doi.org/10.1038/nm.3595>.
- David C Swinney. Biochemical mechanisms of drug action: what does it take for success? *Nature reviews Drug discovery*, 3(9):801–808, 2004.
- 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.
- Huu Ngoc Tran Tran, J Joshua Thomas, and Nurul Hashimah Ahamed Hassain Malim. DeepNC: a framework for drug-target interaction prediction with graph neural networks. *PeerJ*, 10:e13163, 2022.
- Masashi Tsubaki, Kentaro Tomii, and Jun Sese. Compoundprotein interaction prediction with end-to-end learning of neural networks for graphs and sequences. *Bioinformatics*, 35(2):309–318, 07 2018. ISSN 1367-4803. doi: 10.1093/bioinformatics/bty535. URL <https://doi.org/10.1093/bioinformatics/bty535>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, pp. 60006010, Red Hook, NY, USA, 2017a. Curran Associates Inc. ISBN 9781510860964.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017b.
- Petar Velikovi, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks, 2018.
- Lu Wang, Yifeng Zhou, and Qu Chen. AMMVF-DTI: A Novel Model Predicting DrugTarget Interactions Based on Attention Mechanism and Multi-View Fusion. *International Journal of Molecular Sciences*, 24(18), 2023. ISSN 1422-0067. doi: 10.3390/ijms241814142. URL <https://www.mdpi.com/1422-0067/24/18/14142>.
- Minjie Wang, Da Zheng, Zihao Ye, Quan Gan, Mufei Li, Xiang Song, Jinjing Zhou, Chao Ma, Lingfan Yu, Yu Gai, et al. Deep graph library: A graph-centric, highly-performant package for graph neural networks. *arXiv preprint arXiv:1909.01315*, 2019.
- Shuyu Wang, Peng Shan, Yuliang Zhao, and Lei Zuo. GanDTI: A multi-task neural network for drug-target interaction prediction. *Computational Biology and Chemistry*, 92:107476, 2021. ISSN 1476-9271. doi: <https://doi.org/10.1016/j.compbiolchem.2021.107476>. URL <https://www.sciencedirect.com/science/article/pii/S1476927121000438>.
- David Weininger. SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *Journal of Chemical Information and Computer Sciences*, 28(1):31–36, 1988a. doi: 10.1021/ci00057a005. URL <https://doi.org/10.1021/ci00057a005>.
- David Weininger. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36, 1988b.
- Ming Wen, Zhimin Zhang, Shaoyu Niu, Haozhi Sha, Ruihan Yang, Yonghuan Yun, and Hongmei Lu. Deep-learning-based drug–target interaction prediction. *Journal of proteome research*, 16(4):1401–1409, 2017.
- David S Wishart, Craig Knox, An Chi Guo, Dean Cheng, Savita Shrivastava, Dan Tzur, Bijaya Gautam, and Murtaza Hassanali. Drugbank: a knowledgebase for drugs, drug actions and drug targets. *Nucleic acids research*, 36(suppl\_1):D901–D906, 2008.

- Yifan Wu, Min Gao, Min Zeng, Jie Zhang, and Min Li. BridgeDPI: a novel Graph Neural Network for predicting drugprotein interactions. *Bioinformatics*, 38(9):2571–2578, 03 2022. ISSN 1367-4803. doi: 10.1093/bioinformatics/btac155. URL <https://doi.org/10.1093/bioinformatics/btac155>.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2018.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks?, 2019.
- Minghao Xu, Zuobai Zhang, Jiarui Lu, Zhaocheng Zhu, Yangtian Zhang, Ma Chang, Runcheng Liu, and Jian Tang. Peer: a comprehensive and multi-task benchmark for protein sequence understanding. *Advances in Neural Information Processing Systems*, 35:35156–35173, 2022.
- Ziduo Yang, Weihe Zhong, Lu Zhao, and Calvin Yu-Chian Chen. MGraphDTA: deep multiscale graph neural network for explainable drugtarget binding affinity prediction. *Chemical Science*, 13(3):816–833, 2022. doi: 10.1039/D1SC05180F. URL <http://dx.doi.org/10.1039/D1SC05180F>.
- Zeyu Yin, Yu Chen, Yajie Hao, Sanjeevi Pandiyan, Jinsong Shao, and Li Wang. FOTF-CPI: A compound-protein interaction prediction transformer based on the fusion of optimal transport fragments. *Isience*, 27(1), 2024.
- Jing Zhang, Zhi Liu, Yaohua Pan, Hongfei Lin, and Yijia Zhang. IMAEN: An interpretable molecular augmentation model for drugtarget interaction prediction. *Expert Systems with Applications*, 238:121882, 2024. ISSN 0957-4174. doi: <https://doi.org/10.1016/j.eswa.2023.121882>. URL <https://www.sciencedirect.com/science/article/pii/S0957417423023849>.
- Li Zhang, Chun-Chun Wang, and Xing Chen. Predicting drugtarget binding affinity through molecule representation block based on multi-head attention and skip connection. *Briefings in Bioinformatics*, 23(6):bbac468, 10 2022. ISSN 1477-4054. doi: 10.1093/bib/bbac468. URL <https://doi.org/10.1093/bib/bbac468>.
- Shugang Zhang, Mingjian Jiang, Shuang Wang, Xiaofeng Wang, Zhiqiang Wei, and Zhen Li. SAG-DTA: Prediction of DrugTarget Affinity Using Self-Attention Graph Network. *International Journal of Molecular Sciences*, 22(16), 2021. ISSN 1422-0067. doi: 10.3390/ijms22168993. URL <https://www.mdpi.com/1422-0067/22/16/8993>.
- Qichang Zhao, Haochen Zhao, Kai Zheng, and Jianxin Wang. HyperAttentionDTI: improving drugprotein interaction prediction by sequence-based deep learning with attention mechanism. *Bioinformatics*, 38(3): 655–662, 10 2021. ISSN 1367-4803. doi: 10.1093/bioinformatics/btab715. URL <https://doi.org/10.1093/bioinformatics/btab715>.
- Zhiqin Zhu, Zheng Yao, Xin Zheng, Guanqiu Qi, Yuanyuan Li, Neal Mazur, Xinbo Gao, Yifei Gong, and Baisen Cong. Drugtarget affinity prediction method based on multi-scale information interaction and graph optimization. *Computers in Biology and Medicine*, 167:107621, 2023. ISSN 0010-4825. doi: <https://doi.org/10.1016/j.combiomed.2023.107621>. URL <https://www.sciencedirect.com/science/article/pii/S0010482523010867>.



## A Related Works

**GNN-based Methods** GNNs play a crucial role in mining the intricate features of drug molecules for drug-target prediction. Numerous models, including Graph Convolutional Network (GCN), Graph Isomorphism Network (GIN), and Graph Attention Network (GAT), have been utilized (Nguyen et al., 2020; Quan et al., 2019; Wang et al., 2023; Lin et al., 2020; Jin et al., 2021) to process and enhance drug features. Additionally, MGraphDTA (Yang et al., 2022) employs a multi-scale GNN architecture, while DeepGLSTM (Mukherjee et al., 2022) leverages parallel GNN structures for drug representation. DeepNC integrates advanced techniques from generalized aggregation networks (Li et al., 2020) and hypergraph convolution (Bai et al., 2021) to improve feature extraction. BACPI (Li et al., 2022) develops a bi-directional attention network to integrate the representations of drug molecules and proteins, enhancing their mutual interaction. Besides, BridgeDPI (Wu et al., 2022) innovates by incorporating bridging nodes between proteins and drugs, utilizing a three-layer GNN for graph embeddings.

**Transformer-based Methods** Transformers, known for their efficacy in handling sequence data, are extensively applied in drug and protein feature processing. For instance, models like MolTrans (Huang et al., 2020b) and FOTFCPI (Yin et al., 2024) employ self-attention mechanisms to refine embeddings by focusing on drug and protein substructures. MRBDTA (Zhang et al., 2022) uses multi-head attention and skip connection to enhance drug and protein representation. Additionally, a cross-attention mechanism (Pan et al., 2023b; Kurata & Tsukiyama, 2022) is employed to facilitate the integration of drug and protein features, enabling effective mutual querying. TDGraphDTA (Zhu et al., 2023) captures contextual relationships between molecular substructures by using a multi-head cross-attention mechanism and graph optimization. Lastly, DrugormerDTI (Hu et al., 2023) incorporates degree centrality with positional information to highlight the positional relevance of amino acids in proteins.

**Input and Featurization** Structural information is crucial at the input stage for models such as BridgeDPI (Wu et al., 2022). Various libraries, such as DGLGraph (Wang et al., 2019), DGL-lifeSci (Li et al., 2021), and RDKit (Landrum et al., 2006), are employed to process input SMILES of drugs, with RDKit (Landrum et al., 2006) being pivotal for converting SMILE strings into molecular graphs and extracting diverse chemical properties, including chemical bonds, hydrogen presence, electron properties, and so on. Additionally, some approaches (Wang et al., 2023; Lin et al., 2020; Li et al., 2022; Wang et al., 2021) incorporate molecular fingerprints (Rogers & Hahn, 2010) to capture local chemical information. For protein sequences, typical preprocessing involves converting amino acid sequences into N-grams (Pan et al., 2023a; Dong et al., 2005) or integer (Nguyen et al., 2020) sequences. To enhance the expressiveness of embeddings, some models leverage pre-trained Word2Vec (Mikolov et al., 2013; Quan et al., 2019; Wang et al., 2023; Li et al., 2022; Tsubaki et al., 2018; Lin et al., 2020; Cheng et al., 2022) or pre-trained protein language models (Bal et al., 2024).

## B Model Descriptions

This section provides a comprehensive overview of 31 DTI methods, which are classified into GNN-based and Transformer-based approaches. The DTI framework can be simplified as using two encoders to process drugs and proteins separately, followed by an MLP to handle the integrated representations.

### B.1 GNN-based Methods

#### B.1.1 GCN

★ *GraphDTA-GCN* (Nguyen et al., 2020): GraphDTA-GCN uses GCN to process the molecular graph, which is derived from SMILES using the RDkit tool, and a simple CNN with integer encoding to handle protein sequences.

★ *GraphCPI-GCN* (Quan et al., 2019): Similar to GraphDTA, GraphCPI-GCN employs 3-gram encoding with pretrained Word2Vec to process protein sequences, followed by a CNN to handle the protein embeddings.

★ *MGraphDTA* (Yang et al., 2022): MGraphDTA utilizes a multiscale GCN, inspired by dense connections, and a multiscale CNN to process drug graphs and protein sequences, respectively.

★ *SAGDTA* (Zhang et al., 2021): Similar to GraphDTA, SAGDTA introduces global or hierarchical pooling after GCN to aggregate node representations weightedly.

★ *EmbedDTI* (Jin et al., 2021): For protein sequences, EmbedDTI leverages GloVe for pretraining amino acid feature embeddings, which are then fed into a CNN. For drugs, it constructs both an atom graph and a substructure graph to capture structural information at different levels, processed by GCN.

★ *DeepGLSTM* (Mukherjee et al., 2022): DeepGLSTM processes molecular graphs using a parallel GCN module composed of three GCNs with different layers. For protein sequences, it adopts a bi-LSTM.

★ *CPI* (Tsubaki et al., 2018): CPI processes drug graphs using GCN. The protein sequence is handled via n-gram with integer encoding, followed by a CNN.

★ *DeepNC* (Tran et al., 2022): DeepNC adopts advanced techniques from generalized aggregation networks and hypergraph convolution, two variants of GCN, to capture the representations of drugs. For protein sequences, it uses a simple CNN.

★ *DrugBAN* (Zhang et al., 2022): DrugBAN employs GCN and CNN blocks to encode molecular graphs and proteins, respectively. Then they use a bilinear attention network module to learn local interactions between the representations of drugs and proteins.

★ *BridgeDPI* (Wu et al., 2022): BridgeDPI innovates by constructing a learnable drugprotein association network, which is processed using a three-layer GNN for graph embeddings. The learned representations for drug and protein pairs are then concatenated for further processing.

★ *ColdDTA* (Fang et al., 2023): ColdDTA removes the subgraphs of drugs. For the model, they adopt the dense GCN and multiscale CNN from MGraphDTA as the encoders for drugs and proteins, respectively. Additionally, an attention-based method is developed to integrate representations for improved prediction.

★ *IMAEN* (Zhang et al., 2024): IMAEN employs a molecular augmentation mechanism to enhance molecular structures by fully aggregating molecular node neighborhood information. It then uses multiscale GCN and CNN for drug and protein processing, respectively.

★ *GanDTI* (Wang et al., 2021): Inspired by residual networks, GanDTI adds the input drug fingerprints to the output of three GCN layers as graph node features and uses summation to get the final drug representation.

#### B.1.2 GAT

★ *GraphDTA-GAT* (Nguyen et al., 2020): GraphDTA-GAT adopts a GAT as the encoder for drugs, while other components remain the same as in GraphDTA-GCN.

- ★ *GraphDTA-GATGCN* (Nguyen et al., 2020): GraphDTA-GATGCN adopts a combination of GAT and GCN as the encoder for drugs, while other components remain the same as in GraphDTA-GCN.
- ★ *GraphCPI-GAT* (Quan et al., 2019): GraphDTA-CPI adopts a GAT as the encoder for drugs, while other components remain the same as in GraphCPI-GCN.
- ★ *GraphCPI-GATGCN* (Quan et al., 2019): GraphCPI-GATGCN adopts a combination of GAT and GCN as the encoder for drugs, while other components remain the same as in GraphCPI-GCN.
- ★ *BACPI* (Li et al., 2022): BACPI adopts a GAT and a CNN for the features of the fingerprints and protein sequence, respectively. These features are then fed into a bidirectional attention neural network to obtain integrated representations.
- ★ *PGraphDTA-CNN* (Bal et al., 2024): PGraphDTA-CNN is a straightforward method that utilizes GAT for drug feature extraction and CNN for protein sequences.

## B.2 GIN

- ★ *GraphDTA-GIN* (Nguyen et al., 2020): GraphDTA-GAT adopts a GAT as the encoder for drugs, while other components remain the same as in GraphDTA-GCN.
- ★ *GraphCPI-GIN* (Quan et al., 2019): GraphDTA-GAT adopts a GAT as the encoder for drugs, while other components remain the same as in GraphDTA-GCN.
- ★ *SubMDTA* (Pan et al., 2023a): SubMDTA utilizes a pretrained GIN encoder obtained through contrastive learning for the molecular graph. For protein sequences, it employs N-gram embedding with different N to extract features at various scales, which are then processed by a BiLSTM.

## B.3 Transformer-based Methods

### B.3.1 Self-attention

- ★ *AMMVF* (Wang et al., 2023): AWMVF introduces the multi-head mechanism to GAT to learn features in different spaces, and the update function is obtained through the concatenation of different heads’ outputs.
- ★ *IIFDTI* (Cheng et al., 2022): IIFDTI model attains the drug matrix and protein matrix and inputs them to the bi-directional encoder-decoder block, which considers both the drug and target directions. The decoder is mainly composed of multi-head attention.
- ★ *MolTrans* (Huang et al., 2020b): MolTrans uses Transformer encoder layers to augment the embedding of substructure sequences of proteins and drugs.
- ★ *FOTFCPI* (Yin et al., 2024): Similar to MolTrans, FOTFCPI uses Transformer encoder layers to extract the features of protein and drug fragments after the embedding layers.
- ★ *TransformerCPI* (Chen et al., 2020): TransformerCPI uses the decoder module of Transformer, which takes in the atom sequence embedding processed by GCN and the protein sequence embedding processed by word2vec and 1D CNN.
- ★ *MRBDTA* (Zhang et al., 2022): In MRBDTA, after the embedding layer, drug sequences are directly fed into a block consisting of three Transformer encoders. The first encoder has a linear layer before it and the following two encoders are parallel. The protein sequence is also processed by a block with a similar structure.

### B.3.2 Cross attention

- ★ *CSDTI* (Pan et al., 2023b): CSDTI uses cross-attention to fuse the deep representations of drugs and proteins. Specifically, the different projections of protein features are used as keys and values, respectively, while the projection of drug features is used as a query.

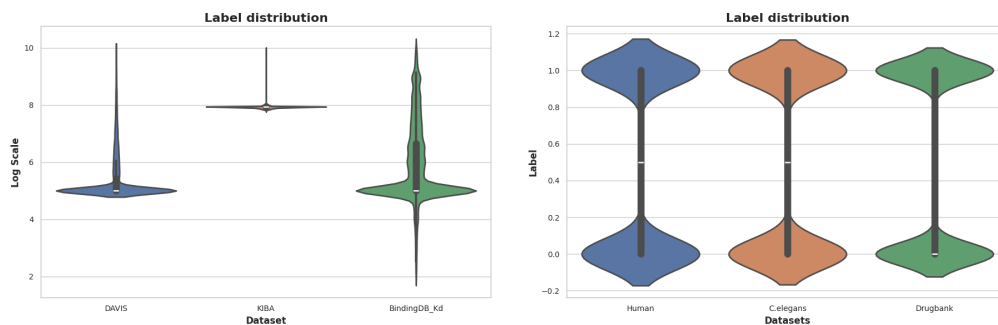
★ *TDGraphDTA*(Zhu et al., 2023): TDGraphDTA uses a multi-head cross-attention mechanism with two attention heads. Both drug and protein features are linearly transformed into query, key, and value matrices. One cross-attention layer uses a drug query matrix, a protein key matrix, and a protein value matrix, while its parallel counterparts use the rest of the matrices. The outputs of these two layers are concatenated and fed into MLP to get the final output.

## C Datasets Descriptions

In this subsection, we provide a detailed description of the datasets for both the regression task and classification task. The statistical characteristics of the datasets are summarized in Table 3. Here we present the statistics after we cleaned the data as described in Section 3.

Table 3: Statistics of the benchmark dataset for two tasks.

	Regression			Classification		
	Davis	KIBA	BindingDB_Kd	Human	<i>C. elegans</i>	DrugBank
Number of drugs	68	2068	10661	2726	1767	6645
Number of target proteins	379	229	1413	2001	1876	4256
Number of total samples	25772	117657	52284	5997	6552	34740



Label distribution of DVAIS, KIBA and Binding\_Kd for regression tasks. Label distribution of Human, *C. elegans* and Drugbank for classification tasks.

Figure 6: Label distribution of different datasets for two tasks.

## D Evaluation Metrics

We adopt distinct sets of metrics to evaluate the classification and regression tasks. In particular, considering the classification task, we utilize the common metrics, including Area Under Receiver Operating Characteristic Curve (ROC-AUC), Precision-Recall Area Under Curve (PR-AUC), LogAUC, accuracy, precision, recall, and F1 score. For the continuous binding affinity regression, we benchmark the models using metrics of mean squared error (MSE), mean absolute error (MAE), coefficient of determination ( $R^2$ ), Pearson correlation coefficient, concordance index (CI), and Spearman correlation coefficient. Each of these metrics offers unique insights into different aspects of model performance, allowing us to assess predictive accuracy, correlation with observed values, and consistency in ranking predictions.

## E Original Hyperparameter

To have a basic understanding of hyperparameters before greedy search and to find the optimized setting for each model, we summarize the hyperparameters reported in the corresponding paper or codes in Table 4.

Table 4: Configurations of basic hyperparameters adopted to implement different approaches.

Category	Models	Batch size	Total epoch	Learning rate & Decay & Decay epoch	Weight decay	Dropout	Optimizer
GNN	GraphDTA-GCN (Nguyen et al., 2020)	512	1000	0.0005	-	0.2	Adam
	GraphDTA-GAT (Nguyen et al., 2020)	512	1000	0.0005	-	0.2	Adam
	GraphDTA-GATGCN (Nguyen et al., 2020)	512	1000	0.0005	-	0.2	Adam
	GraphDTA-GIN (Nguyen et al., 2020)	512	1000	0.0005	-	0.2	Adam
	GraphCPI-GCN (Quan et al., 2019)	512	1000	0.0005	-	0.5	Adam
	GraphCPI-GAT (Quan et al., 2019)	512	1000	0.0005	-	0.6	Adam
	GraphCPI-GATGCN (Quan et al., 2019)	512	1000	0.0005	-	-	Adam
	GraphCPI-GIN (Quan et al., 2019)	512	1000	0.0005	-	0.6	Adam
	MGraphDTA (Yang et al., 2022)	512	3000	0.0005	-	0.1	Adam
	SAGDTA (Zhang et al., 2021)	512	2000	0.001	-	0.1	Adam
	EmbedDTI (Jin et al., 2021)	512	1500	0.0005	-	0.2	Adam
	DeepGLSTM (Mukherjee et al., 2022)	512/128	1000	0.0005	-	0.2	Adam
	CPI (Tsubaki et al., 2018)	1	100	0.001, 0.5, 10	1e-6	0	Adam
	BACPI (Li et al., 2022)	16	20	0.0005, 0.5, 10	-	0.1	Adam
	DeepNC-HGC (Tran et al., 2022)	256	1000	0.0005	-	0.2	Adam
	DeepNC-GEN (Tran et al., 2022)	256	1000	0.0005	-	0.2	Adam
	DrugBAN (Bai et al., 2023)	64	100	0.00005	-	0	Adam
	GANDTI (Wang et al., 2021)	1	30/15	0.001	1e-6	0.5	Adam
	PGraphDTA-CNN (Bal et al., 2024)	512	1500	0.0005	-	0.2	Adam
	BridgeDPI (Wu et al., 2022)	512	100	0.001	-	0.5	Adam
	ColdDTA (Fang et al., 2023)	128	700/300	0.0003	-	0	Adam
	SubMDTA (Pan et al., 2023a)	512	1200	0.0005	-	0.2	Adam
	IMAEN (Zhang et al., 2024)	128	1000	0.0005	-	0.2	Adam
Transformer	CSDTI (Pan et al., 2023b)	256	3000	0.0005	-	0.2	Adam
	AMMVF (Wang et al., 2023)	32	40	0.001, 0.5, 5	1e-4	0.1	Adam
	TDGraphDTA (Zhu et al., 2023)	1024	3000	0.0005	-	0.1	Adam
	IIFDTI (Cheng et al., 2022)	64	200	0.001	1e-6	0.2	AdamW
	ICAN (Kurata & Tsukiyama, 2022)	128	50	0.001	-	0.1	Adam
	MolTrans (Huang et al., 2020b)	64	30	0.00001	-	0.1	Adam
	TransformerCPI (Chen et al., 2020)	8	40	0.0001, 0.5, 5	1e-4	0.2	RAdam
	MRBDTA (Zhang et al., 2022)	1024/256	600/300	0.001	-	0.1	Adam
	FOTFCPI (Yin et al., 2024)	64	100	0.0001	-	0.1	Adam

## F Full experiment

To evaluate the model’s best performance, based on the hyperparameters given in its paper or codes, we found the optimized hyperparameters for each model. On top of the mean value, we also provide the standard deviation across five-fold. The complete result on the regression task is shown in Table 5, and the complete result on the classification task is shown in Table 6.

Table 5: Regression task benchmark on DAIVS, KIBA, and BindingDB datasets, respectively.

Category	Models	DAIVS										KIBA										BindingDB									
		MSE	MAE	R2	PCC	CI	Spearman	MSE ( $\times 10^{-3}$ )	MAE	R2	PCC	CI	Spearman	MSE	MAE	R2	PCC	CI	Spearman												
GNN	GraphDTA-GCN	0.315 ± 0.019	0.332 ± 0.017	0.531 ± 0.028	0.734 ± 0.018	0.840 ± 0.006	0.605 ± 0.010	0.279 ± 0.059	0.041 ± 0.005	-2.114 ± 0.062	0.112 ± 0.026	0.555 ± 0.014	0.156 ± 0.039	0.609 ± 0.033	0.505 ± 0.028	0.669 ± 0.018	0.820 ± 0.011	0.837 ± 0.007	0.744 ± 0.011												
	GraphDTA-GAT	0.382 ± 0.043	0.380 ± 0.025	0.431 ± 0.064	0.671 ± 0.039	0.828 ± 0.014	0.588 ± 0.024	0.428 ± 0.154	0.051 ± 0.011	-3.769 ± 1.712	0.073 ± 0.061	0.534 ± 0.026	0.098 ± 0.074	1.029 ± 0.032	0.676 ± 0.005	0.445 ± 0.038	0.707 ± 0.035	0.786 ± 0.017	0.653 ± 0.011												
	GraphDTA-GATCON	0.306 ± 0.011	0.325 ± 0.013	0.544 ± 0.017	0.741 ± 0.013	0.847 ± 0.003	0.617 ± 0.006	0.311 ± 0.234	0.041 ± 0.018	-2.467 ± 2.604	0.107 ± 0.054	0.578 ± 0.004	0.210 ± 0.249	0.571 ± 0.032	0.478 ± 0.027	0.687 ± 0.017	0.831 ± 0.011	0.844 ± 0.007	0.756 ± 0.012												
	GraphDTA-GIN	0.253 ± 0.010	0.295 ± 0.014	0.623 ± 0.015	0.791 ± 0.008	0.861 ± 0.006	0.638 ± 0.009	0.255 ± 0.007	0.039 ± 0.006	-1.840 ± 0.779	0.124 ± 0.037	0.553 ± 0.019	0.149 ± 0.052	0.563 ± 0.038	0.494 ± 0.023	0.693 ± 0.021	0.836 ± 0.012	0.842 ± 0.007	0.756 ± 0.011												
	GraphCAP-GCN	0.612 ± 0.038	0.501 ± 0.021	0.089 ± 0.056	0.419 ± 0.135	0.718 ± 0.060	0.396 ± 0.123	4.508 ± 2.136	0.169 ± 0.048	-49.762 ± 23.563	-0.021 ± 0.023	0.494 ± 0.013	0.011 ± 0.024	1.199 ± 0.040	0.756 ± 0.019	0.347 ± 0.022	0.606 ± 0.026	0.727 ± 0.015	0.533 ± 0.033												
	GraphCPI-GAT	0.338 ± 0.013	0.364 ± 0.007	0.496 ± 0.019	0.708 ± 0.012	0.838 ± 0.004	0.604 ± 0.007	0.445 ± 0.081	0.051 ± 0.005	-3.957 ± 0.907	0.063 ± 0.050	0.522 ± 0.025	0.081 ± 0.075	0.629 ± 0.012	0.520 ± 0.011	0.657 ± 0.007	0.813 ± 0.005	0.834 ± 0.004	0.739 ± 0.007												
	GraphCPI-GIN	0.274 ± 0.009	0.321 ± 0.007	0.593 ± 0.013	0.773 ± 0.008	0.851 ± 0.008	0.622 ± 0.013	1.681 ± 0.936	0.091 ± 0.042	-17.724 ± 10.533	0.142 ± 0.220	0.553 ± 0.004	0.149 ± 0.246	0.557 ± 0.017	0.473 ± 0.016	0.696 ± 0.009	0.838 ± 0.006	0.847 ± 0.003	0.763 ± 0.006												
	MGraphDTA	0.232 ± 0.012	0.268 ± 0.008	0.655 ± 0.018	0.812 ± 0.011	0.869 ± 0.007	0.650 ± 0.011	0.032 ± 0.012	0.011 ± 0.002	0.642 ± 0.133	0.083 ± 0.079	0.832 ± 0.004	0.793 ± 0.070	0.529 ± 0.011	0.444 ± 0.025	0.712 ± 0.006	0.847 ± 0.005	0.852 ± 0.005	0.769 ± 0.008												
	SAGITA	0.324 ± 0.064	0.329 ± 0.041	0.518 ± 0.096	0.723 ± 0.065	0.833 ± 0.027	0.594 ± 0.044	0.965 ± 0.088	0.017 ± 0.002	0.217 ± 0.885	0.541 ± 0.085	0.713 ± 0.032	0.561 ± 0.075	0.529 ± 0.011	0.444 ± 0.025	0.712 ± 0.006	0.847 ± 0.005	0.852 ± 0.005	0.769 ± 0.008												
	EnsembleDT	0.280 ± 0.024	0.310 ± 0.028	0.583 ± 0.036	0.764 ± 0.023	0.851 ± 0.009	0.623 ± 0.013	0.269 ± 0.142	0.041 ± 0.012	-2.217 ± 1.579	0.131 ± 0.090	0.558 ± 0.028	0.164 ± 0.166	0.542 ± 0.019	0.446 ± 0.023	0.705 ± 0.010	0.843 ± 0.006	0.850 ± 0.004	0.767 ± 0.007												
	DeepGLSTM	0.316 ± 0.023	0.322 ± 0.024	0.529 ± 0.035	0.732 ± 0.022	0.841 ± 0.007	0.609 ± 0.011	0.859 ± 0.749	0.024 ± 0.155	-94.109 ± 83.400	0.040 ± 0.083	0.514 ± 0.036	0.071 ± 0.075	0.594 ± 0.061	0.474 ± 0.046	0.677 ± 0.037	0.843 ± 0.021	0.840 ± 0.013	0.750 ± 0.021												
	CPI	0.402 ± 0.082	0.393 ± 0.054	0.401 ± 0.122	0.629 ± 0.101	0.811 ± 0.033	0.560 ± 0.055	0.032 ± 0.003	0.110 ± 0.008	0.416 ± 0.036	0.654 ± 0.032	0.734 ± 0.037	0.665 ± 0.088	0.762 ± 0.105	0.565 ± 0.088	0.585 ± 0.090	0.768 ± 0.063	0.815 ± 0.028	0.704 ± 0.054												
	BACPI	0.334 ± 0.015	0.323 ± 0.034	0.502 ± 0.023	0.717 ± 0.014	0.837 ± 0.006	0.584 ± 0.010	0.031 ± 0.004	0.011 ± 0.001	0.658 ± 0.043	0.820 ± 0.019	0.831 ± 0.020	0.738 ± 0.032	0.550 ± 0.010	0.436 ± 0.005	0.700 ± 0.006	0.839 ± 0.003	0.845 ± 0.002	0.759 ± 0.003												
	DeepNC-BGC	0.309 ± 0.025	0.331 ± 0.022	0.541 ± 0.037	0.738 ± 0.025	0.841 ± 0.005	0.608 ± 0.008	0.089 ± 0.003	0.019 ± 0.001	0.110 ± 0.035	0.342 ± 0.041	0.667 ± 0.022	0.448 ± 0.048	0.572 ± 0.011	0.486 ± 0.010	0.689 ± 0.006	0.833 ± 0.005	0.844 ± 0.003	0.757 ± 0.005												
	DeepNC-GEN	0.270 ± 0.012	0.298 ± 0.012	0.597 ± 0.017	0.776 ± 0.012	0.832 ± 0.009	0.624 ± 0.014	0.135 ± 0.045	0.027 ± 0.006	-0.509 ± 0.035	0.266 ± 0.073	0.608 ± 0.037	0.301 ± 0.097	0.578 ± 0.029	0.474 ± 0.012	0.685 ± 0.011	0.830 ± 0.006	0.840 ± 0.003	0.749 ± 0.005												
	DrugBAN	0.242 ± 0.007	0.272 ± 0.007	0.649 ± 0.010	0.861 ± 0.007	0.869 ± 0.003	0.631 ± 0.005	0.079 ± 0.003	0.011 ± 0.000	0.676 ± 0.032	0.258 ± 0.020	0.582 ± 0.013	0.800 ± 0.022	0.465 ± 0.018	0.420 ± 0.016	0.717 ± 0.010	0.869 ± 0.006	0.862 ± 0.003	0.785 ± 0.004												
	GANDIT	0.318 ± 0.018	0.301 ± 0.021	0.527 ± 0.027	0.732 ± 0.016	0.844 ± 0.006	0.616 ± 0.013	0.030 ± 0.002	0.011 ± 0.000	0.682 ± 0.026	0.516 ± 0.016	0.831 ± 0.007	0.800 ± 0.011	0.621 ± 0.012	0.489 ± 0.007	0.662 ± 0.006	0.815 ± 0.003	0.836 ± 0.002	0.741 ± 0.005												
BridgeDT	0.141 ± 0.142	0.705 ± 0.060	-0.848 ± 2.133	0.657 ± 0.209	0.827 ± 0.078	0.581 ± 0.128	0.225 ± 0.109	0.010 ± 0.000	0.638 ± 0.121	0.821 ± 0.038	0.857 ± 0.001	0.839 ± 0.002	0.514 ± 0.011	0.413 ± 0.006	0.720 ± 0.006	0.853 ± 0.003	0.861 ± 0.002	0.783 ± 0.003													
ColdDT	0.220 ± 0.009	0.229 ± 0.007	0.672 ± 0.014	0.820 ± 0.009	0.880 ± 0.004	0.666 ± 0.006	0.110 ± 0.029	0.026 ± 0.003	-0.224 ± 0.229	0.414 ± 0.217	0.673 ± 0.079	0.459 ± 0.191	0.463 ± 0.008	0.391 ± 0.007	0.678 ± 0.004	0.866 ± 0.002	0.866 ± 0.001	0.789 ± 0.002													
SubMDTA	0.289 ± 0.012	0.353 ± 0.020	0.570 ± 0.018	0.766 ± 0.012	0.841 ± 0.007	0.604 ± 0.012	0.029 ± 0.002	0.011 ± 0.001	0.675 ± 0.025	0.825 ± 0.015	0.836 ± 0.011	0.865 ± 0.018	0.532 ± 0.022	0.476 ± 0.026	0.710 ± 0.017	0.845 ± 0.010	0.852 ± 0.006	0.772 ± 0.011													
DMAN	0.290 ± 0.009	0.245 ± 0.004	0.674 ± 0.011	0.847 ± 0.004	0.858 ± 0.007	0.606 ± 0.003	0.041 ± 0.002	0.011 ± 0.000	0.484 ± 0.136	0.684 ± 0.176	0.781 ± 0.026	0.713 ± 0.121	0.479 ± 0.015	0.399 ± 0.018	0.739 ± 0.006	0.861 ± 0.001	0.861 ± 0.001	0.788 ± 0.005													
Transformer	CSDT1	0.331 ± 0.012	0.339 ± 0.020	0.508 ± 0.017	0.720 ± 0.009	0.832 ± 0.005	0.593 ± 0.008	0.088 ± 0.004	0.020 ± 0.001	0.014 ± 0.041	0.273 ± 0.084	0.628 ± 0.047	0.350 ± 0.124	0.768 ± 0.021	0.572 ± 0.012	0.582 ± 0.012	0.765 ± 0.008	0.805 ± 0.004	0.689 ± 0.006												
	TDgraphDTA	0.222 ± 0.005	0.265 ± 0.007	0.669 ± 0.008	0.820 ± 0.004	0.853 ± 0.011	0.871 ± 0.007	0.091 ± 0.019	0.022 ± 0.003	-0.009 ± 0.209	0.330 ± 0.117	0.327 ± 0.125	0.619 ± 0.046	0.497 ± 0.016	0.418 ± 0.010	0.729 ± 0.009	0.855 ± 0.005	0.777 ± 0.005	0.687 ± 0.003												
	AMVP	0.377 ± 0.030	0.385 ± 0.016	0.439 ± 0.044	0.660 ± 0.038	0.815 ± 0.005	0.586 ± 0.024	0.075 ± 0.019	0.019 ± 0.003	0.141 ± 0.222	0.079 ± 0.103	0.603 ± 0.141	0.610 ± 0.015	0.517 ± 0.046	0.628 ± 0.038	0.796 ± 0.014	0.825 ± 0.002	0.721 ± 0.004													
	IFTDT	0.313 ± 0.018	0.378 ± 0.039	0.534 ± 0.027	0.754 ± 0.008	0.836 ± 0.008	0.598 ± 0.013	0.054 ± 0.012	0.015 ± 0.001	0.398 ± 0.143	0.691 ± 0.059	0.767 ± 0.009	0.647 ± 0.021	0.634 ± 0.024	0.527 ± 0.020	0.655 ± 0.013	0.820 ± 0.009	0.832 ± 0.006	0.737 ± 0.009												
	ICAN	0.371 ± 0.013	0.359 ± 0.007	0.448 ± 0.020	0.681 ± 0.018	0.818 ± 0.006	0.582 ± 0.009	0.089 ± 0.006	0.022 ± 0.000	-2.052 ± 0.000	0.500 ± 0.000	-	0.747 ± 0.031	0.580 ± 0.018	0.593 ± 0.017	0.785 ± 0.006	0.813 ± 0.004	0.772 ± 0.005													
	McTTrans	0.149 ± 0.136	0.382 ± 0.045	0.200 ± 0.102	0.670 ± 0.107	0.812 ± 0.039	0.591 ± 0.042	4.314 ± 2.290	0.160 ± 0.053	-4.745 ± 25.315	0.032 ± 0.053	0.549 ± 0.021	0.112 ± 0.058	0.695 ± 0.183	0.523 ± 0.031	0.621 ± 0.100	0.809 ± 0.002	0.822 ± 0.009	0.745 ± 0.030												
	TransformerCPI	0.293 ± 0.022	0.445 ± 0.043	0.415 ± 0.032	0.695 ± 0.018	0.802 ± 0.008	0.542 ± 0.015	0.070 ± 0.003	0.019 ± 0.001	0.217 ± 0.033	0.779 ± 0.006	0.800 ± 0.002	0.742 ± 0.004	0.659 ± 0.040	0.548 ± 0.024	0.641 ± 0.022	0.825 ± 0.017	0.829 ± 0.013	0.727 ± 0.022												
	MRBETA	0.241 ± 0.005	0.265 ± 0.006	0.640 ± 0.008	0.802 ± 0.003	0.870 ± 0.007	0.651 ± 0.011	0.050 ± 0.005	0.016 ± 0.001	0.304 ± 0.058	0.600 ± 0.050	0.735 ± 0.015	0.613 ± 0.031	0.507 ± 0.006	0.411 ± 0.006	0.724 ± 0.003	0.853 ± 0.002	0.862 ± 0.002	0.786 ± 0.003												
	FOTFCPI	0.305 ± 0.012	0.302 ± 0.019	0.546 ± 0.018	0.749 ± 0.011	0.839 ± 0.009	0.604 ± 0.015	0.229 ± 0.140	0.034 ± 0.016	-1.554 ± 2.003	0.235 ± 0.264	0.587 ± 0.086	0.414 ± 0.262	0.567 ± 0.008	0.432 ± 0.012	0.605 ± 0.004	0.829 ± 0.004	0.848 ± 0.006	0.763 ± 0.008												
	Our combos	0.211 ± 0.007	0.251 ± 0.008	0.685 ± 0.011	0.829 ± 0.006	0.886 ± 0.004	0.676 ± 0.007	0.036 ± 0.004	0.010 ± 0.001	0.710 ± 0.051	0.845 ± 0.031	0.849 ± 0.023	0.827 ± 0.037	0.461 ± 0.006	0.389 ± 0.007	0.749 ± 0.003	0.867 ± 0.002	0.869 ± 0.002	0.796 ± 0.003												



## G Comparison of different featurization

In this section, we present the summarized featurization methods in Table 7, the detailed description of all properties is shown in Table 8. Besides, an ablation study on featurization strategies is in Table 9.

Table 7: Summary of the featurization of the GNN-based model. Mol. Graphs means molecular graphs, and both mean using molecular graphs and fingerprints.

Model Information		Atomic Properties								Hydrogen Information			Electron Properties		Stereochemistry		Structure
Models	Graph	Atom Type	Degree	Implicit Valence	Explicit Valence	Hybridization	Aromaticity	Formal Charge	# Atom	# Hs	# Explicit Hs	# Implicit Hs	# Radical Electrons	Electron Affinity	CIP	Chirality	Ring
GraphDTA	Mol. Graphs	✓	✓	✓			✓			✓							
GraphCPI	Mol. Graphs	✓	✓	✓			✓		✓	✓	✓		✓	✓			
MGraphDTA	Mol. Graphs	✓	✓	✓	✓	✓		✓					✓				
SAGDTA	Mol. Graphs	✓	✓	✓													
EmbedDTI	Mol. Graphs	✓	✓	✓	✓	✓	✓	✓		✓							✓
DeepGLSTM	Mol. Graphs	✓	✓	✓			✓			✓							
CPI	Fingerprints	✓															
BACPI	Fingerprints	✓															
DeepNC	Mol. Graphs	✓	✓	✓						✓							
DrugBAN	Mol. Graphs	✓	✓	✓		✓	✓	✓		✓			✓				✓
GANDTI	Fingerprints	✓															
PGraphDTA-CNN	Mol. Graphs	✓	✓	✓						✓			✓				✓
BridgeDPI	Mol. Graphs	✓	✓	✓		✓	✓	✓									
ColdDTA	Mol. Graphs	✓	✓	✓			✓			✓					✓	✓	
SubMDTA	Mol. Graphs	✓	✓	✓			✓			✓							
IMAXEN	Mol. Graphs	✓	✓	✓			✓										
CSDTI	Mol. Graphs	✓	✓	✓	✓	✓		✓	✓		✓		✓	✓			
TDGraphDTA	Mol. Graphs	✓	✓	✓						✓							
AMMVF	Both	✓	✓	✓		✓	✓	✓		✓			✓		✓	✓	
TransformerCPI	Mol. Graphs	✓	✓	✓		✓	✓	✓		✓			✓		✓	✓	

Table 8: Description of atomic and molecular properties for node featurization

Name	Description
<b>Atomic Properties</b>	
Atom Type	Type of the atom (e.g., C, N, O, H)
Degree	Number of directly bonded neighbors
Implicit Valence	Number of implicit valence of the atom
Explicit Valence	Number of explicit valence of the atom
Hybridization	The state of hybridization (e.g., sp3, sp2)
Aromaticity	Whether the atom is part of an aromatic system
Formal Charge	The charge assigned to an atom
# Atom	Total number of atoms
<b>Hydrogen Information</b>	
# Hs	Total number of hydrogens
# Explicit Hs	Number of explicit hydrogens on the atom
# Implicit Hs	Number of implicit hydrogens on the atom
<b>Electron Properties</b>	
# Radical Electrons	Number of radical electrons
Electron Affinity	Tendency of an atom to accept electrons
<b>Stereochemistry</b>	
CIP	The CIP code (R or S) of the atom
Chirality	If an atom is a possible chiral center
<b>Structure</b>	
Ring	Whether the atom is part of a ring structure

Table 9: Extra Graph embedding feature exploration. Here Basic: {Atom Type, Degree, Implicit Valence, Aromaticity, # Hs}

Models	Initial Feature	Regression						Classification						
		MSE	MAE	R2	PCC	CI	Spearman	ROC-AUC	PR-AUC	Range-AUC	Acc.	Precision	Recall	F1
GraphDTA	Basic	0.2771	0.2947	0.5873	0.7695	0.8521	0.6241	0.9170	0.8873	0.4830	0.9171	0.9204	0.9189	0.9194
	Basic+AP	0.2772	0.2978	0.5873	0.7671	0.8496	0.6200	0.9153	0.8817	0.4517	0.9157	0.9099	0.9290	0.9191
	Basic+HI	0.2783	0.2983	0.5855	0.7663	0.8483	0.6185	0.9211	0.8903	0.4815	0.9214	0.9188	0.9296	0.9241
	Basic+EP	0.2775	0.3068	0.5868	0.7682	0.8499	0.6205	0.9165	0.8862	0.4795	0.9166	0.9185	0.9198	0.9191
	Basic+Ste	0.2838	0.3030	0.5773	0.7624	0.8523	0.6254	0.9200	0.8905	0.4869	0.9200	0.9216	0.9235	0.9224
	Basic+Str	0.2783	0.2991	0.5857	0.7668	0.8505	0.6228	0.9198	0.8865	0.4649	0.9202	0.9124	0.9351	0.9235
	Basic+AP+HI	0.2851	0.3029	0.5755	0.7610	0.8504	0.6222	0.9163	0.8822	0.4629	0.9168	0.9094	0.9313	0.9201
	Basic+AP+HI+EP	0.2845	0.2917	0.5763	0.7620	0.8510	0.6227	0.9140	0.8811	0.4580	0.9143	0.9115	0.9232	0.9173
	Basic+AP+HI+EP+Ste	0.2811	0.3099	0.5814	0.7640	0.8500	0.6212	0.9192	0.8899	0.4853	0.9193	0.9218	0.9215	0.9216
Basic+AP+HI+EP+Ste+Str	0.2801	0.2916	0.5829	0.7659	0.8538	0.6278	0.9217	0.8905	0.4794	0.9220	0.9180	0.9319	0.9248	
GraphCPI	Basic	0.3291	0.3388	0.5100	0.7265	0.8294	0.5885	0.9060	0.8706	0.4385	0.9064	0.9029	0.9169	0.9098
	Basic+AP	0.3331	0.3389	0.5040	0.7198	0.8223	0.5761	0.9038	0.8657	0.4103	0.9043	0.8955	0.9218	0.9084
	Basic+HI	0.3402	0.3457	0.4934	0.7157	0.8228	0.5769	0.9051	0.8713	0.4495	0.9052	0.9058	0.9094	0.9080
	Basic+EP	0.3408	0.3505	0.4926	0.7123	0.8211	0.5749	0.9053	0.8713	0.4442	0.9055	0.9060	0.9111	0.9085
	Basic+Ste	0.3398	0.3634	0.4940	0.7119	0.8274	0.5855	0.9061	0.8692	0.4261	0.9065	0.8992	0.9221	0.9104
	Basic+Str	0.3419	0.3562	0.4909	0.7113	0.8226	0.5766	0.9066	0.8683	0.4079	0.9073	0.8957	0.9281	0.9115
	Basic+AP+HI	0.3326	0.3471	0.5048	0.7212	0.8210	0.5734	0.9010	0.8659	0.4288	0.9012	0.9018	0.9071	0.9043
	Basic+AP+HI+EP	0.3404	0.3476	0.4931	0.7150	0.8212	0.5748	0.9015	0.8612	0.3821	0.9022	0.8890	0.9258	0.9070
	Basic+AP+HI+EP+Ste	0.3403	0.3445	0.4932	0.7111	0.8169	0.5671	0.9109	0.8763	0.4511	0.9113	0.9065	0.9229	0.9146
Basic+AP+HI+EP+Ste+Str	0.3469	0.3550	0.4834	0.7073	0.8228	0.5775	0.9134	0.8772	0.4440	0.9140	0.9033	0.9328	0.9178	

## H Memory and Parameter Comparison

Table 10: Training time per epoch (s) and the max allocated memory (MB) for representative datasets on both regression (Davis) and classification (Human) tasks when BS is 32.

Categories	Models	Regression			Classification		
		Model parameter	Memory Usage (MB)	Time(s)	Model parameter	Memory Usage (MB)	Run Time (s)
Graph	GraphDTA-GCN	7.87	86.45	8.92	7.87	86.33	2.43
	GraphDTA-GAT	6.58	104.71	9.62	6.58	99.40	2.43
	GraphDTA-GATGCN	18.12	148.25	8.37	18.12	145.13	2.35
	GraphDTA-GIN	5.97	78.00	12.33	5.95	77.47	3.13
	GraphCPI-GCN	10.46	98.13	7.02	10.48	63.37	1.92
	GraphCPI-GAT	9.16	116.19	9.38	9.18	112.34	2.48
	GraphCPI-GATGCN	20.70	158.21	9.47	20.73	156.22	2.20
	GraphCPI-GIN	8.55	88.55	12.54	8.56	88.02	2.92
	MGraphDTA	11.75	235.97	69.84	11.43	217.15	17.59
	SAGDTA	7.45	88.31	20.87	7.44	87.54	4.34
	EmbedDTI	16.97	152.55	17.80	16.97	-	-
	DeepGLSTM	131.92	1287.92	20.69	131.93	1287.16	11.22
	CPI	0.37	14.00	11.29	0.6	14.82	2.69
	BACPI	4.05	1051.91	43.27	6.13	1058.95	12.38
	DeepNC-HGC	16.61	123.70	9.85	16.60	123.65	3.46
	DeepNC-GEN	18.84	174.00	11.35	18.84	166.55	3.46
	DrugBAN	4.10	940.22	30.06	4.10	940.23	7.84
	GANDTI	1.48	35.89	6.01	2.43	39.95	1.54
	PGraphDTA-CNN	9.03	102.85	13.71	9.03	-	-
	BridgeDPI	39.32	232.53	16.27	39.32	232.53	4.36
Transformer	ColdDTA	13.14	282.74	72.98	13.14	262.91	18.56
	SubMDTA	169.37	992.61	35.12	195.50	1095.73	8.49
	IMAEN	10.43	174.34	35.77	10.43	172.86	4.41
	CSDTI	9.67	281.23	17.66	9.66	281.02	4.35
	TDGraphDTA	8.62	247.23	116.38	8.62	236.02	28.43
	AMMVF	6.68	17847.62	216.20	7.49	17850.79	57.99
	HFDTI	10.75	7946.92	141.12	10.75	11890.79	56.95
	ICAN	63.89	649.55	12.44	63.89	648.67	2.84
	MolTrans	239.73	10624.55	70.19	239.74	10624.55	25.06
	TransformerCPI	4.44	1219.58	28.98	4.45	1219.60	7.17
	MRBDTA	17.83	3893.76	66.47	17.84	3893.78	16.13
	FOTFCPI	189.15	6780.35	58.75	189.15	6780.35	14.80
	Our	19.02	1081.99	94.71	19.02	1082.68	13.68