

DOCKEDAC: A DATASET WITH COMPREHENSIVE 3D PROTEIN-LIGAND COMPLEXES FOR ACTIVITY CLIFF ANALYSIS

Anonymous authors

Paper under double-blind review

ABSTRACT

Artificial intelligence has become a crucial tool in drug discovery, excelling in tasks such as molecular property prediction. However, an *activity cliff*—a phenomenon where a minor structural modification to a molecule leads to a significant change in its biological activity—poses a challenge in predictive modeling. The activity cliff depends on the interaction between the target and the ligand, which is largely overlooked by previous ligand-centric studies. However, the limited availability of activity cliff data for target-ligand 3D complexes constrains the predictive power of modern deep learning models. In this paper, we introduce `DockedAC`, a new dataset incorporating the protein target and 3D complex structure information for studying the problem of activity cliffs. By matching protein binding information and ligand bioactivity, we employ molecular docking to generate the complex structure for each activity value. The `DockedAC` dataset contains 82,836 activity data across 52 protein targets annotated with activity cliff information. This dataset represents a significant step toward large-scale activity cliff research using 3D complex structures. We benchmark the dataset with traditional machine learning and deep learning approaches. Our data and benchmark codes are available [here](#).

1 INTRODUCTION

Artificial intelligence (AI) is revolutionizing the drug discovery process as it is capable of large-scale data analysis, pattern recognition, and making accurate predictions (Vamathevan et al., 2019). One important application of AI models is to predict the biological activity of candidate compounds, thereby reducing labor-intensive tasks. A foundational concept in many AI algorithms is the similarity principle, which states that similar objects are likely to share similar features and predictions. However, in drug discovery, a phenomenon known as activity cliffs challenges this concept and poses difficulties for AI models. An **activity cliff** (AC) occurs when structurally similar compounds exhibit significant differences in their biological activity against the same target (Maggiora, 2006), as illustrated in [Figure 1](#) (a).

AC plays a crucial role in drug discovery, as it complicates the optimization of drug candidates by confounding the human experts in the understanding of traditional structure-activity relationships (SARs) (Vogt et al., 2011). On the other hand, knowledge about ACs can be highly beneficial when designing or optimizing compounds to enhance the bioactivity of a given target (Cruz-Monteagudo et al., 2014; Stumpfe et al., 2014). For example, replacing a single atom or adding a methyl group can result in more than 100-fold improvement in bioactivity (Leung et al., 2012; Pennington & Moustakas, 2017). However, the mechanisms underlying ACs in individual drug development programs can be different, making it challenging to process such information and derive transferable experiences. Consequently, various efforts have been made to computationally predict ACs (Stumpfe et al., 2019).

Compared to quantitative structure-activity relationship (QSAR) modeling for other molecular properties, AC prediction is particularly challenging due to the instability that ACs introduce to the models (Cruz-Monteagudo et al., 2016). Early attempts use machine learning methods such as random forest (RF) and support vector machine (SVM) to predict the AC of a compound pair (Guha, 2012; Heikamp et al., 2012). To further improve AC predictions, the matched molecular pair (MMP) kernel (Tamura et al., 2021) and condensed graphs of reaction representations (Horvath

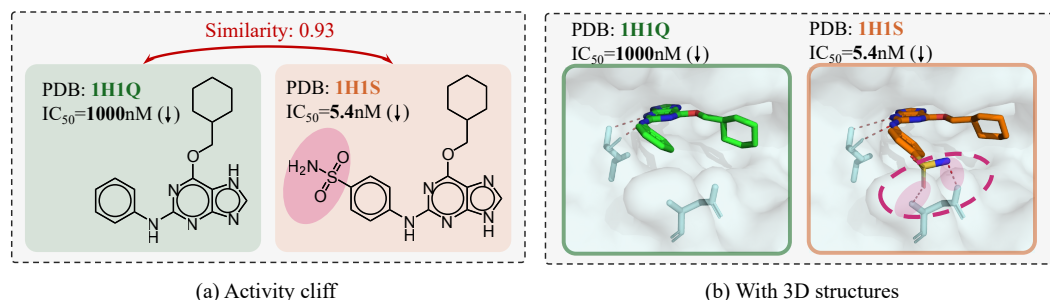


Figure 1: Illustration of activity cliffs. (a) An activity cliff example: two similar molecules with a significant bioactivity difference. (b) From the 3D structure, it is easier to see that the bioactivity of the right ligand is improved due to the formation of two new hydrogen bonds (pink dashed lines).

et al., 2016) have been integrated into various machine learning methods. More recently, deep neural networks-based algorithms, such as convolutional neural networks (Iqbal et al., 2021), graph neural networks (Park et al., 2022) and transformers (Chen et al., 2022), have been applied to predict ACs.

In most previous works, the study of ACs has been ligand-centric and lacked 3D structure consideration, failing to account for interactions between the ligand and the protein target (Husby et al., 2015; Tamura et al., 2023). Many mechanisms of ACs can be analyzed from the structural perspective, such as hydrogen bonding, ionic interactions, hydrophobic or aromatic group interactions (Hu et al., 2012) (e.g. Figure 1 (b)). It is therefore natural to incorporate the information of structures into the modeling of ACs. However, the available structural data for ACs is very limited, with only 215 pairs of AC ligands (Husby et al., 2015). This data scarcity makes it challenging to train deep learning models effectively.

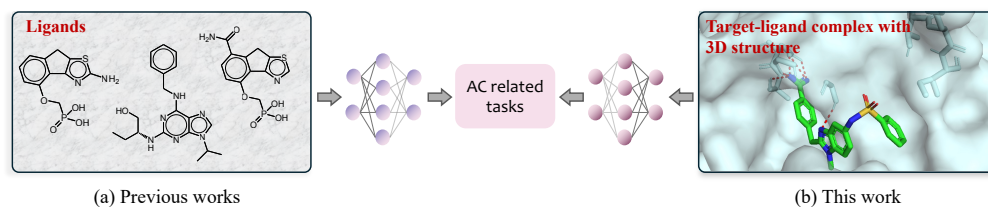


Figure 2: Settings of previous studies and our work about ACs. (a) Previous works mostly consider AC prediction from a ligand-centric view and overlook the target information and 3D complex structure. (b) We construct a dataset with target-ligand complex structures for AC prediction.

In this paper, we present *DockedAC*, a new dataset designed to tackle the challenge of ACs from a structural perspective, enabling large-scale AC modeling with modern AI algorithms. Unlike previous studies, our dataset includes not only the protein target information but also the target-ligand complex structures built using molecular docking (Figure 2). We collect the bioactivity data of more than 80,000 ligands across over 50 protein targets. These protein targets are mapped to their corresponding structures in the RCSB Protein Data Bank (PDB) (Berman et al., 2000), with the ligand binding sites identified for docking. In addition, we provide a benchmarking framework to evaluate the performance of traditional machine learning and deep learning methods on AC prediction and to analyze the impact of ACs on model performance. Our dataset enhances model interpretability, inspires the development of advanced algorithms for AC prediction, and fosters the advancement of more effective 3D feature extraction methods.

2 RELATED WORK

Models on AC prediction. As a crucial phenomenon in drug discovery, ACs have garnered significant attention not only in medicinal chemistry but also in computer science and the intelligence community. Various methods of machine learning and deep learning have been applied to the prediction of ACs (Yu et al., 2025; Iqbal et al., 2021; Chen et al., 2022; Park et al., 2022). In

108 addition, recent research has explored ACs from several different perspectives, such as transfer
109 learning (Ibragimova et al., 2025), QSAR modeling (Dablander et al., 2023), drug design (Hu et al.,
110 2025), and benchmarking of different approaches (van Tilborg et al., 2022). However, due to the
111 limited availability of data, almost all existing works focus on the ligand-centric view of ACs, where
112 the ligand is modeled with a 2D molecular graph or 1D SMILES sequence (Weininger, 1988), without
113 incorporating the 3D structure and protein target information. The 3D activity cliff (3DAC) database,
114 used in a study on structure-based AC prediction, contains only 219 3DAC pairs (Hu et al., 2012;
115 Husby et al., 2015). This motivates us to construct a larger dataset for structure-based ACs.

116 **Existing AC datasets.** Although efforts have been put into AC prediction, few good benchmarking
117 datasets are available. Several works rely on self-collected datasets and are not well documented, or
118 have little information provided about the protein targets (Li et al., 2025; Jiménez-Luna et al., 2022;
119 Dablander et al., 2023; Tamura et al., 2023). Two recent works on AC datasets collect data from the
120 ChEMBL database (Mendez et al., 2019), either for the classification of a pair of AC ligands (Zhang
121 et al., 2023c) or the regression of the bioactivity value of individual AC ligands (van Tilborg et al.,
122 2022). These datasets do not consider modeling the 3D structure of the binding complex, rendering
123 them less appropriate for accurate AC prediction. In our work, we match the obtained bioactivity
124 data to the corresponding protein structures in PDB and generate target-ligand binding structures.

125 **3D protein-ligand binding affinity prediction.** In this work, we consider the regression problem
126 and train different models to predict the binding affinity in the presence of the AC Passaro et al.
127 (2025). Many methods for structure-based binding affinity prediction employ the PDBbind dataset,
128 including convolutional neural networks, graph neural networks, and attention-based models (Zhang
129 et al., 2023a; Jiang et al., 2021; Jiménez et al., 2018; Tan et al., 2024). A comprehensive review
130 of the drug-target interaction prediction can be found in Zeng et al. (2024). In molecular property
131 prediction, activity cliffs can significantly impact model predictions (Deng et al., 2023). We evaluate
132 the performance of 3D target-ligand affinity prediction models with our dataset and compare them
133 with other machine learning or deep neural network models with ligand-only inputs.

134 3 THE DOCKEDAC DATASET

135
136 In summary, the construction of DockedAC involves several key steps: data collection, AC identi-
137 fication, target structure annotation, and target-ligand complex generation. The following section
138 provides a detailed explanation of each step in this process.

139 3.1 DATA COLLECTION

140
141 We first collect bioactivity data from ChEMBL v33 (Mendez et al., 2019) using the ChEMBL web
142 resource client (Davies et al., 2015) for 64 protein targets. The data includes Inhibitory Constant (K_i),
143 Half-Maximal Effective Concentration (EC_{50}), and Half-Maximal Inhibitory Concentration (IC_{50}),
144 all measured in nanomolar (nM). To eliminate significant sources of error, the obtained raw data is
145 checked for validity and reliability. In particular, when a ligand-target pair has multiple entries of the
146 bioactivity data, the ligand is removed if the standard deviation of the activities is larger than 10. The
147 mean value of the activities is used as the ligand-target activity label. A ligand is also removed if it
148 fails the sanitization and standardization by RDKit (Bento et al., 2020). To ensure enough samples of
149 a target for model training, the targets with fewer than 500 ligands are dropped. Finally, the negative
150 logarithm p is applied to the bioactivity values as the regression target (denoted as pK_i ; pEC_{50} ;
151 pIC_{50} in [log units]) (Stewart & Watson, 1983). After this process, we have the ChEMBL id of
152 the target and the corresponding ligands with bioactivity values (the first step in Figure 3 (a)). The
153 resulting dataset has 54 protein targets.

154 3.2 ACTIVITY CLIFF IDENTIFICATION

155
156 An activity cliff is defined as a pair of structurally similar compounds exhibiting a large difference
157 in bioactivities against a given target. To identify similar ligand pairs, we use a consensus of
158 three similarity measures to define the activity cliff pairs following van Tilborg et al. (2022): (a)
159 substructure similarity, calculated using the Tanimoto coefficient on the extended connectivity
160 fingerprint (ECFP) (Tanimoto, 1958; Rogers & Hahn, 2010); (b) scaffold similarity, determined by
161 the Tanimoto coefficient on the ECFP of generic Murcko scaffolds (Bemis & Murcko, 1996); (c)

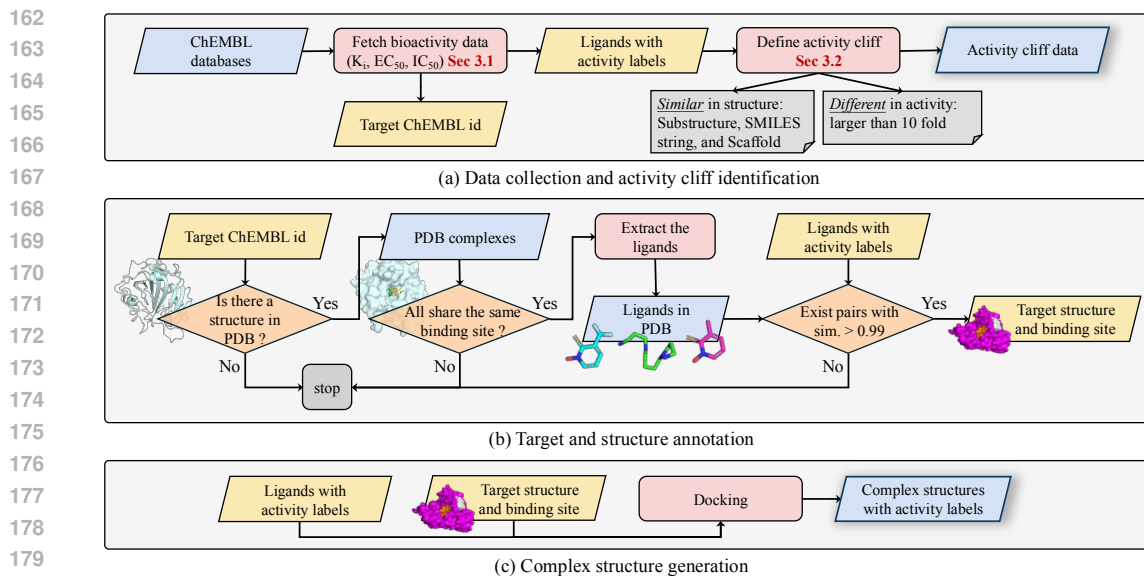


Figure 3: The whole process of building DockedAC with: (a) initial data collection from ChEMBL (Sec. 3.1) and activity cliff identification (Sec. 3.2), (b) mapping targets to 3D structures and identifying binding sites (Sec. 3.3), and (c) generation of target-ligand complex structures (Sec. 3.4).

SMILES similarity, computed as one minus the scaled Levenshtein distance between the canonical SMILES representations (Levenshtein et al., 1966). If any of these three similarity measures is equal to or greater than 0.9, the pair of ligands is further evaluated for differences in bioactivity. There are no widely accepted quantitative definitions of ACs (Stumpfe et al., 2020). Following previous studies (Jiménez-Luna et al., 2022; Hu & Bajorath, 2012), we define an activity cliff as a bioactivity difference exceeding one order of magnitude (10 \times), as illustrated in the second step of Figure 3 (a).

3.3 TARGET AND STRUCTURE ANNOTATION

To generate the target-ligand complex, it is essential to identify the 3D structure of the target protein and its binding site. This mapping process is illustrated in Figure 3 (b). Given a target ChEMBL id, the first step is to map the target protein to its UniProt id (Consortium, 2023) and find all the structures corresponding to the UniProt id in the PDB. We utilize the PDBbind database for the initial search (Wang et al., 2004). If the target is not found in PDBbind, we then search for it in the entire PDB. The retrieved structures containing a small molecule ligand are chosen and aligned to verify whether the ligands bind to the same site. If the binding site is not unique, the target is discarded (see Figure 10 (a)(b)). After alignment, ligands sharing the same binding site are extracted and compared with the ligands that have activity labels from ChEMBL. If a pair of ligands—one from the PDB database and one from ChEMBL—has a fingerprint similarity (Tanimoto coefficient) greater than 0.99, the target structure and the binding site are used. Otherwise, the target is removed from the dataset. When multiple structures satisfy this condition, the structure with the highest resolution is selected. This procedure ensures the correspondence between the bioactivity values and the target binding site. As a result of this structure mapping process, two targets were removed, resulting in a final dataset of 52 protein targets.

3.4 COMPLEX STRUCTURE GENERATION

Given a target protein and its corresponding binding site for a ligand, molecular docking is employed to generate the target-ligand complex, as illustrated in Figure 3 (c). The docking tool DSDP is used, which combines the AutoDock Vina’s pose sampling algorithm with GPU acceleration (Huang et al., 2023; Trott & Olson, 2010). Since the binding site information of the target is already known, local docking is performed within a 25 Å wide box around the given binding region. A docking score (in kcal/mol) greater than zero indicates an inaccurate docking conformation (e.g. Figure 10 (c)), and

the corresponding ligand is removed from the dataset. To further improve the quality of the docking complex, for ligands with high structural similarity to known references, a template-based docking approach is employed (Yang et al., 2022). In this approach, when a query ligand shares significant structural features with a crystallographically resolved reference ligand, the known binding pose is used as a template to guide the placement of similar substructures. This template-based protocol constrains the conformational search space, leading to more accurate pose predictions. To validate the docking result, we have compared the docking poses of compounds with known crystal structures from Husby et al. (2015). The observed root-mean-square deviation (RMSD) values (median = 2.55 Å) fall well within accepted standards for reliable pose prediction.

To enhance the comprehensiveness and diversity of our protein-ligand complex dataset, we incorporate two additional docking methods. First, we employ KarmaDock (Zhang et al., 2023b), a state-of-the-art machine learning-based docking method that leverages deep learning algorithms to predict binding poses. Second, to account for protein flexibility in ligand binding, we utilize DSDPflex (Dong et al., 2024), which accounts for side-chain movements in the binding site region. For DSDPflex, we allow the 10 residues closest to the reference ligand to have flexible side chains and record the top 5 scoring poses for each ligand-target docking simulation. This complementary approach provides insights into the dynamic nature of protein-ligand interactions. While we provide the complexes generated by all three methods in our dataset, our subsequent analyses primarily focused on the results obtained from DSDP, as it offered the most reliable and consistent predictions for our system.

3.5 DATASET SPLITTING

The preparation of datasets for benchmarking machine learning models requires careful data-splitting strategies. For ligand-based methods, separate models are developed for individual protein targets. To assess the ligand-based method, the ligands of each target are split into a training and test set using a double-stratified sampling strategy (van Tilborg et al., 2022). In particular, the ligands of each target are first clustered into 5 groups based on their substructural similarity (Tanimoto similarity of the ECFP). A two-stage stratified splitting (80%/20%) is then performed on the cluster label and the AC label. This procedure ensures that the training and test set have similar ligand distributions.

Our dataset contains 3D structural information and target-specific data, which can be used to train cross-target 3D protein-ligand binding affinity prediction models. This approach allows for making predictions on novel targets that do not exist in the training data. In this case, it is appropriate to use the data with the same activity label types (K_i , EC_{50} or IC_{50}) and separate the dataset by target to evaluate the cross-target modeling capabilities.

3.6 DATASET DESCRIPTION

The final dataset contains 82,836 target-ligand activity values and their corresponding generated complex structures. A brief overview of the dataset is provided in Table 1, while detailed information on each target can be found in Appendix Table 3. The dataset includes popular target families in drug discovery (such as G-protein-coupled receptors (GPCR), kinases, proteases, and nuclear receptors) as well as targets with critical roles in biology (like chaperone and kinesin). In terms of size, the target Carbonic anhydrase II has the most ligands with bioactivity values (5794 unique molecules), while the target with the least ligands (533 unique molecules) is Matrix metalloproteinase 8. As an intensively studied drug target, the GPCR family has the most ligands on average. For all the targets, around 37% of the ligands are annotated as ACs, with percentages ranging from 15.7% to 43.2%.

Table 1: Brief dataset statistics by the target type.

| Target type | # Targets | Avg. # ligands | %AC |
|-----------------------------|-----------|----------------|------|
| G protein-coupled receptor | 12 | 2091 | 41.7 |
| Kinase | 11 | 1234 | 27.5 |
| Protease | 8 | 1667 | 38.0 |
| Nuclear receptor | 8 | 1299 | 35.7 |
| Phosphodiesterase | 3 | 1328 | 34.1 |
| Phosphatase | 2 | 1581 | 18.0 |
| Transporter | 1 | 1051 | 25.3 |
| Transferase | 1 | 960 | 41.8 |
| Oxidoreductase | 1 | 739 | 38.0 |
| Other membrane receptor | 1 | 1328 | 38.2 |
| Lyase | 1 | 5796 | 42.2 |
| Kinesin | 1 | 719 | 43.2 |
| Electrochemical transporter | 1 | 1702 | 37.5 |
| Chaperones | 1 | 999 | 15.7 |

4 BENCHMARK

In addition to the `DockedAC` dataset, we provide a framework to benchmark the performance of various machine learning and deep learning methods on AC prediction. In this section, we briefly introduce the benchmark setup. The detailed experimental results are provided in Section 5.

4.1 MODEL DESCRIPTIONS

In general, three types of learning models are considered:

- Four classic machine learning algorithms for structure-activity relationship prediction using hand-crafted molecular descriptors: K-nearest neighbor (KNN) (Cover & Hart, 1967), random forest (RF) (Breiman, 1996), gradient boosting machine (GBM) (Friedman, 2001), and support vector regression (SVM) (Hearst et al., 1998).
- Deep learning models that only leverage the 1D or 2D ligand information, including (1) three 1D sequential models: transformer (Vaswani et al., 2017), long short-term memory (LSTM) networks (Hochreiter & Schmidhuber, 1997), and 1D CNN (Kimber et al., 2021), and (2) four 2D structural graph neural network (GNN) models: message passing neural network (MPNN) (Gilmer et al., 2017), graph convolutional network (GCN) (Kipf & Welling, 2016), graph attention network (GAT) (Vaswani et al., 2017), and attentive fingerprint (AFP) (Xiong et al., 2019).
- Two 3D structural GNN models: IGN (Jiang et al., 2021) and SS-GNN (Zhang et al., 2023a) are included to study the effect of 3D structures, as our dataset contains 3D structural information.

4.2 FEATURE DESCRIPTIONS

For machine learning algorithms, following previous work van Tilborg et al. (2022), we consider four types of molecule descriptors from several levels of complexity as follows. (1) Extended Connectivity Fingerprints (ECFPs) (Rogers & Hahn, 2010): circular topological fingerprints used for molecular characterization, capturing structural features of molecules. (2) Molecular ACCess System (MACCS) keys (Durant et al., 2002): a set of structural keys utilized for substructure searching and similarity analysis, encoding specific chemical substructures or patterns. (3) Physicochemical (PhysChem) descriptors (Walters & Murcko, 2002): a set of 11 properties indicative of drug-likeness, providing insights into the physical and chemical properties of molecules. (4) Weighted Holistic Invariant Molecular (WHIM) descriptors (Todeschini et al., 1998): capturing three-dimensional geometrical and electronic properties of molecules, invariant to rotation and translation.

Deep learning methods eliminate the need for handcrafted descriptors, allowing direct learning from “unstructured” data representations. For sequential methods, the Simplified Molecular Input Line Entry System (SMILES) (Weininger, 1988) string is used, which is popular for its ability to describe the structure of chemical species in a text format that sequential methods can naturally process. For 2D GNN models, we adopt molecular graphs, which represent the structural formula where nodes represent atoms and edges represent bonds. For 3D GNN models, we employ the target-ligand complexes we have processed that incorporate detailed 3D structure information. Additional descriptions of the features and their corresponding models are available in Appendix B.4 and Table 4.

4.3 METRICS AND IMPLEMENTATIONS

For each target, we train separate regression models on the bioactivity values ($pK_i/pEC_{50}/pIC_{50}$ in [log units]). The regression setting makes it possible to compare the AC and non-AC tasks. The root-mean-square error (RMSE) is employed as the evaluation metric to quantify the performance. The RMSE represents the error calculated across all ligands, whereas $RMSE_{cliff}$ specifically denotes the error computed for AC ligands. For model implementation, we conduct hyperparameter tuning through grid search and report the results from five-fold cross-validation. Further details on these methods and their implementations are provided in Appendix B.3 and B.5.

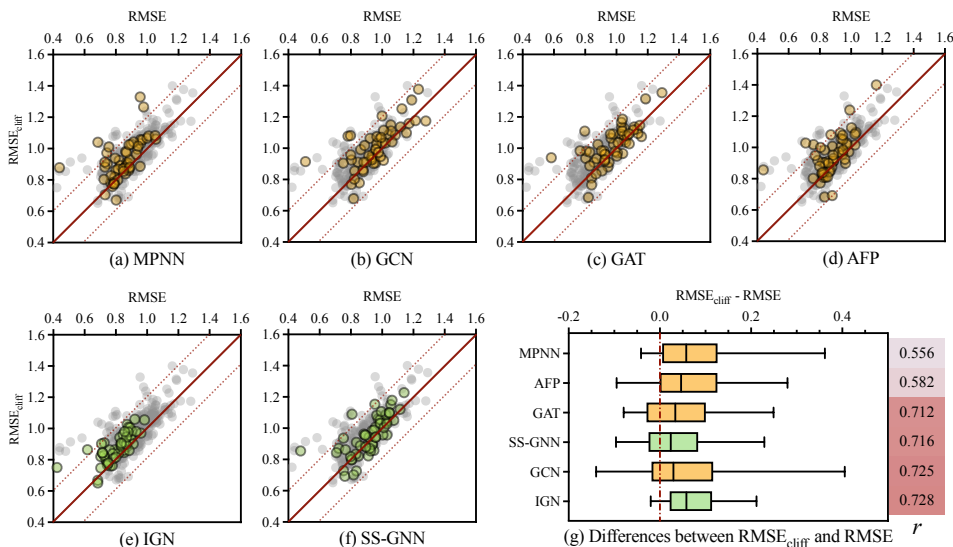


Figure 4: Performance comparison for GNN models. (a)-(f) Comparison between RMSE and RMSE_{cliff} of GNN models across 52 targets. The 2D GNN models are colored in yellow, while the 3D GNN models are colored in green. Gray nodes depict all nodes in these six subgraphs for a clear comparison. Red solid lines show RMSE = RMSE_{cliff}, while red dashed lines indicate a ± 0.2 log units difference. (g) Target-wise differences between overall RMSE and RMSE_{cliff} for all GNN models ordered by Pearson correlation r of RMSE and RMSE_{cliff}.

5 EXPERIMENTAL RESULTS AND ANALYSES

This section provides a structured evaluation of model performance on our DockedAC. It begins with a comparative analysis of 2D and 3D GNN models, followed by an investigation into the target-dependent nature of AC prediction. The impact of the ratio of AC ligands is examined, alongside a benchmarking of machine learning and deep learning methods. Finally, multidimensional scaling is employed to assess the performance positioning of 3D GNN models.

5.1 PERFORMANCE COMPARISON FOR GNN MODELS

To investigate the effect of 3D structure information, we first evaluate 2D GNN models and 3D GNN models across 52 targets. To study AC, scatter plots with RMSE on the x-axis and RMSE_{cliff} on the y-axis are utilized, as shown in Figure 4 (a) to (f).

We have the following empirical observations: (1) The majority of the points are distributed above the line RMSE = RMSE_{cliff}, indicating higher prediction errors on ACs due to their unusual structure-activity relationships. (2) Despite a general correlation between RMSE and RMSE_{cliff}, notable outliers indicate that models with overall high prediction accuracy do not necessarily perform well on ACs. Among these models, SS-GNN exhibits the closest distribution around line RMSE = RMSE_{cliff}, with only two targets deviating by more than 0.2 log units. (3) The distribution of IGNN is primarily clustered in the lower-left corner of the plots, indicating superior performance in both RMSE and RMSE_{cliff}. This suggests that incorporating 3D structural information enhances the prediction of ACs and improves the model’s understanding of standard structure-activity relationships. (4) Figure 4 (g) further presents the target-wise differences between RMSE and RMSE_{cliff} for GNN models, sorted by the Pearson correlation coefficient r of RMSE and RMSE_{cliff}. 3D structure GNN models ranked first and third in terms of r . SS-GNN exhibits the smallest RMSE - RMSE_{cliff} differences, while IGNN has the most concentrated distribution across targets. Its 5%-95% coverage range is only 0.58 times that of MPNN and 0.71 times that of GAT. These findings demonstrate the benefit of incorporating 3D structural information, which leads to a higher degree of correlation between performance on overall ligands and AC ligands, ultimately improving the understanding of structure-activity relationships and aiding in the prediction of ACs.

5.2 THE AC PREDICTION IS TARGET-DEPENDENT

The AC effect is determined by the interaction between the ligand and the target. We hypothesize that the target type may also influence the model performance. Table 2 shows the average $\text{RMSE}_{\text{cliff}}$ of the top four target families: GPCR, kinase, protease, and nuclear receptor. The rankings, represented by color coding, reveal consistent trends across both deep learning and machine learning methods. Protease has the worst $\text{RMSE}_{\text{cliff}}$ for all the methods while kinase is the target family with the best $\text{RMSE}_{\text{cliff}}$ for most methods. Deep learning methods generally perform better on nuclear receptors than GPCRs, while machine learning methods exhibit the opposite trend.

Table 2: The $\text{RMSE}_{\text{cliff}}$ evaluated using GNN models and machine learning algorithms with ECFP featurization across the top four target families. For each method, the colors show the ranking of the target, i.e., first, second, third, fourth.

| Target type | # Target | MPNN | GCN | GAT | AFP | IGN | SS-GNN | KNN | RF | GBM | SVM |
|------------------|----------|-------|-------|-------|-------|-------|--------|-------|-------|-------|-------|
| GPCR | 12 | 0.927 | 0.995 | 1.018 | 0.907 | 0.877 | 0.977 | 0.814 | 0.785 | 0.791 | 0.752 |
| Kinase | 11 | 0.902 | 0.942 | 0.970 | 0.917 | 0.865 | 0.896 | 0.802 | 0.765 | 0.747 | 0.707 |
| Protease | 8 | 0.979 | 1.071 | 1.069 | 1.025 | 0.904 | 1.006 | 0.867 | 0.827 | 0.828 | 0.810 |
| Nuclear receptor | 8 | 0.893 | 0.972 | 0.978 | 0.932 | 0.865 | 0.906 | 0.822 | 0.799 | 0.800 | 0.781 |

5.3 THE PERCENTAGE OF AC MATTERS

In general, machine learning models tend to perform better with more training data. Here, we investigate the factors influencing AC (Activity Cliff) prediction performance. Surprisingly, our analysis reveals that the number of training samples does not exhibit a significant correlation with RMSE , $\text{RMSE}_{\text{cliff}}$, or their difference, i.e., $\text{RMSE}_{\text{cliff}} - \text{RMSE}$ (see Appendix Figure 11). This suggests that simply increasing the size of the training data is not sufficient to improve AC prediction accuracy. However, as shown in Figure 5 (see more results for other models in Figure 9), the ratio of AC ligands in the training set is a significant factor affecting $\text{RMSE}_{\text{cliff}} - \text{RMSE}$, with a p-value of $1.0e-4$. A higher percentage of AC ligands in the training set means more information directly relevant to AC, thereby improving the AC predictive power. Our finding indicates that the knowledge about general bioactivity prediction is different from the knowledge benefiting AC prediction. This underscores the need for new datasets and methodologies tailored specifically for AC prediction, as relying solely on general bioactivity data is insufficient to achieve optimal performance in this domain.

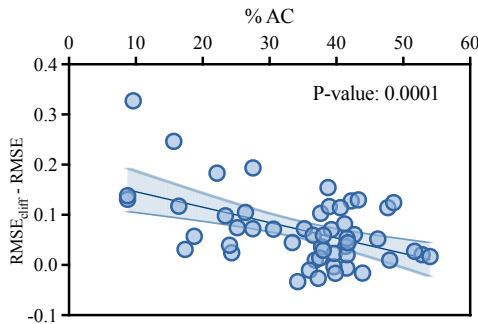


Figure 5: Relationship between the ratio of the AC and $\text{RMSE}_{\text{cliff}} - \text{RMSE}$ of IGN.

5.4 PERFORMANCE COMPARISON WITH MACHINE LEARNING ALGORITHMS

We benchmark the ability of all methods to predict bioactivity in the presence of the AC (measured by $\text{RMSE}_{\text{cliff}}$), as shown in Figure 6 (detailed results in Appendix Figure 12). We have the following empirical observations: (1) Significant performance differences can be observed among targets in the handling of AC compounds, with $\text{RMSE}_{\text{cliff}}$ values spanning from 0.52 to 1.59 log units, which is consistent with previous works (van Tilborg et al., 2022; Sheridan, 2012). (2) Among the four machine learning algorithms, performance disparities primarily stem from the molecule descriptors rather than the learning methods. ECFPs, which are designed specifically for structure-activity modeling by encoding detailed information about each atom’s local environment, yield the lowest prediction error of all methods. Their strong discriminative capability effectively differentiates molecules, even with minor structural differences. (3) For deep learning methods, IGN coupled with 3D structure information achieves the best performance on ACs. This approach benefits from the interaction information between the ligand and the protein target captured within the 3D structure.

Surprisingly, some machine learning models outperform deep learning approaches, which can be primarily attributed to their use of handcrafted features, especially ECFP. To validate this observation,

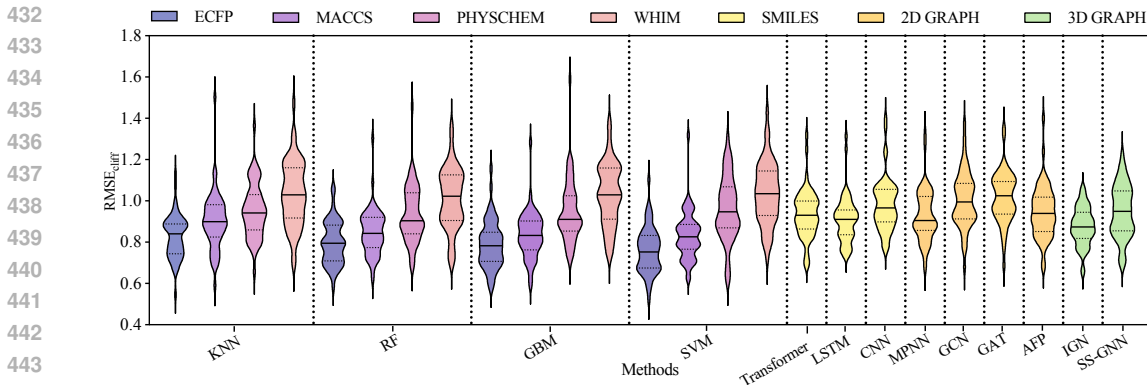


Figure 6: The $RMSE_{cliff}$ evaluated using different methods and features across 52 targets.

we implement a hybrid approach combining the ECFP features with the features extracted from the last layer of the 3D IGN model. These concatenated features are then fed into an MLP for prediction (as illustrated in Appendix Figure 15). The promising results across ten targets (see Appendix Table 10) demonstrate the effectiveness of ECFP in structure-activity relationship learning. On the other hand, this experiment underscores the value of integrating traditional cheminformatics techniques with advanced deep-learning methods in molecular property prediction tasks. Future research could explore optimizing this hybrid approach and investigating its applicability to a broader range of molecular targets and properties.

5.5 3D COMPLEX ENABLES CROSS-TARGET PREDICTION

The incorporation of 3D structural information also enables cross-target modeling capabilities. To investigate this potential, we conduct additional experiments by combining the targets with the same labels (e.g., K_i) and training the IGN model under two scenarios: (i) in-domain setting under all targets with the same labels, and (ii) out-of-distribution (OOD) setting excluding the family of testing targets. Analysis of the Protease and Nuclear receptor targets (Appendix Table 6, 7, 8, and 9) reveals that the absence of the testing family targets in training leads to performance degradation. In Table 6, the average $RMSE_{cliff}$ increasing from 0.9 to 1.4. While multi-target training achieves comparable performance to target-specific training across all targets (Appendix Figure 13), these results indicate that there is still a long way to go to fully exploit the multi-target 3D data and make the model generalize to new targets.

6 CONCLUSION

In this paper, we introduce *DockedAC*, a novel dataset for ACs with 3D complex structures. The dataset contains over 80k ligands from 52 protein targets, with the 3D structure of each target annotated by a unique known binding site. For each target, we generate protein-ligand complexes for at least 500 ligands using molecular docking. Benchmarking with various machine learning and deep learning approaches reveals that graph neural network (GNN)-based methods particularly benefit from 3D structural information, which enhances AC prediction accuracy and narrows the performance gap between general and AC-specific activity prediction.

Our experimental results demonstrate two key findings: (1) the absolute error in AC prediction exhibits significant target dependence, and (2) the proportion of AC ligands in the training set critically influences the disparity between general and AC activity prediction. Notably, current deep learning methods underperform compared to traditional machine learning approaches using molecular fingerprints, underscoring the urgent need for developing next-generation 3D-QSAR algorithms. *DockedAC* represents a crucial first step toward this goal by providing comprehensive 3D complex structures and target-ligand interaction data.

While *DockedAC* offers valuable structural insights, its reliance on molecular docking introduces potential inaccuracies in the generated complex structures. These limitations could be addressed through more advanced computational techniques, such as molecular dynamics simulations, for structural refinement.

486 STATEMENT
487

488 **Ethics Statement.** This research fully adheres to the ICLR Code of Ethics (<https://iclr.cc/public/CodeOfEthics>). We are committed to contributing to society and human well-being
489 by developing AI methods that enhance drug safety assessment, ultimately benefiting public health.
490 We have carefully considered potential harm and implemented safeguards to minimize negative
491 consequences, particularly with regard to patient privacy and data protection. All datasets used are
492 properly licensed (CC-BY-SA 4.0), and we have ensured appropriate attribution to original creators.

493 **Reproducibility Statement.** To ensure full reproducibility of our results, we provide compre-
494 hensive implementation details throughout the paper and appendices. Section B.5 contains com-
495 plete specifications of all experimental details, including data splits, hyperparameters, and op-
496 timizer configurations. We report error bars and statistical significance measures for all ex-
497 perimental results to ensure robust evaluation. Computational resource requirements, includ-
498 ing hardware specifications, memory usage, and execution times, are detailed in Section B.5
499 to facilitate reproduction on similar systems. The code for the benchmark is available here:
500 <https://anonymous.4open.science/r/DockedAC>. The DockedAC dataset and its fu-
501 ture updates can be found here: <https://doi.org/10.5281/zenodo.11485279>.
502
503

504 REFERENCES
505

506 Jinheon Baek, Minki Kang, and Sung Ju Hwang. Accurate learning of graph representations with
507 graph multiset pooling. In *International Conference on Learning Representations*, 2021. URL
508 <https://openreview.net/forum?id=JHcqXGaqiGn>.

509 Guy W Bemis and Mark A Murcko. The properties of known drugs. 1. Molecular frameworks.
510 *Journal of Medicinal Chemistry*, 39(15):2887–2893, 1996.
511

512 A Patrícia Bento, Anne Hersey, Eloy Félix, Greg Landrum, Anna Gaulton, Francis Atkinson, Louisa J
513 Bellis, Marleen De Veij, and Andrew R Leach. An open source chemical structure curation pipeline
514 using RDKit. *Journal of Cheminformatics*, 12:1–16, 2020.
515

516 Helen M Berman, John Westbrook, Zukang Feng, Gary Gilliland, Talapady N Bhat, Helge Weissig,
517 Ilya N Shindyalov, and Philip E Bourne. The protein data bank. *Nucleic Acids Research*, 28(1):
518 235–242, 2000.

519 Leo Breiman. Bagging predictors. *Machine Learning*, 24:123–140, 1996.
520

521 Duanhua Cao, Geng Chen, Jiaxin Jiang, Jie Yu, Runze Zhang, Mingan Chen, Wei Zhang, Lifan Chen,
522 Feisheng Zhong, Yingying Zhang, et al. Generic protein–ligand interaction scoring by integrating
523 physical prior knowledge and data augmentation modelling. *Nature Machine Intelligence*, pp.
524 1–13, 2024.
525

526 Hengwei Chen, Martin Vogt, and Jürgen Bajorath. DeepAC–conditional transformer-based chemical
527 language model for the prediction of activity cliffs formed by bioactive compounds. *Digital
528 Discovery*, 1(6):898–909, 2022.

529 Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. ChemBERTa: large-scale self-
530 supervised pretraining for molecular property prediction. *arXiv preprint arXiv:2010.09885*, 2020.
531

532 The UniProt Consortium. UniProt: the universal protein knowledgebase in 2023. *Nucleic Acids
533 Research*, 51(D1):D523–D531, 2023.
534

535 Thomas Cover and Peter Hart. Nearest neighbor pattern classification. *IEEE Transactions on
536 Information Theory*, 13(1):21–27, 1967.
537

538 Maykel Cruz-Monteagudo, José L Medina-Franco, Yunierkis Pérez-Castillo, Orazio Nicolotti,
539 M Natália DS Cordeiro, and Fernanda Borges. Activity cliffs in drug discovery: Dr Jekyll
or Mr Hyde? *Drug Discovery Today*, 19(8):1069–1080, 2014.

- 540 Maykel Cruz-Monteagudo, José L Medina-Franco, Yunier Perera-Sardiña, Fernanda Borges, Ed-
541 uardo Tejera, Cesar Paz-y Mino, Yunierkis Pérez-Castillo, Aminaél Sánchez-Rodríguez, Zuleidys
542 Contreras-Posada, Natália DS Cordeiro, et al. Probing the hypothesis of sar continuity restoration
543 by the removal of activity cliffs generators in qsar. *Current Pharmaceutical Design*, 22(33):
544 5043–5056, 2016.
- 545 Markus Dablander, Thierry Hanser, Renaud Lambiotte, and Garrett M Morris. Exploring QSAR
546 models for activity-cliff prediction. *Journal of Cheminformatics*, 15(1):47, 2023.
- 548 Mark Davies, Michał Nowotka, George Papadatos, Nathan Dedman, Anna Gaulton, Francis Atkinson,
549 Louisa Bellis, and John P Overington. ChEMBL web services: streamlining access to drug
550 discovery data and utilities. *Nucleic Acids Research*, 43(W1):W612–W620, 2015.
- 552 Jianyuan Deng, Zhibo Yang, Hehe Wang, Iwao Ojima, Dimitris Samaras, and Fusheng Wang. A sys-
553 tematic study of key elements underlying molecular property prediction. *Nature Communications*,
554 14(1):6395, 2023.
- 555 Chengwei Dong, Yu-Peng Huang, Xiaohan Lin, Hong Zhang, and Yi Qin Gao. Dsdpflex: Flexible-
556 receptor docking with gpu acceleration. *Journal of Chemical Information and Modeling*, 2024.
- 558 Joseph L Durant, Burton A Leland, Douglas R Henry, and James G Nourse. Reoptimization of MDL
559 keys for use in drug discovery. *Journal of Chemical Information and Computer Sciences*, 42(6):
560 1273–1280, 2002.
- 561 Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of*
562 *Statistics*, pp. 1189–1232, 2001.
- 564 Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
565 message passing for quantum chemistry. In *International Conference on Machine Learning*, pp.
566 1263–1272. PMLR, 2017.
- 567 Rajarshi Guha. Exploring uncharted territories: Predicting activity cliffs in structure–activity land-
568 scapes. *Journal of Chemical Information and Modeling*, 52(8):2181–2191, 2012.
- 570 Marti A. Hearst, Susan T Dumais, Edgar Osuna, John Platt, and Bernhard Scholkopf. Support vector
571 machines. *IEEE Intelligent Systems and their Applications*, 13(4):18–28, 1998.
- 572 Kathrin Heikamp, Xiaoying Hu, Aixia Yan, and Jürgen Bajorath. Prediction of activity cliffs using
573 support vector machines. *Journal of Chemical Information and Modeling*, 52(9):2354–2365, 2012.
- 575 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):
576 1735–1780, 1997.
- 578 Dragos Horvath, Gilles Marcou, Alexandre Varnek, Shilva Kayastha, Antonio de la Vega de León, and
579 Jürgen Bajorath. Prediction of activity cliffs using condensed graphs of reaction representations,
580 descriptor recombination, support vector machine classification, and support vector regression.
581 *Journal of Chemical Information and Modeling*, 56(9):1631–1640, 2016.
- 582 Xiuyuan Hu, Guoqing Liu, Yang Zhao, and Hao Zhang. Activity cliff-aware reinforcement learning
583 for de novo drug design. *Journal of Cheminformatics*, 17(1):54, 2025.
- 584 Ye Hu and Jürgen Bajorath. Exploration of 3d activity cliffs on the basis of compound binding
585 modes and comparison of 2d and 3d cliffs. *Journal of Chemical Information and Modeling*, 52(3):
586 670–677, 2012.
- 588 Ye Hu, Norbert Furtmann, Michael Gütschow, and Jürgen Bajorath. Systematic identification and
589 classification of three-dimensional activity cliffs. *Journal of Chemical Information and Modeling*,
590 52(6):1490–1498, 2012.
- 591 YuPeng Huang, Hong Zhang, Siyuan Jiang, Dajiong Yue, Xiaohan Lin, Jun Zhang, and Yi Qin Gao.
592 DSDP: A blind docking strategy accelerated by GPUs. *Journal of Chemical Information and*
593 *Modeling*, 63(14):4355–4363, 2023.

- 594 Jarmila Husby, Giovanni Bottegoni, Irina Kufareva, Ruben Abagyan, and Andrea Cavalli. Structure-
595 based predictions of activity cliffs. *Journal of Chemical Information and Modeling*, 55(5):1062–
596 1076, 2015.
- 597 Regina Ibragimova, Dimitrios Iliadis, and Willem Waegeman. Enhancing drug-target interaction
598 prediction through transfer learning from activity cliff prediction tasks. *Journal of Chemical*
599 *Information and Modeling*, 65(13):6558–6567, 2025.
- 600 Javed Iqbal, Martin Vogt, and Jürgen Bajorath. Prediction of activity cliffs on the basis of images
601 using convolutional neural networks. *Journal of Computer-Aided Molecular Design*, pp. 1–8,
602 2021.
- 603 Dejun Jiang, Chang-Yu Hsieh, Zhenxing Wu, Yu Kang, Jike Wang, Ercheng Wang, Ben Liao, Chao
604 Shen, Lei Xu, Jian Wu, et al. InteractionGraphNet: A novel and efficient deep graph representation
605 learning framework for accurate protein–ligand interaction predictions. *Journal of Medicinal*
606 *Chemistry*, 64(24):18209–18232, 2021.
- 607 José Jiménez, Miha Skalic, Gerard Martínez-Rosell, and Gianni De Fabritiis. K deep: protein–ligand
608 absolute binding affinity prediction via 3d-convolutional neural networks. *Journal of Chemical*
609 *Information and Modeling*, 58(2):287–296, 2018.
- 610 José Jiménez-Luna, Miha Skalic, and Nils Weskamp. Benchmarking molecular feature attribution
611 methods with activity cliffs. *Journal of Chemical Information and Modeling*, 62(2):274–283, 2022.
- 612 Wengong Jin, Caroline Uhler, and Nir Hacohen. Se (3) denoising score matching for unsupervised
613 binding energy prediction and nanobody design. In *NeurIPS 2023 Generative AI and Biology*
614 *(GenBio) Workshop*, 2023.
- 615 Talia B Kimber, Maxime Gagnebin, and Andrea Volkamer. Maxsmi: maximizing molecular property
616 prediction performance with confidence estimation using smiles augmentation and deep learning.
617 *Artificial Intelligence in the Life Sciences*, 1:100014, 2021.
- 618 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
619 *arXiv preprint arXiv:1609.02907*, 2016.
- 620 Greg Landrum et al. RDKit: A software suite for cheminformatics, computational chemistry, and
621 predictive modeling. *Greg Landrum*, 8(31.10):5281, 2013.
- 622 Cheryl S Leung, Siegfried SF Leung, Julian Tirado-Rives, and William L Jorgensen. Methyl effects
623 on protein–ligand binding. *Journal of Medicinal Chemistry*, 55(9):4489–4500, 2012.
- 624 Vladimir I Levenshtein et al. Binary codes capable of correcting deletions, insertions, and reversals.
625 In *Soviet physics doklady*, volume 10, pp. 707–710. Soviet Union, 1966.
- 626 Kwei Li, Yuqian Wu, Yinheng Li, Yutong Guo, Yanwen Kong, Yan Wang, Yiyang Liang, Yusi Fan,
627 Lan Huang, Ruochi Zhang, et al. AMPCliff: quantitative definition and benchmarking of activity
628 cliffs in antimicrobial peptides. *Journal of Advanced Research*, 2025.
- 629 Wei Lu, Qifeng Wu, Jixian Zhang, Jiahua Rao, Chengtao Li, and Shuangjia Zheng. Tankbind:
630 Trigonometry-aware neural networks for drug-protein binding structure prediction. *Advances in*
631 *neural information processing systems*, 35:7236–7249, 2022.
- 632 Gerald M Maggiora. On outliers and activity cliffs why QSAR often disappoints, 2006.
- 633 David Mendez, Anna Gaulton, A Patrícia Bento, Jon Chambers, Marleen De Veij, Eloy Félix,
634 María Paula Magariños, Juan F Mosquera, Prudence Mutowo, Michał Nowotka, et al. ChEMBL:
635 towards direct deposition of bioassay data. *Nucleic Acids Research*, 47(D1):D930–D940, 2019.
- 636 Seokhyun Moon, Wonho Zhung, Soojung Yang, Jaechang Lim, and Woo Youn Kim. Pignet: a
637 physics-informed deep learning model toward generalized drug–target interaction predictions.
638 *Chemical Science*, 13(13):3661–3673, 2022.
- 639 Junhui Park, Gaeun Sung, SeungHyun Lee, SeungHo Kang, and ChunKyun Park. ACGCN: graph
640 convolutional networks for activity cliff prediction between matched molecular pairs. *Journal of*
641 *Chemical Information and Modeling*, 62(10):2341–2351, 2022.

- 648 Saro Passaro, Gabriele Corso, Jeremy Wohlwend, Mateo Reveiz, Stephan Thaler, Vignesh Ram
649 Somnath, Noah Getz, Tally Portnoi, Julien Roy, Hannes Stark, et al. Boltz-2: Towards accurate
650 and efficient binding affinity prediction. *BioRxiv*, pp. 2025–06, 2025.
- 651
652 Lewis D Pennington and Demetri T Moustakas. The necessary nitrogen atom: a versatile high-
653 impact design element for multiparameter optimization. *Journal of Medicinal Chemistry*, 60(9):
654 3552–3579, 2017.
- 655 David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of Chemical Informa-
656 tion and Modeling*, 50(5):742–754, 2010.
- 657
658 Chao Shen, Xujun Zhang, Yafeng Deng, Junbo Gao, Dong Wang, Lei Xu, Peichen Pan, Tingjun Hou,
659 and Yu Kang. Boosting protein–ligand binding pose prediction and virtual screening based on
660 residue–atom distance likelihood potential and graph transformer. *Journal of Medicinal Chemistry*,
661 65(15):10691–10706, 2022.
- 662 Robert P Sheridan. Three useful dimensions for domain applicability in QSAR models using random
663 forest. *Journal of Chemical Information and Modeling*, 52(3):814–823, 2012.
- 664
665 MJ Stewart and ID Watson. Standard units for expressing drug concentrations in biological fluids.
666 *British Journal of Clinical Pharmacology*, 16(1):3, 1983.
- 667 Dagmar Stumpfe, Ye Hu, Dilyana Dimova, and Jürgen Bajorath. Recent progress in understanding
668 activity cliffs and their utility in medicinal chemistry: miniperspective. *Journal of medicinal
669 chemistry*, 57(1):18–28, 2014.
- 670 Dagmar Stumpfe, Huabin Hu, and Jürgen Bajorath. Evolving concept of activity cliffs. *Acs Omega*,
671 4(11):14360–14368, 2019.
- 672
673 Dagmar Stumpfe, Huabin Hu, and Jürgen Bajorath. Advances in exploring activity cliffs. *Journal of
674 Computer-Aided Molecular Design*, 34:929–942, 2020.
- 675
676 Shunsuke Tamura, Swarit Jasial, Tomoyuki Miyao, and Kimito Funatsu. Interpretation of ligand-
677 based activity cliff prediction models using the matched molecular pair kernel. *Molecules*, 26(16):
678 4916, 2021.
- 679 Shunsuke Tamura, Tomoyuki Miyao, and Jürgen Bajorath. Large-scale prediction of activity cliffs
680 using machine and deep learning methods of increasing complexity. *Journal of Cheminformatics*,
681 15(1):4, 2023.
- 682 Huishuang Tan, Zhixin Wang, and Guang Hu. GAABind: a geometry-aware attention-based network
683 for accurate protein–ligand binding pose and binding affinity prediction. *Briefings in Bioinformatics*,
684 25(1):bbad462, 2024.
- 685
686 Taffee T Tanimoto. Elementary mathematical theory of classification and prediction. 1958.
- 687 Roberto Todeschini, Paola Gramatica, et al. New 3D molecular descriptors: the WHIM theory and
688 QSAR applications. *Perspectives in Drug Discovery and Design*, 9(0):355–380, 1998.
- 689
690 Oleg Trott and Arthur J Olson. AutoDock Vina: improving the speed and accuracy of docking with
691 a new scoring function, efficient optimization, and multithreading. *Journal of Computational
692 Chemistry*, 31(2):455–461, 2010.
- 693 Jessica Vamathevan, Dominic Clark, Paul Czodrowski, Ian Dunham, Edgardo Ferran, George Lee,
694 Bin Li, Anant Madabhushi, Parantu Shah, Michaela Spitzer, et al. Applications of machine learning
695 in drug discovery and development. *Nature Reviews Drug Discovery*, 18(6):463–477, 2019.
- 696
697 Derek van Tilborg, Alisa Alenicheva, and Francesca Grisoni. Exposing the limitations of molecular
698 machine learning with activity cliffs. *Journal of Chemical Information and Modeling*, 62(23):
699 5938–5951, 2022.
- 700 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
701 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing
Systems*, 30, 2017.

- 702 Martin Vogt, Yun Huang, and Jürgen Bajorath. From activity cliffs to activity ridges: informative
703 data structures for sar analysis. *Journal of Chemical Information and Modeling*, 51(8):1848–1856,
704 2011.
- 705 W Patrick Walters and Mark A Murcko. Prediction of ‘drug-likeness’. *Advanced Drug Delivery*
706 *Reviews*, 54(3):255–271, 2002.
- 707
708 Renxiao Wang, Xueliang Fang, Yipin Lu, and Shaomeng Wang. The PDBbind database: Collection of
709 binding affinities for protein-ligand complexes with known three-dimensional structures. *Journal*
710 *of Medicinal Chemistry*, 47(12):2977–2980, 2004.
- 711 Renxiao Wang, Xueliang Fang, Yipin Lu, Chao-Yie Yang, and Shaomeng Wang. The pddbnd
712 database: methodologies and updates. *Journal of medicinal chemistry*, 48(12):4111–4119, 2005.
- 713
714 David Weininger. SMILES, a chemical language and information system. 1. Introduction to method-
715 ology and encoding rules. *Journal of Chemical Information and Computer Sciences*, 28(1):31–36,
716 1988.
- 717 Zhaoping Xiong, Dingyan Wang, Xiaohong Liu, Feisheng Zhong, Xiaozhe Wan, Xutong Li, Zhaojun
718 Li, Xiaomin Luo, Kaixian Chen, Hualiang Jiang, et al. Pushing the boundaries of molecular
719 representation for drug discovery with the graph attention mechanism. *Journal of Medicinal*
720 *Chemistry*, 63(16):8749–8760, 2019.
- 721
722 Xiaocong Yang, Yang Liu, Jianhong Gan, Zhi-Xiong Xiao, and Yang Cao. FitDock: Protein–ligand
723 docking by template fitting. *Briefings in Bioinformatics*, 23(3):bbac087, 2022.
- 724
725 Xinxin Yu, Yimeng Wang, Long Chen, Weihua Li, Yun Tang, and Guixia Liu. ACtriplet: An improved
726 deep learning model for activity cliffs prediction by integrating triplet loss and pre-training. *Journal*
727 *of Pharmaceutical Analysis*, pp. 101317, 2025.
- 728
729 Xin Zeng, Shu-Juan Li, Shuang-Qing Lv, Meng-Liang Wen, and Yi Li. A comprehensive review
730 of the recent advances on predicting drug-target affinity based on deep learning. *Frontiers in*
731 *Pharmacology*, 15:1375522, 2024.
- 732
733 Shuke Zhang, Yanzhao Jin, Tianmeng Liu, Qi Wang, Zhaohui Zhang, Shuliang Zhao, and Bo Shan.
734 SS-GNN: a simple-structured graph neural network for affinity prediction. *ACS Omega*, 8(25):
735 22496–22507, 2023a.
- 736
737 Xujun Zhang, Odin Zhang, Chao Shen, Wanglin Qu, Shicheng Chen, Hanqun Cao, Yu Kang, Zhe
738 Wang, Ercheng Wang, Jintu Zhang, et al. Efficient and accurate large library ligand docking with
739 karmadock. *Nature Computational Science*, 3(9):789–804, 2023b.
- 740
741 Ziqiao Zhang, Bangyi Zhao, Ailin Xie, Yatao Bian, and Shuigeng Zhou. Activity cliff prediction:
742 Dataset and benchmark. *arXiv preprint arXiv:2302.07541*, 2023c.
- 743
744
745
746
747
748
749
750
751
752
753
754
755

APPENDIX

A LLM USAGE

A large language model (Claude Sonnet 4) was used as a writing assistance tool to improve the linguistic quality and readability of this manuscript. Specifically, we employed it for tasks such as sentence rephrasing, grammar correction, and enhancing the clarity and flow of the text presentation.

It is important to note that the LLM was **not** involved in any aspect of research ideation, methodology development, experimental design, or scientific reasoning. All research concepts, hypotheses, experimental procedures, and results interpretation were independently conceived and conducted by the authors without LLM assistance.

The authors take full responsibility for all content in this manuscript, including any LLM-assisted text, and have verified that all LLM usage adheres to ethical guidelines and does not contribute to plagiarism or scientific misconduct.

Table 3: Dataset overview. n (where $n_{\text{train}}/n_{\text{test}}$, *resp.*) represents the total number of compounds, divided into training and test sets. n^{AC} (where $n_{\text{train}}^{AC}/n_{\text{test}}^{AC}$, *resp.*) denotes the total number of activity cliff compounds within the dataset, also divided into training and test sets.

| Target Name | ChEMBL ID | PDB | Type | n ($n_{\text{train}} / n_{\text{test}}$) | n^{AC} ($n_{\text{train}}^{AC} / n_{\text{test}}^{AC}$) |
|--|------------|--------|-----------|--|---|
| Androgen Receptor | ChEMBL1871 | 2ama | K_i | 617 (492/125) | 135 (109/26) |
| Cannabinoid CB1 receptor | ChEMBL218 | 6kqi | EC_{50} | 1004 (802/202) | 369 (293/76) |
| Coagulation factor X | ChEMBL244 | 2p93 | K_i | 3093 (2474/619) | 1476 (1180/296) |
| Delta opioid receptor | ChEMBL236 | 6pt3 | K_i | 2580 (2060/520) | 1005 (802/203) |
| Dopamine D3 receptor | ChEMBL234 | 3pbl_A | K_i | 3657 (2924/733) | 1604 (1284/320) |
| Dopamine D4 receptor | ChEMBL219 | 5wiu_A | K_i | 1865 (1491/374) | 740 (592/148) |
| Dopamine transporter | ChEMBL238 | 2q6h_A | K_i | 1051 (838/213) | 266 (211/55) |
| Dual specificity protein kinase CLK4 | ChEMBL4203 | 6fyv | K_i | 731 (582/149) | 64 (51/13) |
| Bile acid receptor FXR | ChEMBL2047 | 5q0u | EC_{50} | 631 (503/128) | 245 (195/50) |
| Ghrelin receptor | ChEMBL4616 | 6ko5_A | EC_{50} | 673 (534/139) | 355 (282/73) |
| Glucocorticoid receptor | ChEMBL2034 | 4lsj | K_i | 684 (551/133) | 243 (194/49) |
| Glycogen synthase kinase-3 beta | ChEMBL262 | 6hk3 | K_i | 855 (683/172) | 160 (128/32) |
| Histamine H1 receptor | ChEMBL231 | 3rze_A | K_i | 972 (776/196) | 237 (189/48) |
| Histamine H3 receptor | ChEMBL264 | 7f61_A | K_i | 2862 (2288/574) | 1191 (952/239) |
| Tyrosine-protein kinase JAK1 | ChEMBL2835 | 4k77 | K_i | 615 (489/126) | 60 (47/13) |
| Tyrosine-protein kinase JAK2 | ChEMBL2971 | 4jja | K_i | 976 (779/197) | 162 (128/34) |
| Kappa opioid receptor | ChEMBL237 | 4djh | EC_{50} | 953 (761/192) | 456 (365/91) |
| Kappa opioid receptor | ChEMBL237 | 4djh | K_i | 2599 (2078/521) | 1109 (887/222) |
| Orexin receptor 2 | ChEMBL4792 | 5wqc | K_i | 1471 (1174/297) | 794 (634/160) |
| Peroxisome proliferator-activated receptor alpha | ChEMBL239 | 3kdu | EC_{50} | 1721 (1374/347) | 699 (558/141) |
| Peroxisome proliferator-activated receptor delta | ChEMBL3979 | 5xmx | EC_{50} | 1125 (899/226) | 468 (374/94) |
| Peroxisome proliferator-activated receptor gamma | ChEMBL235 | 2yfe | EC_{50} | 2349 (1877/472) | 885 (707/178) |
| PI3-kinase p110-alpha subunit | ChEMBL4005 | 6gvf | K_i | 960 (767/193) | 401 (320/81) |
| Serine/threonine-protein kinase PIM1 | ChEMBL2147 | 2j2i | K_i | 1456 (1162/294) | 572 (456/116) |
| Serotonin 1a (5-HT1a) receptor | ChEMBL214 | 7e2x_R | K_i | 3317 (2651/666) | 1222 (977/245) |
| Serotonin transporter | ChEMBL228 | 6awo_A | K_i | 1702 (1362/340) | 638 (511/127) |
| Sigma opioid receptor | ChEMBL287 | 6dk1 | K_i | 1328 (1061/267) | 507 (404/103) |
| Thrombin | ChEMBL204 | 1mu8 | K_i | 2747 (2195/552) | 1089 (870/219) |
| Tyrosine-protein kinase ABL | ChEMBL1862 | 2hzi | K_i | 794 (633/161) | 330 (263/67) |
| Mu opioid receptor | ChEMBL233 | 8feo_R | K_i | 3141 (2511/630) | 1294 (1035/259) |
| Cyclin-dependent kinase 2 | ChEMBL301 | 1h1q | IC_{50} | 1454 (1161/293) | 350 (279/71) |
| Serine/threonine-protein kinase Chk1 | ChEMBL4630 | 2brb | IC_{50} | 1701 (1359/342) | 826 (660/166) |
| 3-phosphoinositide dependent protein kinase-1 | ChEMBL2534 | 1uu3 | IC_{50} | 705 (562/143) | 282 (224/58) |
| Phosphodiesterase 5A | ChEMBL1827 | 4ia0 | IC_{50} | 1609 (1285/324) | 667 (532/135) |
| Dihydrofolate reductase | ChEMBL202 | 1u71 | IC_{50} | 739 (590/149) | 281 (223/58) |
| Urokinase-type plasminogen activator | ChEMBL3286 | 1owe | K_i | 718 (572/146) | 191 (151/40) |
| Carbonic anhydrase II | ChEMBL205 | 5sz6 | K_i | 5796 (4636/1160) | 2444 (1957/487) |
| Estrogen receptor alpha | ChEMBL206 | 1qkt | IC_{50} | 2094 (1674/420) | 700 (559/141) |
| Heat shock protein HSP 90-alpha | ChEMBL3880 | 4o0b | IC_{50} | 999 (797/202) | 157 (125/32) |
| Fructose-1,6-bisphosphatase | ChEMBL3975 | 2jjk | IC_{50} | 556 (443/113) | 153 (122/31) |
| Protein-tyrosine phosphatase 1B | ChEMBL335 | 1nny | IC_{50} | 2607 (2084/523) | 229 (183/46) |
| Matrix metalloproteinase 8 | ChEMBL4588 | 3dng | IC_{50} | 533 (425/108) | 163 (130/33) |
| Dipeptidyl peptidase IV | ChEMBL284 | 2ole | IC_{50} | 2507 (2003/504) | 691 (551/140) |
| Vascular endothelial growth factor receptor 2 | ChEMBL279 | 3vhk | K_i | 780 (622/158) | 135 (108/27) |
| Matrix metalloproteinase 13 | ChEMBL280 | 4jpa | IC_{50} | 2112 (1688/424) | 976 (780/196) |
| Methionine aminopeptidase 2 | ChEMBL3922 | 6qef | IC_{50} | 565 (450/115) | 193 (154/39) |
| Kinesin-like protein 1 | ChEMBL4581 | 5zo8 | IC_{50} | 719 (573/146) | 311 (248/63) |
| Beta-secretase 1 | ChEMBL4822 | 4h3j | K_i | 1061 (847/214) | 549 (438/111) |
| Phosphodiesterase 4B | ChEMBL275 | 3w5e | IC_{50} | 1432 (1143/289) | 535 (426/109) |
| Phosphodiesterase 4D | ChEMBL288 | 2qyn | IC_{50} | 942 (752/190) | 220 (176/44) |
| MAP kinase p38 alpha | ChEMBL260 | 2zbl | IC_{50} | 3502 (2799/703) | 1333 (1065/268) |
| Estrogen receptor beta | ChEMBL242 | lzaf | IC_{50} | 1176 (937/239) | 425 (337/88) |

B DATASETS AND BASELINE MODELS

B.1 LICENSE AND AVAILABILITY

The code for benchmark is available here: <https://anonymous.4open.science/r/DockedAC>. The DockedAC dataset and its future updates can be found here: <https://doi.org/10.5281/zenodo.11485279>.

The DockedAC dataset is licensed under the *Creative Commons Attribution-ShareAlike 4.0 International License*. For details, please see <https://creativecommons.org/licenses/by-sa/4.0/>. The content of DockedAC includes data from the following sources: **RCSB PDB**, which is available under the *CC0 1.0 Universal (CC0 1.0) Public Domain Dedication* (more information at <https://creativecommons.org/publicdomain/zero/1.0/>), and **ChEMBL**, which is licensed under the *Creative Commons Attribution-ShareAlike 3.0 Unported License* (see <https://creativecommons.org/licenses/by-sa/3.0/>).

B.2 DATASETS

Our introduced dataset DockedAC¹ comprises 82,836 ligands from 52 protein targets, which is meticulously curated to support various machine learning and deep learning studies related to activity cliff (AC) prediction. Table 3 provides detailed statistics of DockedAC.

B.3 BASELINE MODELS

In this work, we integrate 13 recent baselines commonly used for structure-activity relationship prediction, including four traditional machine learning algorithms: KNN, RF, GBM, and SVM; three sequential models: LSTM, Transformer, and 1D CNN; four 2D GNN models: GCN, GAT, MPNN, and AFP; and two 3D structure GNN models: IGN and SS-GNN.

Besides, we also evaluate 6 state-of-the-art binding affinity prediction models (PIGNet, RTMScore, TANKBind, DSMBind, KarmaDock, and EquiScore) pre-trained on PDBBind to assess their zero-shot generalization ability on our benchmark. The detailed descriptions of these approaches are listed in the following:

- **KNN** (Cover & Hart, 1967). K-Nearest Neighbor (KNN) is a simple, non-parametric method that predicts the target molecule’s response by averaging the response of the k-nearest neighbors from the training set.
- **RF** (Breiman, 1996). Random Forest (RF) is an ensemble method that combines the outputs of multiple decision trees to improve accuracy and reduce over-fitting. Each decision tree is built upon a subset of the training set, and the final prediction is obtained by averaging the results from these individual trees.
- **GBM** (Friedman, 2001). Similar to RF, Gradient Boosting Machine (GBM) also combines the predictions of multiple decision trees. However, in GBM, these trees are built sequentially, with each subsequent tree specially designed to correct the errors of its predecessors.
- **SVM** (Hearst et al., 1998). Support Vector Machine (SVM) aims to identify a linear regression plane in a higher-dimensional space created by applying a designated kernel function. In this work, the Radial Basis Function (RBF) kernel is used.
- **Transformer** (Vaswani et al., 2017). The Transformer model leverages self-attention mechanisms to capture dependencies across different positions in the input sequence. In our work, we employed the pretrained ChemBERTa (Chithrananda et al., 2020) architecture, which has been trained on 10 million compounds.
- **LSTM** (Hochreiter & Schmidhuber, 1997). Long Short-Term Memory (LSTM) can capture temporal dependencies and patterns in sequential data by maintaining long-term memory through their gated structure. In this work, we employ SMILES strings as the input for the model.

¹<https://anonymous.4open.science/r/DockedAC>

Table 4: Featurization and corresponding baseline models.

| Featurization | Baseline Models | Aug. |
|---------------------|---------------------------|--------|
| ECFP Descriptor | KNN, RF, GBM, SVM, | ✗ |
| MACCS Descriptor | KNN, RF, GBM, SVM, | ✗ |
| PHYSCHEM Descriptor | KNN, RF, GBM, SVM, | ✗ |
| WHIM Descriptor | KNN, RF, GBM, SVM, | ✗ |
| SMILES string | LSTM, Transformer, 1D CNN | ✓ × 10 |
| 2D GRAPH | MPNN, GCN, GAT, AFP | ✗ |
| 3D GRAPH | IGN, SS-GNN | ✗ |

- **1D CNN** (Kimber et al., 2021). Convolutional Neural Network (CNN) uses convolutional filters to aggregate spatial information from adjacent positions. For processing sequential SMILES string data, we employ 1D CNNs that perform convolutional operations along a single dimension.
- **MPNN** (Gilmer et al., 2017). Message Passing Neural Network (MPNN) operates by iteratively passing messages between nodes and updating their representations based on neighboring nodes.
- **GCN** (Kipf & Welling, 2016). Graph Convolutional Network (GCN) performs convolution operations on graphs.
- **GAT** (Vaswani et al., 2017). Graph Attention Network (GAT) introduces attention mechanisms to GNN to weigh the importance of different neighbors.
- **AFP** (Xiong et al., 2019). Attentive Fingerprint (AFP) uses attention mechanisms at both the atom and molecule levels to learn local and nonlocal properties, enabling it to capture substructural details effectively.
- **IGN** (Jiang et al., 2021). IGN models the molecular interactions in 3D space. In IGN, two graph convolution modules are layered to learn intramolecular interactions and then sequentially intermolecular interactions.
- **SS-GNN** (Zhang et al., 2023a). Like IGN, SS-GNN is also a 3D structure GNN model tailored for affinity prediction. It constructs a 3D structure graph for protein-ligand interactions based on a distance threshold, reducing both the graph data scale and computational cost by omitting covalent bonds in proteins.
- **PIGNet** (Moon et al., 2022). Physics-Informed Graph Neural Network (PIGNet) predicts binding affinity as a sum of atom-atom pairwise interactions computed through physics-informed equations parameterized with neural networks, enhanced by data augmentation with computationally generated binding poses.
- **RTMScore** (Shen et al., 2022). RTMScore employs a tailored residue-based graph representation strategy with graph transformer layers to learn protein and ligand representations, followed by a mixture density network to obtain residue-atom distance likelihood potential for robust protein-ligand scoring.
- **TANKBind** (Lu et al., 2022). TANKBind uses a trigonometry-aware neural network with kinematics (TANK) to model protein-ligand binding by incorporating geometric constraints and molecular dynamics principles.
- **DSMBind** (Jin et al., 2023). DSMBind leverages SE(3)-equivariant GNN to maintain rotational and translational invariance while predicting protein-ligand binding affinities.
- **KarmaDock** (Zhang et al., 2023b). KarmaDock combines molecular docking with deep learning to predict binding poses and affinities efficiently by incorporating both geometric and chemical features.
- **EquiScore** (Cao et al., 2024). EquiScore employs a heterogeneous GNN to integrate physical prior knowledge and characterize protein-ligand interactions in equivariant geometric space.

B.4 MODEL FEATURES

In addition to the molecular descriptor used for machine learning algorithms (introduced in Sec. 4.2), we further delve into the featurization for deep learning models. Detailed information on all featurizations and the corresponding models used can be found in Table 4.

Table 5: Hyperparameter search space.

| Methods | Hyperparameters | Search Space |
|--|--|---|
| KNN | The number of nearest neighbors, k | $k = [3, 5, 11, 21]$ |
| RF | The number of trees, n_t | $n_t = [100, 250, 500, 1000]$ |
| GBM | The number of boosting stages, n_b | $n_b = [100, 200, 400]$ |
| | The maximum depth of the model, n_d | $n_d = [5, 6, 7]$ |
| SVM | The regularization parameter, C | $C = [1, 10, 100, 1000, 10,000]$ |
| | The kernel coefficient for <i>rbf</i> , γ | $\gamma = [1 \times 10^{-5}, 1 \times 10^{-4}, 1 \times 10^{-3}, 1 \times 10^{-2}, 1 \times 10^{-1}]$ |
| <i>Shared hyperparameters for all deep learning models</i> | | |
| Common | The learning rate, lr | $lr = [5 \times 10^{-4}, 1 \times 10^{-4}, 5 \times 10^{-5}, 1 \times 10^{-5}]$ |
| | The batch size, bs | $bs = [10, 32, 64, 128]$ |
| | The epoch, γ | $\gamma = 500$ |
| <i>Specific hyperparameters for each model</i> | | |
| GCN | The dimension of hidden node features, h_n | $h_n = [64, 128, 256]$ |
| | The dimension of hidden transformer nodes, h_t | $h_t = [64, 128, 256]$ |
| | The dimension of predictor, h_p | $h_p = [128, 256, 512]$ |
| GAT | The dimension of hidden node features, h_n | $h_n = [64, 128, 256]$ |
| | The dimension of hidden transformer nodes, h_t | $h_t = [64, 128, 256]$ |
| | The dimension of predictor, h_p | $h_p = [128, 256, 512]$ |
| MPNN | The dimension of hidden node features, h_n | $h_n = [64, 128, 256]$ |
| | The dimension of hidden edge features, h_e | $h_e = [64, 128, 256]$ |
| | The dimension of hidden transformer nodes, h_t | $h_t = [64, 128, 256]$ |
| AFP | The dimension of hidden node features, h_n | $h_n = [64, 128, 256]$ |
| | The number of iterations for readout, n_r | $n_r = [1, 2, 3, 4, 5]$ |
| LSTM | - <i>pretrained</i> | - <i>pretrained</i> |
| Transformer | - <i>pretrained</i> | - <i>pretrained</i> |
| 1D CNN | The size of convolution kernel, h_c | $h_c = [4, 8, 10]$ |
| | The dimension of hidden features, h_t | $h_t = [64, 128, 256, 512, 1024]$ |
| IGN | The dimension of hidden features, h_t | $h_t = [64, 128, 256]$ |
| SS-GNN | - <i>pretrained</i> | - <i>pretrained</i> |

For sequential methods, SMILES strings were encoded as one-hot vectors, with truncation applied to strings exceeding 200 characters. To enhance model robustness, tenfold data augmentation was applied using up to nine additional noncanonical SMILES strings for each SMILES string in the dataset, generated via RDKit (Landrum et al., 2013).

For 2D GNN methods, the node has the following features: atom type (one-hot), atomic vertex degree (one-hot), orbital hybridization (one-hot), aromaticity (one-hot), atomic weight (float), formal charge (integer), number of radical electrons (integer), and number of connected hydrogens (integer). For MPNN and AFP, two one-hot bond features are used for the edges, i.e., the bond type and conjugation.

For SS-GNN, there are 11 node features, including atom type, formal charge, hybridization, atom valence, atom degree, number of hydrogens, chirality, atomic mass, aromatic, atom coordinates, and whether belonging to the protein. The edge features include covalent bond type, aromatic, bond length, bond direction, bond stereochemistry, and edge type. The atom coordinates and bond length are extracted from the 3D structures. Further details can be found in Zhang et al. (2023a).

For IGN, it uses similar 2D node and edge features. In addition, IGN uses four new edge features from the 3D structures, including bond length, angle statistics, area statistics, and distance statistics. For detailed descriptions of the features, see Jiang et al. (2021).

B.5 ADDITIONAL EXPERIMENTAL DETAILS

Hardware Specifications. All our experiments were carried out on an NVIDIA RTX3090 GPU with 24G memory. The training time of a target for MPNN, GAT, GCN, and AFP is around 0.5 hours. Training of one target takes around 1 hour and 4 hours for SS-GNN and IGN, respectively.

Implementation Details. Traditional machine learning algorithms including KNN, SVM, GBM, and RF regression models were implemented using the Scikit-Learn library².

Deep learning algorithms were trained for 500 epochs with early stopping, set with patience of 10 epochs. Four GNN models are implemented using the PyTorch Geometric package³. For the MPNN, GCN, and GAT, global pooling was enabled using a graph multiset Transformer (Baek et al., 2021) with eight attention heads, followed by a fully connected prediction head. Each of these models utilized two graph layers. The Transformer model was based on the ChemBERTa (Chithrananda et al., 2020) architecture, using weights derived from 10M compounds in PubChem. Fine-tuning was conducted by freezing the original model weights and substituting the final pooling layer with a regression head. Following van Tilborg et al. (2022), the LSTM model is pretrained on the SMILES strings with the next token prediction objective. For the SS-GNN model, we conducted a pretraining phase on the original dataset, PDBbind V2019 (Wang et al., 2004; 2005). In contrast, the IGN model was not fine-tuned using the original dataset due to a mismatch in the model dimensions caused by the varying types of atoms in the dataset. Consequently, we opted to train the IGN model from scratch.

Hyperparameter Optimization. Hyperparameter optimization was conducted through grid search. Hyperparameter combinations were evaluated for all models using five-fold cross-validation. Table 5 shows the detailed hyperparameter search space.

C ADDITIONAL RESULTS AND FIGURES

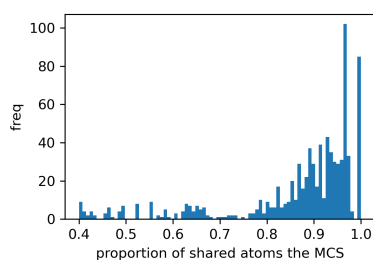


Figure 7: Proportion of shared atoms in AC data pairs with MCS for Target ChEMBL218 EC_{50} .

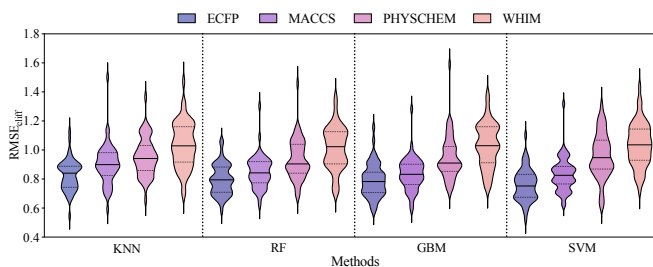


Figure 8: The $RMSE_{cliff}$ evaluated using ML methods under MACCS split across 52 targets.

More dataset features. Figure 10 illustrates three examples of removed targets and ligands. Figure 7 analyzes the proportion of shared atoms between the AC pairs in the target ChEMBL218 EC_{50} using Maximum Common Substructure (MCS). The average proportion of shared atoms (86.78%) in the identified AC pairs confirms high structural similarity in common substructures.

Dataset split. We split the dataset using the Tanimoto similarity of the ECFP. To assess potential bias from ECFP-based data splitting, Figure 8 evaluates ML methods using four molecular descriptors on an alternative MACCS-based split. ECFP maintains superior performance, confirming its inherent descriptive power.

Protein flexibility. Using DSDPFlex (Dong et al., 2024), we investigate protein flexibility by allowing flexible side chains for 10 amino acids nearest to the crystal ligand. Figure 14 shows that performance metrics on 8 K_i targets distribute evenly around the $y = x$ line, suggesting comparable effectiveness between fixed and flexible docking approaches.

Train cross-target models with 3D data. Table 6, 7, 8, and 9 explore the cross-target applicability of 3D models on combined targets under two settings: out-of-distribution (OOD) excluding Protease or Nuclear receptor targets, and in-domain, using all targets with the same labels. Figure 13 shows that multi-target training performs comparably to single-target training, complementing the analysis in § 5.4.

²<https://scikit-learn.org/>

³<https://www.pyg.org/>

Combine the 3D information and ECFP features. To explore the integration of 3D structural information with handcrafted ECFP features, we utilize a 3D model as a feature extractor, combining its output with ECFP descriptors, followed by MLP for affinity prediction (architecture shown [Figure 15](#)). The evaluation across ten targets (shown in [Table 10](#)) highlights two key findings. First, models with 3D information consistently outperform or match those without 3D information across most targets, achieving notable improvements in overall RMSE and $RMSE_{\text{cliff}}$. Second, the integration of 3D features significantly enhances the model’s ability to handle activity cliffs, as evidenced by greater improvements in $RMSE_{\text{cliff}}$ (avg. imp. of 5.61%) compared to overall RMSE (avg. imp. of 3.48%).

Benchmarking the zero-shot ability of more 3D models. To explore the generalization ability of recent 3D binding affinity prediction models, we evaluate six SOTA methods (PIGNet ([Moon et al., 2022](#)), RTMScore ([Shen et al., 2022](#)), TANKBind ([Lu et al., 2022](#)), DSMBind ([Jin et al., 2023](#)), KarmaDock ([Zhang et al., 2023b](#)), and EquiScore ([Cao et al., 2024](#))) trained on PDBBind. [Figure 16](#) presents their Pearson correlation on the complete dataset and activity cliff cases across each target. All these methods perform worse on the AC samples, which is consistent with the result of our benchmark. Additionally, these methods show decreased performance compared to the PDBBind test set, with effectiveness correlating with the presence of homologous proteins in the PDBBind training data. For instance, targets with numerous homologous samples in PDBBind demonstrate superior results: ChEMBL2147 K_i achieves a Pearson correlation of 0.688 (DSMBind, PDB ID: 2j2i) with 103 homologous samples, while ChEMBL2971 K_i reaches 0.671 (DSMBind, PDB ID: 4jia) with 61 homologous samples in PDBBind. In contrast, targets lacking homologous proteins in PDBBind (ChEMBL219 K_i , ChEMBL228 K_i , and ChEMBL233 K_i) show very small correlation (DSMBind, Pearson=-0.021, -0.087, and 0.033 respectively).

Limitations and future work. While our dataset contains a variety of protein targets, the distribution of different types of targets is imbalanced, with several popular drug targets dominating. Increasing the diversity of target types would be beneficial for enhancing the generalization of successive models. Furthermore, the mapping between the target and the unique binding site may introduce bias, as some targets have unknown binding sites. We plan to conduct routine validity checks to update the dataset as more protein structures are deposited into the PDB, ensuring its relevance and accuracy over time. Lastly, the complex structures generated by molecular docking may be inaccurate, and more advanced approaches such as molecular simulation can be employed to refine the complex structures.

DockedAC provides the foundation for studying ACs from a structural perspective, and we anticipate that it will inspire the development of novel 3D QSAR algorithms. Future research could focus on designing advanced deep learning architectures capable of capturing and leveraging 3D structural information to improve AC prediction accuracy. Additionally, the dataset could be expanded to include more diverse targets and ligands, as well as refined complex structures, thereby increasing its utility for AI-driven drug discovery. By enabling a deeper understanding of structure-activity relationships and promoting the integration of 3D molecular data, we believe that our DockedAC will foster the development of innovative computational methods and contribute to the advancement of rational drug design and precision medicine.

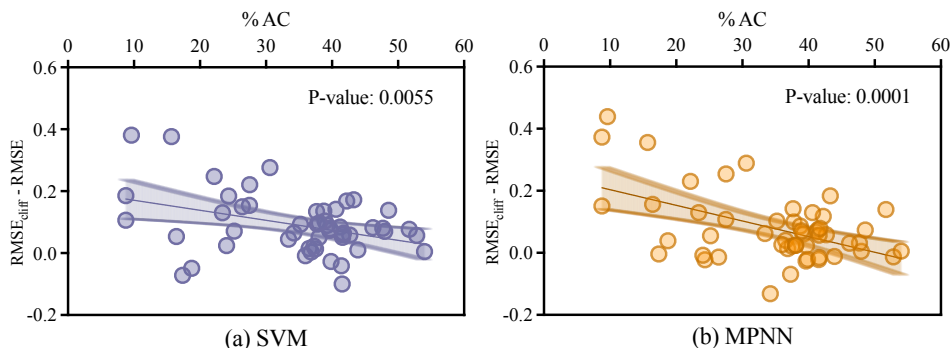
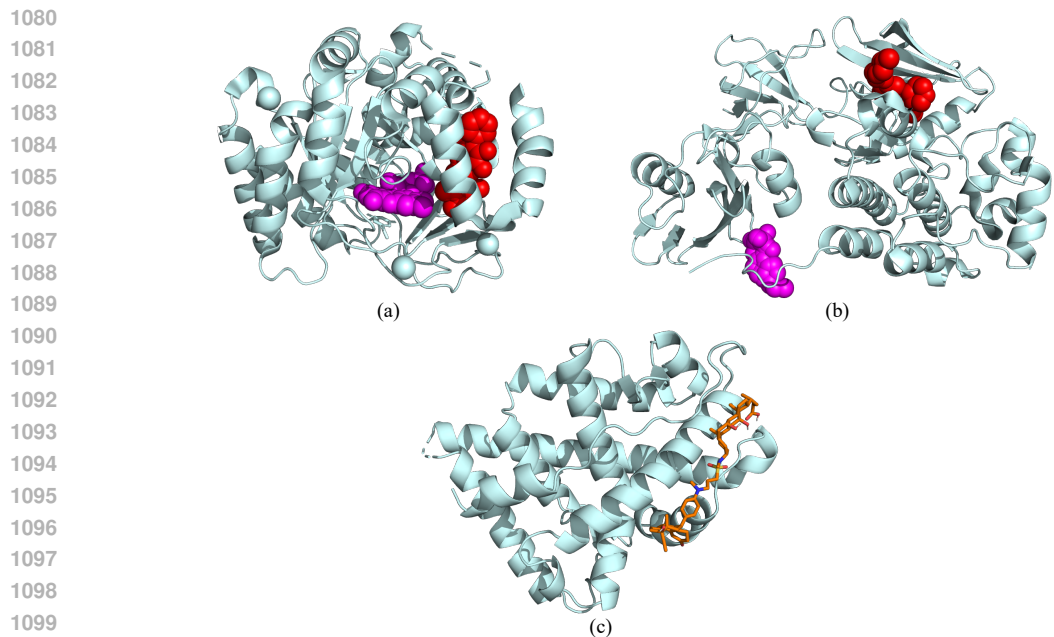
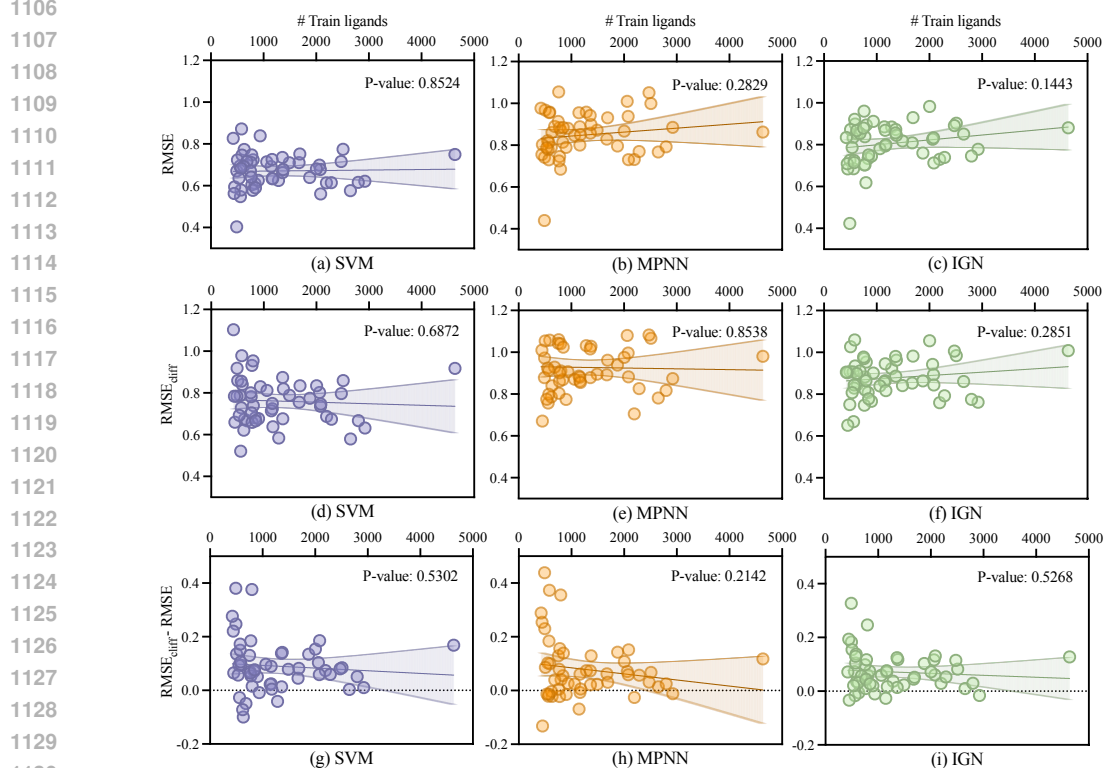


Figure 9: Relationship between the ratio of the AC and $RMSE - RMSE_{\text{cliff}}$ of SVM and MPNN.



1101 Figure 10: Three examples of the removed targets and ligands. (a) The target structure has two ligand
1102 binding sites (PDB: 5mvd). (b) Two structures of the same target have different binding sites (PDB:
1103 2h8h and 1o4j). The two structures are aligned. (c) The ligand docking score is larger than zero
1104 (Target: ChEMBL1871, PDB: 2ama, ligand: ChEMBL406027).



1131 Figure 11: Relationship between the number of training ligands and (a)-(c) RMSE, (d)-(f) $RMSE_{cliff}$
1132 and (g)-(i) their difference on SVM, MPNN, and IGN.
1133

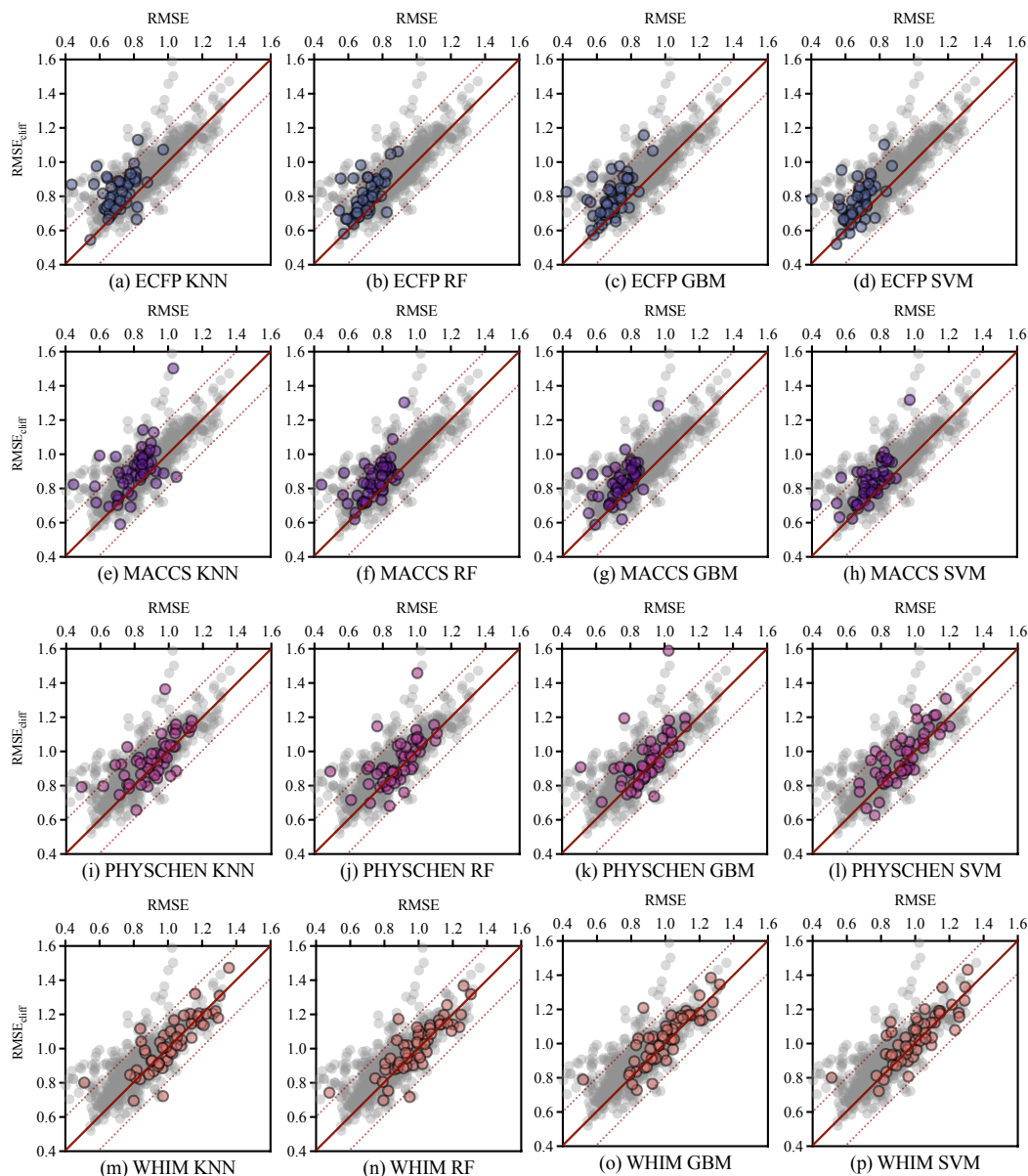


Figure 12: Performance comparison between RMSE and $RMSE_{cliff}$ for classic ML algorithms across 52 targets.

Table 6: The results of Protease family targets with K_i data in the setting of training with in-domain and out-of-distribution (OOD) targets.

| Model | ChEMBL204 K_i | | ChEMBL244 K_i | | ChEMBL3286 K_i | | ChEMBL4822 K_i | |
|---------|-----------------|----------------|-----------------|----------------|------------------|----------------|------------------|----------------|
| | RMSE | $RMSE_{cliff}$ | RMSE | $RMSE_{cliff}$ | RMSE | $RMSE_{cliff}$ | RMSE | $RMSE_{cliff}$ |
| IGN | 0.873 | 1.027 | 0.891 | 1.006 | 0.724 | 0.829 | 0.751 | 0.778 |
| IGN OOD | 1.612 | 1.788 | 1.647 | 1.643 | 1.183 | 1.149 | 1.153 | 1.197 |

Table 7: The results of Protease family targets with IC_{50} data in the setting of training with in-domain and out-of-distribution (OOD) targets.

| Model | ChEMBL280 IC_{50} | | ChEMBL284 IC_{50} | | ChEMBL3922 IC_{50} | | ChEMBL4588 IC_{50} | |
|---------|---------------------|-----------------------|---------------------|-----------------------|----------------------|-----------------------|----------------------|-----------------------|
| | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} |
| IGN | 0.930 | 0.982 | 0.983 | 1.055 | 0.685 | 0.651 | 0.834 | 0.906 |
| IGN OOD | 1.828 | 1.706 | 1.691 | 1.793 | 1.138 | 1.068 | 1.757 | 1.464 |

Table 8: The results of Nuclear receptor family targets with K_i data in the setting of training with in-domain and out-of-distribution (OOD) targets.

| Model | ChEMBL1871 K_i | | ChEMBL2034 K_i | |
|---------|------------------|-----------------------|------------------|-----------------------|
| | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} |
| IGN | 0.725 | 0.908 | 0.834 | 0.907 |
| IGN OOD | 1.005 | 1.337 | 1.567 | 1.505 |

Table 9: The results of Nuclear receptor family targets with EC_{50} data in the setting of training with in-domain and out-of-distribution (OOD) targets.

| Model | ChEMBL2047 EC_{50} | | ChEMBL239 EC_{50} | | ChEMBL3979 EC_{50} | | ChEMBL235 EC_{50} | |
|---------|----------------------|-----------------------|---------------------|-----------------------|----------------------|-----------------------|---------------------|-----------------------|
| | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} |
| IGN | 0.730 | 0.760 | 0.845 | 0.962 | 0.745 | 0.766 | 0.760 | 0.863 |
| IGN OOD | 1.022 | 0.826 | 1.368 | 1.351 | 1.478 | 1.299 | 1.405 | 1.229 |

Table 10: The performance of MLP and IGN using the handcrafted molecule descriptor ECFP.

| Model | ChEMBL205 K_i | | ChEMBL214 K_i | | ChEMBL233 K_i | | ChEMBL237 K_i | | ChEMBL264 K_i | |
|---------|-----------------|-----------------------|-----------------|-----------------------|-----------------|-----------------------|-----------------|-----------------------|-----------------|-----------------------|
| | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} |
| MLP | 0.795 | 0.929 | 0.683 | 0.770 | 0.846 | 0.917 | 0.720 | 0.767 | 0.669 | 0.730 |
| IGN | 0.781 | 0.904 | 0.683 | 0.792 | 0.814 | 0.878 | 0.728 | 0.764 | 0.637 | 0.691 |
| Imp (%) | 1.76 | 2.69 | 0.00 | - | 3.78 | 4.25 | - | 0.39 | 4.78 | 5.34 |

| Model | ChEMBL287 K_i | | ChEMBL1871 K_i | | ChEMBL2047 EC_{50} | | ChEMBL3979 EC_{50} | | ChEMBL4203 K_i | |
|---------|-----------------|-----------------------|------------------|-----------------------|----------------------|-----------------------|----------------------|-----------------------|------------------|-----------------------|
| | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} | RMSE | RMSE _{cliff} |
| MLP | 0.746 | 0.855 | 0.730 | 0.991 | 0.673 | 0.714 | 0.664 | 0.729 | 0.943 | 0.988 |
| IGN | 0.759 | 0.855 | 0.686 | 0.860 | 0.594 | 0.599 | 0.667 | 0.723 | 0.880 | 0.857 |
| Imp (%) | - | 0.00 | 6.03 | 13.22 | 11.74 | 16.11 | - | 0.82 | 6.68 | 13.26 |

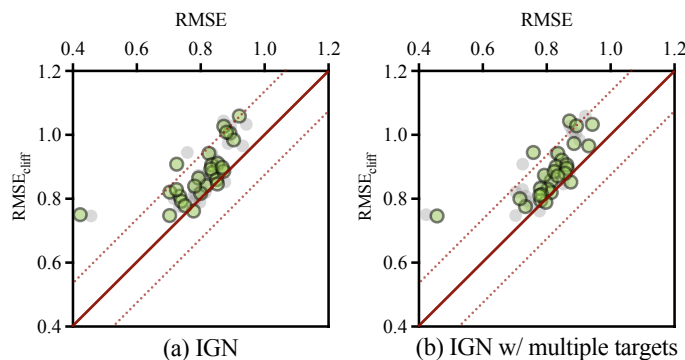


Figure 13: The results of IGN on all the targets of K_i labels when trained separately on (a) each target or (b) the data of multiple targets.

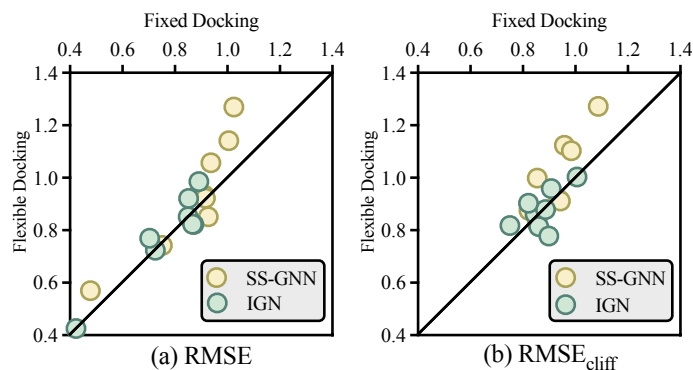
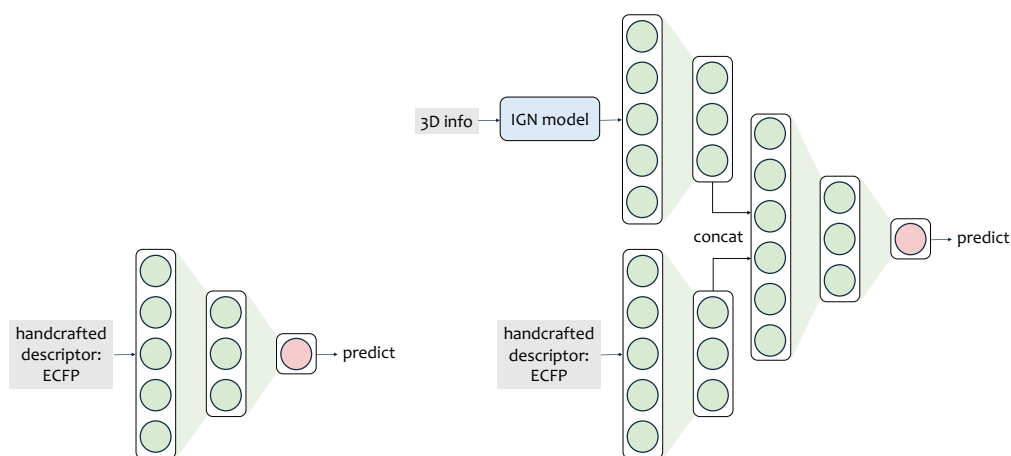


Figure 14: The (a) RMSE and (b) $RMSE_{cliff}$ metric on fixed docking v.s. flexible docking.



(a) The model illustration of MLP with ECFP descriptor (b) The model illustration of IGN combined with ECFP descriptor

Figure 15: The model illustration of MLP and IGN using the handcrafted molecule descriptor ECFP.

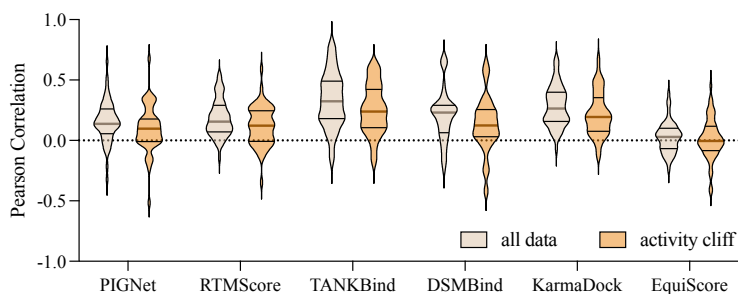


Figure 16: The Pearson and $Pearson_{cliff}$ evaluated on our DockedAC benchmark across 52 targets using PIGNet (Moon et al., 2022), RTMScore (Shen et al., 2022), TANKBind (Lu et al., 2022), DSMBind (Jin et al., 2023), KarmaDock (Zhang et al., 2023b), and EquiScore (Cao et al., 2024), all of which were trained on general binding affinity datasets.