
DockedAC: A Dataset with Comprehensive 3D Protein-ligand Complexes for Activity Cliff Analysis

Zijing Liu^{1*}, Xinni Zhang^{2*}, Yankai Chen³, Bin Feng¹, Mingjun Yang⁴, Zenglin Xu⁵, Yu Li¹, Philip S. Yu³, Irwin King²
¹ International Digital Economy Academy (IDEA) ² The Chinese University of Hong Kong
³ University of Illinois Chicago ⁴ XtalPi Inc. ⁵ Fudan University

Abstract

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

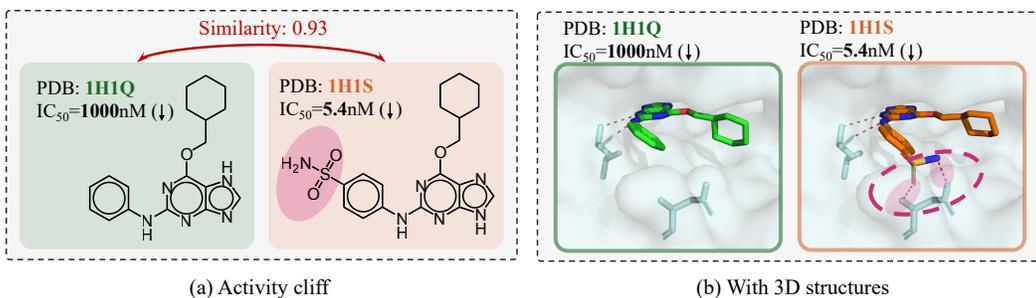
17 1 Introduction

18 Artificial intelligence (AI) is revolutionizing the drug discovery process as it is capable of large-scale
19 data analysis, pattern recognition, and making accurate predictions [61]. One important application
20 of AI models is to predict the biological activity of candidate compounds, thereby reducing labor-
21 intensive tasks. A foundational concept in many AI algorithms is the similarity principle, which states
22 that similar objects are likely to share similar features and predictions. However, in drug discovery, a
23 phenomenon known as activity cliffs challenges this concept and poses difficulties for AI models. An
24 **activity cliff** (AC) occurs when structurally similar compounds exhibit significant differences in their
25 biological activity against the same target [42], as illustrated in [Figure 1](#) (a).

26 AC plays a crucial role in drug discovery, as it complicates the optimization of drug candidates by
27 confounding the human experts in the understanding of traditional structure-activity relationships
28 (SARs) [64]. On the other hand, knowledge about ACs can be highly beneficial when designing or
29 optimizing compounds to enhance the bioactivity of a given target [11, 52]. For example, replacing a
30 single atom or adding a methyl group can result in more than 100-fold improvement in bioactivity [39,
31 47]. However, the mechanisms underlying ACs in individual drug development programs can be
32 different, making it challenging for humans to process such information and derive transferable
33 experiences. Consequently, various efforts have been made to computationally predict ACs [53].

34 Compared to quantitative structure-activity relationship (QSAR) modeling for other molecular
35 properties, AC prediction is particularly challenging due to the instability that ACs introduce to
36 the models [12]. Early attempts use machine learning methods such as random forest (RF) and
37 support vector machine (SVM) to predict the AC of a compound pair [20, 22]. To further improve

*Equal Contribution

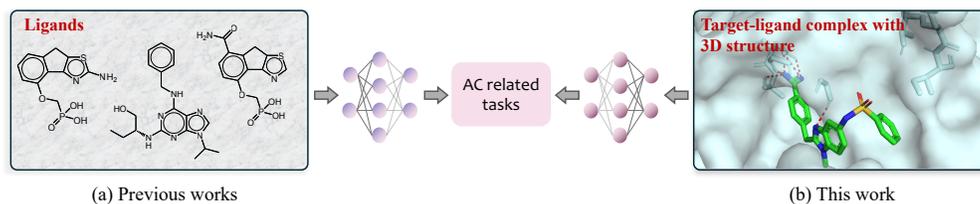


(a) Activity cliff (b) With 3D structures

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).

38 AC predictions, the matched molecular pair (MMP) kernel [55] and condensed graphs of reaction
 39 representations [24] have been integrated into various machine learning methods. More recently,
 40 deep neural networks-based algorithms, such as convolutional neural networks [29], graph neural
 41 networks [46] and transformers [7], have been applied to predict ACs.

42 In most previous works, the study of ACs has been ligand-centric and lacked 3D structure consid-
 43 eration, failing to account for interactions between the ligand and the protein target [28, 56]. Many
 44 mechanisms of ACs can be analyzed from the structural perspective, such as hydrogen bonding,
 45 ionic interactions, hydrophobic or aromatic group interactions [26] (e.g. Figure 1 (b)). It is therefore
 46 natural to incorporate the information of structures into the modeling of ACs. However, the available
 47 structural data for ACs is very limited, with only 215 pairs of AC ligands [28]. This data scarcity
 48 makes it challenging to train deep learning models effectively.



(a) Previous works (b) This work

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.

49 In this paper, we present DockedAC, a new dataset designed to tackle the challenge of ACs from
 50 a structural perspective, enabling large-scale AC modeling with modern AI algorithms. Unlike
 51 previous studies, our dataset includes not only the protein target information but also the target-
 52 ligand complex structures built using molecular docking (Figure 2). We collect the bioactivity data
 53 of more than 80,000 ligands across over 50 protein targets. These protein targets are mapped to
 54 their corresponding structures in the RCSB Protein Data Bank (PDB) [4], with the ligand binding
 55 sites identified for docking. In addition, we provide a benchmarking framework to evaluate the
 56 performance of traditional machine learning and deep learning methods on AC prediction and to
 57 analyze the impact of ACs on model performance. Our dataset enhances model interpretability,
 58 inspires the development of advanced algorithms for AC prediction, and fosters the advancement of
 59 more effective 3D feature extraction methods.

60 2 Related Work

61 **Previous works on AC prediction.** As a crucial phenomenon in drug discovery, ACs have garnered
 62 significant attention not only in medicinal chemistry but also in the computer science and intelligence
 63 community. Various methods of machine learning and deep learning have been applied to the
 64 prediction of ACs [20, 22, 29, 7, 46]. In addition, recent research has explored ACs from several
 65 different perspectives, such as QSAR modeling [13], the complexity of the learning methods [56],
 66 and benchmarking of different approaches [62]. However, due to the limited availability of data,
 67 almost all existing works focus on the ligand-centric view of ACs, where the ligand is modeled with

68 a 2D molecular graph or 1D SMILES sequence [68], without incorporating the 3D structure and the
69 protein target information. The 3D activity cliff (3DAC) database, used in a study on structure-based
70 AC prediction, contains only 219 3DAC pairs [26, 28]. This motivates us to construct a larger dataset
71 for structure-based ACs.

72 **Existing AC datasets.** Although there are several works on AC prediction, few good benchmarking
73 datasets exist. Several works rely on self-collected datasets and are not well documented, or have little
74 information provided about the protein targets [34, 13, 56]. Two recent works on AC datasets both
75 collect data from the ChEMBL database [44], either for the classification of a pair of AC ligands [74]
76 or the regression of the bioactivity value of individual AC ligands [62]. These datasets do not consider
77 modeling the 3D structure of the binding complex, rendering them less appropriate for accurate
78 AC prediction. In our work, we match the obtained bioactivity data to the corresponding protein
79 structures in PDB and generate target-ligand binding structures.

80 **3D protein-ligand binding affinity prediction.** In this work, we consider the regression problem
81 and train different models to predict the bioactivity in the presence of the AC. Given the target-ligand
82 complex structures, nearly all the models for binding affinity prediction use the PDBbind dataset,
83 including convolutional neural networks, graph neural networks, and attention-based models [72, 31,
84 33, 57]. A comprehensive review of the drug-target interaction prediction can be found in [71]. In
85 molecular property prediction, activity cliffs can significantly impact model predictions [15]. We
86 evaluate the performance of 3D target-ligand affinity prediction models with our dataset and compare
87 them with other machine learning or deep neural network models with ligand-only inputs.

88 3 The DockedAC Dataset

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

92 3.1 Data Collection

93 We first collect bioactivity data from ChEMBL v33 [44] using the ChEMBL web resource client [14]
94 for 64 protein targets. The data includes Inhibitory Constant (K_i), Half-Maximal Effective Con-
95 centration (EC_{50}), and Half-Maximal Inhibitory Concentration (IC_{50}), all measured in nanomolar
96 (nM). To eliminate significant sources of error, the obtained raw data is checked for validity and
97 reliability. In particular, when a ligand-target pair has multiple entries of the bioactivity data, the
98 ligand is removed if the standard deviation of the activities is larger than 10. The mean value of the
99 activities is used as the ligand-target activity label. A ligand is also removed if it fails the sanitization
100 and standardization by RDKit [3]. To ensure enough samples of a target for model training, the targets
101 with fewer than 500 ligands are dropped. Finally, the negative logarithm p is applied to the bioactivity
102 values as the regression target (denoted as pK_i ; pEC_{50} ; pIC_{50} in [log units]) [51]. After this process,
103 we have the ChEMBL id of the target and the corresponding ligands with bioactivity values (the first
104 step in Figure 3 (a)). The resulting dataset has 54 protein targets.

105 3.2 Activity Cliff Identification

106 An activity cliff is defined as a pair of structurally similar compounds exhibiting a large difference
107 in bioactivities against a given target. To identify similar ligand pairs, we use a consensus of
108 three similarity measures to define the activity cliff pairs following van Tilborg et al. [62]: (a)
109 substructure similarity, calculated using the Tanimoto coefficient on the extended connectivity
110 fingerprint (ECFP) [58, 48]; (b) scaffold similarity, determined by the Tanimoto coefficient on
111 the ECFP of generic Murcko scaffolds [2]; (c) SMILES similarity, computed as one minus the
112 scaled Levenshtein distance between the canonical SMILES representations [40]. If any of these
113 three similarity measures is equal to or greater than 0.9, the pair of ligands is further evaluated for
114 differences in bioactivity. Currently, there are no widely accepted quantitative definitions of ACs [54].
115 Following previous studies [34, 25], we define an activity cliff as a bioactivity difference exceeding
116 one order of magnitude ($10\times$), as illustrated in the second step of Figure 3 (a)).

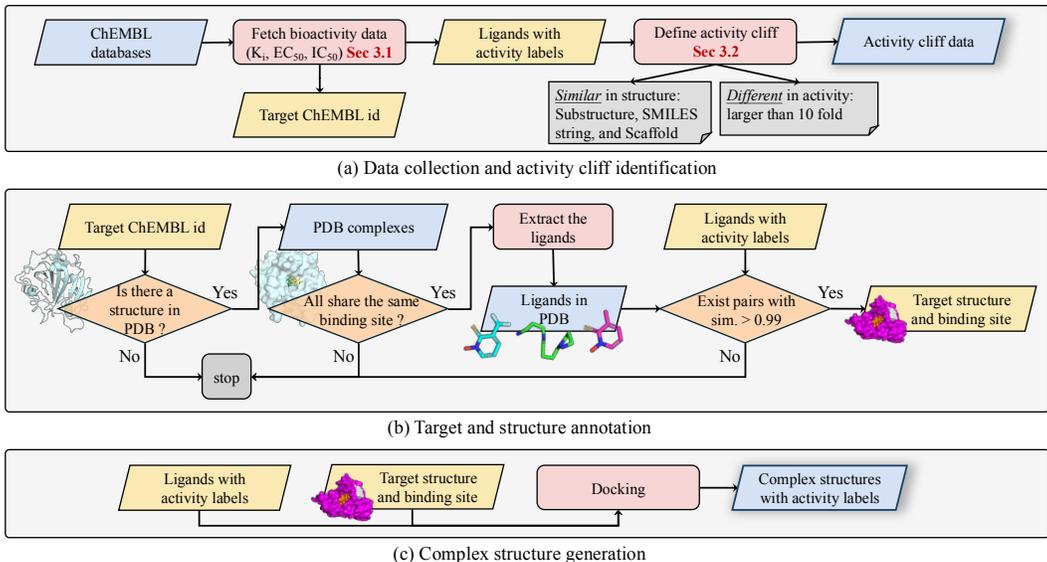


Figure 3: The whole process of building DockedAC with: (a) initial data collection from ChEMBL (Sec. 3.1) and activity cliff identification (Sec. 3.2), (b) mapping targets to 3D structures and identifying binding sites (Sec. 3.3), and (c) generation of target-ligand complex structures (Sec. 3.4).

117 3.3 Target and Structure Annotation

118 To generate the target-ligand complex, it is essential to identify the 3D structure of the target protein
 119 and its binding site. This mapping process is illustrated in Figure 3 (b). Given a target ChEMBL id,
 120 the first step is to map the target protein to its UniProt id [9] and find all the structures corresponding
 121 to the UniProt id in the PDB. We utilize the PDBbind database for the initial search [66]. If the
 122 target is not found in PDBbind, we then search for it in the entire PDB. The retrieved structures
 123 containing a small molecule ligand are chosen and aligned to verify whether the ligands bind to
 124 the same site. If the binding site is not unique, the target is discarded (see Figure 11 (a)(b)). After
 125 alignment, ligands sharing the same binding site are extracted and compared with the ligands that
 126 have activity labels from ChEMBL. If a pair of ligands—one from the PDB database and one from
 127 ChEMBL—has a fingerprint similarity (Tanimoto coefficient) greater than 0.99, the target structure
 128 and the binding site are used. Otherwise, the target is removed from the dataset. When multiple
 129 structures satisfy this condition, the structure with the highest resolution is selected. This procedure
 130 ensures the correspondence between the bioactivity values and the target binding site. As a result of
 131 this structure mapping process, two targets were removed, resulting in a final dataset of 52 protein
 132 targets.

133 3.4 Complex Structure Generation

134 Given a target protein and its corresponding binding site for a ligand, molecular docking is employed
 135 to generate the target-ligand complex, as illustrated in Figure 3 (c). The docking tool DSDP is
 136 used, which combines the AutoDock Vina’s pose sampling algorithm with GPU acceleration [27, 60].
 137 Since the binding site information of the target is already known, local docking is performed within
 138 a 25 Å wide box around the given binding region. A docking score (in kcal/mol) greater than zero
 139 indicates an inaccurate docking conformation (e.g. Figure 11 (c)), and the corresponding ligand is
 140 removed from the dataset. To further improve the quality of the docking complex, for ligands with
 141 high structural similarity to known references, a template-based docking approach is employed [70]. In
 142 this approach, when a query ligand shares significant structural features with a crystallographically
 143 resolved reference ligand, the known binding pose is used as a template to guide the placement
 144 of similar substructures. This template-based protocol constrains the conformational search space,
 145 leading to more accurate pose predictions. To validate the docking result, we have compared the
 146 docking poses of compounds with known crystal structures from [28]. The observed root-mean-
 147 square deviation (RMSD) values (median = 2.55 Å) fall well within accepted standards for reliable
 148 pose prediction.

149 To enhance the comprehensiveness and diversity of our protein-ligand complex dataset, we incorporate
 150 two additional docking methods. First, we employ KarmaDock [73], a state-of-the-art machine
 151 learning-based docking method that leverages deep learning algorithms to predict binding poses.
 152 Second, to account for protein flexibility in ligand binding, we utilize DSDPflex [16], which accounts
 153 for side-chain movements in the binding site region. For DSDPflex, we allow the 10 residues closest
 154 to the reference ligand to have flexible side chains and record the top 5 scoring poses for each
 155 ligand-target docking simulation. This complementary approach provides insights into the dynamic
 156 nature of protein-ligand interactions. While we provide the complexes generated by all three methods
 157 in our dataset, our subsequent analyses primarily focused on the results obtained from DSDP, as it
 158 offered the most reliable and consistent predictions for our system.

159 3.5 Dataset Splitting

160 The preparation of datasets for benchmarking machine learning models requires careful data-splitting
 161 strategies. For ligand-based methods, separate models are developed for individual protein targets.
 162 To assess the ligand-based method, the ligands of each target are split into a training and test set using
 163 a double-stratified sampling strategy [62]. In particular, the ligands of each target are first clustered
 164 into 5 groups based on their substructural similarity (Tanimoto similarity of the ECFP). A two-stage
 165 stratified splitting (80%/20%) is then performed on the cluster label and the AC label. This procedure
 166 ensures that the training and test set have similar ligand distributions.

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

172 3.6 Dataset Description

173 The final dataset contains 82,836
 174 target-ligand activity values and
 175 their corresponding generated com-
 176 plex structures. A brief overview
 177 of the dataset is provided in **Table 1**, while detailed information
 178 on each target can be found in
 179 Appendix **Table 3**. The dataset
 180 includes popular target families
 181 in drug discovery (such as G-
 182 protein-coupled receptors (GPCR),
 183 kinases, proteases, and nuclear re-
 184 ceptors) as well as targets with criti-
 185 cal roles in biology (like chaperone
 186 and kinesin). In terms of size, the
 187 target Carbonic anhydrase II has
 188 the most ligands with bioactivity values (5794 unique molecules), while the target with the least
 189 ligands (533 unique molecules) is Matrix metalloproteinase 8. As an intensively studied drug target,
 190 the GPCR family has the most ligands on average. For all the targets, around 37% of the ligands are
 191 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

193 4 Benchmark

194 In addition to the DockedAC dataset, we also provide a framework to benchmark the performance of
 195 various machine learning and deep learning methods on AC prediction. This section briefly introduces
 196 our benchmark setup. A detailed presentation and analysis of the experimental results are provided in
 197 Section 5.

198 4.1 Model Descriptions

199 In general, three types of learning models are considered:

- 200 • Four classic machine learning algorithms for structure-activity relationship prediction using hand-
201 crafted molecular descriptors: K-nearest neighbor (KNN) [10], random forest (RF) [5], gradient
202 boosting machine (GBM) [18], and support vector regression (SVM) [21].
- 203 • Deep learning models that only leverage the 1D or 2D ligand information, including (1) three
204 1D sequential models: transformer [63], long short-term memory (LSTM) networks [23], and
205 1D CNN [36], and (2) four 2D structural graph neural network (GNN) models: message passing
206 neural network (MPNN) [19], graph convolutional network (GCN) [37], graph attention network
207 (GAT) [63], and attentive fingerprint (AFP) [69].
- 208 • Two 3D structural GNN models: IGN [31] and SS-GNN [72] are included to study the effect of 3D
209 structures, as our dataset contains 3D structural information.

210 4.2 Feature Descriptions

211 For machine learning algorithms, following previous work [62], we consider four types of molecule
212 descriptors from several levels of complexity as follows. (1) Extended Connectivity Fingerprints
213 (ECFPs) [48]: circular topological fingerprints used for molecular characterization, capturing struc-
214 tural features of molecules. (2) Molecular ACCess System (MACCS) keys [17]: a set of structural
215 keys utilized for substructure searching and similarity analysis, encoding specific chemical sub-
216 structures or patterns. (3) Physicochemical (PhysChem) descriptors [65]: a set of 11 properties
217 indicative of drug-likeness, providing insights into the physical and chemical properties of molecules.
218 (4) Weighted Holistic Invariant Molecular (WHIM) descriptors [59]: capturing three-dimensional
219 geometrical and electronic properties of molecules, invariant to rotation and translation.

220 Deep learning methods eliminate the need for handcrafted descriptors, allowing direct learning from
221 “unstructured” data representations. For sequential methods, the Simplified Molecular Input Line
222 Entry System (SMILES) [68] string is used, which is popular for its ability to describe the structure
223 of chemical species in a text format that sequential methods can naturally process. For 2D GNN
224 models, we adopt molecular graphs, which represent the structural formula where nodes represent
225 atoms and edges represent bonds. For 3D GNN models, we employ the target-ligand complexes we
226 have processed that incorporate detailed 3D structure information. Additional descriptions of the
227 features and their corresponding models are available in Appendix A.4 and Table 5.

228 4.3 Metrics and Implementations

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

236 5 Experimental Results and Analyses

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

242 5.1 Performance Comparison for GNN Models

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

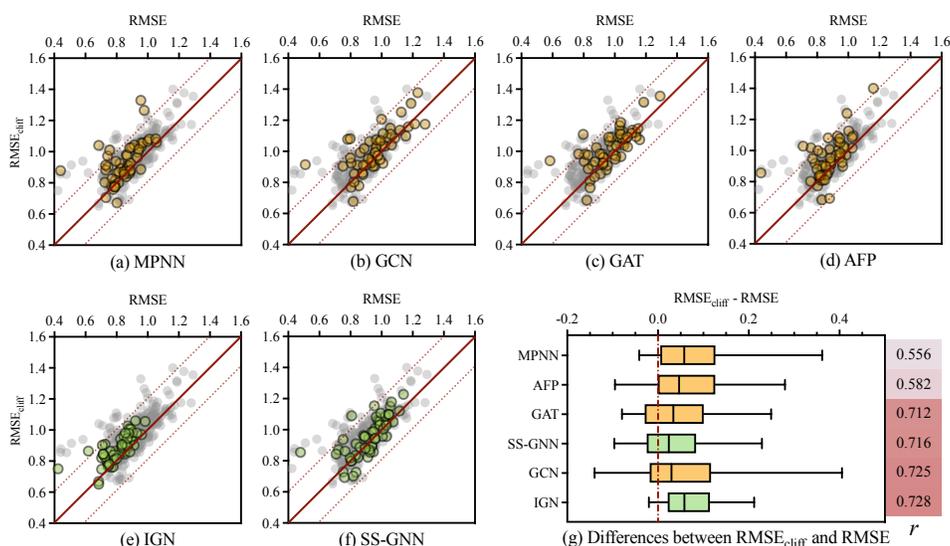


Figure 4: Performance comparison for GNN models. (a)-(f) Comparison between RMSE and RMSE_{cliff} of GNN models across 52 targets. The 2D GNN models are colored in yellow, while the 3D GNN models are colored in green. Gray nodes depict all nodes in these six subgraphs for a clear comparison. Red solid 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}.

246 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.

263 5.2 The AC Prediction is Target-dependent

264 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 RMSE_{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. 268 Protease has the worst RMSE_{cliff} for all the methods while kinase is the target family with the best RMSE_{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 RMSE_{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	IGNN	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

271 5.3 The Percentage of AC Matters

272 In general, machine learning models tend to perform better with more training data. Here, we
273 investigate the factors influencing AC (Activity Cliff) prediction performance. Surprisingly, our
274 analysis reveals that the number of training samples does not exhibit a significant correlation
275 with RMSE, $\text{RMSE}_{\text{cliff}}$, or their difference, i.e., $\text{RMSE}_{\text{cliff}} - \text{RMSE}$ (see Appendix Figure 12).
276 This suggests that simply increasing the size of the training data is not sufficient to improve AC
277 prediction accuracy. However, as shown in Figure 5 (see more results for other models in Fig-
278 ure 10), the ratio of AC ligands in the training set is a significant factor affecting $\text{RMSE}_{\text{cliff}} -$
279 RMSE , with a p-value of $1.0\text{e-}4$. A higher per-
280 centage of AC ligands in the training set means more information directly relevant to AC, thereby
281 improving the AC predictive power. Our finding indicates that the knowledge about general bioactiv-
282 ity prediction is different from the knowledge benefiting AC prediction. This underscores the need for
283 new datasets and methodologies tailored specifically for AC prediction, as relying solely on general
284 bioactivity data is insufficient to achieve optimal performance in this domain.

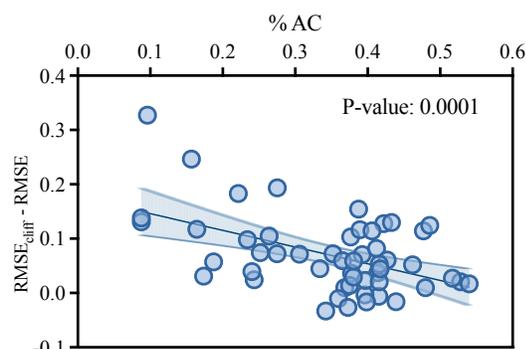


Figure 5: Relationship between the ratio of the AC and $\text{RMSE}_{\text{cliff}} - \text{RMSE}$ of IGN.

285
286
287
288
289
290
291

292 5.4 Performance Comparison with Machine Learning Algorithms

293 We benchmark the ability of all methods to predict bioactivity in the presence of the AC (measured
294 by $\text{RMSE}_{\text{cliff}}$), as shown in Figure 6 (detailed results in Appendix Figure 13). We have the following
295 empirical observations: (1) Significant performance differences can be observed among targets
296 in the handling of AC compounds, with $\text{RMSE}_{\text{cliff}}$ values spanning from 0.52 to 1.59 log units,
297 which is consistent with previous works [62, 50]. (2) Among the four machine learning algorithms,
298 performance disparities primarily stem from the molecule descriptors rather than the learning meth-
299 ods. ECFPs, which are designed specifically for structure-activity modeling by encoding detailed
300 information about each atom’s local environment, yield the lowest prediction error of all methods.
301 Their strong discriminative capability effectively differentiates molecules, even with minor structural
302 differences. (3) For deep learning methods, IGN coupled with 3D structure information achieves
303 the best performance on ACs. This approach benefits from the interaction information between the
304 ligand and the protein target captured within the 3D structure.

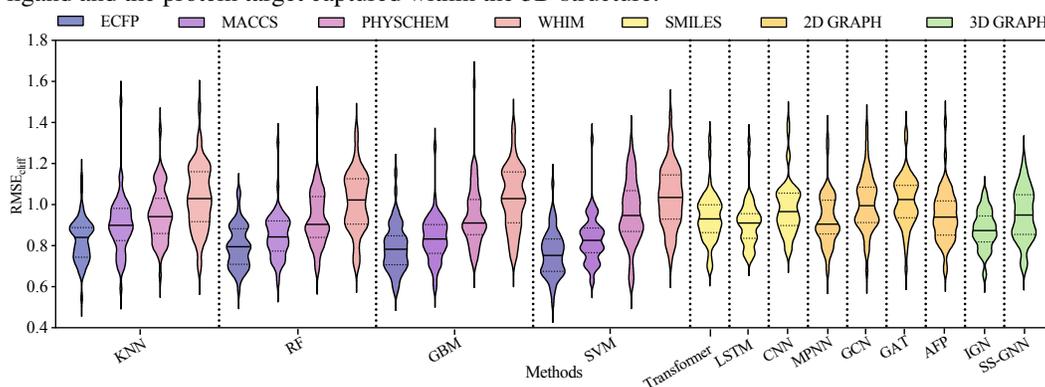


Figure 6: The $\text{RMSE}_{\text{cliff}}$ evaluated using different methods and features across 52 targets.

305 Surprisingly, some machine learning models outperform deep learning approaches, which can be
306 primarily attributed to their use of handcrafted features, especially ECFP. To validate this observation,
307 we implement a hybrid approach combining the ECFP features with the features extracted from
308 the last layer of the 3D IGN model. These concatenated features are then fed into an MLP for
309 prediction (as illustrated in Appendix Figure 16). The promising results across ten targets (see
310 Appendix Table 7) demonstrate the effectiveness of ECFP in structure-activity relationship learning.
311 On the other hand, this experiment underscores the value of integrating traditional cheminformatics

312 techniques with advanced deep-learning methods in molecular property prediction tasks. Future
313 research could explore optimizing this hybrid approach and investigating its applicability to a broader
314 range of molecular targets and properties.

315 5.5 Performance Positioning of 3D GNN Methods

316 Our experimental results show that machine learning
317 methods significantly outperform deep learning meth-
318 ods, especially with the ECFP featurization. This finding
319 aligns with previous studies on molecular property pre-
320 diction [32, 30]. To provide a global assessment of the
321 methods and demonstrate the effect of the target on 3D
322 GNN methods, we take the $RSME_{cliff}$ values of the 52
323 targets as features and compute the Pearson correlation
324 between the methods. The correlation serves as a similar-
325 ity measure in multidimensional scaling (MDS) [43] to
326 visualize the methods in a 2D plane (Figure 7). We then
327 identify the direction that determines the performance
328 of the methods. Although the average performance of
329 SS-GNN and IGN does not surpass machine learning
330 methods, there are specific targets where IGN and SS-
331 GNN outperform SVM. In contrast, GAT and GCN con-
332 sistentlly have larger $RSME_{cliff}$ values than SVM across
333 all targets. Handcraft features, such as ECFP, have been
334 optimized for QSAR over decades. Our
335 analysis indicates that the models with 3D structures can offer insights that handcrafted features do
336 not capture. Therefore, in practice, models based on 3D structures can be important complements to
the machine learning methods.

337 The incorporation of 3D structural information also enables cross-target modeling capabilities. To
338 investigate this potential, we conduct additional experiments by combining K_i targets and training the
339 IGN model under two scenarios: (i) in-domain setting under all K_i targets, and (ii) out-of-distribution
340 (OOD) setting excluding Protease-type targets. Analysis of four Protease targets (Appendix Table 6)
341 reveals that the absence of Protease targets in training leads to performance degradation, with average
342 $RMSE_{cliff}$ increasing from 0.9 to 1.4. While multi-target training achieves comparable performance
343 to target-specific training across all targets (Appendix Figure 14), these results indicate that there is
344 still a long way to go to fully exploit the multi-target 3D data and make the model generalize to new
345 targets.

346 6 Conclusion

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

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

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

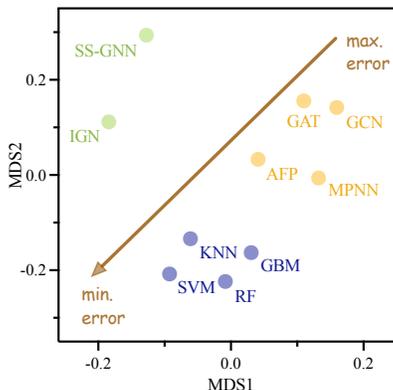


Figure 7: MDS visualization of the model performances by $RMSE_{cliff}$. The machine learning algorithms are with ECFP featurization.

365 References

- 366 [1] Jinheon Baek, Minki Kang, and Sung Ju Hwang. Accurate learning of graph representations with graph
367 multiset pooling. In *International Conference on Learning Representations*, 2021. URL [https://openre](https://openreview.net/forum?id=JHcqXGaqiGn)
368 [view.net/forum?id=JHcqXGaqiGn](https://openreview.net/forum?id=JHcqXGaqiGn).
- 369 [2] Guy W Bemis and Mark A Murcko. The properties of known drugs. 1. Molecular frameworks. *Journal of*
370 *Medicinal Chemistry*, 39(15):2887–2893, 1996.
- 371 [3] A Patrícia Bento, Anne Hersey, Eloy Félix, Greg Landrum, Anna Gaulton, Francis Atkinson, Louisa J Bellis,
372 Marleen De Veij, and Andrew R Leach. An open source chemical structure curation pipeline using RDKit.
373 *Journal of Cheminformatics*, 12:1–16, 2020.
- 374 [4] Helen M Berman, John Westbrook, Zukang Feng, Gary Gilliland, Talapady N Bhat, Helge Weissig, Ilya N
375 Shindyalov, and Philip E Bourne. The protein data bank. *Nucleic Acids Research*, 28(1):235–242, 2000.
- 376 [5] Leo Breiman. Bagging predictors. *Machine Learning*, 24:123–140, 1996.
- 377 [6] Duanhua Cao, Geng Chen, Jiabin Jiang, Jie Yu, Runze Zhang, Mingan Chen, Wei Zhang, Lifan Chen,
378 Feisheng Zhong, Yingying Zhang, et al. Generic protein–ligand interaction scoring by integrating physical
379 prior knowledge and data augmentation modelling. *Nature Machine Intelligence*, pp. 1–13, 2024.
- 380 [7] Hengwei Chen, Martin Vogt, and Jürgen Bajorath. DeepAC–conditional transformer-based chemical
381 language model for the prediction of activity cliffs formed by bioactive compounds. *Digital Discovery*, 1(6):
382 898–909, 2022.
- 383 [8] Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. ChemBERTa: large-scale self-supervised
384 pretraining for molecular property prediction. *arXiv preprint arXiv:2010.09885*, 2020.
- 385 [9] The UniProt Consortium. UniProt: the universal protein knowledgebase in 2023. *Nucleic Acids Research*,
386 51(D1):D523–D531, 2023.
- 387 [10] Thomas Cover and Peter Hart. Nearest neighbor pattern classification. *IEEE Transactions on Information*
388 *Theory*, 13(1):21–27, 1967.
- 389 [11] Maykel Cruz-Monteagudo, José L Medina-Franco, Yunierkis Pérez-Castillo, Orazio Nicolotti,
390 M Natália DS Cordeiro, and Fernanda Borges. Activity cliffs in drug discovery: Dr Jekyll or Mr Hyde?
391 *Drug Discovery Today*, 19(8):1069–1080, 2014.
- 392 [12] Maykel Cruz-Monteagudo, José L Medina-Franco, Yunier Perera-Sardiña, Fernanda Borges, Eduardo
393 Tejera, Cesar Paz-y Mino, Yunierkis Pérez-Castillo, Aminael Sánchez-Rodríguez, Zuleidys Contreras-
394 Posada, Natália DS Cordeiro, et al. Probing the hypothesis of sar continuity restoration by the removal of
395 activity cliffs generators in qsar. *Current Pharmaceutical Design*, 22(33):5043–5056, 2016.
- 396 [13] Markus Dablander, Thierry Hanser, Renaud Lambiotte, and Garrett M Morris. Exploring QSAR models
397 for activity-cliff prediction. *Journal of Cheminformatics*, 15(1):47, 2023.
- 398 [14] Mark Davies, Michał Nowotka, George Papadatos, Nathan Dedman, Anna Gaulton, Francis Atkinson,
399 Louisa Bellis, and John P Overington. ChEMBL web services: streamlining access to drug discovery data
400 and utilities. *Nucleic Acids Research*, 43(W1):W612–W620, 2015.
- 401 [15] Jianyuan Deng, Zhibo Yang, Hehe Wang, Iwao Ojima, Dimitris Samaras, and Fusheng Wang. A systematic
402 study of key elements underlying molecular property prediction. *Nature Communications*, 14(1):6395,
403 2023.
- 404 [16] Chengwei Dong, Yu-Peng Huang, Xiaohan Lin, Hong Zhang, and Yi Qin Gao. Dsdpflex: Flexible-receptor
405 docking with gpu acceleration. *Journal of Chemical Information and Modeling*, 2024.
- 406 [17] Joseph L Durant, Burton A Leland, Douglas R Henry, and James G Nourse. Reoptimization of MDL keys
407 for use in drug discovery. *Journal of Chemical Information and Computer Sciences*, 42(6):1273–1280,
408 2002.
- 409 [18] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of Statistics*,
410 pp. 1189–1232, 2001.
- 411 [19] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message
412 passing for quantum chemistry. In *International Conference on Machine Learning*, pp. 1263–1272. PMLR,
413 2017.

- 414 [20] Rajarshi Guha. Exploring uncharted territories: Predicting activity cliffs in structure–activity landscapes.
415 *Journal of Chemical Information and Modeling*, 52(8):2181–2191, 2012.
- 416 [21] Marti A. Hearst, Susan T Dumais, Edgar Osuna, John Platt, and Bernhard Scholkopf. Support vector
417 machines. *IEEE Intelligent Systems and their Applications*, 13(4):18–28, 1998.
- 418 [22] Kathrin Heikamp, Xiaoying Hu, Aixia Yan, and Jürgen Bajorath. Prediction of activity cliffs using support
419 vector machines. *Journal of Chemical Information and Modeling*, 52(9):2354–2365, 2012.
- 420 [23] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780,
421 1997.
- 422 [24] Dragos Horvath, Gilles Marcou, Alexandre Varnek, Shilva Kayastha, Antonio de la Vega de León, and
423 Jürgen Bajorath. Prediction of activity cliffs using condensed graphs of reaction representations, descriptor
424 recombination, support vector machine classification, and support vector regression. *Journal of Chemical
425 Information and Modeling*, 56(9):1631–1640, 2016.
- 426 [25] Ye Hu and Jürgen Bajorath. Exploration of 3d activity cliffs on the basis of compound binding modes and
427 comparison of 2d and 3d cliffs. *Journal of Chemical Information and Modeling*, 52(3):670–677, 2012.
- 428 [26] Ye Hu, Norbert Furtmann, Michael Gütschow, and Jürgen Bajorath. Systematic identification and classifica-
429 tion of three-dimensional activity cliffs. *Journal of Chemical Information and Modeling*, 52(6):1490–1498,
430 2012.
- 431 [27] YuPeng Huang, Hong Zhang, Siyuan Jiang, Dajiong Yue, Xiaohan Lin, Jun Zhang, and Yi Qin Gao. DSDP:
432 A blind docking strategy accelerated by GPUs. *Journal of Chemical Information and Modeling*, 63(14):
433 4355–4363, 2023.
- 434 [28] Jarmila Husby, Giovanni Bottegoni, Irina Kufareva, Ruben Abagyan, and Andrea Cavalli. Structure-based
435 predictions of activity cliffs. *Journal of Chemical Information and Modeling*, 55(5):1062–1076, 2015.
- 436 [29] Javed Iqbal, Martin Vogt, and Jürgen Bajorath. Prediction of activity cliffs on the basis of images using
437 convolutional neural networks. *Journal of Computer-Aided Molecular Design*, pp. 1–8, 2021.
- 438 [30] Tiago Janela and Jürgen Bajorath. Simple nearest-neighbour analysis meets the accuracy of compound
439 potency predictions using complex machine learning models. *Nature Machine Intelligence*, 4(12):1246–
440 1255, 2022.
- 441 [31] Dejun Jiang, Chang-Yu Hsieh, Zhenxing Wu, Yu Kang, Jike Wang, Ercheng Wang, Ben Liao, Chao Shen,
442 Lei Xu, Jian Wu, et al. InteractionGraphNet: A novel and efficient deep graph representation learning
443 framework for accurate protein–ligand interaction predictions. *Journal of Medicinal Chemistry*, 64(24):
444 18209–18232, 2021.
- 445 [32] Dejun Jiang, Zhenxing Wu, Chang-Yu Hsieh, Guangyong Chen, Ben Liao, Zhe Wang, Chao Shen,
446 Dongsheng Cao, Jian Wu, and Tingjun Hou. Could graph neural networks learn better molecular represen-
447 tation for drug discovery? A comparison study of descriptor-based and graph-based models. *Journal of
448 Cheminformatics*, 13:1–23, 2021.
- 449 [33] José Jiménez, Miha Skalic, Gerard Martínez-Rosell, and Gianni De Fabritiis. K deep: protein–ligand
450 absolute binding affinity prediction via 3d-convolutional neural networks. *Journal of Chemical Information
451 and Modeling*, 58(2):287–296, 2018.
- 452 [34] José Jiménez-Luna, Miha Skalic, and Nils Weskamp. Benchmarking molecular feature attribution methods
453 with activity cliffs. *Journal of Chemical Information and Modeling*, 62(2):274–283, 2022.
- 454 [35] Wengong Jin, Caroline Uhler, and Nir Hacohen. Se (3) denoising score matching for unsupervised binding
455 energy prediction and nanobody design. In *NeurIPS 2023 Generative AI and Biology (GenBio) Workshop*,
456 2023.
- 457 [36] Talia B Kimber, Maxime Gagnebin, and Andrea Volkamer. Maxsmi: maximizing molecular property
458 prediction performance with confidence estimation using smiles augmentation and deep learning. *Artificial
459 Intelligence in the Life Sciences*, 1:100014, 2021.
- 460 [37] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
461 *arXiv preprint arXiv:1609.02907*, 2016.
- 462 [38] Greg Landrum et al. RDKit: A software suite for cheminformatics, computational chemistry, and predictive
463 modeling. *Greg Landrum*, 8(31.10):5281, 2013.

- 464 [39] Cheryl S Leung, Siegfried SF Leung, Julian Tirado-Rives, and William L Jorgensen. Methyl effects on
465 protein–ligand binding. *Journal of Medicinal Chemistry*, 55(9):4489–4500, 2012.
- 466 [40] Vladimir I Levenshtein et al. Binary codes capable of correcting deletions, insertions, and reversals. In
467 *Soviet physics doklady*, volume 10, pp. 707–710. Soviet Union, 1966.
- 468 [41] Wei Lu, Qifeng Wu, Jixian Zhang, Jiahua Rao, Chengtao Li, and Shuangjia Zheng. Tankbind:
469 Trigonometry-aware neural networks for drug-protein binding structure prediction. *Advances in neural*
470 *information processing systems*, 35:7236–7249, 2022.
- 471 [42] Gerald M Maggiora. On outliers and activity cliffs why QSAR often disappoints, 2006.
- 472 [43] Al Mead. Review of the development of multidimensional scaling methods. *Journal of the Royal Statistical*
473 *Society: Series D (The Statistician)*, 41(1):27–39, 1992.
- 474 [44] David Mendez, Anna Gaulton, A Patrícia Bento, Jon Chambers, Marleen De Veij, Eloy Félix, María Paula
475 Magariños, Juan F Mosquera, Prudence Mutowo, Michał Nowotka, et al. ChEMBL: towards direct
476 deposition of bioassay data. *Nucleic Acids Research*, 47(D1):D930–D940, 2019.
- 477 [45] Seokhyun Moon, Wonho Zhung, Soojung Yang, Jaechang Lim, and Woo Youn Kim. Pignet: a physics-
478 informed deep learning model toward generalized drug–target interaction predictions. *Chemical Science*, 13
479 (13):3661–3673, 2022.
- 480 [46] Junhui Park, Gaeun Sung, SeungHyun Lee, SeungHo Kang, and ChunKyun Park. ACGCN: graph
481 convolutional networks for activity cliff prediction between matched molecular pairs. *Journal of Chemical*
482 *Information and Modeling*, 62(10):2341–2351, 2022.
- 483 [47] Lewis D Pennington and Demetri T Moustakas. The necessary nitrogen atom: a versatile high-impact
484 design element for multiparameter optimization. *Journal of Medicinal Chemistry*, 60(9):3552–3579, 2017.
- 485 [48] David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of Chemical Information*
486 *and Modeling*, 50(5):742–754, 2010.
- 487 [49] Chao Shen, Xujun Zhang, Yafeng Deng, Junbo Gao, Dong Wang, Lei Xu, Peichen Pan, Tingjun Hou, and
488 Yu Kang. Boosting protein–ligand binding pose prediction and virtual screening based on residue–atom
489 distance likelihood potential and graph transformer. *Journal of Medicinal Chemistry*, 65(15):10691–10706,
490 2022.
- 491 [50] Robert P Sheridan. Three useful dimensions for domain applicability in QSAR models using random
492 forest. *Journal of Chemical Information and Modeling*, 52(3):814–823, 2012.
- 493 [51] MJ Stewart and ID Watson. Standard units for expressing drug concentrations in biological fluids. *British*
494 *Journal of Clinical Pharmacology*, 16(1):3, 1983.
- 495 [52] Dagmar Stumpfe, Ye Hu, Dilyana Dimova, and Jürgen Bajorath. Recent progress in understanding activity
496 cliffs and their utility in medicinal chemistry: miniperspective. *Journal of medicinal chemistry*, 57(1):18–28,
497 2014.
- 498 [53] Dagmar Stumpfe, Huabin Hu, and Jürgen Bajorath. Evolving concept of activity cliffs. *Acs Omega*, 4(11):
499 14360–14368, 2019.
- 500 [54] Dagmar Stumpfe, Huabin Hu, and Jürgen Bajorath. Advances in exploring activity cliffs. *Journal of*
501 *Computer-Aided Molecular Design*, 34:929–942, 2020.
- 502 [55] Shunsuke Tamura, Swarit Jasial, Tomoyuki Miyao, and Kimito Funatsu. Interpretation of ligand-based
503 activity cliff prediction models using the matched molecular pair kernel. *Molecules*, 26(16):4916, 2021.
- 504 [56] Shunsuke Tamura, Tomoyuki Miyao, and Jürgen Bajorath. Large-scale prediction of activity cliffs using
505 machine and deep learning methods of increasing complexity. *Journal of Cheminformatics*, 15(1):4, 2023.
- 506 [57] Huishuang Tan, Zhixin Wang, and Guang Hu. GAABind: a geometry-aware attention-based network for
507 accurate protein–ligand binding pose and binding affinity prediction. *Briefings in Bioinformatics*, 25(1):
508 bbad462, 2024.
- 509 [58] Taffee T Tanimoto. Elementary mathematical theory of classification and prediction. 1958.
- 510 [59] Roberto Todeschini, Paola Gramatica, et al. New 3D molecular descriptors: the WHIM theory and QSAR
511 applications. *Perspectives in Drug Discovery and Design*, 9(0):355–380, 1998.

- 512 [60] Oleg Trott and Arthur J Olson. AutoDock Vina: improving the speed and accuracy of docking with a new
513 scoring function, efficient optimization, and multithreading. *Journal of Computational Chemistry*, 31(2):
514 455–461, 2010.
- 515 [61] Jessica Vamathevan, Dominic Clark, Paul Czodrowski, Ian Dunham, Edgardo Ferran, George Lee, Bin
516 Li, Anant Madabhushi, Parantu Shah, Michaela Spitzer, et al. Applications of machine learning in drug
517 discovery and development. *Nature Reviews Drug Discovery*, 18(6):463–477, 2019.
- 518 [62] Derek van Tilborg, Alisa Alenicheva, and Francesca Grisoni. Exposing the limitations of molecular
519 machine learning with activity cliffs. *Journal of Chemical Information and Modeling*, 62(23):5938–5951,
520 2022.
- 521 [63] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
522 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*,
523 30, 2017.
- 524 [64] Martin Vogt, Yun Huang, and Jürgen Bajorath. From activity cliffs to activity ridges: informative data
525 structures for sar analysis. *Journal of Chemical Information and Modeling*, 51(8):1848–1856, 2011.
- 526 [65] W Patrick Walters and Mark A Murcko. Prediction of ‘drug-likeness’. *Advanced Drug Delivery Reviews*,
527 54(3):255–271, 2002.
- 528 [66] Renxiao Wang, Xueliang Fang, Yipin Lu, and Shaomeng Wang. The PDBbind database: Collection
529 of binding affinities for protein-ligand complexes with known three-dimensional structures. *Journal of*
530 *Medicinal Chemistry*, 47(12):2977–2980, 2004.
- 531 [67] Renxiao Wang, Xueliang Fang, Yipin Lu, Chao-Yie Yang, and Shaomeng Wang. The pddbnd database:
532 methodologies and updates. *Journal of medicinal chemistry*, 48(12):4111–4119, 2005.
- 533 [68] David Weininger. SMILES, a chemical language and information system. 1. Introduction to methodology
534 and encoding rules. *Journal of Chemical Information and Computer Sciences*, 28(1):31–36, 1988.
- 535 [69] Zhaoping Xiong, Dingyan Wang, Xiaohong Liu, Feisheng Zhong, Xiaozhe Wan, Xutong Li, Zhaojun Li,
536 Xiaomin Luo, Kaixian Chen, Hualiang Jiang, et al. Pushing the boundaries of molecular representation for
537 drug discovery with the graph attention mechanism. *Journal of Medicinal Chemistry*, 63(16):8749–8760,
538 2019.
- 539 [70] Xiaocong Yang, Yang Liu, Jianhong Gan, Zhi-Xiong Xiao, and Yang Cao. FitDock: Protein–ligand
540 docking by template fitting. *Briefings in Bioinformatics*, 23(3):bbac087, 2022.
- 541 [71] Xin Zeng, Shu-Juan Li, Shuang-Qing Lv, Meng-Liang Wen, and Yi Li. A comprehensive review of the
542 recent advances on predicting drug-target affinity based on deep learning. *Frontiers in Pharmacology*, 15:
543 1375522, 2024.
- 544 [72] Shuke Zhang, Yanzhao Jin, Tianmeng Liu, Qi Wang, Zhaohui Zhang, Shuliang Zhao, and Bo Shan.
545 SS-GNN: a simple-structured graph neural network for affinity prediction. *ACS Omega*, 8(25):22496–22507,
546 2023.
- 547 [73] Xujun Zhang, Odin Zhang, Chao Shen, Wanglin Qu, Shicheng Chen, Hanqun Cao, Yu Kang, Zhe Wang,
548 Ercheng Wang, Jintu Zhang, et al. Efficient and accurate large library ligand docking with karmadock.
549 *Nature Computational Science*, 3(9):789–804, 2023.
- 550 [74] Ziqiao Zhang, Bangyi Zhao, Ailin Xie, Yatao Bian, and Shuigeng Zhou. Activity cliff prediction: Dataset
551 and benchmark. *arXiv preprint arXiv:2302.07541*, 2023.

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	7f6l_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	4jia	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	5mxm	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	1u7l	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	2ijk	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)

552 Appendix A Datasets and Baseline Models

553 A.1 License and Availability

554 The code for benchmark is available here: <https://anonymous.4open.science/r/DockedAC>.
 555 The DockedAC dataset and its future updates can be found here: <https://doi.org/10.5281/zenodo.11485280>.
 556

557 The DockedAC dataet is licensed under the *Creative Commons Attribution-ShareAlike 4.0 International License*. For details, please see <https://creativecommons.org/licenses/by-sa/4.0/>.
 558 The content of DockedAC includes data from the following sources: **RCSB PDB**, which is
 559 available under the *CC0 1.0 Universal (CC0 1.0) Public Domain Dedication* (more informa-
 560 tion at <https://creativecommons.org/publicdomain/zero/1.0/>), and **ChEMBL**,
 561 which is licensed under the *Creative Commons Attribution-ShareAlike 3.0 Unported License* (see
 562 <https://creativecommons.org/licenses/by-sa/3.0/>).
 563

Table 4: 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 The maximum depth of the model, n_d	$n_b = [100, 200, 400]$ $n_d = [5, 6, 7]$
SVM	The regularization parameter, C The kernel coefficient for <i>rbf</i> , γ	$C = [1, 10, 100, 1000, 10,000]$ $\gamma = [1 \times 10^{-5}, 1 \times 10^{-4}, 1 \times 10^{-3}, 1 \times 10^{-2}, 1 \times 10^{-1}]$
<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>

564 A.2 Datasets

565 Our introduced dataset DockedAC² comprises 82,836 ligands from 52 protein targets, which is
 566 meticulously curated to support various machine learning and deep learning studies related to activity
 567 cliff (AC) prediction. Table 3 provides detailed statistics of DockedAC.

568 A.3 Baseline Models

569 In this work, we integrate 13 recent baselines commonly used for structure-activity relationship
 570 prediction, including four traditional machine learning algorithms: KNN, RF, GBM, and SVM; three
 571 sequential models: LSTM, Transformer, and 1D CNN; four 2D GNN models: GCN, GAT, MPNN,
 572 and AFP; and two 3D structure GNN models: IGN and SS-GNN. The detailed descriptions of these
 573 approaches are listed in the following:

- 574 • **KNN** [10]. K-Nearest Neighbor (KNN) is a simple, non-parametric method that predicts the target
 575 molecule’s response by averaging the response of the k-nearest neighbors from the training set.
- 576 • **RF** [5]. Random Forest (RF) is an ensemble method that combines the outputs of multiple decision
 577 trees to improve accuracy and reduce over-fitting. Each decision tree is built upon a subset of the
 578 training set, and the final prediction is obtained by averaging the results from these individual trees.
- 579 • **GBM** [18]. Similar to RF, Gradient Boosting Machine (GBM) also combines the predictions of
 580 multiple decision trees. However, in GBM, these trees are built sequentially, with each subsequent
 581 tree specially designed to correct the errors of its predecessors.
- 582 • **SVM** [21]. Support Vector Machine (SVM) aims to identify a linear regression plane in a higher-
 583 dimensional space created by applying a designated kernel function. In this work, the Radial Basis
 584 Function (RBF) kernel is used.
- 585 • **Transformer** [63]. The Transformer model leverages self-attention mechanisms to capture depen-
 586 dencies across different positions in the input sequence. In our work, we employed the pretrained
 587 ChemBERTa [8] architecture, which has been trained on 10 million compounds.

²<https://anonymous.4open.science/r/DockedAC>

Table 5: 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	✗

- 588 • **LSTM** [23]. Long Short-Term Memory (LSTM) can capture temporal dependencies and patterns
589 in sequential data by maintaining long-term memory through their gated structure. In this work, we
590 employ SMILES strings as the input for the model.
- 591 • **1D CNN** [36]. Convolutional Neural Network (CNN) uses convolutional filters to aggregate spatial
592 information from adjacent positions. For processing sequential SMILES string data, we employ
593 1D CNNs that perform convolutional operations along a single dimension.
- 594 • **MPNN** [19]. Message Passing Neural Network (MPNN) operates by iteratively passing messages
595 between nodes and updating their representations based on neighboring nodes.
- 596 • **GCN** [37]. Graph Convolutional Network (GCN) performs convolution operations on graphs.
- 597 • **GAT** [63]. Graph Attention Network (GAT) introduces attention mechanisms to GNN to weigh the
598 importance of different neighbors.
- 599 • **AFP** [69]. Attentive Fingerprint (AFP) uses attention mechanisms at both the atom and molecule
600 levels to learn local and nonlocal properties, enabling it to capture substructural details effectively.
- 601 • **IGN** [31]. IGN models the molecular interactions in 3D space. In IGN, two graph convolution
602 modules are layered to learn intramolecular interactions and then sequentially intermolecular
603 interactions.
- 604 • **SS-GNN** [72]. Like IGN, SS-GNN is also a 3D structure GNN model tailored for affinity prediction.
605 It constructs a 3D structure graph for protein-ligand interactions based on a distance threshold,
606 reducing both the graph data scale and computational cost by omitting covalent bonds in proteins.

607 A.4 Model Features

608 In addition to the molecular descriptor used for machine learning algorithms (introduced in Sec. 4.2),
609 we further delve into the featurization for deep learning models. Detailed information on all featur-
610 izations and the corresponding models used can be found in Table 5.

611 For sequential methods, SMILES strings were encoded as one-hot vectors, with truncation applied
612 to strings exceeding 200 characters. To enhance model robustness, tenfold data augmentation was
613 applied using up to nine additional noncanonical SMILES strings for each SMILES string in the
614 dataset, generated via RDKit [38].

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

619 For SS-GNN, there are 11 node features, including atom type, formal charge, hybridization, atom
620 valence, atom degree, number of hydrogens, chirality, atomic mass, aromatic, atom coordinates, and
621 whether belonging to the protein. The edge features include covalent bond type, aromatic, bond
622 length, bond direction, bond stereochemistry, and edge type. The atom coordinates and bond length
623 are extracted from the 3D structures. Further details can be found in Zhang et al. [72].

624 For IGN, it uses similar 2D node and edge features. In addition, IGN uses four new edge features
625 from the 3D structures, including bond length, angle statistics, area statistics, and distance statistics.
626 For detailed descriptions of the features, see Jiang et al. [31].

627 A.5 Additional Experimental Details

628 **Hardware Specifications.** All our experiments were carried out on an NVIDIA RTX3090 GPU with
629 24G memory. The training time of a target for MPNN, GAT, GCN, and AFP is around 0.5 hours.
630 Training of one target takes around 1 hour and 4 hours for SS-GNN and IGN, respectively.

631 **Implementation Details.** Traditional machine learning algorithms including KNN, SVM, GBM, and
632 RF regression models were implemented using the Scikit-Learn library³.

633 Deep learning algorithms were trained for 500 epochs with early stopping, set with patience of 10
634 epochs. Four GNN models are implemented using the PyTorch Geometric package⁴. For the MPNN,
635 GCN, and GAT, global pooling was enabled using a graph multiset Transformer [1] with eight
636 attention heads, followed by a fully connected prediction head. Each of these models utilized two
637 graph layers. The Transformer model was based on the ChemBERTa [8] architecture, using weights
638 derived from 10M compounds in PubChem. Fine-tuning was conducted by freezing the original
639 model weights and substituting the final pooling layer with a regression head. Following van Tilborg
640 et al. [62], the LSTM model is pretrained on the SMILES strings with the next token prediction
641 objective. For the SS-GNN model, we conducted a pretraining phase on the original dataset, PDBbind
642 V2019 [66, 67]. In contrast, the IGN model was not fine-tuned using the original dataset due to a
643 mismatch in the model dimensions caused by the varying types of atoms in the dataset. Consequently,
644 we opted to train the IGN model from scratch.

645 **Hyperparameter Optimization.** Hyperparameter optimization was conducted through grid search.
646 Hyperparameter combinations were evaluated for all models using five-fold cross-validation. Table 4
647 shows the detailed hyperparameter search space.

648 Appendix B Additional Results and Figures

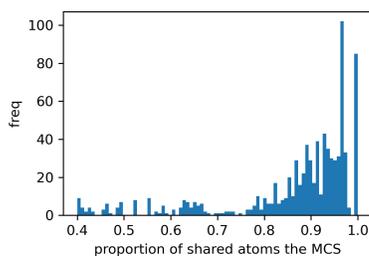


Figure 8: A histogram showing the proportion of shared atoms in AC data pairs with MCS for Target ChEMBL218 EC_{50} .

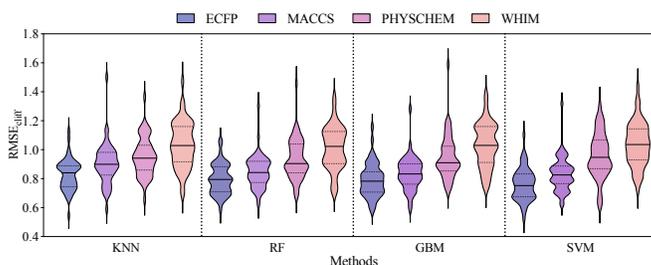


Figure 9: The $RMSE_{cliff}$ evaluated using ML methods under MACCS split across 52 targets.

649 **More dataset features.** Figure 11 illustrates three examples of removed targets and ligands. Figure 8
650 analyzes the proportion of shared atoms between the AC pairs in the target ChEMBL218 EC_{50} using
651 Maximum Common Substructure (MCS). The average proportion of shared atoms (86.78%) in the
652 identified AC pairs confirms high structural similarity in common substructures.

653 **Dataset split.** We split the dataset using the Tanimoto similarity of the ECFP. To assess potential bias
654 from ECFP-based data splitting, Figure 9 evaluates ML methods using four molecular descriptors on
655 an alternative MACCS-based split. ECFP maintains superior performance, confirming its inherent
656 descriptive power.

657 **Protein flexibility.** Using DSDPFlex [16], we investigate protein flexibility by allowing flexible side
658 chains for 10 amino acids nearest to the crystal ligand. Figure 15 shows that performance metrics on
659 8 K_i targets distribute evenly around the $y = x$ line, suggesting comparable effectiveness between
660 fixed and flexible docking approaches.

661 **Train cross-target models with 3D data.** Table 6 explore the cross-target applicability of 3D models
662 on combined K_i targets under two settings: out-of-distribution (OOD) excluding Protease targets,

³<https://scikit-learn.org/>

⁴<https://www.pyg.org/>

663 and in-domain, using all K_i targets. Figure 14 shows multi-target training performs comparable to
664 single-target training, complementing the analysis in Section 5.4.

665 **Combine the 3D information and ECFP features.** To explore the integration of 3D structural infor-
666 mation with handcrafted ECFP features, we utilize a 3D model as a feature extractor, combining its
667 output with ECFP descriptors, followed by MLP for affinity prediction (architecture shown Figure 16)
668 The evaluation across ten targets (shown in Table 7) highlights two key findings. First, models with
669 3D information consistently outperform or match those without 3D information across most targets,
670 achieving notable improvements in overall RMSE and $RMSE_{cliff}$. Second, the integration of 3D
671 features significantly enhances the model’s ability to handle activity cliffs, as evidenced by greater
672 improvements in $RMSE_{cliff}$ (avg. imp. of 5.61%) compared to overall RMSE (avg. imp. of 3.48%).

673 **Benchmarking the zero-shot ability of more 3D models.** To explore the generalization ability
674 of recent 3D binding affinity prediction models, we evaluate six SOTA methods (PIGNet [45],
675 RTMScore [49], TANKBind [41], DSMBind [35], KarmaDock [73], and EquiScore [6]) trained on
676 PDBBind. Figure 17 presents their Pearson correlation on the complete dataset and activity cliff cases
677 across each target. All these methods perform worse on the AC samples, which is consistent with
678 the result of our benchmark. Additionally, these methods show decreased performance compared
679 to the PDBBind test set, with effectiveness correlating with the presence of homologous proteins in
680 the PDBBind training data. For instance, targets with numerous homologous samples in PDBBind
681 demonstrate superior results: ChEMBL2147 K_i achieves a Pearson correlation of 0.688 (DSMBind,
682 PDB ID: 2j2i) with 103 homologous samples, while ChEMBL2971 K_i reaches 0.671 (DSMBind,
683 PDB ID: 4jia) with 61 homologous samples in PDBBind. In contrast, targets lacking homologous
684 proteins in PDBBind (ChEMBL219 K_i , ChEMBL228 K_i , and ChEMBL233 K_i) show very small
685 correlation (DSMBind, Pearson=-0.021, -0.087, and 0.033 respectively).

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

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

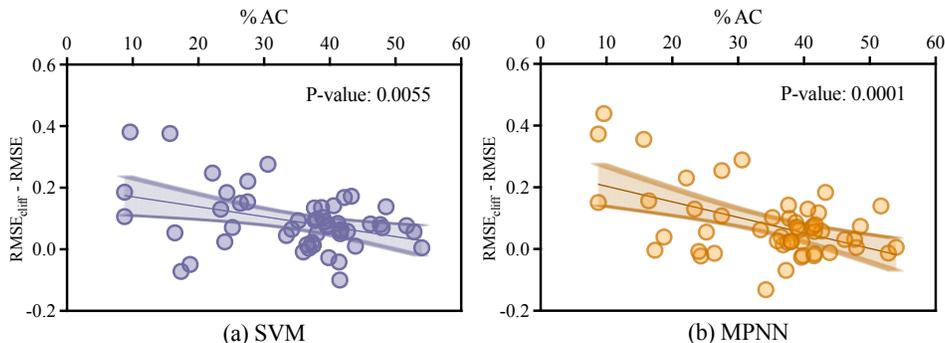


Figure 10: Relationship between the ratio of the AC and $RMSE - RMSE_{cliff}$ of SVM and MPNN.

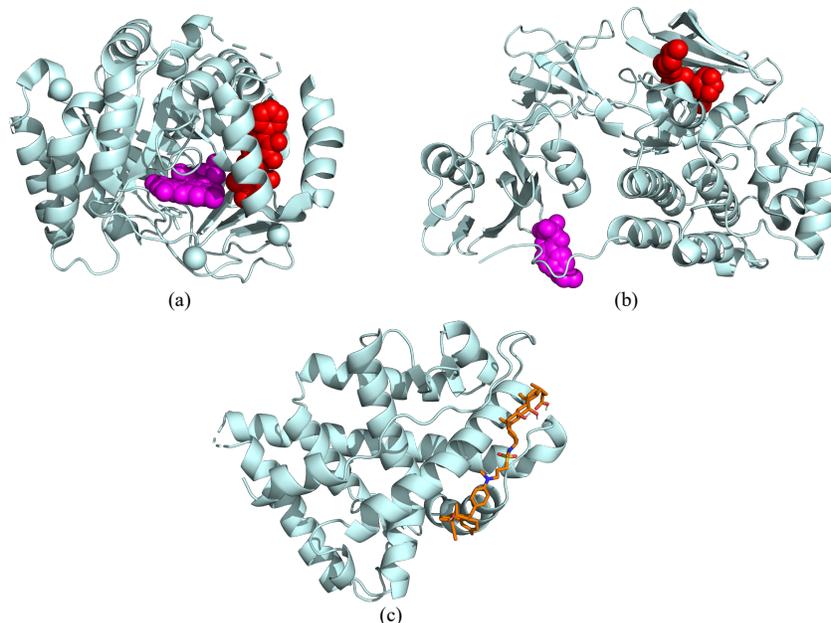


Figure 11: Three examples of the removed targets and ligands. (a) The target structure has two ligand binding sites (PDB: 5mvd). (b) Two structures of the same target have different binding sites (PDB: 2h8h and 1o4j). The two structures are aligned. (c) The ligand docking score is larger than zero (Target: ChEMBL1871, PDB: 2ama, ligand: ChEMBL406027).

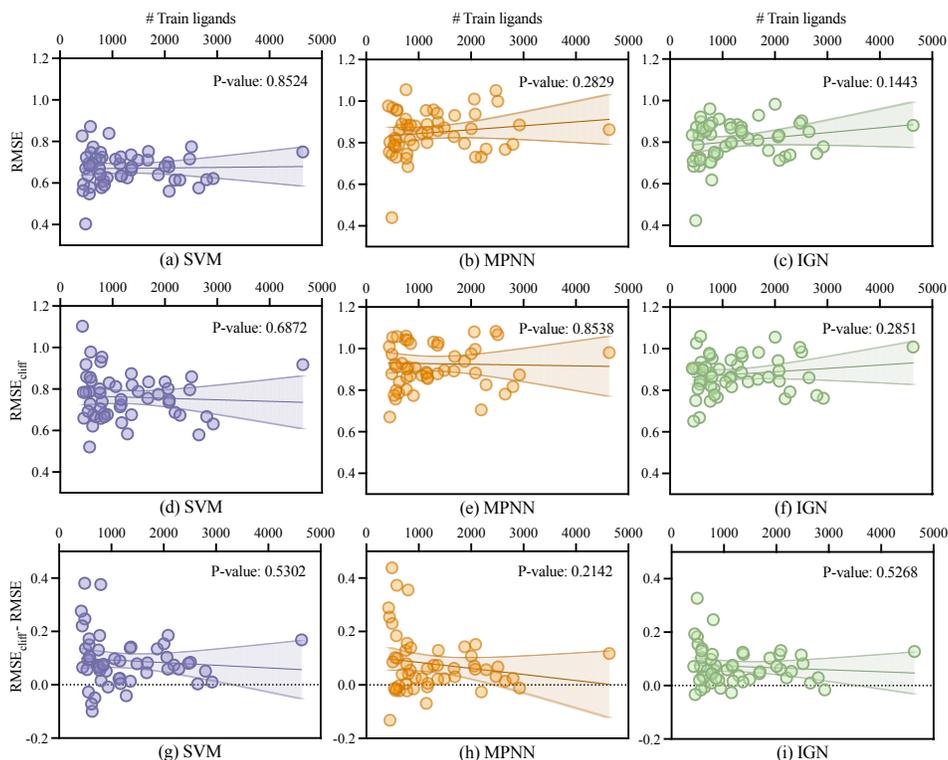


Figure 12: Relationship between the number of training ligands and (a)-(c) RMSE, (d)-(f) $RMSE_{cliff}$ and (g)-(i) their difference on SVM, MPNN, and IGN.

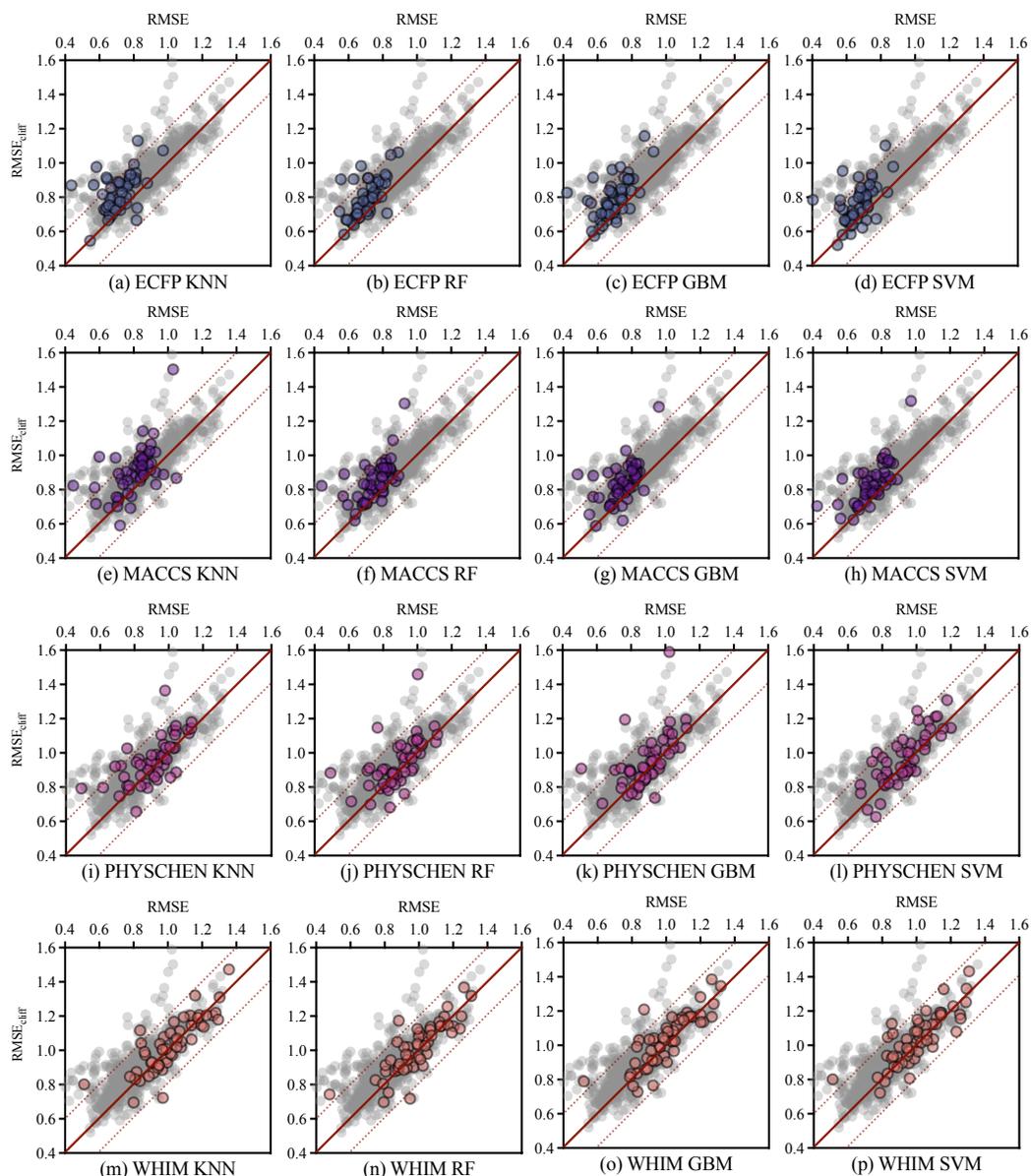


Figure 13: Performance comparison between RMSE and $RMSE_{cliff}$ for classic ML algorithms across 52 targets.

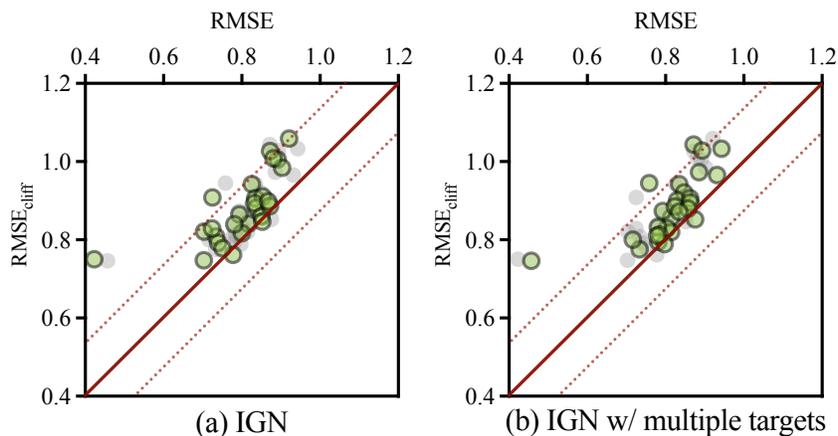


Figure 14: The results of IGN on all the targets of K_i labels when trained separately on (a) each target or (b) the data of multiple targets.

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

Model	ChEMBL204 Ki		ChEMBL244 Ki		ChEMBL3286 Ki		ChEMBL4822 Ki	
	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 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}								
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 EC50		ChEMBL3979 EC50		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

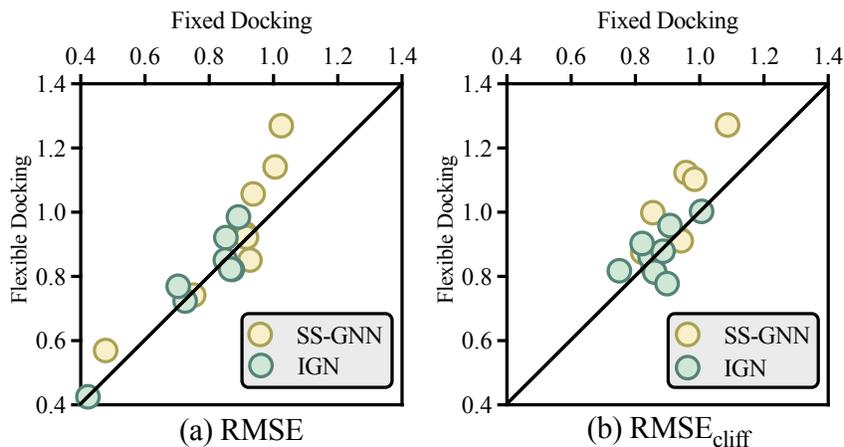
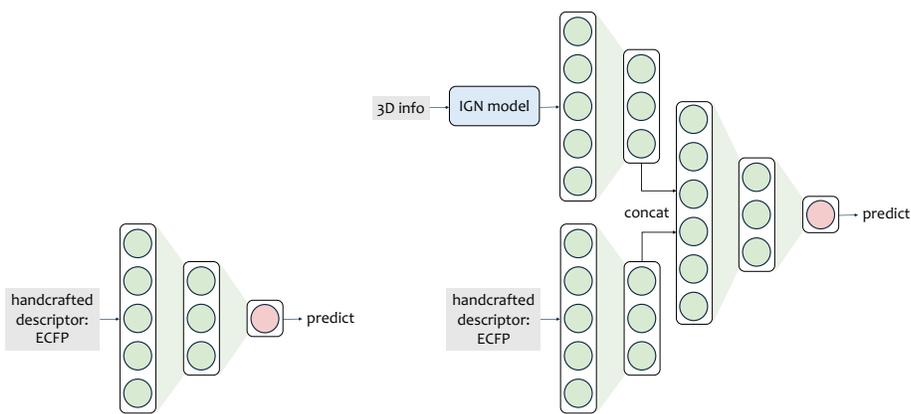


Figure 15: The (a) RMSE and (b) RMSE_{cliff} metric on fixed docking v.s. flexible docking.



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

Figure 16: The model illustration of MLP and IGN using the handcrafted molecule descriptor ECFP.

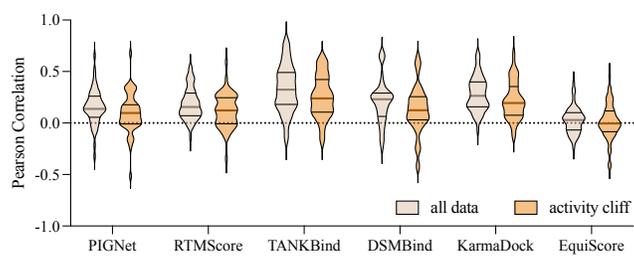


Figure 17: The Pearson and $\text{Pearson}_{\text{cliff}}$ evaluated on our DockedAC benchmark across 52 targets using PIGNet [45], RTMScore [49], TANKBind [41], DSMBind [35], KarmaDock [73], and EquiScore [6], all of which were trained on general binding affinity datasets.

704 **NeurIPS Paper Checklist**

705 **1. Claims**

706 Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s
707 contributions and scope?

708 Answer: [Yes]

709 **2. Limitations**

710 Question: Does the paper discuss the limitations of the work performed by the authors?

711 Answer: [Yes]

712 Justification: the limitations are briefly stated in Sec. 6 and a detailed discussion appears in
713 Appendix B.

714 **3. Theory assumptions and proofs**

715 Question: For each theoretical result, does the paper provide the full set of assumptions and a
716 complete (and correct) proof?

717 Answer: [NA]

718 **4. Experimental result reproducibility**

719 Question: Does the paper fully disclose all the information needed to reproduce the main experi-
720 mental results of the paper to the extent that it affects the main claims and/or conclusions of the
721 paper (regardless of whether the code and data are provided or not)?

722 Answer: [Yes]

723 **5. Open access to data and code**

724 Question: Does the paper provide open access to the data and code, with sufficient instructions to
725 faithfully reproduce the main experimental results, as described in supplemental material?

726 Answer: [Yes]

727 Justification: The code for benchmark is available here: [https://anonymous.4open.scie
728 nce/r/DockedAC](https://anonymous.4open.science/r/DockedAC). The DockedAC dataset and its future updates can be found here: [https:
729 //doi.org/10.5281/zenodo.11485280](https://doi.org/10.5281/zenodo.11485280).

730 **6. Experimental setting/details**

731 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,
732 how they were chosen, type of optimizer, etc.) necessary to understand the results?

733 Answer: [Yes]

734 Justification: see Sec. A.5.

735 **7. Experiment statistical significance**

736 Question: Does the paper report error bars suitably and correctly defined or other appropriate
737 information about the statistical significance of the experiments?

738 Answer: [Yes]

739 **8. Experiments compute resources**

740 Question: For each experiment, does the paper provide sufficient information on the computer
741 resources (type of compute workers, memory, time of execution) needed to reproduce the experi-
742 ments?

743 Answer: [Yes]

744 Justification: see Sec. A.5.

745 **9. Code of ethics**

746 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS
747 Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

748 Answer: [Yes]

749 Justification: This research conducted in the paper conforms with the NeurIPS Code of Ethics.

750 **10. Broader impacts**

751 Question: Does the paper discuss both potential positive societal impacts and negative societal
752 impacts of the work performed?

753 Answer: [Yes]

754 **11. Safeguards**

755 Question: Does the paper describe safeguards that have been put in place for responsible release of
756 data or models that have a high risk for misuse (e.g., pretrained language models, image generators,
757 or scraped datasets)?

758 Answer: [NA]

759 **12. Licenses for existing assets**

760 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the
761 paper, properly credited and are the license and terms of use explicitly mentioned and properly
762 respected?

763 Answer: [Yes]

764 Justification: The license of the dataset in work, CC-BY-SA 4.0, complies with the licenses of all
765 the cited previous resources.

766 **13. New assets**

767 Question: Are new assets introduced in the paper well documented and is the documentation
768 provided alongside the assets?

769 Answer: [Yes]

770 Justification: see Appendix Sec. A.1.

771 **14. Crowdsourcing and research with human subjects**

772 Question: For crowdsourcing experiments and research with human subjects, does the paper
773 include the full text of instructions given to participants and screenshots, if applicable, as well as
774 details about compensation (if any)?

775 Answer: [NA]

776 **15. Institutional review board (IRB) approvals or equivalent for research with human subjects**

777 Question: Does the paper describe potential risks incurred by study participants, whether such
778 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals
779 (or an equivalent approval/review based on the requirements of your country or institution) were
780 obtained?

781 Answer: [NA]

782 **16. Declaration of LLM usage**

783 Question: Does the paper describe the usage of LLMs if it is an important, original, or non-
784 standard component of the core methods in this research? Note that if the LLM is used only for
785 writing, editing, or formatting purposes and does not impact the core methodology, scientific
786 rigorosity, or originality of the research, declaration is not required.

787 Answer: [No]