PoseX: AI Defeats Physics-based Methods on Protein Ligand Cross-Docking

Yize Jiang^{1,*}, Xinze Li^{1,*}, Yuanyuan Zhang^{2,*}, Jin Han^{3,*}, Youjun Xu^{4,*},

Ayush Pandit⁵, Zaixi Zhang⁶, Mengdi Wang⁶, Mengyang Wang⁷, Chong Liu⁸, Guang Yang⁹,

Yejin Choi^{5,10}, Wu-Jun Li^{3,†}, Tianfan Fu^{3,†}, Fang Wu^{5,†}, Junhong Liu^{1,†,*}

¹MicroCyto, ²Purdue University, ³Nanjing University, ⁴ByteDance,

⁵Stanford University, ⁶Princeton University, ⁷Peking University, ⁸Central South University

⁹Imperial College London, ¹⁰NVIDIA

*Equal Contributions, [†]Correspondence

Abstract

Recently, significant progress has been made in protein-ligand docking, especially in deep learning-based methods, and some benchmarks were proposed, such as PoseBench and PLINDER. However, these studies typically focus on the selfdocking scenario, which is less practical in real-world applications. Moreover, some studies involve heavy frameworks requiring extensive training, posing challenges for convenient and efficient assessment of docking methods. To address these gaps, we introduce PoseX, an open-source benchmark designed to evaluate both self-docking and cross-docking, enabling a practical and comprehensive assessment of algorithmic advances. Specifically, we curated a novel dataset comprising 718 entries for self-docking and 1,312 entries for cross-docking; secondly, we incorporated 23 docking methods in three methodological categories, including physics-based methods (e.g., Schrödinger Glide), AI docking methods (e.g., DiffDock) and AI co-folding methods (e.g., AlphaFold3); thirdly, we developed a relaxation method for post-processing to minimize conformational energy and refine binding poses; fourthly, we established a leaderboard to rank submitted models in real-time. We derived some key insights and conclusions from extensive experiments: (1) AI-based approaches have consistently outperformed physicsbased methods in overall docking success rate. (2) Most intra- and intermolecular clashes of AI-based approaches can be greatly alleviated with relaxation, which means combining AI modeling with physics-based post-processing could achieve excellent performance. (3) AI co-folding methods commonly exhibit ligand chirality issues, except for Boltz-1x, which introduced physics-inspired potentials to fix hallucinations, suggesting the modeling on stereochemistry improves the structural plausibility markedly. (4) Specifying binding pockets significantly promotes docking performance, indicating that pocket information can be leveraged adequately, particularly for AI co-folding methods, in future modeling efforts. The code, dataset, and leaderboard are released at https://github.com/CataAI/PoseX.

1 Introduction

2

3

5

6

7

8

9

10

12

13

14

15

16

17

18

19

20 21

22

23

24

25

26

27

Protein-ligand docking is crucial to drug discovery since it predicts how a ligand interacts with a protein, helping to identify potential drug candidates and accelerate the development of new therapeutics. Learning from known protein-ligand complexes through machine learning, especially deep learning (DL) techniques, the AI-based approaches have the potential to revolutionize proteinligand docking by significantly enhancing the speed and accuracy of prediction, and substantial

- progress has been made recently. In response to the large number of new approaches, recent work has introduced several benchmarks, such as PoseBench [1] and PLINDER [2], with corresponding 35 datasets and metrics focusing on the evaluation of protein-ligand interaction. Despite the rapid 36 progress, existing studies still encounter several challenges, summarized as follows.
- 1. Impractical evaluation scenarios. Most existing benchmarks, such as PoseBuster [3] and 38 PoseBench, focus on the self-docking scenario, which is less practical in real-world applications. 39 For instance, pharmaceutical chemists usually design new drug molecules and dock them with 40 the targets, of which the conformations are extracted from existing complex structures that are 41 co-crystallized with other published compounds. 42
- 2. Heavy framework and low accessibility. Some benchmarks (e.g., PLINDER) suffer from heavy 43 evaluation frameworks that involve data splitting and training, which are hard to use. While 44 studies such as PoseBuster that only concentrate on evaluation rather than together with training 45 are worthy of reference, which are lightweight and user-friendly. 46
- 3. Limited model selection for benchmarking. Existing studies often restrict their comparative 47 scope to a narrow set of models. For instance, PoseBuster evaluated only 5 AI-based approaches 48 and 2 physics-based methods, while PLINDER exclusively benchmarked against DiffDock [4], 49 neglecting other notable algorithms. 50
- Therefore, we propose several solutions to address these issues: 51

37

64

65

68

- 1. Practical evaluation setup. To better evaluate the capacity of various docking methods in a more 52 practical scenario, we incorporate cross-docking, which involves docking various small molecules 53 extracted from distinct complexes of the same protein with all the conformations except the native 54 co-crystalized one. 55
- 2. Construction of new dataset. We curated a new dataset named PoseX that collects newly found 56 crystal structures of protein-ligand complexes in RCSB PDB, which contains 718 entries for 57 self-docking and 1,312 entries for cross-docking. 58
- 3. A wide variety of models. We evaluated 23 docking methods encompassing nearly all relevant 59 models published in peer-reviewed journals and conferences alongside established commercial 60 61 docking software across three different research lines, including 5 physics-based methods such as Schrödinger Glide [5], 11 AI docking methods such as DiffDock, and 7 AI co-folding methods 62 such as AlphaFold3 [6]. 63
 - In addition, we developed a novel relaxation module (also known as energy minimization), which serves as a post-processing method to refine AI-generated binding poses. It is especially helpful to promote structural plausibility. We also established an online leaderboard, which enables researchers to benchmark their models against a standardized dataset, fostering transparency and facilitating easy and fair comparisons for the broader community. The key differences between the existing docking benchmarks and ours are summarized in Table 1.

Table 1: Comparison of existing docking benchmark studies.

Benchmarks	PoseBuster	PoseBench	PLINDER	PoseX (Ours)
Code of dataset pipeline	Х	Х	/	√
Relaxation	×	coarse	×	well-designed
Self-docking evaluation	✓	✓	✓	✓
Cross-docking evaluation	×	X	X	✓
# Open-source docking software	2	1	0	2
# Commercial docking software	0	0	0	3
# Physics-based methods	2	1	0	5
# AI docking methods	5	2	1	11
# AI co-folding methods	0	4	0	7
# Total methods	7	7	1	23
Real-time leaderboard	×	×	×	✓

2 Methods

We categorize all the docking approaches into three distinct classes: (1) *physics-based methods* utilize physics-based scoring functions and sampling algorithms to estimate protein-ligand interactions, including Discovery Studio [7], Schrödinger Glide [5], MOE [8], AutoDock Vina [9, 10], and GNINA [11]; (2) *AI docking methods* produce ligand binding poses based on the three-dimensional structure of proteins, including DeepDock [12], EquiBind [13], TankBind [14], DiffDock [4], Uni-Mol [15], FABind [16], DiffDock-L [17], DiffDock-Pocket [18], DynamicBind [19], Interformer [20], SurfDock [21]; (3) *AI co-folding methods* predict both the ligand's binding conformation and the protein's conformational changes induced by ligand binding, which account for simultaneous structural adaptations of the protein and ligand, enabling more accurate modeling of their interactions; we involve 7 *AI co-folding methods*, including NeuralPLexer [22], RoseTTAFold-All-Atom (RFAA) [23], AlphaFold3 [6], Chai-1 [24], Boltz-1 [25], Boltz-1x [25], Protenix [26]. For comparative analysis, we summarize all the compared methods in Table 2, and the detailed settings of these methods are shown in Appendix B.

Table 2: Comparison of various methods.

Method	Pub. Year	License	Pocket Required	Pocket Changed	Avg. Runtime Per Sample ¹
	Physics-based methods				
Discovery Studio	late 1990s	Commercial	✓	X	14.4 min
Schrödinger Glide	2004	Commercial	✓	X	7.2 min
MOE	2008	Commercial	✓	X	50 sec
AutoDock Vina	2010, 2021	Apache-2.0	✓	X	18 sec
GNINA	2021	Apache-2.0	✓	X	12 sec
		AI docking methods			
DeepDock	2021	MIT	✓	X	2.7 min
EquiBind	2022	MIT	X	X	1.4 sec
TankBind	2022	MIT	X	X	7.8 sec
DiffDock	2022	MIT	X	X	1.2 min
Uni-Mol	2024	MIT	✓	X	24 sec
FABind	2023	MIT	X	X	8.8 sec
DiffDock-L	2024	MIT	X	X	1.5 min
DiffDock-Pocket	2024	MIT	✓	✓	1.7min
DynamicBind	2024	MIT	X	✓	2.4 min
Interformer	2024	Apache-2.0	✓	X	0.6 min
SurfDock	2024	MIT	✓	×	10.8 sec
AI co-folding methods					
NeuralPLexer	2024	BSD	X	✓	1.5 min
RoseTTAFold-All-Atom	2023	BSD	×	✓	9 min
AlphaFold3	2024	CC-BY-NC-SA 4.0	×	✓	16.5 min
Chai-1	2024	Apache-2.0	×	✓	3 min
Boltz-1	2024	MIT	×	✓	3 min
Boltz-1x	2025	MIT	×	✓	3 min
Protenix	2025	Apache-2.0	Х	✓	3.6 min

¹ Regarding computational runtime performance, different methods operate on varied computational environments. Details for each method are provided in Appendix B.

Relaxation as Post-processing Relaxation in molecular docking, also known as *energy minimization*, is a post-processing method used to refine and optimize docked protein-ligand complexes [27, 28]. It involves energy minimization and sometimes short molecular dynamics simulations to resolve steric clashes, improve atomic interactions, and ensure the system reaches a stable, low-energy conformation. This step enhances the physical realism and the accuracy of the docking results, making the predicted binding poses more reliable for further analysis or experimental validation. In this paper, we introduce a novel relaxation module, the novelty of which is summarized as: (1) Implemented an automated relaxation process for complexes based on OpenMM [29]. (2) Established a comprehensive automatic data processing pipeline for proteins and small molecules, including fixing missing chains, capping the N- and C-terminals, adding formal charges to proteins and small molecules, and applying restraints to backbone atoms (CA, C, N, O). (3) Supports small molecule force field parameters from GAFF and OpenFF [30]. (4) Supports partial charge calculation methods for small molecules, including Gasteiger and MMFF94. (5) Effectively alleviates unreasonable

predicted conformations, improving the pass rate of PB-Valid. The technical details of the relaxation
 process are provided in Appendix C.

99 3 Dataset

100 3.1 Self-docking Versus Cross-docking

Self-docking. Self-docking involves docking a ligand back into its native co-crystallized conformation [31]. This is typically used to check if the docking software can accurately reproduce the known binding pose, helping validate the method. Most existing benchmarks only consider the self-docking setup.

Cross-docking. Cross-docking refers to dock molecules extracted from distinct complexes of the same protein with all conformations except the native co-crystalized one. This approach is considered a more versatile evaluation, as it takes into account the fact that the receptor protein may undergo conformational changes and might not be fully optimized for docking with the ligand. The difference between self-docking and cross-docking is illustrated in Figure 1.

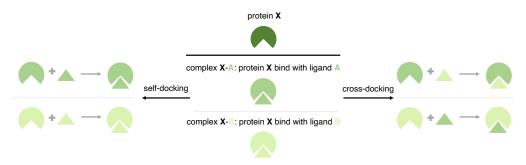


Figure 1: Self-docking vs. Cross-docking.

110 3.2 Astex

116

The Astex Diverse set [32], published in 2007, is a set of hand-picked, relevant, diverse, and highquality protein-ligand complexes from the RCSB PDB. It comprises 85 unique and significant protein-ligand complexes. These complexes have been appropriately formatted for docking purposes and will be made freely accessible to the entire research community via the website (http://www. ccdc.cam.ac.uk). It only supports self-docking evaluation.

3.3 PoseX: Our Curated Dataset

In this paper, we curated a high-quality protein-ligand complex structure dataset designed to evaluate 117 molecular docking methods named PoseX. It consists of carefully selected crystal structures from 118 the RCSB Protein Data Bank (RCSB PDB) [33] with two subsets for evaluating self-docking and 119 cross-docking tasks. The dataset only includes complex structures published from 2022 January 120 1st to 2025 January 1st, ensuring that there is no overlap with the training data of all AI-based 121 approaches that are being evaluated (as shown in Table S3). The construction steps of the two subsets 122 PoseX Self-Docking (PoseX-SD) and PoseX Cross-Docking (PoseX-CD) are shown in Table S1 and 123 Table S2. Ultimately, there are 718 entries for PoseX-SD and 1,312 entries for PoseX-CD, comprising 124 109 protein targets (a total of 371 structures) and 362 small molecules. The distribution of the number 125 of conformation structures per target is shown in Figure S1a, and the distribution of pocket similarity 126 is shown in Figure S1b.

4 Experiments

129 4.1 Evaluation Metrics

- Performance evaluation of protein-ligand docking involves metrics that assess both the quality of the predicted binding pose and the chemical validity as well as the structural plausibility, which are described in detail as follows.
- RMSD. In accordance with most benchmarking studies, we evaluate the quality of the binding poses with Root Mean Square Deviation (RMSD), which measures the distance between the predicted and the ground-truth complex structures. Lower RMSD scores indicate better binding poses.
- PB-Valid. The physicochemical validity and structural plausibility of the generated binding poses are measured with the PoseBuster test suite (*i.e.*, PB-Valid). This suite evaluates whether predicted ligand poses are consistent with known chemical and structural constraints. See Appendix D for more details.
- Success rate. The docking success rate is defined as the percentage of the top-1 ranked predictions satisfying either of the following criteria: (1) RMSD < 2Å, or (2) RMSD < 2Å & PB-Valid. For PoseX-CD, we report the averaged success rate at the target level in view of the uneven distribution of docking sizes per target (as shown in Figure S1a). Higher success rates indicate better performance.

44 4.2 Results

145 4.2.1 Overall Performance Analysis

- Figure 2 and Figure S2 present a comprehensive evaluation of various docking approaches on three benchmarks Astex, PoseX-SD, and PoseX-CD under RMSD < 2Å and PB-Valid criteria. From these results, we highlight several main observations and provide a more granular analysis of these results.
- 1. **AI-based approaches lead in success rate.** The latest AI-based approaches, both *AI docking methods* (*e.g.*, SurfDock) and *AI co-folding methods* (*e.g.*, AlphaFold3) have consistently outperformed physics-based methods in overall docking pose and validity.
- Relaxation mitigates clashing significantly. The intra- and intermolecular clashes of AI-based approaches can be greatly alleviated with relaxation, which means that the force field-based energy minimization step is very crucial to achieve excellent performance in real-world applications, particularly for AI modeling.
- 157 3. **Chirality warrants further improvement.** Most of the *AI co-folding methods* exhibit ligand chirality issues, such as AlphaFold3 and Chai-1, except for Boltz-1x, which introduces an inference time steering technique employing physics-inspired potential to fix hallucinations and enhance structural plausibility.
- Pocket information is crucial to docking. Explicit modeling of binding pocket substantially improves docking performance, as seen by the consistent performance gains of DiffDock-Pocket over its counterpart DiffDock across both self-docking and cross-docking, indicating that pocket information can be leveraged adequately, especially for *AI co-folding methods*, in future modeling efforts.
- **Astex Benchmark.** The Astex benchmark represents an idealized docking scenario with high-quality 166 co-crystal structures. In this setting, AI docking methods outperform all other categories overall. 167 Uni-Mol and SurfDock achieve the highest docking success rates (94.1%) when integrated with our 168 structural relaxation protocol, surpassing physics-based methods, such as Glide and Discovery Studio, 169 by over 25%. DiffDock-Pocket, Interformer, and DiffDock-L also perform strongly, achieving success 170 rates above 83.5%. While AI co-folding methods such as AlphaFold3, Protenix, and Chai-1 deliver 171 competitive results (over 80% success), they are marginally outperformed by docking-specialized 172 architectures. *Physics-based methods* like AutoDock Vina and MOE plateau around 56.5%–67.1%, 173 even with induced-fit docking (e.g., Glide IFD). These results illustrate the substantial performance 174 gains offered by AI modeling tailored specifically for pose prediction. 175
- PoseX-SD Benchmark. For PoseX-SD evaluation, SurfDock (78.4%) achieves the overall state-ofthe-art performance, and Uni-Mol takes the second place. DiffDock-Pocket shows clear advantages

over its pocket-agnostic counterpart, with a success rate of 52.2%. Among AI co-folding methods, AlphaFold3 and Protenix perform well (60.3% and 56.0%, respectively), demonstrating their capacity to model close-range binding interactions. In contrast, earlier AI docking methods such as EquiBind and TankBind perform poorly (below 20%), meanwhile, they exhibit significant issues with structural plausibility. Physics-based methods such as Glide and Discovery Studio remain clustered in the 48–55% range. Most AI-based approaches benefit from the relaxation method we developed, and their intra- and intermolecular validity are significantly improved.

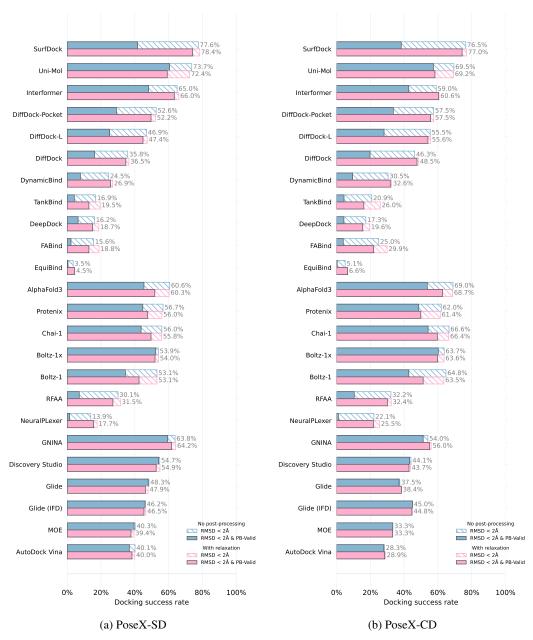


Figure 2: Performance on PoseX-SD and PoseX-CD.

PoseX-CD Benchmark. For PoseX-CD evaluation, SurfDock (77.0%) and Uni-Mol (69.2%) are still the top performers in all three categories of docking methods, as well as AlphaFold3, which achieves competitive performance (68.7%) against Uni-Mol. We observed that *AI docking methods* have developed rapidly in recent years, of which the latest models (such as SurfDock, Uni-Mol, Interformer, and DiffDock-Pocket) demonstrably surpass the earlier models (such as EquiBind, TankBind, and DeepDock). For *AI co-folding methods*, AlphaFold3 defeats other models (such as

Chai-1, Boltz-1, Boltz-1x and Proteinx) by a narrow margin (2.3% - 7.3%). Notably, physics-based 191 methods struggle significantly in this scenario. For example, in the PoseX-SD task, only 3 AI 192 docking methods outperform the leading physics-based method, GNINA, in terms of the percentage 193 of RMSD < 2Å with relaxation. However, in the PoseX-CD task, 9 AI-based approaches (including 194 4 AI docking methods and 5 AI co-folding methods) surpass GNINA (56.0%). This underscores 195 a significant advantage of AI-based approaches over physics-based methods in the cross-docking 196 scenario. Figure S10 and Figure S11 depict an illustrative example of the superior performance of 197 AI-based methods. Relaxation yields consistent improvements across most approaches, emphasizing 198 its role in resolving steric or geometric inconsistencies. 199

4.2.2 Pocket Similarity based Generalizability Analysis

200

To further understand the generalization capacity of various docking approaches, we analyze the 201 relationship between pocket similarity and ligand RMSD across different scenarios. In view of the 202 cut-off time of the training data for each method (as shown in Table S3), the pocket similarity is 203 calculated as the maximum TM-score compared to pockets extracted from crystal structures released 204 before 2022 on RCSB PDB, where the pocket is defined as the residues within 10.0Å of the ligand. 205 Figure S3 and Figure S4 present per-sample scatter plots of pocket similarity versus docking RMSD 206 for self-docking and cross-docking, respectively. Each plot reports Pearson's correlation coefficient 207 to quantify the strength and direction of the relationship. Figure 3 complements these results by 208 summarizing the average ligand RMSD separately for test cases with similar and dissimilar pockets. 209 Figure S5 and Figure S6 illustrate the relationship between the ligand RMSD and the decreasing binding pocket similarity of AI-based approaches.

Self-Docking Observations. In the self-docking scenario, most AI-based approaches exhibit a moderate negative correlation between pocket similarity and ligand RMSD, indicating that the leakage of pocket information is associated with improved ligand pose accuracy. For example, Protenix and Chai-1 show stronger correlations (r = -0.390 and r = -0.389, respectively), while other models such as AlphaFold3 (r = -0.313) and Boltz-1 (r = -0.276) exhibit similar trends. DiffDock and DiffDock-L display similar correlations (r = -0.283 and r = -0.278, respectively), suggesting that docking-specific models also benefit from the pocket leakage.

In contrast, *physics-based methods* show weaker or near-zero correlations. Glide (r = 0.010), AutoDock Vina (r = -0.009), and Discovery Studio (r = -0.001) exhibit negligible correlations, suggesting consistent docking performance across varying pocket similarities.

Notably, SurfDock (r=-0.091) and Uni-Mol (r=-0.134), which achieve top performance overall, show only weak correlation between pocket similarity and ligand RMSD. These findings suggest that their success likely stems from robust pose prediction mechanisms that have less sensitivity for pocket information leakage. These results highlight the importance of robust pose prediction in achieving high docking performance, even when pocket similarity is limited, in the self-docking scenario.

Cross-Docking Observations. The cross-docking setting reveals an overall stronger correlation between pocket similarity and ligand RMSD, particularly for *AI co-folding methods* and *AI docking methods*. Chai-1 (r=-0.526), Boltz-1 (r=-0.521), and Protenix (r=-0.553) exhibit strong negative correlations, suggesting that successful docking in cross-docking is highly contingent upon correctly modeling the target pocket's conformation. DiffDock and its variants continue to reflect this trend (e.g., DiffDock r=-0.505; DiffDock-L r=-0.498), further confirming the influence of pocket leakage under receptor shift scenarios.

Models such as DynamicBind (r=-0.576) and DiffDock-Pocket (r=-0.425) also show a strong correlation between pocket similarity and ligand RMSD, reinforcing that flexible or dynamic AI docking methods also have constrained generalization. In contrast, physics-based methods such as Glide (r=0.015) and Discovery Studio (r=0.053) again exhibit negligible correlation.

Even high-performing models like SurfDock (r=-0.376) and Uni-Mol (r=-0.280) show stronger correlations in this setting than in self-docking, indicating that pocket modeling becomes more critical in the presence of conformational variance. This further highlights the need for improved pocket-conditioned pose generation in cross-docking scenarios.

Performance Stratified by Pocket Similarity. Figure 3 further stratifies the average ligand RMSD for each method, where the evaluation set is split into two groups, a similar group (Pocket Similarity

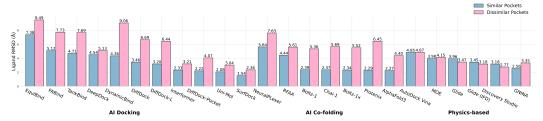


Figure 3: Cross-docking performance difference on "similar" and "dissimilar" binding pockets.

≥ 0.70) and a dissimilar group (Pocket Similarity < 0.70). Across all the AI-based approaches, both *AI docking methods* and *AI co-folding methods*, docking into similar pockets consistently achieve lower RMSD. However, the degradation of different models in dissimilar pockets evaluation varies significantly. *Physics-based methods* such as Glide, MOE, and Discovery Studio consistently demonstrate a very small gap between similar and dissimilar evaluations, which shows excellent generalizability that outperforms most of the AI-based approaches in the dissimilar pocket scenario. Earlier *AI docking methods* (e.g., TankBind) and most *AI co-folding methods* (e.g., Chai-1, Protenix, AlphaFold3) suffer steep performance drops—TankBind degrades from 4.71Å to 7.69Å, Chai-1 degrades from 2.37Å to 5.69Å, and Protenix degrades from 2.29Å to 6.45Å—highlighting their overreliance on pocket leakage and lack of adaptability. The latest *AI docking methods*, particularly SurfDock (1.54Å to 2.36Å) and Uni-Mol (2.08Å to 3.04Å), demonstrate much smaller gaps and showcase robust generalization.

Overall Implications. These analyses collectively suggest that pocket similarity is a key determinant of successful docking, particularly for the cross-docking scenario. *AI co-folding method* and *AI docking methods* reveal a stronger dependence on pocket information, while *physics-based methods* show little sensitivity. Notably, even the state-of-the-art models such as SurfDock and Uni-Mol exhibit varying levels of dependence on pocket fidelity, indicating that future improvements in docking may arise from synergistically enhancing both pocket modeling and pose prediction.

4.2.3 Impact of Relaxation from a Physically-based Validation Perspective

We systematically evaluated the docking performance of various methods using the PoseBuster test suite, comprising 20 physicochemical validation metrics that assess stereochemistry and intra- and intermolecular validity. Figures S7 and S8 illustrate the failure rates of the PB-Valid metric before and after relaxation in self-docking and cross-docking settings, respectively.

Without Relaxation. In the absence of relaxation, most AI docking methods generate ligand poses that violate physicochemical constraints. Notably, models such as EquiBind, FABind, and DeepDock exhibit a high failure rate in intermolecular validity, especially in the minimum distance-to-protein metric, with only approximately 10% of the predictions passing the test. Even SurfDock, which achieves the lowest RMSD, fails in nearly half of its predictions for this metric. Among the AI docking methods, Uni-Mol demonstrates the best performance on PB-Valid, but still exhibits chirality prediction errors. Among AI co-folding methods, NeuralPLexer and RFAA perform poorly in intermolecular validity. AlphaFold3 and similar models show relatively stable performance, but are not immune to chirality errors. In comparison, the recently introduced Boltz-1x model effectively addresses these issues, achieving the highest PB-Valid pass rate among all AI methods. Physics-based methods consistently perform well in structural plausibility, achieving high pass rates.

With Relaxation. Most AI-based approaches benefit significantly from our relaxation protocol, which effectively mitigates intra- and intermolecular clashes. SurfDock emerges as the top-performing method in the benchmark post-relaxation. However, relaxation does not resolve chirality errors, and Uni-Mol shows no performance improvement in this process. Similarly, *AI co-folding methods*, including AlphaFold3, Chai-1, Boltz-1, and Protenix, exhibit limited improvement due to persistent chirality errors. Figure S9 illustrates two representative cases of chirality errors in docking predictions. Chiral errors still exist in the prediction results of AlphaFold3, Chai-1, and Boltz-1, but are resolved in Boltz-1x.

Summary. Integrating relaxation with *AI docking methods* yields the state-of-the-art performance.
Concurrently, advancements in AI for biology are driving progress in docking methodologies. Boltz1x incorporates physical mechanisms to produce docking results that satisfy physical constraints without relying on relaxation. These findings highlight the critical role of combining physically informed generation with refinement procedures in docking pipelines, particularly when applied to drug design scenarios requiring atomic-level accuracy.

293 5 Conclusion

313

314

315

316

317

This paper proposed PoseX, a comprehensive benchmark for protein-ligand docking. Specifically, 294 we curated a new dataset with newly released protein-ligand complex crystal structures focusing 295 on both self-docking and cross-docking, and incorporated 23 docking methods across three main 296 research lines (physics-based methods, AI docking methods, and AI co-folding methods) to make an 297 exhaustive comparison. We also designed a novel relaxation module to refine the AI-generated binding 298 pose through energy minimization. Furthermore, we developed an online leaderboard that fosters 299 transparency and facilitates easy and fair comparisons for protein-ligand docking. By conducting 300 thorough empirical studies, we drew several key conclusions: (1) Both AI docking methods and AI 301 co-folding methods have outperformed physics-based methods in overall docking success rate. (2) 302 Most structural plausibility (except chirality) of AI-based approaches can be enhanced with relaxation, 303 which means combining AI modeling with physics-based post-processing may achieve excellent 304 performance. (3) Almost all the AI co-folding methods are plagued by ligand chirality, except for 305 Boltz-1x, which introduced a new inference time steering technique to fix hallucinations, pointing out the direction of incorporation of advantages of AI and physics. (4) Pocket information can be 307 leveraged adequately, especially for AI co-folding methods, to further promote the performance in 308 real-world applications. 309

310 6 Limitation and Future Work

Here, we briefly summarize the limitations of this work and present some directions for future research.

- 1. Evaluation on downstream tasks with binding affinities. While we focus on pose prediction and structural plausibility, binding affinity prediction remains an underexplored but complementary objective. Joint evaluation of structure and affinity on downstream tasks such as drug-target interaction and enzyme-substrate interaction would enable a more holistic assessment of docking algorithms and also remain an exciting direction for future research.
- 2. **Benchmarking on multi-ligand systems.** So far, most existing benchmarks focus on the evaluation of single-ligand docking, while multi-ligand docking is also practical in real-world applications such as enzyme engineering, where enzymes usually catalyze substrates together with co-factors. Thus, it is worth being assessed exhaustively in the future.
- 322 3. **Taking protein dynamics into account.** To date, existing studies always evaluate docking with rigid protein conformations, while integrating protein dynamics will better reflect the kinetic nature of biomolecular interactions in vivo. Future benchmarks could incorporate conformational ensembles of receptor structures to evaluate various models in a comprehensive way.

References

- [1] Alex Morehead, Nabin Giri, Jian Liu, Pawan Neupane, and Jianlin Cheng. Deep learning for protein-ligand docking: Are we there yet? *ArXiv*, pages arXiv–2405, 2025.
- Janani Durairaj, Yusuf Adeshina, Zhonglin Cao, Xuejin Zhang, Vladas Oleinikovas, Thomas Duignan, Zachary McClure, Xavier Robin, Daniel Kovtun, Emanuele Rossi, et al. Plinder: The protein-ligand interactions dataset and evaluation resource. *bioRxiv*, pages 2024–07, 2024.
- [3] Martin Buttenschoen, Garrett M Morris, and Charlotte M Deane. Posebusters: Ai-based docking methods fail to generate physically valid poses or generalise to novel sequences. *Chemical Science*, 15(9):3130–3139, 2024.
- Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, and Tommi Jaakkola. Diffdock: Diffusion steps, twists, and turns for molecular docking. *arXiv preprint arXiv:2210.01776*, 2022.
- [5] Richard A Friesner, Jay L Banks, Robert B Murphy, Thomas A Halgren, Jasna J Klicic, Daniel T
 Mainz, Matthew P Repasky, Eric H Knoll, Mee Shelley, Jason K Perry, et al. Glide: a new
 approach for rapid, accurate docking and scoring. 1. method and assessment of docking accuracy.
 Journal of medicinal chemistry, 47(7):1739–1749, 2004.
- If Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, pages 1–3, 2024.
- [7] Shravani S Pawar and Sachin H Rohane. Review on discovery studio: An important tool for
 molecular docking. 2021.
- Santiago Vilar, Giorgio Cozza, and Stefano Moro. Medicinal chemistry and the molecular operating environment (moe): application of qsar and molecular docking to drug discovery. *Current topics in medicinal chemistry*, 8(18):1555–1572, 2008.
- 9 Oleg Trott and Arthur J Olson. Autodock vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *Journal of computational chemistry*, 31(2):455–461, 2010.
- [10] Jerome Eberhardt, Diogo Santos-Martins, Andreas F Tillack, and Stefano Forli. AutoDock Vina
 1.2.0: New docking methods, expanded force field, and python bindings. *Journal of chemical information and modeling*, 61(8):3891–3898, 2021.
- Andrew T McNutt, Paul Francoeur, Rishal Aggarwal, Tomohide Masuda, Rocco Meli, Matthew Ragoza, Jocelyn Sunseri, and David Ryan Koes. Gnina 1.0: molecular docking with deep learning. *Journal of cheminformatics*, 13(1):43, 2021.
- Oscar Méndez-Lucio, Mazen Ahmad, Ehecatl Antonio del Rio-Chanona, and Jörg Kurt Wegner.
 A geometric deep learning approach to predict binding conformations of bioactive molecules.

 Nature Machine Intelligence, 3(12):1033–1039, 2021.
- [13] Hannes Stärk, Octavian Ganea, Lagnajit Pattanaik, Regina Barzilay, and Tommi Jaakkola.
 Equibind: Geometric deep learning for drug binding structure prediction. In *International conference on machine learning*, pages 20503–20521. PMLR, 2022.
- [14] Wei Lu, Qifeng Wu, Jixian Zhang, Jiahua Rao, Chengtao Li, and Shuangjia Zheng. Tankbind:
 Trigonometry-aware neural networks for drug-protein binding structure prediction. In *Advances* in neural information processing systems, volume 35, pages 7236–7249, 2022.
- Eric Alcaide, Zhifeng Gao, Guolin Ke, Yaqi Li, Linfeng Zhang, Hang Zheng, and Gengmo Zhou. Uni-mol docking v2: Towards realistic and accurate binding pose prediction. *arXiv* preprint arXiv:2405.11769, 2024.
- [16] Qizhi Pei, Kaiyuan Gao, Lijun Wu, Jinhua Zhu, Yingce Xia, Shufang Xie, Tao Qin, Kun He,
 Tie-Yan Liu, and Rui Yan. Fabind: Fast and accurate protein-ligand binding. Advances in
 Neural Information Processing Systems, 36:55963–55980, 2023.

- Gabriele Corso, Arthur Deng, Benjamin Fry, Nicholas Polizzi, Regina Barzilay, and Tommi Jaakkola. Deep confident steps to new pockets: Strategies for docking generalization. *ArXiv*, pages arXiv–2402, 2024.
- 377 [18] Michael Plainer, Marcella Toth, Simon Dobers, Hannes Stark, Gabriele Corso, Céline Marquet, 378 and Regina Barzilay. Diffdock-pocket: Diffusion for pocket-level docking with sidechain 379 flexibility. 2023.
- Wei Lu, Jixian Zhang, Weifeng Huang, Ziqiao Zhang, Xiangyu Jia, Zhenyu Wang, Leilei
 Shi, Chengtao Li, Peter G Wolynes, and Shuangjia Zheng. Dynamicbind: predicting ligand specific protein-ligand complex structure with a deep equivariant generative model. *Nature Communications*, 15(1):1071, 2024.
- Houtim Lai, Longyue Wang, Ruiyuan Qian, Junhong Huang, Peng Zhou, Geyan Ye, Fandi Wu, Fang Wu, Xiangxiang Zeng, and Wei Liu. Interformer: an interaction-aware model for protein-ligand docking and affinity prediction. *Nature Communications*, 15(1):10223, 2024.
- Duanhua Cao, Mingan Chen, Runze Zhang, Zhaokun Wang, Manlin Huang, Jie Yu, Xinyu Jiang, Zhehuan Fan, Wei Zhang, Hao Zhou, et al. Surfdock is a surface-informed diffusion generative model for reliable and accurate protein–ligand complex prediction. *Nature Methods*, pages 1–13, 2024.
- Zhuoran Qiao, Weili Nie, Arash Vahdat, Thomas F Miller III, and Animashree Anandkumar.
 State-specific protein-ligand complex structure prediction with a multiscale deep generative
 model. *Nature Machine Intelligence*, 6(2):195–208.
- Rohith Krishna, Jue Wang, Woody Ahern, Pascal Sturmfels, Preetham Venkatesh, Indrek Kalvet, Gyu Rie Lee, Felix S Morey-Burrows, Ivan Anishchenko, Ian R Humphreys, et al. Generalized biomolecular modeling and design with rosettafold all-atom. *Science*, 384(6693):eadl2528, 2024.
- Chai Discovery, Jacques Boitreaud, Jack Dent, Matthew McPartlon, Joshua Meier, Vinicius Reis, Alex Rogozhnikov, and Kevin Wu. Chai-1: Decoding the molecular interactions of life.
 bioRxiv, pages 2024–10, 2024.
- [25] Jeremy Wohlwend, Gabriele Corso, Saro Passaro, Mateo Reveiz, Ken Leidal, Wojtek Swiderski,
 Tally Portnoi, Itamar Chinn, Jacob Silterra, Tommi Jaakkola, et al. Boltz-1: Democratizing
 biomolecular interaction modeling. *bioRxiv*, pages 2024–11, 2024.
- 404 [26] ByteDance AML AI4Science Team, Xinshi Chen, Yuxuan Zhang, Chan Lu, Wenzhi Ma, Jiaqi
 405 Guan, Chengyue Gong, Jincai Yang, Hanyu Zhang, Ke Zhang, et al. Protenix-advancing
 406 structure prediction through a comprehensive alphafold3 reproduction. bioRxiv, pages 2025–01,
 407 2025.
- Isabella A Guedes, Camila S de Magalhães, and Laurent E Dardenne. Receptor–ligand molecular
 docking. *Biophysical reviews*, 6:75–87, 2014.
- Rommie E Amaro, Riccardo Baron, and J Andrew McCammon. An improved relaxed complex scheme for receptor flexibility in computer-aided drug design. *Journal of computer-aided molecular design*, 22:693–705, 2008.
- [29] Peter Eastman, Jason Swails, John D. Chodera, Robert T. McGibbon, Yutong Zhao, et al.
 Openmm 7: Rapid development of high performance algorithms for molecular dynamics. *PLOS Computational Biology*, 13(7):e1005659, 2017. doi: 10.1371/journal.pcbi.1005659.
- 416 [30] Open Force Field Consortium. Open force field initiative: Openff toolkit 2.1.0, December 2024.
 417 URL https://github.com/openforceff-toolkit.
- In Sameer Kawatkar, Hongming Wang, Ryszard Czerminski, and Diane Joseph-McCarthy. Virtual fragment screening: an exploration of various docking and scoring protocols for fragments using glide. *Journal of computer-aided molecular design*, 23(8):527–539, 2009.

- 421 [32] Michael J Hartshorn, Marcel L Verdonk, Gianni Chessari, Suzanne C Brewerton, Wijnand TM
 422 Mooij, Paul N Mortenson, and Christopher W Murray. Diverse, high-quality test set for the
 423 validation of protein- ligand docking performance. *Journal of medicinal chemistry*, 50(4):
 424 726–741, 2007.
- [33] Peter W Rose, Andreas Prlić, Ali Altunkaya, Chunxiao Bi, Anthony R Bradley, Cole H Christie,
 Luigi Di Costanzo, Jose M Duarte, Shuchismita Dutta, Zukang Feng, et al. The RCSB protein
 data bank: integrative view of protein, gene and 3D structural information. *Nucleic acids* research, page gkw1000, 2016.
- 429 [34] Jas Bhachoo and Thijs Beuming. Investigating protein–peptide interactions using the schrödinger computational suite. *Modeling peptide-protein interactions: methods and protocols*, pages 235–254, 2017.
- [35] Andrew T McNutt, Yanjing Li, Rocco Meli, Rishal Aggarwal, and David Ryan Koes. Gnina
 1.3: the next increment in molecular docking with deep learning. *Journal of Cheminformatics*,
 17(1):28, 2025.
- [36] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al.
 Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- Ignormal factories and Francisco and Francis
- 440 [38] Peter Eastman et al. Pdbfixer: Automated protein structure repair tool. https://openmm. 441 org/pdbfixer, 2012-2025. URL https://github.com/openmm/pdbfixer. Accessed: 442 2025-04-12.
- 443 [39] RDKit Contributors. Rdkit: Open-source cheminformatics software. https://www.rdkit.org, 2006-2024. Accessed: 2025-04-12.
- [40] Junmei Wang, Romain M. Wolf, James W. Caldwell, Peter A. Kollman, and David A. Case.
 Development and testing of a general amber force field. *Journal of Computational Chemistry*,
 25(9):1157–1174, 2004. doi: 10.1002/jcc.20035.

NeurIPS Paper Checklist

1. Claims

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

Answer: [Yes]

Justification: The claims can be found in Abstract and Introduction.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
 contributions made in the paper and important assumptions and limitations. A No or
 NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

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

Answer: [Yes]

Justification: The limitations can be found in Limitation and Future Work.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was
 only tested on a few datasets or with a few runs. In general, empirical results often
 depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

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

Answer: [NA]

Justification: This paper does not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

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

Answer: [Yes]

Justification: All the correlated information can be found in Appendix B, C and the github repository mentioned in Abstract.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

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

554 Answer: [Yes]

Justification: We have released the dataset and code through github repository that can be found in Abstract.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be
 possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
 including code, unless this is central to the contribution (e.g., for a new open-source
 benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

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

Answer: [Yes]

Justification: This paper focuses on benchmarking of various docking approaches, thus the work has nothing to do with training, but only testing or evaluation. The setup details of evaluation can be found in Appendix B, C and the github repository mentioned in Abstract.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

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

Answer: [No]

Justification: For all docking methods, we selected the best predictions given by the models for evaluation.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how
 they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

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

Answer: [Yes]

605

606

607

608 609

610

611

612

613

616

617

618

619

620 621

622

623 624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

646

647

649

650

651

652

653

654

655

656

Justification: The information of computation resources can be found in Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

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

Answer: [Yes]

Justification: We follow close to the line of the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
 deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

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

Answer: [Yes]

Justification: The description can be found in Introduction.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

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

Answer: [NA]

Justification: This paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
 not require this, but we encourage authors to take this into account and make a best
 faith effort.

12. Licenses for existing assets

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

Answer: [Yes]

Justification: The related information can be found in Table 2.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

 If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

709

710

711

712

713

714

715

716 717

718

720

721

722 723

724

725

726

727

729

730

731

732

733

734

735

736

737

738

739 740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

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

Answer: [Yes]

Justification: We have released the dataset and code through github repository that can be found in Abstract alongside the instructions.

Guidelines

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

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

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

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

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

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

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

765 Answer: [NA]

Justification: We did not use LLMs for studies in this paper.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

A Dataset Construction and Statistical Analysis

A.1 Dataset Construction Process

Table S1: Selection process of the PDB entries and ligands for the PoseX Self-Docking (PoseX-SD).

Selection Step	# proteins (unique PDB IDs)	# ligands (unique CCD IDs)
PDB entries released from January 1st, 2022 to January 1st, 2025	13207	6877
feature a refinement resolution of 2 Å or better and include at least one		
protein and one ligand		
Remove unknown ligands (e.g., UNX, UNL)	13202	6875
Remove proteins with a sequence length greater than 2000	11771	6442
Ligands weighing from 100 Da to 900 Da	9768	6196
Ligands with at least 3 heavy atoms	9706	6163
Ligands containing only H, C, O, N, P, S, F, Cl atoms	9030	5741
Ligands that are not covalently bound to protein	8383	5185
Structures with no unknown atoms (e.g., element X)	8349	5166
Ligand real space R-factor is at most 0.2	7521	4476
Ligand real space correlation coefficient is at least 0.95	5734	3426
Ligand model completeness is 100%	5645	3358
Ligand starting conformation could be generated with ETKDGv3	5638	3351
All ligand SDF files can be loaded with RDKit and pass its sanitization	5634	3345
PDB ligand report does not list stereochemical errors	5600	3317
PDB ligand report does not list any atomic clashes	3971	2541
Select single protein-ligand conformation ¹	3971	2541
Intermolecular distance between the ligand(s) and the protein is at least 0.2 Å	3945	2527
Intermolecular distance between ligand(s) and other small organic molecules is at least 0.2 Å	3889	2477
Intermolecular distance between ligand(s) and ion metals in complex is at least $0.2~\textrm{Å}$	3889	2477
Remove ligands which are within 5.0 Å of any protein symmetry mate	2451	1598
Get a set with unique pdbs and unique ccds by Hopcroft–Karp matching algorithm	1587	1587
Select representative PDB entries by clustering protein sequences	718	718

¹ The first conformation is chosen when multiple conformations are available in the PDB entry.

² Clustering with MMseqs2 is done with a sequence identity threshold of 0% and a minimum coverage of 100%.

Table S2: Selection process of the PDB entries and ligands for the PoseX Cross-Docking (PoseX-CD).

Selection Step	# proteins (unique PDB IDs)	# ligands (unique CCD IDs)
PDB entries released from January 1st, 2022 to January 1st, 2025	13207	6877
feature a refinement resolution of 2 Å or better and include at least one		
protein and one ligand		
Remove unknown ligands (e.g., UNX, UNL)	13202	6875
Remove proteins with a sequence length greater than 2000	11771	6442
Ligands weighing from 100 Da to 900 Da	9768	6196
Ligands with at least 3 heavy atoms	9706	6163
Ligands containing only H, C, O, N, P, S, F, Cl atoms	9030	5741
Ligands that are not covalently bound to protein	8383	5185
Structures with no unknown atoms (e.g., element X)	8349	5166
Ligand real space R-factor is at most 0.2	7521	4476
Ligand real space correlation coefficient is at least 0.95	5734	3426
Ligand model completeness is 100%	5645	3358
Ligand starting conformation could be generated with ETKDGv3	5638	3351
All ligand SDF files can be loaded with RDKit and pass its sanitization	5634	3345
PDB ligand report does not list stereochemical errors	5600	3317
PDB ligand report does not list any atomic clashes	3971	2541
Select single protein-ligand conformation ¹	3971	2541
Intermolecular distance between the ligand(s) and the protein is at least $0.2~\textrm{Å}$	3945	2527
Intermolecular distance between the ligand(s) and the other ligands is at least 5.0 $\hbox{\normalfont\AA}$	2232	1536
Remove ligands which are within 5.0 Å of any protein symmetry mate	1240	908
Cluster proteins that have at least 90% sequence identity ²	890	708
Structures can be successfully aligned to the reference structure in each cluster ³	371	362

¹ The first conformation is chosen when multiple conformations are available in the PDB entry.

774 A.2 Statistical Characteristics

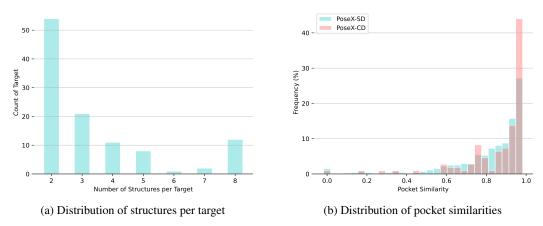


Figure S1: (a) The distribution of structures per target shows that every protein adopts at least two distinct conformations, and about half of the targets are represented by just two. (b) The distribution of pocket similarities.

² Clustering with MMseqs2 is done with a sequence identity threshold of 90% and a minimum coverage of 80%.

³ Each candidate protein is structurally aligned to the reference protein via the superposition of C_{α} atom of amino acid residues using PyMOL. A candidate PDB entry is removed if the RMSD of the protein alignment is greater than 2.0 Å and a candidate ligand is removed if it is 4.0 Å away from the reference ligand.

775 B Docking Methods and Evaluation Settings

This section presents the docking methods employed in our evaluation and illustrates the corresponding setups.

778 B.1 Physics-based Methods

Physics-based methods employ physical forces and geometric complementarity to model molecular 779 interactions, predicting ligand binding to the target protein. Usually, the atomic coordinates of the 780 protein's binding site remain fixed, while the ligand undergoes flexible conformational changes. This 781 schema reduces the computational complexity of docking simulations by neglecting the dynamic 782 flexibility of the protein structure. However, although computationally efficient, this method may 783 fail to fully account for the inherent flexibility of proteins, as biological systems often exhibit 784 protein conformational changes upon ligand binding. We include 5 physics-based methods in this 785 paper, including Discovery Studio [7], Schrödinger Glide [5], MOE [8], AutoDock Vina [9, 10] and 786 GNINA [11]. 787

788 B.1.1 Schrödinger Glide

Schrödinger Glide is a leading provider of biomolecular simulation software, and Glide is one of its flagship products, focusing on precise molecular docking simulations [5, 34]. Glide adopts a unique hierarchical docking approach, starting with coarse screening and then performing fine optimization on high-scoring results to improve prediction accuracy.

793 **Software Version**: Schrödinger Suite 2022-1, Build 141

794 Docking Workflow

795

796

797

798

799

800

801 802

803

804

814

815

- 1. Use **PrepWizard** to preprocess the protein files by adding hydrogens and optimizing with the OPLS3 force field at pH 7.4.
- Use LigPrep to preprocess small molecules, preserving the chirality of the input ligand. Use Epik to predict the pKa and protonation states of small molecules at pH 7.0. Optimize the small-molecule conformations using the S-OPLS force field, and output one small-molecule conformation as the input for docking.
- 3. Define the INNERBOX dimensions as $10 \times 10 \times 10$ Å, and the OUTERBOX dimensions as:

$$\begin{pmatrix} \mathsf{Size}_x \\ \mathsf{Size}_y \\ \mathsf{Size}_z \end{pmatrix} = \begin{pmatrix} x_{\max} - x_{\min} + 20 \\ y_{\max} - y_{\min} + 20 \\ z_{\max} - z_{\min} + 20 \end{pmatrix}$$

The force field is set to OPLS3, and all other parameters are set by default. Generate a grid file.

4. Perform molecular docking using Glide SP (Standard Precision), and output one small molecule pose as the docking result.

805 **Runtime Environment**: Run on an Intel i9-10920X CPU using 16 cores.

806 B.1.2 Discovery Studio

Discovery Studio [7], developed by Dassault Systèmes BIOVIA, is a comprehensive life sciences research platform that covers molecular modeling, virtual screening, and more. For protein-ligand binding, Discovery Studio performs conformational sampling around a given binding site and ranks potential poses using physics-based scoring functions like CDOCKER (which combines grid-based molecular dynamics and CHARMM force fields).

812 **Software Version**: v2021.1.0.20298.

B Docking Workflow:

1. Use the **Proteins Preparation** components in Discovery Studio to process the protein files. The protein was protonated at pH 7.4 with a solvent ionic strength of 0.145 M. Minimization

- was performed using the **CHARMm** force field to optimize the protein structure, and all other parameters are set by default.
 - 2. Use the **Ligands Preparation** components in Discovery Studio to process the ligand files. Enumerate ionization states for each ligand within a pH range of 6.5-8.5. Enumerate automeric forms for each ligand with a maximum of 10 tautomers per ligand. Fix the bad valencies by adjusting formal charges, and all other parameters are set by default.
 - Dock the prepared proteins and the corresponding prepared ligands using the CDOCKER components in Discovery Studio. The docking site was centered at:

$$\begin{pmatrix} x_c \\ y_c \\ z_c \end{pmatrix} = \begin{pmatrix} \frac{x_{\text{max}} + x_{\text{min}}}{2} \\ \frac{y_{\text{max}} + y_{\text{min}}}{2} \\ \frac{z_{\text{max}} + z_{\text{min}}}{2} \end{pmatrix}$$

Define the binding sphere radius as:

$$R = \max\{(x_{\max} - x_{\min}), (y_{\max} - y_{\min}) - (z_{\max} - z_{\min})\} + 20$$

The docking simulations were performed using the **CHARMm** force field. Assign the partial charges to the ligands via the **Momany-Rone** method, and all other parameters are set by default. 10 top docking poses output each docking run, and the best-scored pose was selected as the final docking result.

Runtime Environment: Run on an Intel Ultra 5 125H CPU using 14 cores.

827 B.1.3 Molecular Operating Environment (MOE)

- Molecular Operating Environment (MOE) [8], developed by the Canadian company Chemical Computing Group, is a commercial drug discovery software platform that combines visualization, modeling, simulations, and methodology development into a single, unified package.
- 831 Software Version: MOE 2024.06.

832 Docking Workflow

818

819

820

821

822

823

824

825

833

834

835

836 837

838

839

840

- 1. An SVL script automates the docking pipeline.
- 2. The **StructurePreparation** function is employed to preprocess protein structures.
 - 3. The binding site is defined by reference ligands.
- 4. The **Triangle Matcher** algorithm is utilized to generate initial ligand poses.
 - 5. The scoring function is configured as **London dG**, with a maximum of 30 poses generated.
 - Poses are refined using a fixed receptor, optimizing only the ligand's position and conformation, with the re-scoring function configured as GBVI/WSA dG and a maximum of 5 poses retained.
- Runtime Environment: Run on an AMD EPYC 9554 CPU.

842 B.1.4 AutoDock Vina

- AutoDock Vina [10] is one of the fastest and most widely used **open-source** molecule docking programs. It combines global search (to identify potential binding modes) with local optimization (to refine these modes).
- 846 Software Versions
- AutoDock-Vina: 1.2.6
- MGLTools: 1.5.7
- Reduce: 4.14.230914
- OpenBabel: 3.1.0
- Meeko: 0.6.1

852 **Docking Workflow**

- 1. Use **Reduce** to add polar hydrogens to the protein structure.
- Use OpenBabel to add non-polar hydrogens and normalize atom names, exporting the protein in a format recognizable by MGLTools.
- Use the receptor_prepare4.py script from MGLTools to convert the hydrogen-added protein PDB file into a PDBQT file.
- 4. Use **OpenBabel** to add hydrogens to the ligand molecule at pH 7.4.
- Use the mk_prepare_ligand.py script from Meeko to convert the hydrogen-added ligand SDF file into a PDBQT file.
- 6. Define the docking box center and size as follows:

$$\begin{pmatrix} x_c \\ y_c \\ z_c \end{pmatrix} = \begin{pmatrix} \frac{x_{\max} + x_{\min}}{2} \\ \frac{y_{\max} + y_{\min}}{2} \\ \frac{z_{\max} + z_{\min}}{2} \end{pmatrix}$$

$$\begin{pmatrix} \operatorname{Size}_x \\ \operatorname{Size}_y \\ \operatorname{Size}_z \end{pmatrix} = \begin{pmatrix} x_{\max} - x_{\min} + 20 \\ y_{\max} - y_{\min} + 20 \\ z_{\max} - z_{\min} + 20 \end{pmatrix}$$

- 7. Perform molecular docking using the prepared protein and ligand PDBQT files.
- 8. Use **vina_split** to split the output file, extract the best-scored pose for each ligand, and convert the resulting PDBQT file into an SDF file using Meeko for the final output.

Runtime Environment: Run on an AMD EPYC 9554 CPU, with no specified core limit and up to 256 cores available.

866 B.1.5 GNINA

853

854

855

856

857

858

859 860

861

862

863

GNINA [11, 35] is a relatively new project that introduces DL techniques into the field of molecular docking, particularly leveraging convolutional neural networks (CNNs) as scoring functions to improve docking scoring. It is an **open-source** software.

B70 Docker Image: https://hub.docker.com/layers/gnina/gnina/latest/images

Running Parameters: The command used is:

```
gnina -r rec.pdb -l lig.sdf -autobox_ref.sdf -o out.sdf,
```

where lig.sdf is PDB_CCD_ligand_start_conf.sdf and ref.sdf is PDB_CCD_ligand.sdf.

Runtime Environment: Run on Nvidia A6000 GPU.

873 B.2 AI Docking Methods

Al docking methods utilize SMILES strings of ligands and three-dimensional structures of protein 874 targets as input to predict energetically favorable ligand conformations bound to target proteins. 875 These methods systematically explore the conformational space of small molecules to identify 876 low-energy configurations that optimize the binding affinity to proteins. By sampling diverse 877 ligand conformations, AI docking methods enhance the optimization of spatial arrangements to 878 maximize interactions with protein active sites, including hydrogen bonds, hydrophobic interactions, 879 and electrostatic complementarity. We involve 11 AI docking methods in this paper, including 880 DeepDock [12], EquiBind [13], TankBind [14], DiffDock [4], Uni-Mol [15], FABind [16], DiffDock-881 L [17], DiffDock-Pocket [18], DynamicBind [19], Interformer [20] and SurfDock [21]. 882

883 B.2.1 DeepDock

DeepDock [12] is a geometric DL model that learns a statistical potential based on the distance likelihood.

886 GitHub Repository: https://github.com/OptiMaL-PSE-Lab/DeepDock

887 **GitHub Commit Hash**: ab1e45044c5e0a69105b48d09ea984c6a5ebc26c

888 **Running Parameters**: Default parameters are used in evaluation.

Runtime Environment: Run on Intel(R) Xeon(R) CPU E5-2620 v4.

890 B.2.2 EquiBind

- 891 **EquiBind** [13] is an SE(3)-equivariant geometric DL model designed for direct-shot prediction of
- both i) the receptor binding site (blind docking) and ii) the ligand's bound pose and orientation.
- 893 GitHub Repository: https://github.com/HannesStark/EquiBind
- 894 **GitHub Commit Hash**: 41bd00fd6801b95d2cf6c4d300cd76ae5e6dab5e
- 895 **Running Parameters**: Default parameters are used in evaluation.
- 896 **Runtime Environment**: Run on Nvidia A6000 GPU.

897 B.2.3 TankBind

- 898 TankBind [14] incorporates trigonometric constraints as a robust inductive bias into the model, and
- explicitly examines all potential binding sites for each protein by dividing the entire protein into
- functional blocks. establishes an efficient diffusion process within this space.
- 901 GitHub Repository: https://github.com/luwei0917/TankBind
- 902 **GitHub Commit Hash**: ff85f511db11d7a3e648d2e01cd6fdb4f9823483
- 903 **Running Parameters**: Use the structure of the entire protein as input for prediction, rather than
- ochains within 10Å of the ligand in the default setting.
- 905 **Runtime Environment**: Run on an AMD EPYC 9554 CPU.

906 B.2.4 DiffDock

- 907 **DiffDock** [4] is a diffusion-based generative model defined on the non-Euclidean manifold of ligand
- poses. It maps this manifold to the product space of the degrees of freedom (translational, rotational,
- and torsional) relevant to docking and establishes an efficient diffusion process within this space.
- 910 GitHub Repository: https://github.com/gcorso/DiffDock
- 911 **GitHub Commit Hash**: bc6b5151457ea5304ee69779d92de0fded599a2c
- 912 **Running Parameters**: Default parameters are used in evaluation.
- 913 **Runtime Environment**: Run on Nvidia A800 GPU.

914 B.2.5 DiffDock-L

- 915 **DiffDock-L** [17] is a variant of DiffDock that scales up data and model size by integrating synthetic
- 916 data strategies.
- 917 GitHub Repository: https://github.com/gcorso/DiffDock
- 918 **GitHub Commit Hash**: b4704d94de74d8cb2acbe7ec84ad234c09e78009
- Running Parameters: samples_per_complex is changed from the default value of 10 to 40.
- 920 Runtime Environment: Run on Nvidia A800 GPU.

921 B.2.6 DiffDock-Pocket

- 922 **DiffDock-Pocket** [18] is a variant of DiffDock with additional binding pocket specification.
- 923 GitHub Repository: https://github.com/plainerman/DiffDock-Pocket
- 924 **GitHub Commit Hash**: 3902bdd4d42ee5254d37aa694d005a992c92ad93
- 925 **Running Parameters**: Default parameters are used in evaluation.
- Runtime Environment: Run on Nvidia A6000 GPU.

7 B.2.7 DynamicBind

- 928 **DynamicBind** [19] utilizes equivariant geometric diffusion networks to generate a smooth energy
- 929 landscape, facilitating efficient transitions between various equilibrium states. DynamicBind accu-
- 930 rately identifies ligand-specific conformations from unbound protein structures, eliminating the need
- 931 for holo-structures or extensive sampling.
- 932 GitHub Repository: https://github.com/luwei0917/DynamicBind
- 933 **GitHub Commit Hash**: abdcd83f313cd20d50c3917e04615e989a8f63e5
- 934 **Running Parameters**: Default parameters are used in evaluation.
- 935 **Runtime Environment**: Run on Nvidia A800 GPU.

936 B.2.8 FABind

- 937 **FABind** [16] is an end-to-end model that integrates pocket prediction and docking to achieve precise
- and efficient protein-ligand binding predictions. It involves a ligand-informed pocket prediction
- module, which is also utilized to enhance the accuracy of docking pose estimation.
- 940 **GitHub Repository**: https://github.com/gcorso/DiffDock
- 941 **GitHub Commit Hash**: bc6b5151457ea5304ee69779d92de0fded599a2c
- **Running Parameters**: Default parameters are used in evaluation.
- 943 **Runtime Environment**: Run on Nvidia A800 GPU.

944 B.2.9 Uni-Mol

- 945 Uni-Mol [15] represents Uni-Mol Docking v2. It combines the pretrained molecular and pocket
- models to learn the distance matrix, and then uses a coordinate model to predict the final coordinates
- 947 of the molecule.
- 948 GitHub Repository: https://github.com/deepmodeling/Uni-Mol/tree/main/unimol_
- 949 docking_v2
- 950 **GitHub Commit Hash**: c0365df6535b90197246399417a9b21250268352
- 951 Running Parameters: Default parameters are used in prediction. About one-fifth of the molecules in
- the model output will encounter RDKit's sanitization check errors. This issue is resolved by reading
- in the correct molecular topology and then assigning the coordinates predicted by Uni-Mol to the
- 954 molecules with the new topology
- 955 **Runtime Environment**: Run on Nvidia A6000 GPU.

956 B.2.10 Interformer

- 957 **Interformer** [20], a unified model based on the Graph-Transformer architecture, is specifically
- 958 designed to capture non-covalent interactions using an interaction-aware mixture density network.
- 959 Furthermore, it implements a negative sampling strategy to effectively adjust the interaction distribu-
- 960 tion, enhancing affinity prediction accuracy.
- 961 GitHub Repository: https://github.com/tencent-ailab/Interformer
- 962 **GitHub Commit Hash**: 8cced9b8a5d8c887787a8c8731d9f087563d4c7e
- 963 Running Parameters: Use PDB_CCD_ligand.sdf to obtain the pocket, perform UFF opti-
- 964 mization on PDB_CCD_ligand_start_conf.sdf and replace it in the uff folder, and use the
- 965 -uff_as_ligand option during prediction.
- 966 **Runtime Environment**: Run on Nvidia A6000 GPU.

967 B.2.11 SurfDock

- 968 SurfDock [21] combines protein sequences, three-dimensional structural graphs, and surface-level
- features within an equivariant architecture. It leverages a generative diffusion model on a non-

- 970 Euclidean manifold to optimize molecular translations, rotations, and torsions, producing accurate
- 971 and reliable binding poses.
- 972 **GitHub Repository**: https://github.com/CAODH/SurfDock
- 973 GitHub Commit Hash: 2f0422f6ddcfdfefc3fa61ef12a1d6406a589bce
- 974 **Running Parameters**: Default parameters are used in evaluation.
- 975 **Runtime Environment**: Run on Nvidia A6000 GPU.

976 B.3 AI Co-folding Methods

- 977 AI co-folding methods represent a significant advance in computational biology by simultaneously
- 978 predicting the conformation of both the protein and its associated ligand, which sets them apart from
- 979 physics-based methods and AI docking methods. In contrast to physics-based methods, which typically
- assume a fixed protein structure and focus on optimizing ligand placement, or AI docking methods
- 981 that may still rely on predefined protein conformations, AI co-folding methods adopt a more holistic
- 982 strategy-taking only the protein's amino acid sequence and ligand's SMILES strings as input.
- These methods aim to capture the dynamic interaction between proteins and ligands by predicting
- their structures in tandem, enabling a more accurate representation of how these molecules interact in
- biological systems. In this paper, we involve 7 AI co-folding methods, including NeuralPLexer [22],
- 986 RoseTTAFold-All-Atom (RFAA) [23], AlphaFold3 [6], Chai-1 [24], Boltz-1 [25], Boltz-1x [25] and
- Protenix [26]. It should be noted that in our evaluation of AI co-folding methods, we did not consider
- 988 post-translational modifications and used unmodified protein sequences as input.

989 B.3.1 NeuralPLexer

- 990 **NeuralPLexer** [22] is a physics-inspired flow-based generative model for biomolecular complex
- 991 structure prediction based on sequences only. NeuralPLexer combines a protein language model to
- 992 learn sequence information and graph encoding to represent 3D molecular structure and bioactivity
- 993 information.
- 994 GitHub Repository:https://github.com/zrqiao/NeuralPLexer
- 995 **GitHub Commit Hash**: 2c52b10d3094e836661dfecfa3be76f47dcdea7e
- 996 **Running Parameters**: Default parameters are used in evaluation.
- 997 **Runtime Environment**: Run on Nvidia A6000 GPU.

998 B.3.2 RoseTTAFold-All-Atom

- 999 RoseTTAFold-All-Atom (RFAA) [23] is a generalized foundation model for all-atom biomolec-
- ular structure prediction and design, including protein, nucleic acid, and other small molecules.
- 1001 RoseTTAFold-All-Atom is a 3-track based architecture incorporating equivariant neural networks for
- all atomic structure prediction. Meanwhile, it integrates with RFDiffusion for molecular design.
- 1003 GitHub Repository: https://github.com/baker-laboratory/RoseTTAFold-All-Atom
- 1004 **GitHub Commit Hash**: 6c8514053acf76da0f9edde2aa51b40abff68fa1
- 1005 Running Parameters: Default parameters are used in evaluation.
- 1006 Runtime Environment: Run on Nvidia A800 GPU.

1007 B.3.3 AlphaFold3

- AlphaFold3 [6], developed by DeepMind, represents the latest advancement in protein structure
- prediction technology. Building on the successes of its predecessor AlphaFold 2 [36]), AlphaFold3
- adopts a diffusion model instead of a structure module in AlphaFold2, not only improving the
- accuracy of protein folding but also supporting the structure prediction of complexes (e.g., protein-
- 1012 RNA, protein-ligand), which enables its usage in protein-ligand docking.
- o13 **Software Version**: 3.0.0

- Running Parameters: Except for the number of seeds being set to 1, the rest of the predictions are
- made using the default parameters. We finally select the top 1 result for evaluation.
- 1016 **Runtime Environment**: Run on Nvidia A800 GPU.

1017 B.3.4 Chai-1

- 1018 Chai-1 [37] is a multimodal molecular foundation model that can also predict structures with a single
- 1019 sequence. By leveraging the decoder-only Transformer framework, which is widely used in Large
- Language Models (LLM) like GPT, Chai-1 encodes sequential information without database search.
- Moreover, Chai-1 accepts various chemical or biological constraint features as input to predict more
- 1022 accurate molecular structures.
- 1023 **Software Version**: 0.5.2
- 1024 **Running Parameters**: Use the online MSA server to obtain MSA information, keep the rest as
- default settings, and select the top 1 result for evaluation.
- 1026 **Runtime Environment**: Run on Nyidia A800 GPU.

1027 B.3.5 Boltz-1

- Boltz-1 [25] aims at reproducing AlphaFold3 and releasing all the codes (model architecture, train-
- ing, inference), which achieves competitive performance. Additionally, Boltz-1 introduces several
- architectural innovations, including a novel reverse diffusion process and a revamped confidence
- model, enhancing its predictive accuracy and robustness.
- 1032 **Software Version**: 0.4.0
- 1033 Running Parameters: Use the MSA online server to obtain MSA information, set diffusion samples
- to 5, and select the top 1 result for evaluation.
- 1035 **Runtime Environment**: Run on Nvidia A800 GPU.

1036 B.3.6 Boltz-1x

- Boltz-1x [25] is an advanced version of the Boltz-1 model. It introduces a novel inference-time
- steering technique, which enhances the physical quality of predicted poses by reducing hallucinations
- and non-physical predictions. This ensures more reliable and biologically plausible structures.
- 1040 **Software Version**: 1.0.0
- **Running Parameters:** Use the MSA online server to obtain MSA information, set diffusion samples
- to 5, and select the top 1 result for evaluation.
- 1043 **Runtime Environment**: Run on Nvidia A800 GPU.

1044 B.3.7 Protenix

- Protenix [26] is a comprehensive and open-source reproduction of AlphaFold3, developed by
- 1046 ByteDance. It introduces several architectural innovations, including a modular PyTorch framework
- that facilitates full training and inference, and optimizations such as custom CUDA kernels and BF16
- training to enhance computational efficiency.
- 1049 **Software Version**: 0.4.2
- Running Parameters: Use the MSA online server to obtain MSA information, and the seed is set to
- 1051 101.
- 1052 **Runtime Environment**: Run on Nvidia A6000 GPU.

B.4 Training Data Cutoff Times

Method	Training Data Cutoff Time	
Traditional physics-based methods		
Discovery Studio [7]	N/A	
Schrödinger Glide [5]	N/A	
MOE [8]	N/A	
AutoDock Vina [9, 10]	N/A	
GNINA [11]	2018-12	
AI docking methods		
DeepDock [12]	2018-12	
EquiBind [13]	2019-12	
TankBind [14]	2018-12	
DiffDock [4]	2018-12	
Uni-Mol [15]	2019-12	
FABind [16]	2018-12	
DiffDock-L [17]	2018-12	
DiffDock-Pocket [18]	2019-12	
DynamicBind [19]	2018-12	
Interformer [20]	2019-12	
SurfDock [21]	2019-12	
AI co-fo	olding methods	
NeuralPLexer [22]	2018-12	
RoseTTAFold-All-Atom [
AlphaFold3 [6]	2021-10	
Chai-1 [24]	2021-02	
Boltz-1 [25]	2021-10	
Boltz-1x [25]	2021-10	
Protenix [26]	2021-10	

4 C Technical Details of Relaxation Process

1059

1060

1061

1062

1063

1064

1065

1068

1069

1070

1071

1072

1073

1075

1076

1077

1080

1081

1082

- Our relaxation is based on the following software: OpenMM 7.7 [29], PDBFixer 1.8 [38], RDKit 2023.09 [39], AmberTools 23, and OpenFF 2.1.0 [30]. It contains the following essential steps:
- Structure preprocessing and integrity restoration. Use PDBFixer (v1.8) to handle the initial structure files:
 - Parse complete protein sequence information from CIF files, retaining water molecules and metal ions within a 5 Å range of the ligand in AI-predicted models.
 - Standardize non-canonical amino acids to canonical forms (e.g., SEP to SER), simultaneously correcting the protein sequence database.
 - Detect structural deficiencies using the findMissingResidues/findMissingAtoms algorithms, and apply the AddMissingAtoms module to complete atoms (including N-terminal ACE and C-terminal NME capping).
- Molecular topology construction and validation. To address the lack of bond order information in PDBFixer:
 - Integrate Amber ff14SB force field atom types and topology bond parameters to establish bond order matching rules.
 - Build a molecular graph model with RDKit (v2023.09) and perform SanitizeMol standardization checks (including charge correction and stereochemistry validation).
 - Apply the RDKit AddHs module for protonation, optimizing the spatial arrangement of hydrogen atoms.
- Force field parameterization. Employ a multi-scale force field combination strategy:
 - For protein systems: Generate Amber ff14SB force field parameters using OpenMM 7.7.
 - For ligand systems: Perform GAFF-2.11 [40] parameterization using the OpenFF 2.1.0 toolkit, including mmff94s charge calculations and XML topology generation.
- Constrained molecular dynamics optimization. Implement energy minimization on the OpenMM 7.7 platform [29]:
 - Constraints: Apply additional forces $(0.5 * k * ((x x_0)^2 + (y y_0)^2 + (z z_0)^2)$ (where $k = 10, x_0, y_0, z_0$ are original 3D coordinate) to constrain backbone atomic positions in the protein structure, keeping newly added atoms free.
- Integration parameters: Langevin thermostat (300 K, friction coefficient 1 ps^{-1}), time step 0.004 ps.
- Convergence criteria: Energy gradient convergence threshold ≤ 10 kJ/mol/nm.

D Description of Validity

The validity checks for the structures analyzed in this study were conducted using PoseBuster [3], a tool to ensure the reliability and accuracy of the molecular poses. The validation process encompasses chemical validity and consistency, intramolecular validity, and intermolecular validity, each assessed with specific criteria as detailed below. In this study, we define **structural plausibility** as stereochemical correctness and intra- and intermolecular validity.

D.1 Chemical Validity and Consistency

- File loads: The input molecule can be successfully loaded into a molecule object by RDKit.
- Sanitisation: The input molecule passes RDKit's chemical sanitisation checks, ensuring it adheres to basic chemical rules.
- Molecular formula: The molecular formula of the input molecule is identical to that of the true molecule.
- Bonds: The bonds in the input molecule are the same as in the true molecule.
- **Tetrahedral chirality**: The specified tetrahedral chirality in the input molecule is the same as in the true molecule.
- **Double bond stereochemistry**: The specified double bond stereochemistry in the input molecule is the same as in the true molecule.

D.2 Intramolecular Validity

- **Bond lengths**: The bond lengths in the input molecule are within 0.75 of the lower and 1.25 of the upper bounds determined by distance geometry.
- **Bond angles**: The angles in the input molecule are within 0.75 of the lower and 1.25 of the upper bounds determined by distance geometry.
- Planar aromatic rings: All atoms in aromatic rings with 5 or 6 members are within 0.25 Å
 of the closest shared plane.
- Planar double bonds: The two carbons of aromatic carbon-carbon double bonds and their ring neighbours are within 0.25 Å of the closest shared plane.
- **Internal steric clash**: The interatomic distance between pairs of non-covalently bound atoms is above 0.7 of the lower bound determined by distance geometry.
- Energy ratio: The calculated energy of the input molecule is no more than 100 times the average energy of an ensemble of 50 conformations generated for the input molecule. The energy is calculated using the UFF in RDKit and the conformations are generated with ETKDGv3 followed by force field relaxation using the UFF with up to 200 iterations.

D.3 Intermolecular Validity

- **Minimum protein-ligand distance**: The distance between protein-ligand atom pairs is larger than 0.75 times the sum of the pairs van der Waals radii.
- **Minimum distance to organic cofactors**: The distance between ligand and organic cofactor atoms is larger than 0.75 times the sum of the pairs van der Waals radii.
- **Minimum distance to inorganic cofactors**: The distance between ligand and inorganic cofactor atoms is larger than 0.75 times the sum of the pairs covalent radii.
- **Volume overlap with protein**: The share of ligand volume that intersects with the protein is less than 7.5%. The volumes are defined by the van der Waals radii around the heavy atoms scaled by 0.8.
- Volume overlap with organic cofactors: The share of ligand volume that intersects with organic cofactors is less than 7.5%. The volumes are defined by the van der Waals radii around the heavy atoms scaled by 0.8.
- Volume overlap with inorganic cofactors: The share of ligand volume that intersects with inorganic cofactors is less than 7.5%. The volumes are defined by the van der Waals radii around the heavy atoms scaled by 0.5.

1134 E Additional Figures for Model Evaluation

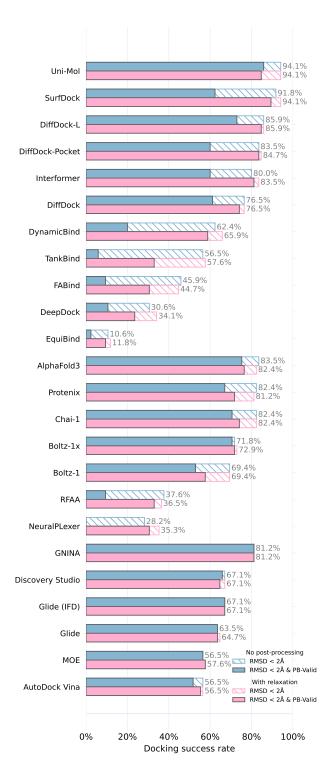


Figure S2: Performance on Astex.

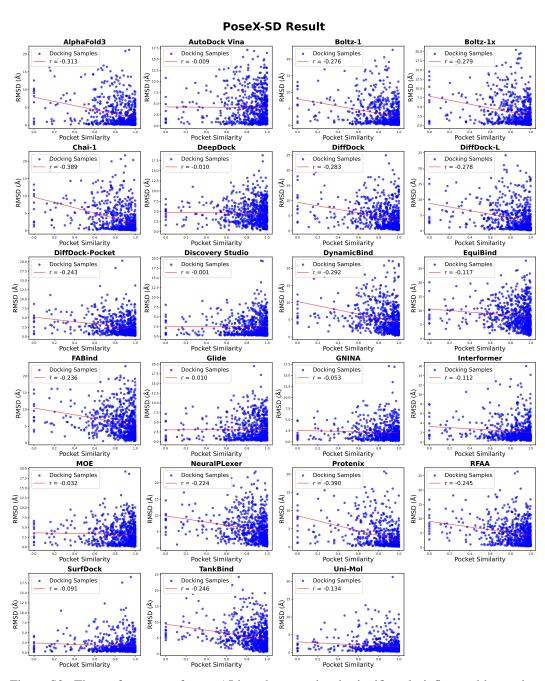


Figure S3: The performance of most AI-based approaches is significantly influenced by pocket similarity under **self-docking** setup. Among them, Protenix [26] exhibits the strongest negative correlation (r = -0.390), whereas SurfDock [21], an AI-based model, demonstrates minimal statistical association. In contrast, *physics-based methods*, such as AutoDock Vina and Glide, are relatively unaffected by protein similarity.

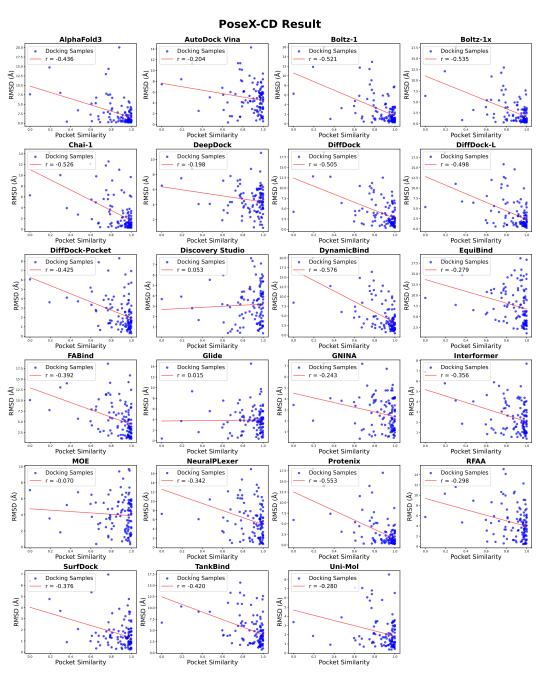


Figure S4: The performance of most AI-based approaches is significantly influenced by the pocket similarity in **cross-docking** scenario, where similar conclusions as the self-docking scenario can be derived.



Figure S5: **Performance on the PoseX-SD dataset.** Samples are sorted by pocket similarity in descending order, and the RMSD results are processed with a moving average (window size: 100). It can be seen that most AI-based approaches degrade as pocket similarity decreases, while *physics-based methods* perform relatively stably.



Figure S6: **Performance on the PoseX-CD dataset.** The results are similar to those on PoseX-SD (window size: 20).

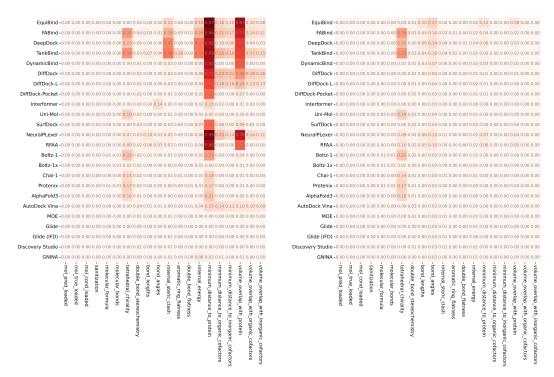


Figure S7: The proportion of models filtered out based on various filtering criteria in PB-Valid (PoseX-SD).

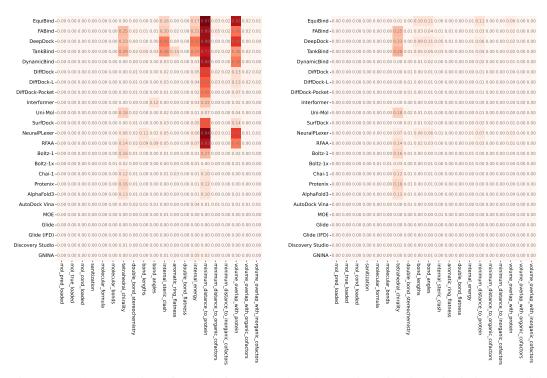


Figure S8: The proportion of models filtered out based on various filtering criteria in PB-Valid (PoseX-CD).

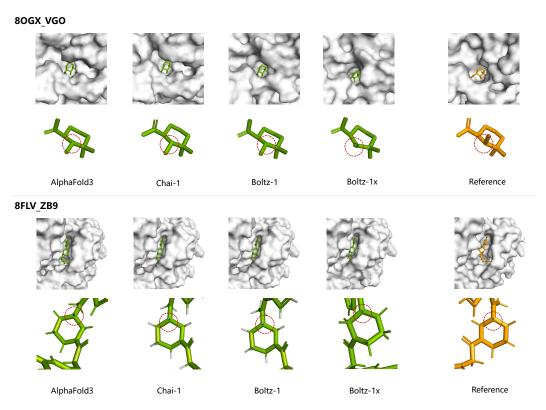


Figure S9: **Case study of** *AI co-folding methods* in chirality validation. We compared AlphaFold3, Chai-1, Boltz-1, and Boltz-1x models on the **8OGX_VGO** and **8FLV_ZB9** complexes. The figure illustrates the docking results, with chiral centers marked by red circles, revealing that all co-folding models except Boltz-1x exhibit chirality errors.

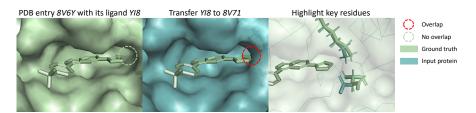


Figure S10: Analysis of **8V71_YI8** in PoseX-CD. When transferring the ligand from its co-crystal structure to the protein structure used for docking through structural alignment, steric clashes arise between the ligand and the protein, underscoring the challenges associated with cross-docking. In this case, all *physics-based methods* failed (RMSD \geq 2Å), while the top-performing *AI docking method* and *AI co-folding method* (SurfDock and AlphaFold3, respectively) accurately predicted the pose. The rightmost column illustrates the conformational variations in residues that overlap with the ligand across the two protein structures.

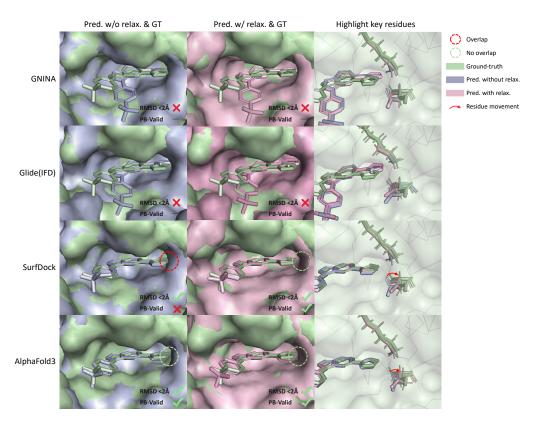


Figure S11: Analysis of docking results for **8V71_YI8**. The *physics-based methods* GNINA and Glide(IFD) generate ligand conformations that substantially deviate from the ground-truth structure. In contrast, SurfDock and AlphaFold3 generate docking poses that closely align with the ground-truth structure. SurfDock's docking poses exhibit steric clashes, which are resolved through relaxation, whereas AlphaFold3's poses are sterically compatible. The rightmost column demonstrates that, for both SurfDock and AlphaFold3, key residues shift toward their corresponding positions in the ground-truth structure after relaxation.