






Protein language models are performant in structure-free virtual screening

Hilbert Yuen In Lam ^{1,2}, Jia Sheng Guan ¹, Xing Er Ong ², Robbe Pincket ³, Yuguang Mu ^{1,2,*}

¹School of Biological Sciences, Nanyang Technological University, 60 Nanyang Dr, Singapore 637551, Singapore, Republic of Singapore

²MagMol Pte. Ltd., 68 Circular Road, #02-01, Singapore 049422, Singapore, Republic of Singapore

³Heliovision, Asstraat 5, 3000 Leuven, Leuven, Kingdom of Belgium

*Corresponding author. Mu Yuguang. Tel.: +65 6316 2885; E-mail: ygmu@ntu.edu.sg

Abstract

Hitherto virtual screening (VS) has been typically performed using a structure-based drug design paradigm. Such methods typically require the use of molecular docking on high-resolution three-dimensional structures of a target protein—a computationally-intensive and time-consuming exercise. This work demonstrates that by employing protein language models and molecular graphs as inputs to a novel graph-to-transformer cross-attention mechanism, a screening power comparable to state-of-the-art structure-based models can be achieved. The implications thereof include highly expedited VS due to the greatly reduced compute required to run this model, and the ability to perform early stages of computer-aided drug design in the complete absence of 3D protein structures.

Keywords: virtual screening; computer-aided drug design; protein language models; cheminformatics

Introduction

One major pillar of rational drug design has always been the use of virtual screening (VS) through docking in what is commonly known as structure-based drug design (SBDD) [1]. In docking, molecular ligands are typically conformationally explored in a protein pocket either through the use of biophysically defined constraints or machine learning (ML) methods, and a best pose with its corresponding computed binding affinity reported. A quintessential VS pipeline will iteratively perform docking through a library, usually consisting of millions to billions of unique chemical compounds, and rank the ligands based on the derived affinity - the top scored ligands will then proceed onto the next phase of drug development, either through computational means such as molecular dynamics (MD) simulations or through experimental validation [2].

In order to enhance the accuracy of SBDD and VS, significant progress has been made in developing docking tools and rescoring functions—the latter of which are typically deep learning methods that, given a docked ligand, output a correction term or a new score entirely. These re-scoring methods have shown large promise, increasing the screening power in benchmarks [3–6]. Screening power is defined as the ability of a given model to differentiate between what will and will not bind to a target experimentally [7], and is largely considered the ultimate test of any VS model. Screening power can be measured by the metric known as enrichment factor, which is measured on a target-by-target basis (usually then averaged across multiple targets in a benchmark), and reflects the concentration of active ligands among the highest scoring hits by a model compared to the concentration of active ligands amongst the entire dataset [8]. Therefore, by increasing

screening power, a model can more effectively discover leads, and overall improve the efficacy of computer-aided drug design (CADD).

However, docking itself and the addition of a rescoring term greatly increases the computational time and expense required to score each ligand. With very large libraries being used and an estimated chemical space of 10^{63} unique compounds [9], this increased compute in already time-consuming VS means that either cost would have to go up per target or the amount of ligands that can be screened per target would have to go down, inadvertently compromising the comprehensiveness of a VS pipeline on the chemical space explored.

Aggravatingly, recent work also suggests that many SBDD deep learning models merely memorize the ligands, resulting in poor generalizability and outcomes [10]. In the case of models predicting the binding affinity directly, studies have shown that many of these models consider less of the protein–ligand interaction but more so of quantitative structure analysis relationship (QSAR) of the ligand alone in making their predictions [11]. Intrinsic biases have also been found in common datasets used for training SBDD deep learning models [11]. There is also the issue of limited data, with common datasets such as the PDBind+ dataset [12] consisting of 22 920 ligand–protein pairs, further complicating the issue of training performative and generalizable models due to data scarcity. SBDD also has a pitfall in which only a single snapshot of a pocket is considered during docking and subsequent rescoring - and although these limitations have been addressed through many works throughout the years [13], the fundamental flexible and dynamic nature of the protein is still, unfortunately, largely given little consideration [14]. The performance of

Received: May 24, 2024. Revised: August 17, 2024. Accepted: September 12, 2024

© The Author(s) 2024. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

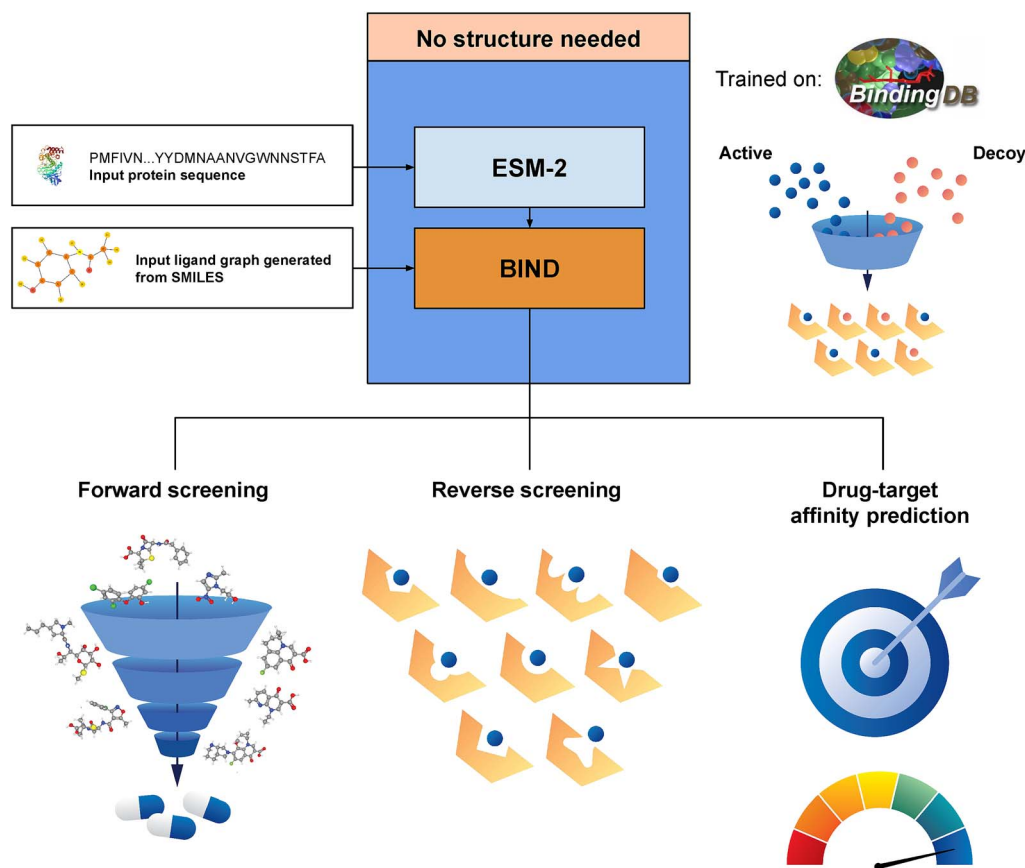


Figure 1. BIND is able to perform forward screening, reverse screening, and drug-target affinity prediction, all without structural input. BIND is trained on the BindingDB dataset only consisting of protein sequences and experimentally determined DTA values. By attaching BIND to a pre-trained ESM-2 protein language model and feeding in a protein sequence and SMILES molecular representation, which are then converted into graphs, the model can effectively discriminate between active and decoy ligands. This allows BIND to be used in forward and reverse screening, while predicting DTA values.

molecular docking can also vary between different conformational states of a protein (i.e., the apo or holo state), and hence docking predominantly suffers from considering the induced fitting of ligands [15].

To address this, some have resorted to instead predicting drug-target affinity (DTA), an experimentally derived figure such as the half-maximal inhibitory concentration (IC_{50}) or dissociation constants. Such work includes SSM-DTA [16], a sequence-based deep learning method that achieves state-of-the-art performance in DTA prediction. However, in this work, the authors demonstrate that DTA-only models do not necessarily have the highest screening power.

Therefore, to tackle the issue of limited data in SBDD, the single snapshot issue, the lack of consideration of the entire protein during SBDD, induced fitting model and different conformations, this work combines protein language models (PLMs), a form of large language model together with graph neural networks to represent ligands to perform VS without inputting structural information (Fig. 1). The fundamental idea is that any structural and dynamic information of proteins is implicit in any PLMs pre-trained through self-supervised means such as Evolutionary Scale Modelling 2 (ESM-2) and this theorem has been further demonstrated in works such as ESMFold [17]. This work further introduces a novel graph-to-transformer cross-attention block, which essentially treats every single node in a graph as a token and allows it to query a separate sequence. The combined model, named Binding Interaction Determination (BIND), achieves comparable performance to state-of-the-art SBDD models with a fraction of the compute and time required and with

only protein sequence and ligand information as inputs. Primarily, this work differs from previous work in its use of a decoy/true binder classifier and application to datasets typically used to benchmark SBDD models such as DUD-AD [18]. Through the use of this model, a state-of-the-art reverse screening performance was also achieved when compared with SBDD models on AlphaFold2-predicted structures. This work also showed that when trained with the multi-objective classification and DTA as used in BIND, the model's DTA prediction also outperforms state-of-the-art DTA-only models in the screening power domain.

Results

Screening power of BIND comparable to top models

The screening power of the model (Fig. 2) is similar to that of top SBDD models for all benchmarks performed (Table 1-5). Overall, a high enrichment factor is observed for all the benchmarks tested. In the DEKOIS 2.0, DUD-AD, DUD-E, and LIT-PCBA datasets, the BIND model appears to have the highest enrichment factor at 1% cutoff with respect to reviewed literature (Fig. 3c-f). However, the model falls short in the CASF-2016 forward screening dataset (Fig. 3a-b), achieving higher enrichment than Glide-SP and Glide-XP but is overall deficient in the success rate and when compared to other very recent ML methods. Overall, in the CASF-2016 dataset, most of the top performing targets when sorted by BEDROC belonged to *Homo sapiens*, whereas the enrichment was poorer in distantly related species (Supplementary Table 1). Furthermore, the model's enrichment factor is higher across all

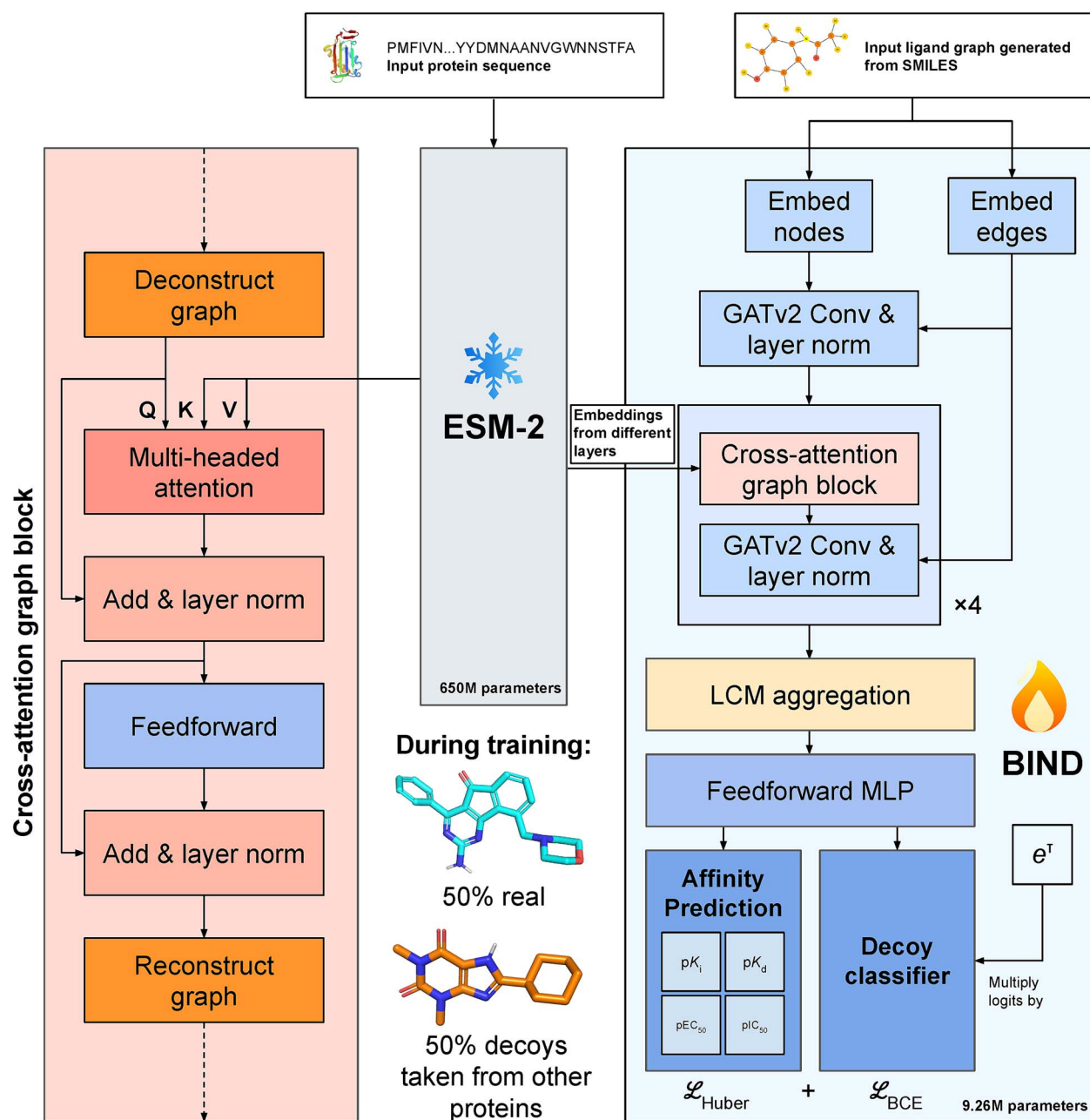


Figure 2. BIND architecture incorporates a proposed cross-attention graph block, and is trained with both true ligands and decoys taken from other proteins in the same dataset. The cross-attention graph block essentially deconstructs a graph and treats each node as a token for cross-attention—this allows the ligand to ‘query’ the protein and its important parts. The loss is a summation of Huber loss for the affinity predictions, and binary cross entropy loss for the decoy classifier. During training, ESM-2’s weights are frozen such that only the 9.26 M parameters from BIND are tuned. Q, K, and V in the diagram represent the query, key and value inputs characteristic of transformers.

datasets tested when compared to using the pIC_{50} prediction from SSM-DTA, a fully DTA prediction model (Fig. 1). The training and validation losses of the evaluated model are described in Supplementary Fig. 1, and individual benchmark scores for each protein are given in Supplementary File 1. Further training experiments to optimize the model and their corresponding results are also available in the supplementary material. Moreover, in zero-shot protein settings where highly homologous protein sequences are removed from training, the zero-shot BIND model still achieves good performance, outperforming some contemporary score functions in the datasets tested. The individual zero-shot BIND results are given in Supplementary File 2. Further side-by-side comparisons in the enrichment factor with another

sequence-based model, TransformerCPI2.0, on the DUD-E and DEKOIS 2.0 datasets, are shown in Supplementary Tables 2 and 3, respectively.

Protein language models can perform reverse screening

Reverse screening was performed on the benchmark as described by Luo et al., 2023 [44]. In summary, 90 selected ligands were scored against 12,195 human proteins representative of the human proteome. From the results, BIND is comparable to that of top models with rescoring such as OnionNet-SFCT with Glide-SP docking on AlphaFold2-predicted structures and PointSite or SiteMap for pocket determination; BIND being able to

Table 1. **Screening power results for top models and commonly used models on CASF-2016, with all scores represented as the mean.** Arrows indicate direction of better results.

Score function / model	Method	Prediction target	EF _{1%} ↑	Forward screening success rate 1% ↑
GenScore [19]	Deep learning, structure-based	Distance likelihood	28.20	71.9%
RTMScore [5]	Deep learning, structure-based	Distance likelihood	28.00	66.7%
PIGNet2 [20]	Deep learning, physics-based, structure-based	Binding affinity	24.90	66.7%
DeepRMSD [21]	Deep learning, structure-based	Δ binding affinity	21.95	47.4%
DeepDock [22]	Deep learning, structure-based	Distance likelihood	16.41	43.9%
OnionNet-SFCT [4]	Deep learning, structure-based	Δ binding affinity	15.5	35.1%
BIND	Deep learning, sequence-based	Binary classification between true binders and decoys	14.91	28.7%
ChemPLP@GOLD [23]	Physics-based, Structure-based	Binding affinity	11.91 [7]	35.1% [7]
ΔVinaRF ₂₀ [24]	Machine learning, structure-based	Δ binding affinity	11.73 [7]	42.1% [7]
Glide-SP [25]	Physics-based, structure-based	Binding affinity	11.44 [7]	36.8% [7]
Glide-XP [25]	Physics-based, structure-based	Binding affinity	8.83 [7]	26.3% [7]
ChemScore@GOLD [26]	Machine learning, structure-based	Binding affinity	8.65 [7]	28.1% [7]
Autodock Vina [27]	Physics-based, structure-based	Binding affinity	7.70 [7]	29.8% [7]
Zero-shot BIND (90% homologous sequences removed)	Deep learning, sequence-based	Binary classification between true binders and decoys	6.98	13.5%
SSM-DTA (pIC ₅₀) [16]	Machine learning, sequence-based	Drug-target affinity	3.67 [7]	8.2% [7]

^abest model's results as reported by respective authors.

Table 2. **Screening power results for DEKOIS 2.0, with all scores represented as the mean.** Arrows indicate direction of better results.

Model	EF _{0.5%} ↑	EF _{1%} ↑	AUROC ↑	BEDROC (α = 80.5) ↑
BIND	25.36	24.46	0.944	0.730
SSM-DTA [16]	21.21	18.49	-	0.504
RTMScore [5]	20.99	18.53	-	0.558
GenScore [19]	20.24	17.87	0.757	0.539
PIGNet2 [20]	20.00	18.60	0.812	0.544
KarmaDock [28]	19.41	16.29	0.783	0.519
Zero-shot BIND (90% homologous sequences removed)	16.84	15.19	0.832	0.453
SCORCH [29]	15.01	13.78	-	-
DyScore [30]	-	7.5	0.702	0.216
Vina [27]	-	4.8 [30]	0.631 [30]	0.140 [30]

^abest model's results as reported by respective authors.

Table 3. **Screening power results for DUD-AD, with all scores represented as the mean.** Arrows indicate direction of better results.

Score function / model	EF _{1%} ↑	AUROC ↑	BEDROC (α = 80.5) ↑
BIND	24.14	0.952	0.698
Zero-shot BIND (90% protein homology sequences removed)	12.03	0.851	0.396
SSM-DTA (pIC ₅₀) [16]	10.19	-	0.336
DeepRMSD [21]	6.04	-	-
OnionNet-SFCT with Vina [4]	5.22	0.549	-
Gnina [31]	2.33 [4]	-	-
Vina [27]	1.90 [4]	-	-

^abest model's results as reported by respective authors.

cumulatively rank the highest number of correct pairs for the top 2–1000 pairs (Fig. 4a-b). Furthermore, the reverse screening is expedient, taking only approximately eight hours on a Ryzen Threadripper Pro 5975WX CPU to run through all 1 097 550 protein–ligand combinations. The reverse docking performance of BIND was also evaluated on the Astex Diver dataset, of which in the top 1 scoring outperformed DOCK, Glide, and Autodock Vina, with comparable performance on the top 5 metrics (Fig. 4c).

Model can perform DTA prediction

The model is able to perform on the DTA prediction datasets DAVIS, despite not being specifically trained on this dataset but

only BindingDB (Table 6). Furthermore, on the BindingDB test split, it achieved state-of-the-art RMSE in pK_d prediction specifically, but falls short in the other metrics of pK_i, pIC₅₀ and pEC₅₀ (Table 7).

Drug-target affinity prediction alone has lower screening power compared to classification

When ranked using pK_i, pK_d, pIC₅₀, or pEC₅₀, the enrichment factors across benchmark datasets dropped compared to classification alone (Fig. 5a–e). Furthermore, the classifier also has the strongest enriching effects in the AlphaFold2 reverse docking benchmark (Fig. 5f). Lastly, in spite of state-of-the-art

Table 4. Screening power results for DUD-E, with all scores represented as the mean.

Score function / model	EF _{0.5%} ↑	EF _{1%} ↑	AUROC ↑	BEDROC ($\alpha = 80.5$) ↑
BIND	51.88	46.35	0.944	0.733
RTMScore [5]	42.47	35.10	0.831	0.558
GenScore [19]	44.02	33.96	0.828	0.546
DyScore (target-aware split) [30]	-	32.20	0.858	0.502
PIGNet2 [20]	36.80	31.20	0.850	0.515
Zero-shot BIND (90% protein homology sequences removed)	30.52	26.39	0.826	0.430
SSM-DTA (pIC ₅₀) [16]	33.59	24.78	-	0.397
Gnina [31]	-	20.40 [32]	0.795 [32]	-
OnionNet-SFCT with Vina [4]	-	15.54	0.724	-
DeepRMSD with Vina* [21]	-	14.80	-	-
DyScore (target-unaware) [30]	-	12.10	0.755	0.216
Vina [27]	-	8.82 [4]	0.697 [4]	-

^abest model's results as reported by respective authors.

Table 5. Screening power results for LIT-PCBA, with all scores represented as the mean.

Score function / model	EF _{1%} ↑	AUROC ↑	BEDROC ($\alpha = 80.5$) ↑
BIND	10.93	0.597	0.112
IFP [33]	7.46 [34]	-	-
GRIM [35]	6.87 [34]	-	-
GenScore [19]	6.80	-	-
DyScore-MF [30]	5.92	0.594	0.071
Δ VinaRF ₂₀ [24]	5.38 [34]	-	-
Pafnucy [36]	5.32 [34]	-	-
FRAGSITE [37]	4.78	-	-
Zero-shot BIND (90% protein homology sequences removed)	4.45	0.564	0.047
EViS [38]	4.18	-	-
FINDSITE ^{comb2.0} [39]	3.03 [38]	-	-
SSM-DTA (pIC ₅₀) [16]	3.03	-	0.040
Score2 [30]	3.00	0.621	0.047
Gnina [31]	2.58 [32]	0.616 [32]	-
Surflex [40]	2.51 [34]	-	-
Vina [27]	2.33	0.565	0.037
BANANA [41]	1.81	0.580	-
NNScore-v2 [42]	1.70 [30]	0.557 [30]	0.025 [30]
RFScore-4 [3]	1.28 [30]	0.600 [30]	-
DSX [30]	1.00	0.523	0.025
Vinardo [43]	0.99 [32]	0.577 [32]	-
RFScore-VS [3]	0.73 [32]	0.542 [32]	-

^abest model's results as reported by respective authors.

DTA predictors being more accurate in predicting DTA compared to BIND (Table 6), BIND still achieves higher screening power on CASF-2016, DEKOIS 2.0, DUD-AD, DUD-E and LIT-PCBA using its pIC₅₀ prediction compared to SSM-DTA's pIC₅₀ (Supplementary Fig. 2).

BIND is significantly faster than other SBDD models

BIND demonstrates around three orders of magnitude in speed improvement over OnionNet-SFCT and approximately two orders of magnitude in speed improvement over QuickVina2.1 and Gnina (Supplementary Fig. 3).

Discussion

The use of PLMs in predicting DTA is not new. Previous work has delivered very accurate DTA predictors using PLMs such as SSM-DTA [16]. However, there has been far less application of this to determine screening power - arguably the most important

goal in CADD. Previous work by Tsubaki *et al.*, 2019, used a one-dimensional n-gram convolutional neural network and cross-attentions to do similar work, achieving good discriminative power between real and decoy ligands when trained and evaluated on the DUD-E dataset using fivefold cross-validation [54]. This work builds upon that foundation and uses the pre-trained ESM-2 to transfer pre-learned knowledge about the protein language into the VS pipeline.

Primarily, BIND serves as a proof of concept that even in the complete absence of protein structure or binding pocket information, it is possible to achieve screening power similar to that of and sometimes even exceed certain top SBDD models with PLMs by training the model to discriminate between real and decoy ligands selected from the same dataset. The results show that by using a model such as the proposed BIND, it may be possible to perform VS for drug discovery, or to elucidate the identity of protein targets which bind to a specific ligand as per reverse docking. However, from the CASF-2016 results, it appears that BIND is not the best at ranking the very top and most potent binders as apparent

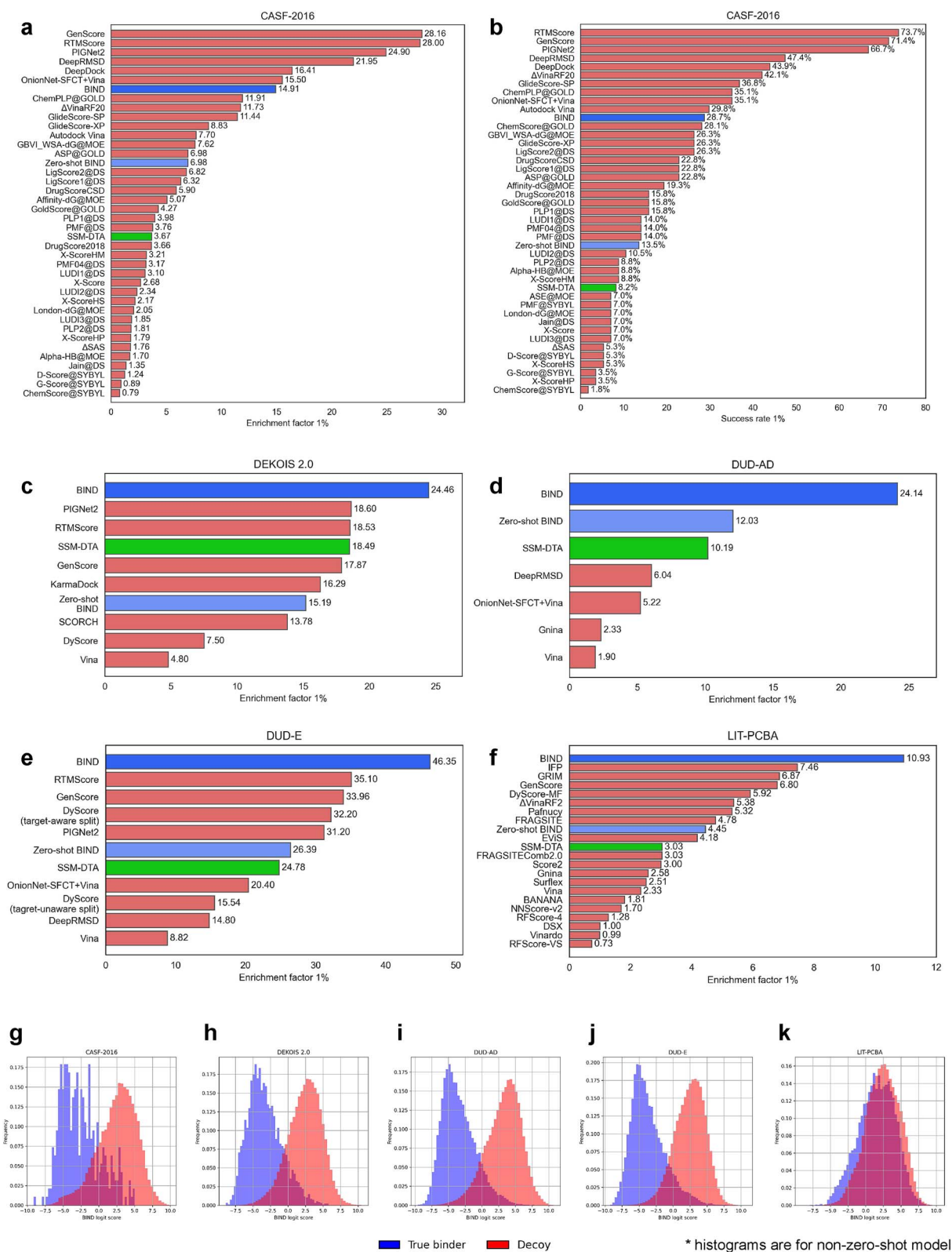


Figure 3. The model is performant on docking benchmarks and outperforms the state-of-the-art DTA only model in enrichment. CASF-2016, DEKOIS 2.0, DUD-AD, DUD-E, and LIT-PCBA were evaluated. (a, b) CASF- 2016 1% enrichment factors and 1% success rate, (c-f) DEKOIS 2.0, DUD-AD, DUD-E and LIT-PCBA 1% enrichment factors, (g-k) probability-normalized histogram of the non-zero-shot BIND classifier logit output showing separation between true binders and decoys. All enrichment factors and success rates are averaged across the entire datasets evaluated, and histograms shown are distributions of logit scores for all proteins and ligands tested. The green bar indicates the enrichment factor of the published SSM-DTA model, which predicts the pC50 DTA. The zero-shot BIND model indicates the model in which proteins with >90% homology comparative to the evaluation datasets are removed during training.

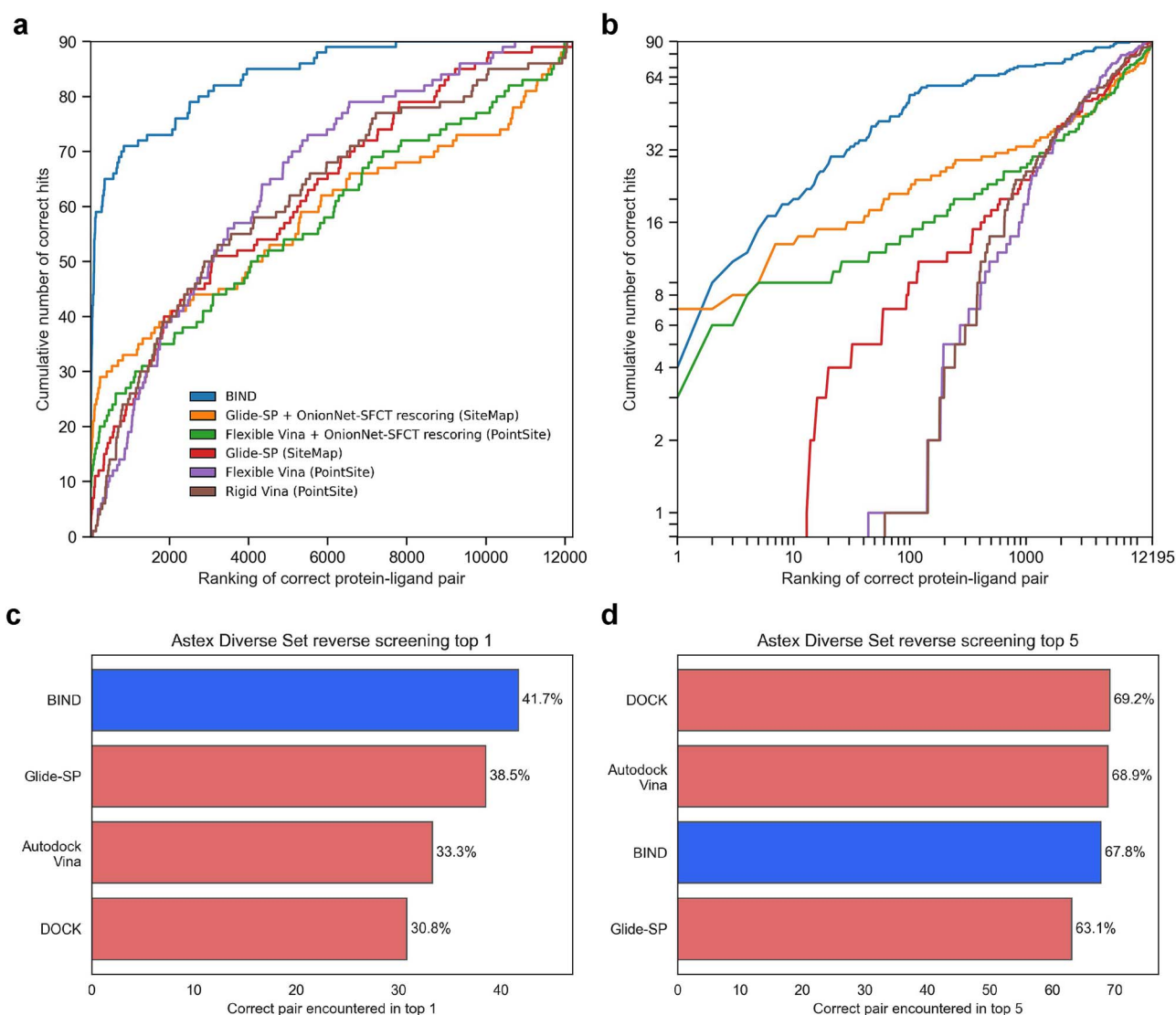


Figure 4. Protein language models can outperform standard reverse docking pipelines using structures from AlphaFold2. A total of 12 195 proteins and 90 ligands from Luo et al., 2023’s AlphaFold2 reverse docking benchmark were individually pairwise scored using BIND and the ligands ranked by classification score for each protein. (a, b) the cumulative ranking of BIND in standard and logarithmic plots respectively. The legend indicates the type of score function used, and the software used to determine the pocket locations in parentheses. (c, d) ranking of 84 ligands against 85 proteins in the reverse docking benchmark on the Astex dataset, with other benchmark scores obtained from Luo et al., 2017.

Table 6. Drug-target affinity on benchmark datasets.

Dataset	DAVIS (pK_d)	
	MSE ↓	CI ↑
SSM-DTA [16]	0.219	0.875
DeepDTA [45]	0.262	0.870
ELECTRA-DTA [46]	0.238	0.897
SubMDTA [47]	0.218	0.894
DeepPurpose [48]	0.242	0.881
AttentionMGT-DTA [49]	0.193	0.891
AttentionDTA [50]	0.195	0.888
BIND (trained on BindingDB only)	0.409	0.747

^abest model’s results as reported by respective authors.

in its poorer success rate in a much smaller set of ligands as seen in CASF-2016. This could also be due to the skew in the BindingDB dataset, with the majority of the targets being from

H. sapiens instead of other species, and that the CASF-2016 dataset contains multiple chains which are not accounted for in the BindingDB training dataset. From this, and with the comparable performance and occasionally state-of-the-art enrichment factor scores in the other larger datasets, it can be said that BIND has rather demonstrated its utility in enrichment of large set; and given the high speed at which BIND can screen compounds, it can be pragmatically deployed in initial screening stages to filter out a large number of non-binders (i.e., from millions to thousands), of which refinement can then be done with other score functions such as GenScore, PIGNet2 or OnionNet-SFCT to pull out only the top binders. The strong enriching effects of BIND in reverse screening also indicates that BIND could be better at finding a target protein from a large number of proteins given a ligand than performing docking on AlphaFold2 structures. This is especially useful as, in many reverse docking pipelines, the experimental structure of all proteins is typically not available, and reverse docking may be performed on non-binding protein conformations [55].

Table 7. Drug-target affinity on BindingDB.

Metric	pK _i		pK _d		pIC ₅₀		pEC ₅₀	
	RMSE ↓	R ↑	RMSE ↓	R ↑	RMSE ↓	R ↑	RMSE ↓	R ↑
DeepAffinity* [51]	0.840	0.840	-	-	0.780	0.840	-	-
SSM-DTA [16]	0.792	0.863	-	-	0.712	0.878	-	-
BACPI [52]	0.800	0.860	1.080	0.730	0.740	0.860	0.780	0.850
BIND	0.965	0.798	0.935	0.798	0.924	0.808	0.891	0.822
MONN [53]	-	-	-	-	0.764	0.858	-	-

^abest model's results as reported by respective authors.

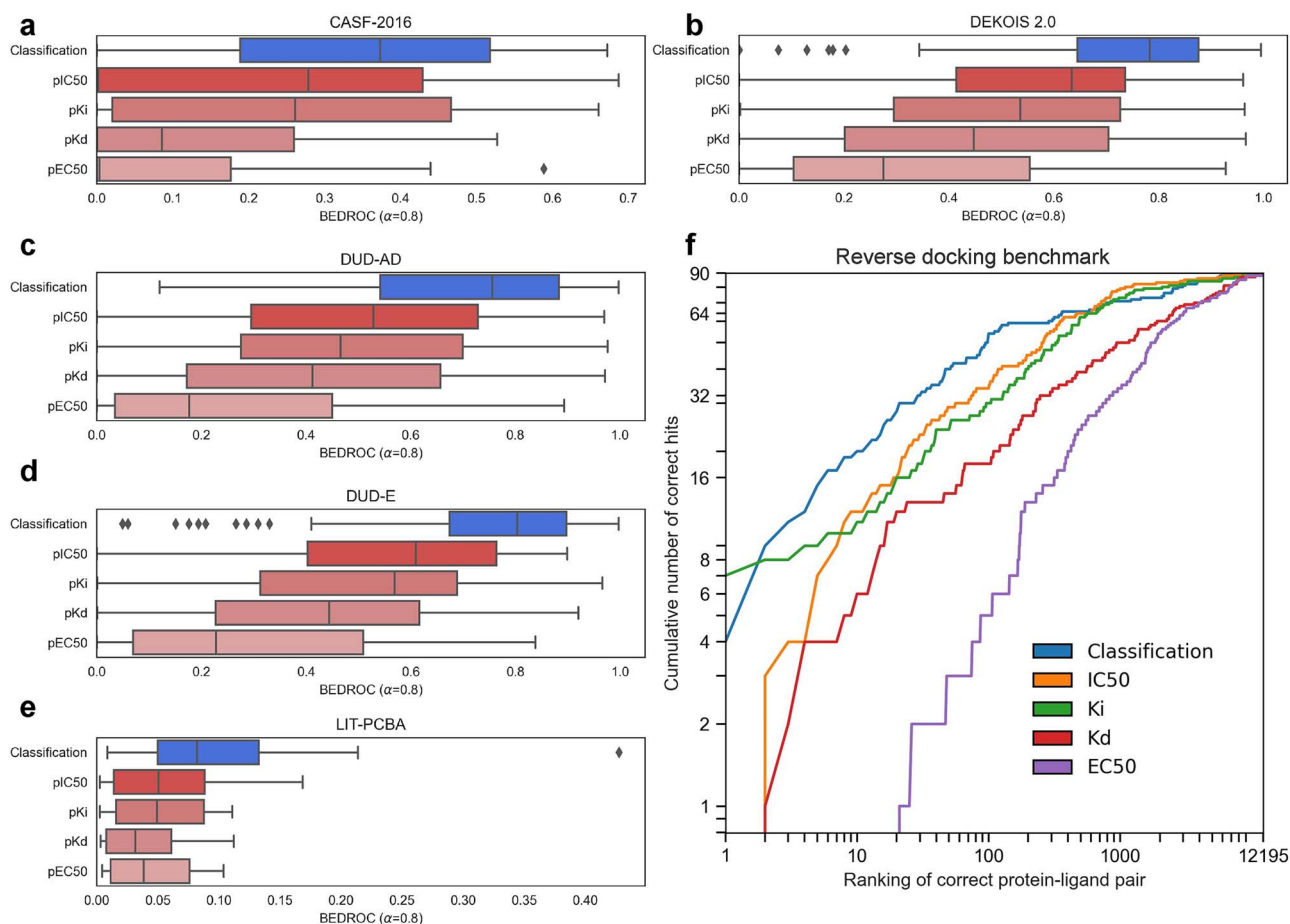


Figure 5. Predicted binding affinity has lower enriching power compared to the classifier in both forward and reverse screening. (a–e) Box-and-whisker plots of BEDROC values across the entire CASF-2016, DEKOIS 2.0, DUD-AD, DUD-E, and LIT-PCBA datasets respectively, with boxes representing the interquartile range and the median demarcated in the box, whiskers showing the range and diamonds showing outliers. (f) Cumulative logarithmic plot of ranking in reverse screening on Luo et al., 2023's AlphaFold2 reverse docking benchmark dataset.

The authors posit from the results that BIND is able to achieve the reported screening power by borrowing the pre-learned context from PLMs such as ESM-2, which are able to understand the implicit structure and underlying language of proteins. This includes understanding biochemical properties, dynamic movement of proteins, and interdomain interactions. The performance of the zero-shot model also shows BIND has sufficient generalization capability in screening unseen proteins, potentially obtaining this capability from its ESM-2 pre-training which allowed it to see a much more diverse set of protein sequences. This premise of stronger performance from transfer learning is supported by that of literature [56], and previous work such as ESMFold have demonstrated that structure can be inferred from sequence alone with the help of PLMs. Therefore, it may not be far-fetched that the

flexibility and even how a protein interacts with other molecules such as ligands be implicitly learnt through the masked language modelling and other self-supervised objectives of PLMs in pre-training. PLMs may be able to better model molecular interactions with disordered regions in proteins and even decipher cryptic pockets which evade the typical SBDD VS pipeline due to their transient nature [57]. There remains much to be explored in this domain.

Furthermore, it is interesting that the BIND model's DTA prediction has lower screening power and enriching effects compared to discriminating between true and false binders (i.e., the classifier in BIND) alone. It is highly plausible that this is due to internal biases of the training dataset - as all of the training protein-ligand pairs are true binders, the algorithm is misled into thinking

that every protein–ligand pair that is given is always a binder and hence predicts a binders' score regardless to prevent heavy loss penalties, as this is reflective of what it has seen and points raised by previous work. This result reinforces the indication that robust augmentation very likely goes hand-in-hand with screening power. The use of sole DTA models in screening may also not be appropriate in all cases, as DTA datasets, including BindingDB, which derives its affinity values from multiple sources, have been shown to be noisy across multiple assays [58]—this underscores that it may be more practical to predict a binary label as to whether or not a drug binds than a quantitative experimental affinity that is subject to many other confounding and unforeseen experimental variables, which could also explain BIND's higher screening power compared to SSM-DTA although BIND achieves significantly worse R^2 and MSE values in IC_{50} and other DTA predictions across the board.

The authors acknowledge several limitations of this work: (i) The BindingDB dataset is mainly obtained from ChEMBL and other sources, and contains overlapping information with the CASF-2016, DEKOIS 2.0, DUD-AD, DUD-E, and LIT-PCBA evaluation datasets, in the process inflating the results for the forward screening for the non-zero-shot model. Although this is an issue, it is difficult to directly and objectively quantify the overlaps and individually eradicate the overlaps. Furthermore, by removing overlaps, a subjective defined stringency is required (i.e., Tanimoto similarity, protein sequence similarity), in the process getting rid of a large, diverse, and representative chunk of protein–ligand binding data in *H. sapiens* and other key protein families (which is the entire rationale of the benchmark datasets). This would prejudice BIND in evaluation as other score functions are typically trained with knowledge of that protein previously (i.e., from the PDBBind dataset [12], which has directly overlapping PDB structures with the evaluation dataset and are typically not removed during model training). By removing direct ligand–protein matches but not proteins during evaluation, the evaluation cannot be fairly performed with other score functions using the same dataset as the dataset would be modified and the enrichment factors cannot be directly compared. To mitigate this limitation and to assess the true screening power of BIND, the zero-shot model was attempted in which sequences containing >90% homology on the proteins were removed from the training set, and the model still showed generalization capability compared to contemporary score functions; (ii) the explainability of BIND is limited—the cross-attentions from BIND are largely uninterpretable in-line with post-hoc explainability analyses performed on transformer models, meaning that there is little explanation as to which residues contribute most to binding; (iii) much like previous work such as SSM-DTA, the training expense of this model is high. To compensate with limited resources, the authors used a large number of gradient accumulation steps, in the process dragging out the training time. The model may also further benefit in the DTA prediction criteria from a scaled-up training regime, which could not be achieved in this work due to limited computational resources; (iv) No wet lab experimental validation was performed for this work—although BIND performs reasonably well on benchmarks as explored in this work, many factors ultimately affect the efficacy of a drug and experimental validation is needed.

Overall, this work proposes the use of PLMs as a potential alternative for SBDD VS—showing that, even without structural information, is comparable to top-of-the-line SBDD models. The authors reinforce that this methodology be termed sequence-based drug discovery (SeBDD) [59]. SeBDD could be useful in work

in which the protein structure cannot be accurately resolved, or when the binding pockets are unknown. In current screening pipelines, SeBDD can be used as an initial filter in a VS campaign, in which the top molecules from SeBDD are selected for standard SBDD docking and MD simulations. Due to the simplicity of SeBDD compared with SBDD (no need for docking, preparing the Gasteiger partial charges of ligands and targets, etc.), pipelines can also be greatly accelerated by its use.

Future work in this domain likely involves simulated annealing of molecules and integrating SeBDD with fragment-based drug discovery, the use of synthons [60], and ML-based chemical spaces to explore very large amounts of chemical spaces quickly and accurately.

Methods

Model training, architecture, and training dataset

The BIND model was written and evaluated in Python 3.11 using PyTorch Geometric 2.4.0 with a PyTorch 2.2.0+cu121 backend and HuggingFace Transformers 4.37.2. The full model, along with its trained weights, is disclosed in the GitHub repository. Three models were trained with the one with the lowest validation loss selected and benchmarked in this study, with each model trained on a single Nvidia RTX A6000 for approximately 28 days. The whole model is described in Fig. 2. The implementation code for the cross-attention graph block is also in the same GitHub repository—in essence, each node in the graph is unwrapped and treated as an individual token, before a standard cross-attention similar to that of the original transformer decoder was used—the graph is then reconstructed afterwards. Five prediction heads are used in the final layer of the model, four of which are for regression objectives—to predict the pK_1 , pK_d , pIC_{50} , and pEC_{50} , respectively. A classification objective is also added to predict which ligands are true and which are decoys. The logits of the classification head are multiplied by the natural exponential of a trainable temperature parameter, τ , initialized at 0.07 as per previous work [61]. After scaling, logits are clipped from -100 to 100 to ensure training stability. The loss functions used are Huber loss with a $\delta = 2.0$ and binary cross entropy for regression and classification heads respectively. Training was performed using the AdamW optimizer with a learning rate of $1e-4$ and a weight decay of $1e-3$. A dropout of 0.1 was applied to the multi-headed attention layer, and a leaky rectified linear unit activation of $\alpha = 0.05$ is used unless otherwise specified. Learning rate was scheduled using cosine decay with no warm up. The 650 million parameter ESM-2 (Facebook/esm2_t33_650M_UR50D) model's weights were retrieved using the HuggingFace Transformers 4.37.1 library and are frozen during training. Latents from the ESM-2 layers 1, 11, 21, 31 are cross-attended to. Convolutions for the molecular graph are GATv2 [62] for all instances, and pooling used for the interaction network is Learnable Commutative Monoid (LCM) Aggregation [63]. Model training was done in automatic mixed precision and FlashAttention-2 [64] was used where applicable. The training dataset used is BindingDB [65], with a total of 2 469 626 (90%) single-chained protein–ligand pairs used for training, 54,902 (2%) for validation, and 224,738 (8%) for testing. A batch size of 1 was used during training together with 256 gradient accumulation steps, for a total of 100 000 full iterations (effectively around 10.37 epochs' worth). To optimize for computation, all protein sequences exceeding 2048 amino acids were removed. In total, 1 202 465 unique SMILES strings as part of the BindingDB dataset were used. Since there are multiple prediction heads, for entries in which one or more of the regression values were unavailable,

the local loss was set to zero, meaning that there is no penalty for that term regardless of what the model predicts. For decoy ligands specifically, the loss for K_i , K_d , IC_{50} , and EC_{50} terms are set to zero. All affinity values were normalized with $-\log_{10}(x/1e9)$, where x is in nM, consistent with all previous work [16]. For the reported training loss, the zeroized terms are disregarded in averaging. Approximately 50% of ligands used during training are decoys as part of a data-balanced regime.

Graph node and edge generation

The construction of the molecular graph is identical to that done in previous work [66], with the exception that the PySmiles library is now set to not reinterpret aromatic bonds due to a software bug. PySmiles 1.1.2 and NetworkX 3.2.1 were used to generate the graphs used by the model.

Decoy ligand determination

Decoys were determined at training by randomly selecting another ligand from the loaded SMILES, and programmatically checking if there is another entry in the BindingDB data that tags that SMILES to the same protein sequence.

Evaluation datasets

The DEKOIS 2.0 [67], DUD-E [68], DUD-AD [18], LIT-PCBA [69] and CASF-2016 [7] datasets were used to determine screening power. In all instances, ligands were converted into SMILES using OpenBabel 3.1.0. For all datasets, the UniProt sequences were obtained which corresponded to each RCSB PDB entry. In instances where multiple sequences relating to the target were present, the sequences were concatenated before inputting into the model - protein sequences that were unrelated to the test protein such as hirudin in thrombin and nuclear cofactors in the estrogen receptors were removed. All protein sequences used in evaluation are available in [Supplementary File 3](#). DTA evaluation was performed on the DAVIS [70] dataset obtained from TDCcommons [71], and binding affinity prediction was also performed on a test set of BindingDB split prior to training. Unlike training, there is no cutoff on protein length used during evaluation/testing.

Dataset splitting and zero-shot protein performance evaluation

As there are overlaps between the BindingDB dataset and the screening power evaluation datasets, sequence-ligand pairs in BindingDB which contained >90% overlap in the protein sequence compared to the UniProt sequences in the evaluating datasets were removed for the zero-shot model. The percentage alignment was calculated using Biopython 1.78's pairwise2 align module, with the UniProt sequence used as the reference length in computing the final homology percentage. The zero-shot model was re-trained on this dataset with the exact same parameters as stated in 5.1.

Evaluation with SSM-DTA

The SSM-DTA [16] pIC₅₀ model (BindingDB_IC50/checkpoint_best_20221021) was used on the same datasets described in 5.4, with the same protein sequence input and SMILES as per [Supplementary File 3](#), after canonicalization and tokenization as provided by the SSM-DTA source code.

Metric calculation

BEDROC [72] and enrichment factors were calculated using RDKit 2023.09.5, whereas AUROC was calculated using the SciPy 1.11.0

library. Reported scores are averaged across all entries in the respective evaluation datasets (57 for CASF-2016, 81 proteins for DEKOIS 2.0, 102 for DUD-AD and DUD-E, 15 for LIT-PCBA). All decoys and ligands are used in benchmarking as per the original datasets. The enrichment factor is as determined by the following formula:

$$\text{Concentration of actives in subset} = \text{Actives}_{\text{subset}} / \text{Total}_{\text{subset}}$$

$$\text{Concentration of actives in dataset} = \text{Actives}_{\text{dataset}} / \text{Total}_{\text{dataset}}$$

$$\text{EF}_{\text{subset}} = \text{Concentration of active ligands in subset} / \text{concentration of ligands in dataset}$$

Where the subset is the top quotient of ligands as selected by the model (i.e., 5% of highest scoring molecules as determined by the model).

The success rate is determined as the concentration of actives in the subset.

Reverse screening

Reverse screening was performed as per Luo et al., 2023 [44]. In essence, 90 ligands were selected and reverse screened on a trimmed down human proteome of 12 195 proteins from UniProt using BIND. For the reverse docking performed using SBDD models done by Luo et al., 2023, structures obtained from the AlphaFold Protein Structure Database [73] were used. PointSite [74] or SiteMap [75] were used to determine potential binding pockets. Reverse screening was also done on the Astex Diverse Set using the same procedure as with other standard docking tools described by Hartshorn et al., 2007 [76]. Summarily, the 84 ligands found in the dataset were docked against the 85 proteins to identify the correct protein target for the ligand and ranked, with the percentage of correct pairs identified in the top 1 and top 5 rankings as per literature [77].

Speed benchmarking

Speed benchmarking on Gnina and QuickVina2.1 [78] was performed in previous work on the *Escherichia coli* EPSP synthase (PDB ID: 2QFT) [66]. For BIND, a batch size of one was used and the sequence P0A6D3 obtained from UniProt, and the ESM-2 embeddings compute time not included in the speed benchmark as only one embedding has to be calculated per protein, and can be reused for any subsequent ligands when screened against the same protein. Benchmarks were performed on an Nvidia RTX 3090 GPU and an AMD Ryzen 5950X CPU, with the same machine used in benchmarking Gnina and QuickVina2.1.

Key Points

- BIND, a fully sequence-based model based on pre-trained protein language models, is able to achieve competitive results compared to structure-based methods in virtual screening, removing the need for elucidating protein structures in the initial stages of computer-aided drug design
- By creating a simple discriminator between decoys and true binders, BIND is able to achieve state-of-the-art performance in some virtual screening benchmarks
- This sequence-only model is also able to competently perform reverse screening
- This work also introduces a cross-attention graph block, in which graph neural networks can cross-attend to transformer models

Acknowledgements

The authors thank Saxena Shikhar for the fruitful discussions as part of this work. We further would like to express our appreciation to Qing Luo and Jingjing Guo of Macao Polytechnic University for supplying their full reverse docking data from their previous work.

Supplementary data

Supplementary data are available online at <http://bib.oxfordjournals.org/>.

Conflict of interest: None declared.

Funding

This work is supported by the Singapore Ministry of Education (MOE) Tier 1 grant RG97/22. The computational work for this article was partially performed on resources of the National Supercomputing Centre, Singapore (<https://www.nsc.sg>) and the HADLEY supercomputer by Singapore Centre for Environmental Life Sciences Engineering (SCELS).

Code availability

The code is available at: <https://github.com/Chokyotager/BIND>.

Data availability

All data utilised in this study are available from the RCSB Protein Data Bank. Additionally, the training data used, sourced from BindingDB, are publicly accessible.

References

1. Ferreira LG, dos Santos R, Oliva G. et al. Molecular docking and structure-based drug design strategies. *Molecules* 2015;**20**: 13384–421. <https://doi.org/10.3390/molecules200713384>.
2. Zhang B, Li H, Yu K. et al. Molecular docking-based computational platform for high-throughput virtual screening. *CCF Trans High Perform Comput* 2022;**4**:63–74. <https://doi.org/10.1007/s42514-021-00086-5>.
3. Wojcikowski M, Ballester PJ, Siedlecki P. Performance of machine-learning scoring functions in structure-based virtual screening. *Sci Rep* 2017;**7**:46710. <https://doi.org/10.1038/srep46710>.
4. Zheng L, Meng J, Jiang K. et al. Improving protein–ligand docking and screening accuracies by incorporating a scoring function correction term. *Brief Bioinform* 2022;**23**. <https://doi.org/10.1093/bib/bbac051>.
5. Shen C, Zhang X, Deng Y. et al. Boosting protein–ligand binding pose prediction and virtual screening based on residue-atom distance likelihood potential and graph transformer. *J Med Chem* 2022;**65**:10691–706. <https://doi.org/10.1021/acs.jmedchem.2c00991>.
6. Morrone JA, Weber JK, Huynh T. et al. Combining docking pose rank and structure with deep learning improves protein–ligand binding mode prediction over a baseline docking approach. *J Chem Inf Model* 2020;**60**:4170–9. <https://doi.org/10.1021/acs.jcim.9b00927>.
7. Su M, Yang Q, du Y. et al. Comparative assessment of scoring functions: the CASF-2016 update. *J Chem Inf Model* 2019;**59**: 895–913. <https://doi.org/10.1021/acs.jcim.8b00545>.
8. Huang N, Shoichet BK, Irwin JJ. Benchmarking sets for molecular docking. *J Med Chem* 2006;**49**:6789–801. <https://doi.org/10.1021/jm0608356>.
9. Bohacek RS, McMartin C, Guida WC. The art and practice of structure-based drug design: a molecular modeling perspective. *Med Res Rev* 1996;**16**:3–50. [https://doi.org/10.1002/\(SICI\)1098-1128\(199601\)16:1<3::AID-MED1>3.0.CO;2-6](https://doi.org/10.1002/(SICI)1098-1128(199601)16:1<3::AID-MED1>3.0.CO;2-6).
10. Powers AS, Yu HH, Suriana P. et al. Geometric deep learning for structure-based ligand design. *ACS Cent Sci* 2023;**9**:2257–67. <https://doi.org/10.1021/acscentsci.3c00572>.
11. Libouban PY, Aci-Sèche S, Gómez-Tamayo JC. et al. The impact of data on structure-based binding affinity predictions using deep neural networks. *Int J Mol Sci* 2023;**24**. <https://doi.org/10.3390/ijms242216120>.
12. Wang R, Fang X, Lu Y. et al. The PDBbind database: methodologies and updates. *J Med Chem* 2005;**48**:4111–9. <https://doi.org/10.1021/jm048957q>.
13. Andrusier N, Mashiach E, Nussinov R. et al. Principles of flexible protein–protein docking. *Proteins* 2008;**73**:271–89. <https://doi.org/10.1002/prot.22170>.
14. Lexa KW, Carlson HA. Protein flexibility in docking and surface mapping. *Q Rev Biophys* 2012;**45**:301–43. <https://doi.org/10.1017/S0033583512000066>.
15. Fan H, Irwin JJ, Webb BM. et al. Molecular docking screens using comparative models of proteins. *J Chem Inf Model* 2009;**49**: 2512–27. <https://doi.org/10.1021/ci9003706>.
16. Pei Q, Wu L, Zhu J. et al. Breaking the barriers of data scarcity in drug-target affinity prediction. *Brief Bioinform* 2023;**24**. <https://doi.org/10.1093/bib/bbad386>.
17. Lin Z, Akin H, Rao R. et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science* 2023;**379**:1123–30. <https://doi.org/10.1126/science.ade2574>.
18. Chen L, Cruz A, Ramsey S. et al. Hidden bias in the DUD-E dataset leads to misleading performance of deep learning in structure-based virtual screening. *PLoS One* 2019;**14**:e0220113. <https://doi.org/10.1371/journal.pone.0220113>.
19. Shen C, Zhang X, Hsieh CY. et al. A generalized protein–ligand scoring framework with balanced scoring, docking, ranking and screening powers. *Chem Sci* 2023;**14**:8129–46. <https://doi.org/10.1039/D3SC02044D>.
20. Moon SH, Hwang SY, Lim J. et al. PIGNet2: a versatile deep learning-based protein–ligand interaction prediction model for binding affinity scoring and virtual screening. *Digital. Discovery* 2023;**3**:287–99. <https://doi.org/10.1039/D3DD00149K>.
21. Wang Z, Zheng L, Wang S. et al. A fully differentiable ligand pose optimization framework guided by deep learning and a traditional scoring function. *Brief Bioinform* 2023;**24**. <https://doi.org/10.1093/bib/bbac520>.
22. Oscar M-L, Mazen A, del Ehecatal AR-C. et al. A geometric deep learning approach to predict binding conformations of bioactive molecules. *Nat Mach Intell* 2021;**3**:1033–9. <https://doi.org/10.1038/s42256-021-00409-9>.
23. Korb O, Stutzle T, Exner TE. Empirical scoring functions for advanced protein–ligand docking with PLANTS. *J Chem Inf Model* 2009;**49**:84–96. <https://doi.org/10.1021/ci800298z>.
24. Wang C, Zhang Y. Improving scoring-docking-screening powers of protein–ligand scoring functions using random forest. *J Comput Chem* 2017;**38**:169–77. <https://doi.org/10.1002/jcc.24667>.
25. Schrödinger Inc., Glide.
26. Verdonk ML, Cole JC, Hartshorn MJ. et al. Improved protein–ligand docking using GOLD. *Proteins* 2003;**52**:609–23. <https://doi.org/10.1002/prot.10465>.

27. Trott O, Olson AJ. AutoDock Vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *J Comput Chem* 2010;**31**:455–61. <https://doi.org/10.1002/jcc.21334>.
28. Zhang X, Zhang O, Shen C. et al. Efficient and accurate large library ligand docking with KarmaDock. *Nat Comput Sci* 2023;**3**:789–804. <https://doi.org/10.1038/s43588-023-00511-5>.
29. McGibbon M, Money-Kyrle S, Blay V. et al. SCORCH: improving structure-based virtual screening with machine learning classifiers, data augmentation, and uncertainty estimation. *J Adv Res* 2023;**46**:135–47. <https://doi.org/10.1016/j.jare.2022.07.001>.
30. Li Y, Zhou D, Zheng G. et al. DyScore: a boosting scoring method with dynamic properties for identifying true binders and non-binders in structure-based drug discovery. *J Chem Inf Model* 2022;**62**:5550–67. <https://doi.org/10.1021/acs.jcim.2c00926>.
31. McNutt AT, Francoeur P, Aggarwal R. et al. GNINA 1.0: molecular docking with deep learning. *J Chem* 2021;**13**:43. <https://doi.org/10.1186/s13321-021-00522-2>.
32. Sunseri J, Koes DR. Virtual screening with Gnina 1.0. *Molecules* 2021;**26**. <https://doi.org/10.3390/molecules26237369>.
33. Marcou G, Rognan D. Optimizing fragment and scaffold docking by use of molecular interaction fingerprints. *J Chem Inf Model* 2007;**47**:195–207. <https://doi.org/10.1021/ci600342e>.
34. Tran-Nguyen VK, Bret G, Rognan D. True accuracy of fast scoring functions to predict high-throughput screening data from docking poses: the simpler the better. *J Chem Inf Model* 2021;**61**:2788–97. <https://doi.org/10.1021/acs.jcim.1c00292>.
35. Desaphy J, Raimbaud E, Ducrot P. et al. Encoding protein–ligand interaction patterns in fingerprints and graphs. *J Chem Inf Model* 2013;**53**:623–37. <https://doi.org/10.1021/ci300566n>.
36. Stepniewska-Dziubinska MM, Zielenkiewicz P, Siedlecki P. Development and evaluation of a deep learning model for protein–ligand binding affinity prediction. *Bioinformatics* 2018;**34**:3666–74. <https://doi.org/10.1093/bioinformatics/bty374>.
37. Zhou H, Cao H, Skolnick J. FRAGSITE: a fragment-based approach for virtual ligand screening. *J Chem Inf Model* 2021;**61**:2074–89. <https://doi.org/10.1021/acs.jcim.0c01160>.
38. Zhang W, Huang J. EVIS: an enhanced virtual screening approach based on pocket–ligand similarity. *J Chem Inf Model* 2022;**62**:498–510. <https://doi.org/10.1021/acs.jcim.1c00944>.
39. Zhou H, Cao H, Skolnick J. FINDSITE(comb2.0): a new approach for virtual ligand screening of proteins and virtual target screening of biomolecules. *J Chem Inf Model* 2018;**58**:2343–54. <https://doi.org/10.1021/acs.jcim.8b00309>.
40. Jain AN. Surfex: fully automatic flexible molecular docking using a molecular similarity-based search engine. *J Med Chem* 2003;**46**:499–511. <https://doi.org/10.1021/jm020406h>.
41. Brocchiacono M, Francoeur P, Aggarwal R. et al. BigBind: learning from nonstructural data for structure-based virtual screening. *J Chem Inf Model* 2024;**64**:2488–95. <https://doi.org/10.1021/acs.jcim.3c01211>.
42. Durrant JD, McCammon JA. NNScore 2.0: a neural-network receptor–ligand scoring function. *J Chem Inf Model* 2011;**51**:2897–903. <https://doi.org/10.1021/ci2003889>.
43. Quiroga R, Villarreal MA, Vinardo: a scoring function based on Autodock Vina improves scoring, docking, and virtual screening. *PLoS One* 2016;**11**:e0155183. <https://doi.org/10.1371/journal.pone.0155183>.
44. Qing, Luo SW, Li HY, Zheng L. et al. Benchmarking reverse docking through AlphaFold2 human proteome. *bioRxiv*, 2023.
45. Ozturk H, Ozgur A, Ozkirimli E. DeepDTA: deep drug–target binding affinity prediction. *Bioinformatics* 2018;**34**:i821–9. <https://doi.org/10.1093/bioinformatics/bty593>.
46. Wang J, Wen NF, Wang C. et al. ELECTRA-DTA: a new compound–protein binding affinity prediction model based on the contextualized sequence encoding. *J Chem* 2022;**14**:14. <https://doi.org/10.1186/s13321-022-00591-x>.
47. Pan S, Xia L, Xu L. et al. SubMDTA: drug target affinity prediction based on substructure extraction and multi-scale features. *BMC Bioinformatics* 2023;**24**:334. <https://doi.org/10.1186/s12859-023-05460-4>.
48. Huang K, Fu T, Glass LM. et al. DeepPurpose: a deep learning library for drug–target interaction prediction. *Bioinformatics* 2021;**36**:5545–7. <https://doi.org/10.1093/bioinformatics/btaa1005>.
49. Wu H, Liu J, Jiang T. et al. AttentionMGT-DTA: a multi-modal drug–target affinity prediction using graph transformer and attention mechanism. *Neural Netw* 2024;**169**:623–36. <https://doi.org/10.1016/j.neunet.2023.11.018>.
50. Zhao Q, Duan G, Yang M. et al. AttentionDTA: drug–target binding affinity prediction by sequence-based deep learning with attention mechanism. *IEEE/ACM Trans Comput Biol Bioinform* 2023;**20**:852–63. <https://doi.org/10.1109/TCBB.2022.3170365>.
51. Karimi M, Wu D, Wang Z. et al. DeepAffinity: interpretable deep learning of compound–protein affinity through unified recurrent and convolutional neural networks. *Bioinformatics* 2019;**35**:3329–38. <https://doi.org/10.1093/bioinformatics/btz111>.
52. Li M, Lu Z, Wu Y. et al. BACPI: a bi-directional attention neural network for compound–protein interaction and binding affinity prediction. *Bioinformatics* 2022;**38**:1995–2002. <https://doi.org/10.1093/bioinformatics/btac035>.
53. Li S, Wan F, Shu H. et al. MONN: a multi-objective neural network for predicting compound–protein interactions and affinities. *Cell Systems* 2020;**10**:308–322.e11. <https://doi.org/10.1016/j.cels.2020.03.002>.
54. Tsubaki M, Tomii K, Sese J. Compound–protein interaction prediction with end-to-end learning of neural networks for graphs and sequences. *Bioinformatics* 2019;**35**:309–18. <https://doi.org/10.1093/bioinformatics/bty535>.
55. Pinzi L, Rastelli G. Molecular docking: shifting paradigms in drug discovery. *Int J Mol Sci* 2019;**20**. <https://doi.org/10.3390/ijms20184331>.
56. Bugnon LA, Fenoy E, Edera AA. et al. Transfer learning: the key to functionally annotate the protein universe. *Patterns (N Y)* 2023;**4**:100691. <https://doi.org/10.1016/j.patter.2023.100691>.
57. Bloore DA, Kim JC, Kapoor K. et al. Protein Language Models Enable Accurate Cryptic Ligand Binding Pocket Prediction. *arXiv*, 2024.
58. Landrum GA, Riniker S. Combining IC(50) or K(i) values from different sources is a source of significant noise. *J Chem Inf Model* 2024;**64**:1560–7. <https://doi.org/10.1021/acs.jcim.4c00049>.
59. Chen L, Fan Z, Chang J. et al. Sequence-based drug design as a concept in computational drug design. *Nat Commun* 2023;**14**:4217. <https://doi.org/10.1038/s41467-023-39856-w>.
60. Sadybekov AA, Sadybekov AV, Liu Y. et al. Synthon-based ligand discovery in virtual libraries of over 11 billion compounds. *Nature* 2022;**601**:452–9. <https://doi.org/10.1038/s41586-021-04220-9>.
61. Wu ZX, Xiong Y, Yu S. et al. Unsupervised Feature Learning Via Non-Parametric Instance-Level Discrimination. *arXiv*, 2018.
62. Brody, S, Alon U, Yahav, E., How Attentive Are Graph Attention Networks?. *arXiv*, 2021.
63. Ong EV, Veličković P. Learnable Commutative Monoids for Graph Neural Networks. *arXiv*, 2022.
64. Dao T. FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning. *arXiv*, 2023.
65. Liu T, Lin Y, Wen X. et al. BindingDB: a web-accessible database of experimentally determined protein–ligand binding affinities.

- Nucleic Acids Res* 2007;**35**:D198–201. <https://doi.org/10.1093/nar/gkl999>.
66. Lam HYI, Pincket R, Han H. et al. Application of variational graph encoders as an effective generalist algorithm in computer-aided drug design. *Nature Machine Intelligence* 2023;**5**:754–64. <https://doi.org/10.1038/s42256-023-00683-9>.
 67. Bauer MR, Ibrahim TM, Vogel SM. et al. Evaluation and optimization of virtual screening workflows with DEKOIS 2.0 – a public library of challenging docking benchmark sets. *J Chem Inf Model* 2013;**53**:1447–62. <https://doi.org/10.1021/ci400115b>.
 68. Mysinger MM, Carchia M, Irwin JJ. et al. Directory of useful decoys, enhanced (DUD-E): better ligands and decoys for better benchmarking. *J Med Chem* 2012;**55**:6582–94. <https://doi.org/10.1021/jm300687e>.
 69. Tran-Nguyen VK, Jacquemard C, Rognan D. LIT-PCBA: an unbiased data set for machine learning and virtual screening. *J Chem Inf Model* 2020;**60**:4263–73. <https://doi.org/10.1021/acs.jcim.0c00155>.
 70. Davis MI, Hunt JP, Herrgard S. et al. Comprehensive analysis of kinase inhibitor selectivity. *Nat Biotechnol* 2011;**29**:1046–51. <https://doi.org/10.1038/nbt.1990>.
 71. Huang K, Fu T, Gao W. et al. Artificial intelligence foundation for therapeutic science. *Nat Chem Biol* 2022;**18**:1033–6. <https://doi.org/10.1038/s41589-022-01131-2>.
 72. Truchon JF, Bayly CI. Evaluating virtual screening methods: good and bad metrics for the “early recognition” problem. *J Chem Inf Model* 2007;**47**:488–508. <https://doi.org/10.1021/ci600426e>.
 73. Varadi M, Anyango S, Deshpande M. et al. AlphaFold protein structure database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Res* 2022;**50**:D439–44. <https://doi.org/10.1093/nar/gkab1061>.
 74. Yan X, Lu Y, Li Z. et al. PointSite: a point cloud segmentation tool for identification of protein ligand binding atoms. *J Chem Inf Model* 2022;**62**:2835–45. <https://doi.org/10.1021/acs.jcim.1c01512>.
 75. Schrödinger Inc. SiteMap.
 76. Hartshorn MJ, Verdonk ML, Chessari G. et al. Diverse, high-quality test set for the validation of protein–ligand docking performance. *J Med Chem* 2007;**50**:726–41. <https://doi.org/10.1021/jm061277y>.
 77. Luo Q, Zhao L, Hu J. et al. The scoring bias in reverse docking and the score normalization strategy to improve success rate of target fishing. *PloS One* 2017;**12**:e0171433. <https://doi.org/10.1371/journal.pone.0171433>.
 78. Alhossary A, Handoko SD, Mu Y. et al. Fast, accurate, and reliable molecular docking with QuickVina 2. *Bioinformatics* 2015;**31**:2214–6. <https://doi.org/10.1093/bioinformatics/btv082>.