# AI-Bind: Improving Binding Predictions for Novel Protein Targets and Ligands

Ayan Chatterjee[1,2],     Robin Walters[3],     Zohair Shafi[3],     Omair Shafi Ahmed [3],

Michael Sebek[1],     Deisy Gysi[1,4],     Rose Yu[5],     Tina Eliassi-Rad[2,3] ,

Albert-László Barabási[1,4,6],     Giulia Menichetti[1,4†]

[1]Center for Complex Network Research, Northeastern University, Boston, USA
[2]Network Science Institute, Northeastern University, Boston, USA
[3]Khoury College of Computer Sciences, Northeastern University, Boston, USA
[4]Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, USA
[5]Department of Computer Science and Engineering, University of California, San Diego, USA
[6]Center for Network Science, Central European University, Budapest, Hungary

## 1   ABSTRACT

Identifying novel drug-target interactions (DTI) is a critical and rate limiting step in drug discovery. While deep learning models have been proposed to accelerate the identification process, we show that state-of-the-art models fail to generalize to novel (i.e., never-before-seen) structures. We first unveil the mechanisms responsible for this shortcoming, demonstrating how models rely on shortcuts that leverage the topology of the protein-ligand bipartite network, rather than learning the node features. Then, we introduce AI-Bind, a pipeline that combines network-based sampling strategies with unsupervised pre-training, allowing us to limit the annotation imbalance and improve binding predictions for novel proteins and ligands. We illustrate the value of AI-Bind by predicting drugs and natural compounds with binding affinity to SARS-CoV-2 viral proteins and the associated human proteins. We also validate these predictions via docking simulations and comparison with recent experimental evidence, and step up the process of interpreting machine learning prediction of protein-ligand binding by identifying potential active binding sites on the amino acid sequence. Overall, AI-Bind offers a powerful high-throughput approach to identify drug-target combinations, with the potential of becoming a powerful tool in drug discovery.

---

[†]Corresponding author. e-mail: menicgiulia@gmail.com

## 2    INTRODUCTION

The accurate prediction of binding interactions between chemicals and proteins is a critical step in drug discovery, necessary to identify new drugs and novel therapeutic targets, and to reduce the failure rate in clinical trials, and to predict the safety of drugs [1–3]. While molecular dynamics and docking simulations [4–6] are frequently employed to identify potential protein-ligand binding, the computational complexity (namely, run-times) of the simulations and the lack of 3D protein structures significantly limit the coverage and the feasibility of large-scale testing. Therefore, machine learning (ML) and artificial intelligence (AI) based models have been proposed to circumvent the computational limitations of the existing approaches [7, 8], leading to the development of models that rely either on deep learning architectures or chemical feature representations [9–15].

Deep learning frameworks formulate the binding prediction problem as either a binary classification task or a regression task. The successful training of a binary classifier requires positive samples, pairs of proteins and ligands that are known to bind to each other, typically extracted from protein-ligand binding databases like DrugBank [16], BindingDB [17], Tox21 [18], ChEMBL [19], Davis [20], or Drug Target Commons [21]. Training also requires negative samples, i.e., pairs that do not interact or only weakly interact. However, the positive and negative annotations associated with different proteins and ligands are not evenly distributed, but some proteins and ligands have disproportionately more positive annotations than negative ones, and vice-versa, an annotation imbalance learned by the ML models, which then predict that some proteins and ligands bind disproportionately more often than the others. In other words, the ML models learn the binding patterns from the degree of the nodes in the protein-ligand interaction network, neglecting relevant node metadata, like the chemical structures of the ligands or the amino-acid sequences of the proteins [9, 22–24]. This annotation imbalance leads to good performance as quantified by the Area Under the Receiver Operating Characteristics (AUROC) and the Area Under the Precision Recall Curve (AUPRC) for the unknown annotations associated with missing links in the protein-ligand interaction network used for training. A key signal of such shortcut learning is the degradation of the performance of an ML model when asked to predict binding between novel (i.e., never-before-seen) protein targets and ligands. This modeling limitation is in-line with the findings of Geirhos et al. [25], who showed that deep learning methods tend to exploit shortcuts in training data to achieve good performance. Laarhoven et al. discuss similar bias in drug–target interaction data and its effect on cross-validation performance [26]. Lee et al. [27] and Wang et al. [28] proposed approaches that partly address shortcut learning, but fail to generalize to unexplored proteins, i.e., proteins that lack

sufficient binding annotations, or originate from organisms with no close relatives in current protein databases. More recently, models such as MolTrans [29], MONN [30], KGE_NFM [31], TransDTI [32], HoTS [33], and DTIHNC [34], explore innovative structural representations of protein and ligand molecules. Though these models better leverage the molecular structures to predict binding, end-to-end training limits their ability to generalize beyond the molecular scaffolds present in the training data.

Here, we introduce AI-Bind, a pipeline for predicting protein-ligand binding which can successfully generalize to unseen proteins and ligands. AI-Bind combines network science methods with unsupervised pre-training to control for the over-fitting and the annotation imbalance of existing libraries. We leverage the notion of shortest path distance on a network to identify distant protein-ligand pairs as negative samples. Combining these network-derived negatives with experimentally validated non-binding protein-ligand pairs, we ensure sufficient positive and negative samples for each node in the training data. Additionally, AI-Bind learns, in an unsupervised fashion, the representation of the node features, i.e., the chemical structures of ligand molecules or the amino-acid sequences of protein targets, helping circumvent the model's dependency on limited binding data. Instead of training the deep neural networks in an end-to-end fashion using binding data, we pre-train the embeddings for proteins and ligands using larger chemical libraries, allowing us to generalize the prediction task to chemical structures, beyond those present in the training data.

## 3  RESULTS

### Limitations of existing ML models

ML models characterize the likelihood of each node (proteins and ligands) to bind to other nodes according to the features and annotations in the training data. While annotations capture known protein-ligand interactions, features refer to the chemical structures of proteins and ligands which determine their physical and chemical properties, and are expressed as amino acid sequences or 3D structures for proteins, and chemical SMILES [35] for ligands. In an ideal scenario, the ML model learns the patterns characterizing the features which drive the protein-ligand interactions, capturing the physical and chemical properties of a protein and of a ligand that determine the mutual binding affinity. Yet, as we show next, multiple state-of-the-art deep learning models, such as DeepPurpose [9], ignore the features and rely largely on annotations, i.e., the degree information for each protein and ligand in the drug-target interaction (DTI) network, as a shortcut to make new binding predictions. A bipartite network represents the binding information as a graph with two different types of *nodes*: one corresponding to pro-

3

teins (also called targets, representing for example a human or a viral protein) and the other corresponding to ligands (representing potential drugs or natural compounds), respectively. A protein-ligand annotation, i.e., evidence that a ligand binds to a protein, is represented as a *link* between the protein and the ligand in the bipartite network [36]. Experimentally validated annotations define the known DTI network. While binding depends only on the detailed chemical characteristics of the nodes (proteins and ligands), as we show here, many current ML models predictions are primarily driven by the topology of the DTI network. We begin by noticing that the number of annotations linked to a protein or a ligand follows a fat-tailed distribution [37], indicating that the vast majority of proteins and ligands have only a small number of annotations, which then coexist with a few *hubs*, nodes with an exceptionally large number of binding records [36]. For example, the number of annotations for proteins follows a power law distribution with degree exponent $\gamma_p = 2.84$ in the BindingDB data used for training and testing DeepPurpose, while the ligands have a degree exponent $\gamma_l = 2.94$ (Figure 1A) [38]. For these degree exponents, the second moment of the distribution diverges for large sample sizes, implying that the expected uncertainty in the binding information is highly significant, limiting our ability to predict the binding between a single protein and ligand [36, 39]. Furthermore, positive and negative annotations are determined by applying a threshold on kinetic constants like the constant of disassociation $K_d$. If the kinetic constant associated with a protein-ligand pair is less than a set threshold, we consider that pair as a positive or binding sample; otherwise, the pair is tagged as negative or non-binding. However, $K_d$ is not randomly distributed across the records, but the number of annotations $k$ and the average $K_d$ per $k$ (i.e., $\langle K_d \rangle$) are anti-correlated, indicating stronger binding propensity for proteins and ligands with more annotations ($r_{Spearman}(k_p, \langle K_d \rangle) = -0.47$ for proteins, $r_{Spearman}(k_l, \langle K_d \rangle) = -0.29$ for ligands in the BindingDB data used by DeepPurpose). Alongside this negative correlation, we observe higher variance for $\langle K_d \rangle$ values given a fixed $k$ for the low degree nodes compared to the hubs. As the annotations follow fat-tailed distributions, the observed anti-correlation drives the hub proteins and ligands to have disproportionately more binding records on average, whereas proteins and ligands with fewer annotations have both binding and non-binding examples. This *annotation imbalance* prompts the ML models to leverage degree information (positive and negative annotations) in making binding prediction instead of learning binding patterns from the molecular structures. We term this phenomenon as the *emergence of topological shortcuts* (see Section S1).

To investigate the emergence of topological shortcuts, for each node $i$ with number of anno-

tations $k_i$, we quantify the balance of the available training information via the *degree ratio*,

$$\rho_i = \frac{k_i^+}{k_i^+ + k_i^-} = \frac{k_i^+}{k_i}, \tag{1}$$

where, $k_i^+$ is the positive degree, corresponding to the number of known binding annotations in the training data, and $k_i^-$ is the negative degree, or the number of known non-binding annotations in the training data (Figure 1C). As most proteins and ligands lack either binding or non-binding annotations for (Table 1), the resulting $\{\rho_i\}$ are close to 1 or 0 (Figs. 2A, B), with extreme $\rho$ values representing the annotation imbalance in the prediction problem. As many state-of-the-art deep learning models, such as DeepPurpose [9], uniformly sample the available positive and negative annotations, they assign higher binding probability to proteins and ligands with higher $\rho$ (Figure 2C, D). Consequently, their binding predictions are driven by topological shortcuts in the protein-ligand network, which are associated with the the positive and negative annotations present in the training data rather than the structural features characterizing proteins and ligands.

The higher binding predictions in DeepPurpose for proteins with large degree ratios (Figure 2C) prompted us to compare the performance of DeepPurpose with network configuration models, algorithms that ignore the features of proteins and ligands and instead predict the likelihood of binding by leveraging only topological constraints derived from the network degree sequence [36, 40–42]. In the configuration model (Figure 3A, Methods), the probability of observing a link is determined only by the the degrees of its end nodes. In a 5-fold cross-validation on the benchmark BindingDB dataset (Table 1), we find that the top-performing DeepPurpose architecture, Transformer-CNN [9], achieves AUROC of 0.85 ($\pm$ 0.005) and AUPRC of 0.65 ($\pm$ 0.008). At the same time, the network configuration model on the same data achieves an AUROC of 0.86 ($\pm$ 0.005) and AUPRC of 0.61 ($\pm$ 0.008) (Figure 3B). A consistent behavior is observed for DeepPurpose under node attribute reshuffling, suggesting that both deep neural networks and the configuration model leverage mainly topological information (see Table 3). In other words, the network configuration model, relying only on annotations, performs just as well as the deep learning model, confirming that the topology of the protein-ligand interaction network drives the prediction task. The major driving factor of the topological shortcuts is the monotone relation between $k$ and $\langle K_d \rangle$, which associates a link type with the degree of its end nodes as the $\langle K_d \rangle$ values are directly associated with the link types after thresholding. Moreover, in BindingDB we observe that hubs encounter less variance for $\langle K_d \rangle$ compared to the low degree nodes. Thus, the configuration model is able to achieve good test performance in predicting the link types associated with the hubs. Since hub nodes are associated with the

majority of the links in the protein-ligand bipartite network, the configuration model achieves excellent test performance by making correct predictions that mainly leverage the degree information of the hubs. To further investigate this hypothesis, we tested three distinct scenarios: (i) unseen edges (Transductive test), when both proteins and ligands from the test dataset are present in the training data; (ii) unseen targets (Semi-inductive test), when only the ligands from the test dataset are present in the training data; (iii) unseen nodes (Inductive test), when both proteins and ligands from the test dataset are absent in the training data.

We find that both DeepPurpose and the configuration model perform well in scenarios (i) and (ii) (Figures 3C, D). However, for the inductive test scenario (iii), when confronted with new proteins and ligands, both performances drop significantly (Table 2). DeepPurpose has an AUROC of 0.60 ($\pm$ 0.066) and AUPRC of 0.42 ($\pm$ 0.063), comparable to the configuration model, for which we have AUROC of 0.50 and AUPRC of 0.30 ($\pm$ 0.034). To offer a final piece of evidence that DeepPurpose disregards node features, we randomly shuffled the chemical SMILES [35] and amino acid sequences in the training set, while keeping the same positive and negative annotations per node, an operation that did not change the test performance (Table 3). These tests confirm that DeepPurpose leverages network topology as a learning shortcut and fails to generalize predictions to proteins and ligands beyond the training data, indicating that we must use inductive testing to evaluate the true performance of ML models.

Beyond DeepPurpose, models such as MolTrans [29] explore different structural representations of protein and ligand molecules. We investigated transductive, semi-inductive, and inductive performances for MolTrans, a state-of-the-art protein-ligand binding prediction model which uses a combination of sub-structural pattern mining algorithm, interaction modeling module, and an augmented transformer encoder to better learn the molecular structures (see Section S8). While the innovative representation of the molecules improves upon DeepPurpose in transductive tests (AUROC of 0.952 ($\pm$ 0.051), AUPRC of 0.872 ($\pm$ 0.131)), the same representation still relies only on the training DTI and fails to generalize to novel molecular structures, as captured by the poor performance in inductive tests (AUROC of 0.575 ($\pm$ 0.059), AUPRC of 0.430 ($\pm$ 0.098)).

## AI-Bind and statistics across models

AI-Bind is a deep learning pipeline that combines network-derived learning strategies with unsupervised pre-trained node features, to optimize the exploration of the binding properties of novel proteins and ligands. Our pipeline is compatible with various neural architectures, three of which we propose here: VecNet, Siamese model, and VAENet. AI-Bind uses two inputs (Figure

4A): For ligands, it takes as input isomeric SMILES, which capture the structures of ligand molecules. AI-Bind considers a search-space consisting of all the drug molecules available in DrugBank and the naturally occurring compounds in the Natural Compounds in Food Database (NCFD) (see Section S4.4), and can be extended by leveraging larger chemical libraries like PubChem [43]. For proteins, AI-Bind uses as input the amino acid sequences retrieved from the protein databases Protein Data Bank (PDB) [44], the Universal Protein knowledgebase (UniProt) [45], and GeneCards [46].

AI-Bind benefits from several novel features compared to the state-of-the-art: (a) It relies on network-derived negatives to balance the number of positive and negative samples for each protein and ligand. To be specific, it uses protein-ligand pairs with shortest path distance $\geq 7$ as negative samples, ensuring that the neural networks observe both binding and non-binding examples for each protein and ligand (see Figure 5, Methods, Section S5). (b) During unsupervised pre-training, AI-Bind uses the node embeddings trained on larger collections of chemical and protein structures, compared to the set with known binding annotations, allowing AI-Bind to learn a wider variety of structural patterns. Indeed, while models like DeepPurpose were trained on 862,337 ligands and 7,504 proteins provided in BindingDB, or 7,307 ligands and 4,762 proteins provided in DrugBank, the unsupervised representation in AI-Bind's VecNet is trained on 19.9 million compounds from ZINC [47] and ChEMBL [48] databases, and on 546,790 proteins from Swiss-Prot [49].

We begin the model's validation by systematically comparing the performance of AI-Bind to DeepPurpose and the configuration model on a 5-fold cross-validation using the network-derived dataset for transductive, semi-inductive, and inductive tests. AI-Bind's VecNet model uses pre-trained `mol2vec` [50] and `protvec` [51] embeddings combined with a simple multi-layer perceptron [52] to learn protein-ligand binding (Figure 4B, see Methods). We find that the configuration model performs poorly in inductive testing (AUROC 0.5, AUPRC $0.469 \pm 0.014$). Due to the network-derived negatives that remove the annotation imbalance, DeepPurpose shows improved performance for novel proteins and ligands (AUROC $0.642 \pm 0.025$, AUPRC $0.583 \pm 0.016$). The best performance on unseen nodes is observed for AI-Bind's VecNet, with AUROC of $0.745 \pm 0.032$ and AUPRC of $0.729 \pm 0.038$ (see Figure 4C and see Table S3 for a summary of the performances). The unsupervised pre-training for ligand embeddings allows us to generalize AI-Bind to naturally occurring compounds, characterized by complex chemical structures and fewer training annotations compared to drugs (see Section S2), obtaining performances comparable to those obtained for drugs (Figure 4D).

Beyond DeepPurpose, AI-Bind's VecNet consistently achieves better inductive performance

(AUROC $0.745 \pm 0.032$, and AUPRC $0.729 \pm 0.038$) compared to MolTrans (AUROC $0.619 \pm 0.021$, and AUPRC $0.480 \pm 0.028$). The comparison between AI-Bind and state-of-the-art models like DeepPrupose and MolTrans validates how unsupervised pre-training of the molecular embeddings improves the generalizability of binding prediction models (see Section S8).

## Validation of AI-Bind predictions on COVID-19 proteins

For a better understanding of the reliability of the AI-Bind predictions, we move beyond standard ML cross-validation and compare our predictions with molecular docking simulations, and in vitro and clinical results on protein-ligand binding. Docking simulations offer a reliable but computationally complex method to predict (or validate) binding between proteins and ligands [53]. Motivated by the need to model rapid response to sudden health crises, we chose as our validation set the 26 SARS-CoV-2 viral proteins and the 332 human proteins targeted by the SARS-CoV-2 viral proteins [54, 55]. These proteins are missing from the training data of AI-Bind, hence represent novel targets and allow us to rely on recent efforts to understand the biology of COVID-19 to validate the AI-Bind predictions. We could retrieve the amino acid sequences in FASTA format for 16 SARS-CoV-2 viral proteins and 330 human proteins from UniProt [45], and use them as input to AI-Bind's VecNet. Binding between viral and human proteins is necessary for the virus to synthesize its own viral proteins and to facilitate its replication. Our goal is to predict drugs in DrugBank or naturally occurring compounds that can bind to any of the 16 SARS-CoV-2 or 330 human proteins associated with COVID-19, potentially disrupting the viral infection. After sorting all protein-ligand pairs based on their binding probability predicted by AI-Bind's VecNet ($p_{ij}^{VecNet}$), we tested the predicted top 100 and bottom 100 binding interactions with blind docking simulations using AutoDock Vina [53], which estimates binding affinity by considering all possible binding locations on the 3D protein structures (see Methods). Of the 54 proteins present in the top 100 and bottom 100 predicted pairs, 23 had 3D structures available in PDB [44] and UniProt [45], and 51 of the 59 involved ligand structures available on PubChem [43], allowing us to perform 128 docking simulations (84 involving the top and 44 involving the bottom predictions). We find that 74 out of 84 top predictions from AI-Bind are indeed validated binding pairs. Furthermore, we find that the median binding affinity for the top VecNet predictions is $-7.65$ kcal/mole, while for the bottom ones is $-3.0$ kcal/mole (Figure 6A), confirming that for AI-Bind, the top predictions show significantly higher binding propensity than the bottom ones (Kruskal-Wallis H-test p-value of $2.5*10^{-5}$) [56,57]. As a second test, we obtained the binary labels (binding or non-binding) from docking and AI-Bind predictions using the threshold of $-1.75$ kcal/mole for binding affinities [58]

and the optimal threshold on $p_{ij}^{VecNet}$ corresponding to the highest F1-Score on the inductive test set (see Section S7, Figure S12). In the derived confusion matrix we observe sensitivity = 0.76, representing the fraction of binding predictions made by AI-Bind that are true binders, i.e., the ratio $True\ Positives/(True\ Positives + False\ Negatives)$, and F1-Score = 0.82. These two numbers confirm that the rank list provided by AI-Bind predictions shows a significant similarity to the rank list obtained by binding affinities compared to a random selection (Figure 6B). We further check the stability of these performance metrics by randomly choosing 20 protein-ligand pairs in a 5-fold bootstrapping set-up and observe F1-Score = $0.90 \pm 0.02$. Additionally, we find that the predictions made by AI-Bind's VecNet ($p_{ij}^{VecNet}$) and the free energy of protein-ligand binding obtained from docking ($\Delta G$) are anti-correlated with $r_{Spearman}(p_{ij}^{VecNet}, \Delta G) = -0.51$. The top 20 VecNet predictions show $r_{Spearman}(p_{ij}^{VecNet}, \Delta G) = -0.17$. As lower binding affinity values correspond to stronger binding, these results document the agreement between AI-Bind predictions and docking simulations.

Among the 50 ligands with the highest average binding probability we find two FDA-approved drugs Anidulafungin (NDA#021948) and Cyclosporine (ANDA#065017). Experimental evidence [59] shows that these drugs have anti-viral activity at very low concentrations in the dose-response curves, and have $IC_{50}$ values of 4.64 $\mu M$ and 5.82 $\mu M$ (see Figure S1), respectively, measured by immunofluorescence analysis with an antibody specific for the viral N protein of SARS-CoV-2. These low $IC_{50}$ values support anti-viral activity, confirming that Anidulafungin and Cyclosporine bind to COVID-19 related proteins [60], and the activity at low concentrations indicate that they are safe to use for treating COVID-19 patients [61]. Anidulafungin binds to the SARS-CoV-2 viral protein nsp12, a key therapeutic target for coronaviruses [62].

AI-Bind also offers several novel predictions with potential therapeutic relevance. For example, it predicts that the naturally occurring compounds Spironolactone, Oleanolic acid, and Echinocystic acid are potential ligands for COVID-19 proteins, all three ligands binding to Tripartite motif-containing protein 59 (TRIM59), a human protein to which the SARS-CoV-2 viral proteins ORF3a and NSP9 bind [63, 64]. AutoDock Vina supports these predictions, offering binding affinities -7.1 kcal/mole, -8.0 kcal/mole, and -7.6 kcal/mole, respectively.

Spironolactone, found in rainbow trout [65], has been suggested to reduce COVID susceptibility [66–69]. Oleanolic acid is present in apple, tomato, strawberry, and peach, and has been proposed as a potential anti-viral agent for COVID-19 [70, 71]. Oleanolic acid, which passed the drug efficacy benchmark ADME (Absorption, Distribution, Metabolism, and Excretion), plays an important role in controlling viral replication of SARS-CoV-2 [72] and is effective in

preventing virus entry at low viral loads [71]. Finally, Echinocystic acid, found in sunflower, basil, and gala apples, is known for its anti-inflammatory [73–75] and anti-viral activity [76,77], but its potential anti-viral role in COVID-19 is yet to be validated.

## Identifying active binding sites

Beyond predicting binding probability, AI-Bind can also be used to identify the probable active binding sites on the amino acid sequence, even in absence of a 3D protein structure. Specifically, we can use AI-Bind to identify which trigrams in the amino acid sequence play the most significant role in binding predictions, indicative of potential protein-ligand binding locations. We perturb each amino acid trigram in the sequence and observe the changes in AI-Bind prediction (see Section S9). Valleys in the obtained binding probability profile represent the trigrams most predictive of binding locations on the amino acid sequence. To validate the AI-Bind predicted binding sites we focus on the human protein TRIM59, a protein for which we have results from multiple docking simulations, using PyMOL [78] and identified the amino acid residues binding to the ligand molecules (Figure 6C). We find that the amino acid residues responsible for binding directly map to the valleys in the binding probability profile identified by AI-Bind. By visualizing the docking results for Pipecuronium, Buprenorphine and Voclosporin, ligands that bind to three different pockets on TRIM59, we mark the valleys corresponding to the respective binding sites on the binding probability profiles (Figure 6C). For example, pocket 1, where Pipecuronium binds, corresponds to four AI-Bind predicted valleys marked by 1A, 1B, 1C and 1D.

Since not all valleys in the binding probability profile map to binding sites that we could match with ligands, we use the protein secondary structure to prioritize the valleys. We predict the secondary structure from the amino acid sequence using S4PRED [79] and identify the regions with $\alpha$-helix, $\beta$-sheet and coil. In particular, $\alpha$-helices prefer non-solvent accessible environments [80], contain non-polar amino acid residues [81], and consist of weaker inter-molecular interactions [82]. Thus, helices show less propensity for protein-ligand binding. In particular, $\alpha$-helices prefer non-solvent accessible environments [80], contain non-polar amino acid residues [81], and consist of weaker inter-molecular interactions [82]. Thus, the presence of alpha-helices reduce the chances of binding between a ligand and a protein. In contrast, $\beta$-sheets and non-regular coil regions (unstructured regions) are preferred by ligands as active binding sites since they provide more binding opportunity to other molecules [83]. Indeed, most of the valleys in Figure 6C where ligands bind map to $\beta$-sheets and coils on TRIM59, associated with pockets 1 and 2 (27 out of 34 ligands validated by docking). Valleys which have large overlap

with the $\beta$-sheets and coils provide most of the predicted binding. By combining the binding probability profile predicted by AI-Bind and the secondary structure predicted by S4PRED, we can create an optimal search grid for the subsequent docking simulations, drastically reducing runtime.

We pursued further validation of AI-Bind predicted binding sites with a gold standard protein binding dataset [84] and with p2rank, another state-of-the-art binding prediction model [85], to extensively assess the reliability of the AI-Bind pipeline (see Section S13).

In summary, ML models often fail in real world settings when making predictions on data that they were not explicitly trained upon despite achieving good test performance based on traditional ML-based metrics [86]. It is therefore necessary to validate the applicability of these models before deploying them. The documented validation of the AI-Bind predictions with molecular dynamic simulations and in vitro experiments offers us confidence AI-Bind is an effective prioritization tool in diverse settings.

## 4  DISCUSSION

The accurate prediction of drug-target interactions is an essential precondition of drug discovery. Here we showed that by taking topological shortcuts, existing deep learning models significantly limit their predictive power. Indeed, a mechanistic and quantitative understanding of the origins of these shortcuts indicates that uniform sampling in presence of annotation imbalance drives ML models to disregard the features of proteins and ligands, limiting their ability to generalize to novel protein targets and ligand structures. To address these shortcomings, we introduced a new pipeline, AI-bind, which mitigates the annotation imbalance of the training data by introducing network-derived negative annotations inferred via shortest path distance, and improves the transferability of the ML models to novel protein and ligand structures by unsupervised pre-training. The proposed unsupervised pre-training of node features also influences the quality of false predictions, removing potential structural biases towards specific protein families (see Section S10). Once we improved the statistical sampling of the training data and generated the node embeddings in an unsupervised fashion, we observed an increase in performance compared to DeepPurpose, resulting in commendable AUROC (24% improvement) and AUPRC (74% improvement) and, most importantly, an ability to predict beyond proteins and ligands present in the training dataset.

A major limitation of using binding predictions in drug discovery is that binding to disease-related protein targets does not always imply a therapeutic treatment. As future work, we plan to extend our implementation by introducing an ML-based classifier to sort the list of potential

ligands according to their pharmaceutical (therapeutic) effects, combining the current node features with additional metrics derived from traditional network medicine approaches [87, 88].

AI-Bind leverages ligands' Morgan fingerprints and proteins' amino acid sequences, which encode relevant properties of the molecules: from the presence of hydrogen donors, hydrogen acceptors, count of different atoms, chirality, and solubility for ligands, to the existence of R groups, N or C terminus in proteins. All these properties influence the mechanisms driving protein-ligand binding (see Section S11) [89]. Yet, the binding phenomenon is largely dependent on the 3D structures of the molecules, which determines the binding pocket structures and the rotation of the bonds. We plan to embed the 3D structures of protein and ligand molecules, which will take into account higher order molecular properties driving protein-ligand binding and refine the predictive power of AI-Bind. To maximize generalization across 3D structure, we will use SE(3) equivariant networks to learn embeddings. Equivariance has proven to be a powerful tool for improving generalization over molecular structure [90–93]. We also plan to explore the performance of AI-Bind over the entire druggable genome [94], allowing us to predict for each protein, which domains are responsible for the binding predictions, revealing binding locations of the ligands and the proteins. Finally, we envision enabling AI-Bind to predict the kinetic constants $K_d$, $K_i$, $IC_{50}$, and $EC_{50}$ by formulating a regression task over these variables.

The existing docking infrastructures allow screening for a specific protein structure against wide chemical libraries. Indeed, VirtualFlow [95], an open-source drug discovery platform offers virtual screening over more than 1.4 billion commercially available ligands. However, running docking simulations over these vast libraries incurs high costs for data preparation and computation time and are often limited to only proteins with 3D structures [44]. For example, in our validation step, only half (23 out of 54) of the 3D structures of the proteins associated with COVID-19 were available. Since AI-Bind only requires the chemical SMILES for ligands [35] and amino acid sequences for proteins, it can offer fast screening for large libraries of targets and molecules without requiring 3D structures, guiding the computationally expensive docking simulations on selected protein-ligand pairs.

## 5 METHODS

### Data Preparation

We use InChIKeys [96] and amino acid sequences as the unique identifiers for ligands and targets, respectively. Positive and negative samples are selected from DrugBank, BindingDB and DTC (see Section S4). We consider samples from BindingDB and DTC to be binding or non-binding based on the kinetic constants $K_i$, $K_d$, $IC_{50}$, and $EC_{50}$. We use thresholds of $\leq 10^3 nM$

and $\geq 10^6 nM$ to obtain positive and (absolute) negative annotations, respectively [58]. We then filter out all samples outside the temperature range 20°C-45°C to remove ambiguous pairs. All amino acid sequences were obtained from UniProt [45].

## Positive Samples

We consider the binding information from DrugBank as positive samples. From these annotations, we removed 53 pairs which are available in BindingDB and have kinetic constants $\geq 10^6 nM$. To obtain additional positive samples for drugs, we searched in BindingDB using their InChIKeys. We obtained 4,330 binding annotations from BindingDB related to the drugs in DrugBank. Overall, we gathered a total of 28,188 positive samples for drugs. We identified naturally occurring/food-borne compounds by leveraging the Natural Compounds in Food Database (NCFD) database (see Section S4.4). We queried BindingDB and DTC with the associated InChIKeys, obtaining a total of 1,555 positive samples.

## Network-Derived Negative Samples

To generate annotation-balanced training data for AI-Bind, we merged the positive annotations derived from DrugBank, BindingDB, and DTC, for a total of 5,104 targets and 8,111 ligands, of which 485 are naturally occurring, and calculated the shortest path distribution. All odd-path lengths in the bipartite network correspond to protein-ligand pairs (Figure 5C). Overall, the longer the shortest path distance separating a protein and a ligand, the higher the kinetic constant observed in BindingDB (Figure 5D). In particular, pairs more than 7 hops apart have, on average, kinetic constants $K_i \geq 10^6 nM$, which is generally considered above the protein-ligand binding threshold [58] (see Section S5). We randomly selected a subset of protein-ligand pairs which are 7 hops apart as negative samples, to create an overall class balance between positive and negative samples in the training data. Finally, we removed all nodes with only positive or only negative samples and obtained the *network-derived negative samples*. We performed testing and validation on $\geq 11$-hop distant pairs. Additionally, we included in testing and validation the absolute non-binding pairs derived from BindingDB by thresholding the kinetic constants ($K_i$, $K_d$, $IC_{50}$, and $EC_{50}$).

## Network Configuration Model

### Overview

Protein-ligand annotations are naturally embedded in a bipartite duplex network, consisting of a set of nodes, comprising all proteins and ligands, interacting in two layers, each reflecting a

distinct type of interaction linking the same pair of nodes [41]. More specifically, one layer (Layer 1) captures the positive or binding annotations, while the second layer (Layer 2) collects the negative or non-binding annotations (Figure 3A). A multilink $\mathbf{m}$ between two nodes encodes the pattern of links connecting these nodes in different layers. In particular, $\mathbf{m} = (1, 0)$ indicates positive interactions, $\mathbf{m} = (0, 1)$ refers to negative interactions, $\mathbf{m} = (0, 0)$ represents the absence of any type of annotations, and $\mathbf{m} = (1, 1)$ is mathematically forbidden, as binding and non-binding cannot coexist for the same pair of protein and ligand.

We developed a canonical bipartite duplex null model that conserves on average the number of positive and negative annotations of each node, while correctly rewiring positive and negative links and avoiding forbidden configurations. By means of entropy maximization with constraints, we derive the analytical formulation of each multilink probability and the conditional probability of observing positive binding once an annotation is reported.

**Mathematical Formulation**

Let $A_{ij}^{\mathbf{m}}$ be the multi-adjacency matrix representing the bipartite duplex of ligands ($\{i\}$) and proteins ($\{j\}$), with elements equal to 1 if there is a multilink $\mathbf{m}$ between $i$ and $j$ and zero otherwise. We define the multidegree of ligand $i$ and target $j$ as

$$k_i^{\mathbf{m}} = \sum_{j=1}^{N_T} A_{ij}^{\mathbf{m}}, \qquad t_j^{\mathbf{m}} = \sum_{i=1}^{N_L} A_{ij}^{\mathbf{m}}, \tag{2}$$

where $N_T$ is the number of targets and $N_L$ is the number of ligands.

A bipartite duplex network ensemble can be defined as the set of all duplexes satisfying a given set of constraints, such as the expected multidegree sequences defined in Eq. 2. We determine the probability of observing a bipartite duplex network $P(\vec{G})$ by entropy maximization with multidegree constraints $\{k_i^{(1,0)}\}$, $\{k_i^{(0,1)}\}$, $\{t_j^{(1,0)}\}$, and $\{t_j^{(0,1)}\}$, and corresponding Lagrangian multipliers $\{\lambda_i^{(1,0)}\}$, $\{\lambda_i^{(0,1)}\}$, $\{\mu_j^{(1,0)}\}$, and $\{\mu_j^{(0,1)}\}$ [41,42,97]. The probability $P(\vec{G})$ factorizes as

$$P(\vec{G}) = \frac{1}{Z} \prod_{ij} \exp\left[ - \sum_{\mathbf{m} \neq (0,0),(1,1)} (\lambda_i^{\mathbf{m}} + \mu_j^{\mathbf{m}}) A_{ij}^{\mathbf{m}} \right], \tag{3}$$

with

$$Z = \prod_{ij} \left[ 1 + \sum_{\mathbf{m} \neq (0,0),(1,1)} e^{-(\lambda_i^{\mathbf{m}} + \mu_j^{\mathbf{m}})} \right]. \tag{4}$$

Multilink probabilities $p_{ij}^{\mathbf{m}}$ are determined by the derivatives of $\log(Z)$ according to $(\lambda_i^{\mathbf{m}} +$

$\mu_j^{\mathbf{m}}$). For instance, the probability of observing a positive annotation is

$$p_{ij}^{(1,0)} = \frac{e^{-(\lambda_i^{(1,0)}+\mu_j^{(1,0)})}}{1 + e^{-(\lambda_i^{(1,0)}+\mu_j^{(1,0)})} + e^{-(\lambda_i^{(0,1)}+\mu_j^{(0,1)})}}, \tag{5}$$

while the probability of observing a negative annotation follows

$$p_{ij}^{(0,1)} = \frac{e^{-(\lambda_i^{(0,1)}+\mu_j^{(0,1)})}}{1 + e^{-(\lambda_i^{(1,0)}+\mu_j^{(1,0)})} + e^{-(\lambda_i^{(0,1)}+\mu_j^{(0,1)})}}, \tag{6}$$

with $p_{ij}^{(1,0)} + p_{ij}^{(0,1)} + p_{ij}^{(0,0)} = 1$.

In this theoretical framework, binding prediction is inherently conditional, as for each ligand $i$ and protein $j$, we test only the presence of positive and negative annotations. Consequently, $p_{ij}^{(1,0)}$ and $p_{ij}^{(0,1)}$ are normalized by the probability of observing a generic annotation $p_{ij}^{(1,0)}+p_{ij}^{(0,1)}$. In case of unseen edges, binding prediction is determined by

$$p_{ij}^{\text{conditional}} = \frac{p_{ij}^{(1,0)}}{p_{ij}^{(1,0)} + p_{ij}^{(0,1)}}, \tag{7}$$

while in case of unseen target $j^*$, the binding probability towards a known compound $i$ follows

$$p_{ij^*}^{\text{conditional}} = \frac{\langle p_{ij}^{(1,0)}\rangle_j}{\langle p_{ij}^{(1,0)}\rangle_j + \langle p_{ij}^{(0,1)}\rangle_j} = \rho_i, \tag{8}$$

where $\langle\cdot\rangle_j$ denotes the average over all known targets, and $\rho_i$ follows from Eq. 1.

In case of unseen ligand $i^*$ and target $j^*$, the binding probability is determined by the overall number of positive ($L^{(1,0)}$) and negative ($L^{(0,1)}$) annotations, i.e.,

$$p_{i^*j^*}^{\text{conditional}} = \frac{\langle p_{ij}^{(1,0)}\rangle_{ij}}{\langle p_{ij}^{(1,0)}\rangle_{ij} + \langle p_{ij}^{(0,1)}\rangle_{ij}} = \frac{L^{(1,0)}}{L^{(1,0)} + L^{(0,1)}}, \tag{9}$$

where $\langle\cdot\rangle_{ij}$ indicates the average over all known pairs of ligands and targets.

## Novel Deep Learning Architectures

### VecNet

VecNet uses the pre-trained `mol2vec` [50] and `protvec` [51] models (Figure 4B). These models create 300- and 100-dimensional embeddings for ligands and proteins, respectively. Based on `word2vec` [98], they treat the Morgan fingerprint [99] and the amino acid sequences as sentences, where words are fingerprint fragments or amino acid trigrams. The training is unsupervised

and independent from the following binding prediction task.

## VAENet

VAENet uses a Variational Auto-Encoder [100], an unsupervised learning technique, to embed ligands onto a latent space. The Morgan fingerprint is directly fed to convolutional layers. The auto-encoder creates latent space embeddings by minimizing the loss of information while reconstructing the molecule from the latent representation. We train the Variational Auto-Encoder on 9.5 million chemicals from ZINC database [47], and all drugs and natural compounds in our binding dataset. Similar to VecNet, we use ProtVec for target embeddings.

## Siamese Model

The Siamese model embeds ligands and proteins into the same space using a one-shot learning approach [101]. We construct triplets of the form ⟨protein target, non-binding ligand, binding ligand⟩ and train the model to find an embedding space that maximizes the Euclidean distances between non-binding pairs, while minimizing it for the binding ones.

## File Preparation for Docking Simulations

The steps to implement docking simulations in AutoDock Vina [53] include:

1. Obtain the 3D ligand structures in SDF format from PubChem and save it in .pdb format with PyMOL for use in AutoDockTools.

2. Download the 3D protein structures in .pdb format and load them into AutoDockTools to remove water molecules from the protein structure, add all hydrogen atoms, and the Kollman charge to the protein.

3. Save both the protein and the ligand structures in .pdbqt format using AutoDockTools.

4. Create the grid for docking that encompasses the whole protein structure. This grid selection ensures a blind docking set-up, so that all locations on the protein are considered for determining the binding affinities. The selected grid sizes are available in gridsizes.txt (see Data and Code availability).

5. Create the configuration files with the grid details for each protein and launch the docking simulation. We consider the protein molecules to be rigid, whereas the ligand molecules are flexible, i.e., we allow rotatable bonds on the ligands.

## Acknowledgement

## Author Contributions

A.C. contributed to writing the manuscript, data curation and preparation, generating the predictions for the network configuration model, performing experiments to identify the emergence of topological shortcuts, implementing negative sample generation, developing and testing of VecNet and VAENet, running docking simulations and developing the method to predict the active binding sites.

R.W. contributed to writing the manuscript, generating the predictions for the network configuration model, designing and training VecNet and VAENet.

Z.S. contributed to training and testing of all the deep learning models, and designing the Siamese model.

O.S.A. contributed to the deep learning literature review, running the DeepPurpose models, implementing negative sample generation, and training VAENet.

M.S. contributed to exploring the optimal representation of molecules and developing the method to predict the active binding sites.

D.G. contributed to the data curation and preparation, and performed the gene phylogeny study.

R.Y., T.E.R., and A.L.B. have provided guidance on designing the experiments and writing the manuscript.

G.M. conceived the project, developed the duplex configuration model, designed experiments to identify the emergence of topological shortcuts, contributed to data preparation, data analysis, and writing the manuscript.

## Competing Interests

A.L.B. is the founder of Scipher Medicine and Naring Health, companies that explore the use of network-based tools in health, and Datapolis, that focuses on urban data.

## Materials & Correspondence

Correspondence and requests for materials should be addressed to G.M.

**Data and Code availability**

The codes that support the findings of this study are openly available at our GitHub page at `https://github.com/Barabasi-Lab/AI-Bind`. The data files are shared via Zenodo at `https://zenodo.org/record/7226641`. Top binding predictions from AI-Bind's VecNet on the COVID-19 related proteins, arranged in the descending order of predicted probabilities and validated by docking, are available at `https://github.com/Barabasi-Lab/AI-Bind/blob/main/Validation/Predictions.csv`.

Table 1: **BindingDB Training Data for DeepPurpose.** Most ligands and proteins in Deep-Purpose training data have either binding or non-binding annotations, which creates imbalance in the degree ratio (see Eq. 1).

| Node Type | Has Only Positive Annotations | Has Only Negative Annotations | Has both annotations | Total Node Count |
|---|---|---|---|---|
| Ligand | $3,084$ | $6,539$ | 793 | $10,416$ |
| Protein | 168 | 556 | 667 | $1,391$ |

Table 2: **DeepPurpose and Duplex Configuration Model Performances on BindingDB dataset.** DeepPurpose and the duplex configuration model perform well in both transductive and inductive tests on the benchmark BindingDB data. Both models fail to achieve good performance in the inductive test, i.e., while predicting over both unseen proteins and ligands.

| Model | Transductive | | Semi-inductive | | Inductive | |
|---|---|---|---|---|---|---|
| | AUROC | AUPRC | AUROC | AUPRC | AUROC | AUPRC |
| DeepPurpose | $0.82\pm0.003$ | $0.48\pm0.004$ | $0.76\pm0.036$ | $0.69\pm0.064$ | $0.60\pm0.066$ | $0.42\pm0.063$ |
| Config. Model | $0.83\pm0.009$ | $0.5\pm0.011$ | $0.77\pm0.048$ | $0.71\pm0.065$ | $0.50\pm0.00$ | $0.30\pm0.034$ |

Table 3: **Assigning SMILES and Amino Acid Sequences Randomly.** A random reshuffle of SMILES and amino acid sequences does not affect the performance of DeepPurpose. This outcome suggests the limitation of DeepPurpose in learning chemical structures.

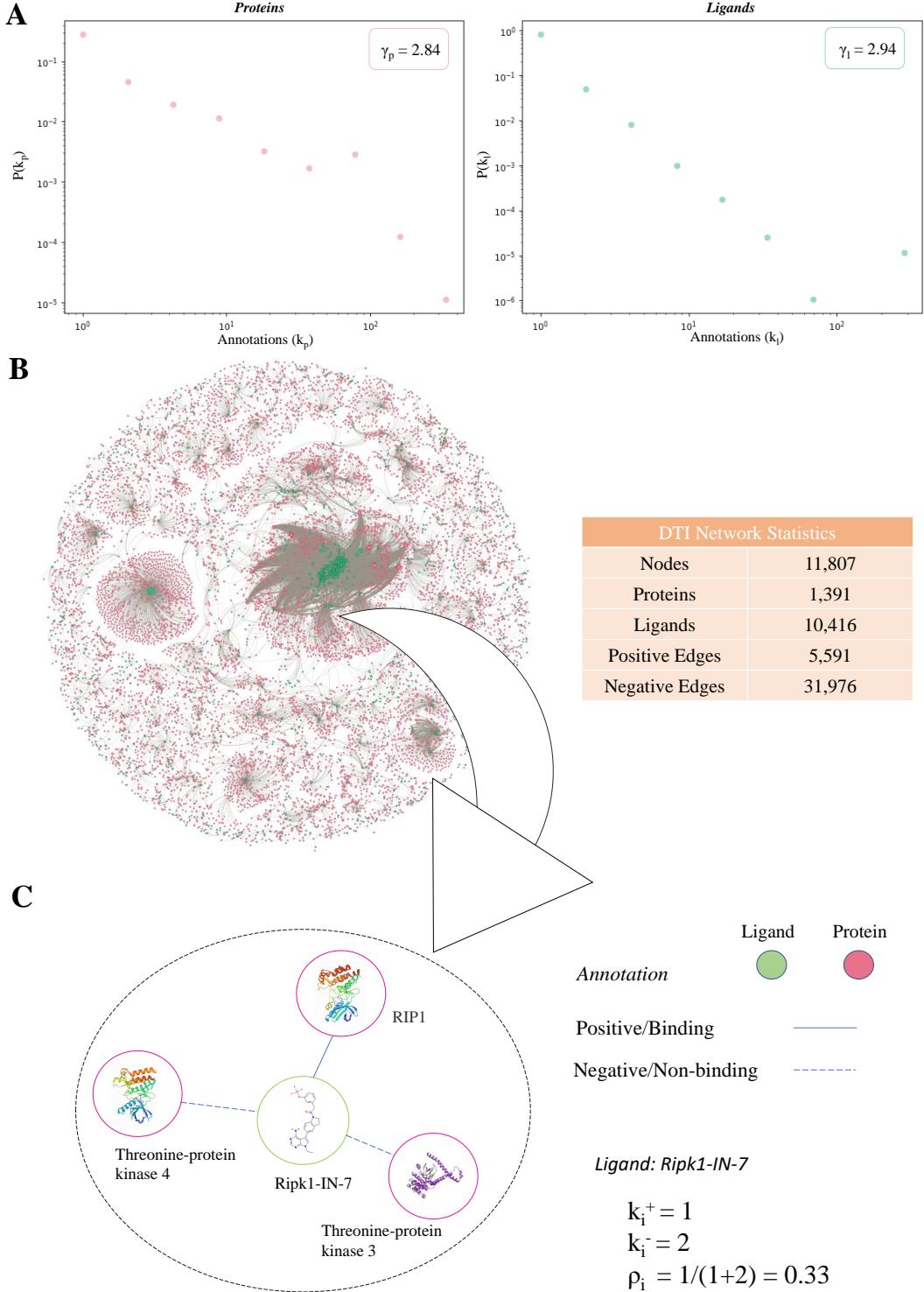| Version | AUROC | AUPRC |
|---|---|---|
| Original | $0.85 \pm 0.005$ | $0.64 \pm 0.008$ |
| Randomized | $0.85 \pm 0.005$ | $0.63 \pm 0.008$ |

Figure 1: **Drug-target Interaction Network. (A)** Distributions of the number of annotations in the benchmark BindingDB data are shown in double logarithmic axes (log-log plot), indicate that $P(k_p)$ and $P(k_l)$ are well approximated by power law for both proteins and ligands, with approximate degree exponents $\gamma_p = 2.84$ and $\gamma_l = 2.94$, respectively. **(B)** The drug-target interaction network used to train the DeepPurpose models, consisting of 10,416 ligands and 1,391 protein targets. Ligands and proteins are represented by green and pink nodes, respectively. **(C)** Network neighborhood of the ligand Ripk1-IN-7. Solid links represent positive or binding annotations, while dashed links refer to negative or non-binding annotations. Ripk1-IN-7 has one positive and two negative annotations in the training data, implying a degree ratio $\rho$ of 0.33.

Figure 2: **Annotation bias in BindingDB training data and DeepPurpose predictions.**
**(A)-(B)** The distribution of degree ratios $\{\rho_p\}$ for the proteins in the original DeepPurpose training set (in a fold from the 5 fold cross-validation). Degree ratio defined in Eq. 1 refers to the ratio of positive annotations to the total annotations for a given node in the protein-ligand interaction network. The average $K_d$ for different degree values $\{k_p\}$ are negatively correlated with $r_{Spearman}(k_p, \langle K_d \rangle) \approx -0.47$. We observe larger variance in $\langle K_d \rangle$ for the low degree nodes. After thresholding $K_d$ values associated with each link to create the binary labels, the hubs get many positive or binding annotations, whereas the low degree nodes get both binding and non-binding annotations. As the hubs are associated with many links in the network, learning the type of binding from degree information helps ML models to achieve good performance leveraging shortcut learning. We observe similar association patterns for the ligands with $r_{Spearman}(k_p, \langle K_d \rangle) \approx -0.29$. **(C)** Protein degree ratios $\{\rho_p\}$ and DeepPurpose predictions are highly correlated with $r_{Spearman} = 0.94$. We observe that the top 100 false positive predictions include the proteins with large $\{\rho_p\}$ represented by the red crosses, whereas the false negatives are contributed by the proteins with small $\{\rho_p\}$ which are represented by the blue dots. **(D)** Examples of proteins and ligands with large degree ratios and contributing to false positive predictions.
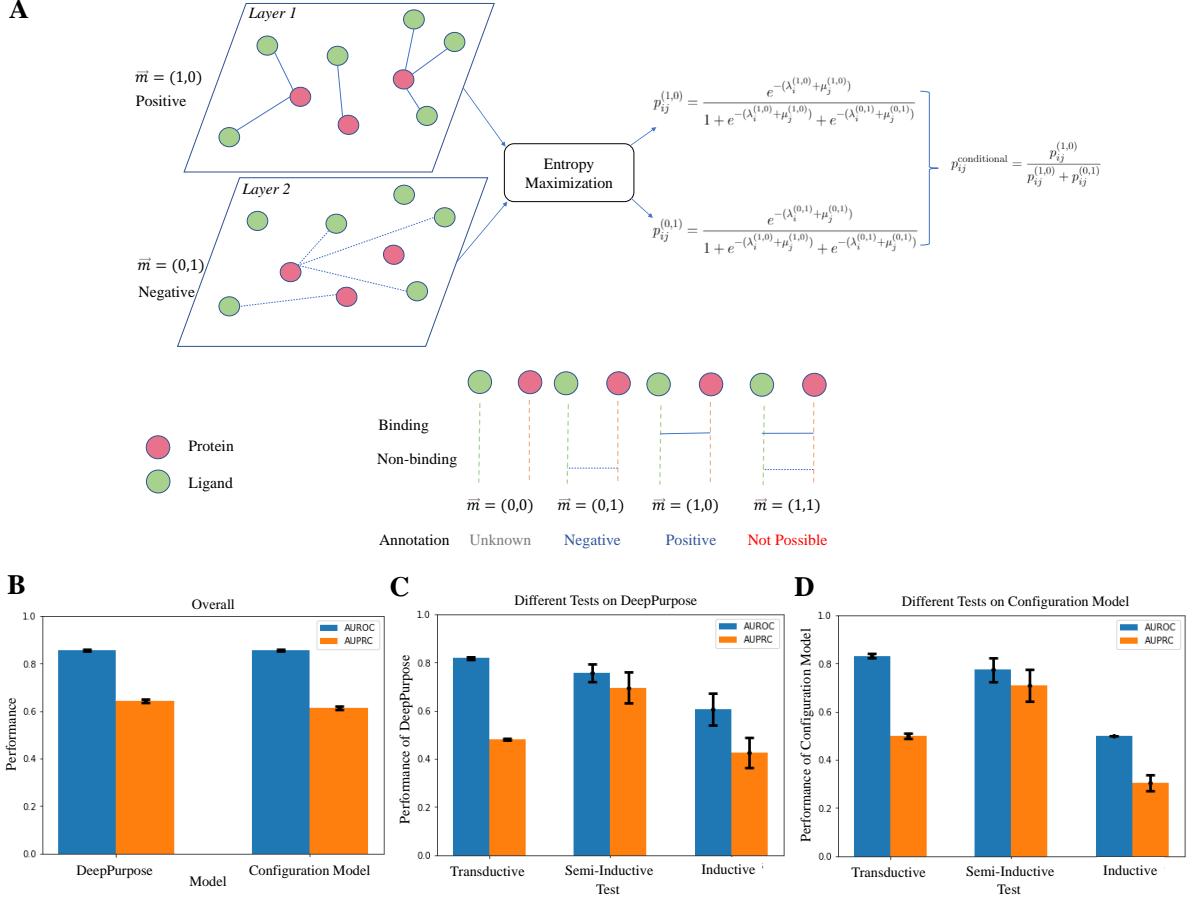
21

Figure 3: **Comparing DeepPurpose and the Duplex Configuration Model. (A)** The duplex configuration model includes two layers corresponding to binding and non-binding annotations. Positive and negative link probabilities are determined by entropy maximization (see Methods), and used to estimate the conditional probability in transductive (Eq. 7), semi-inductive (Eq. 8), and inductive (Eq. 9) scenario. **(B)-(D)** The configuration model achieves similar test performance as DeepPurpose on the the benchmark BindingDB data in a 5-fold cross-validation. Breakdown of performances shows good predictive performance on unseen edges and unseen targets. But the same models have poor predictive performance on unseen nodes.
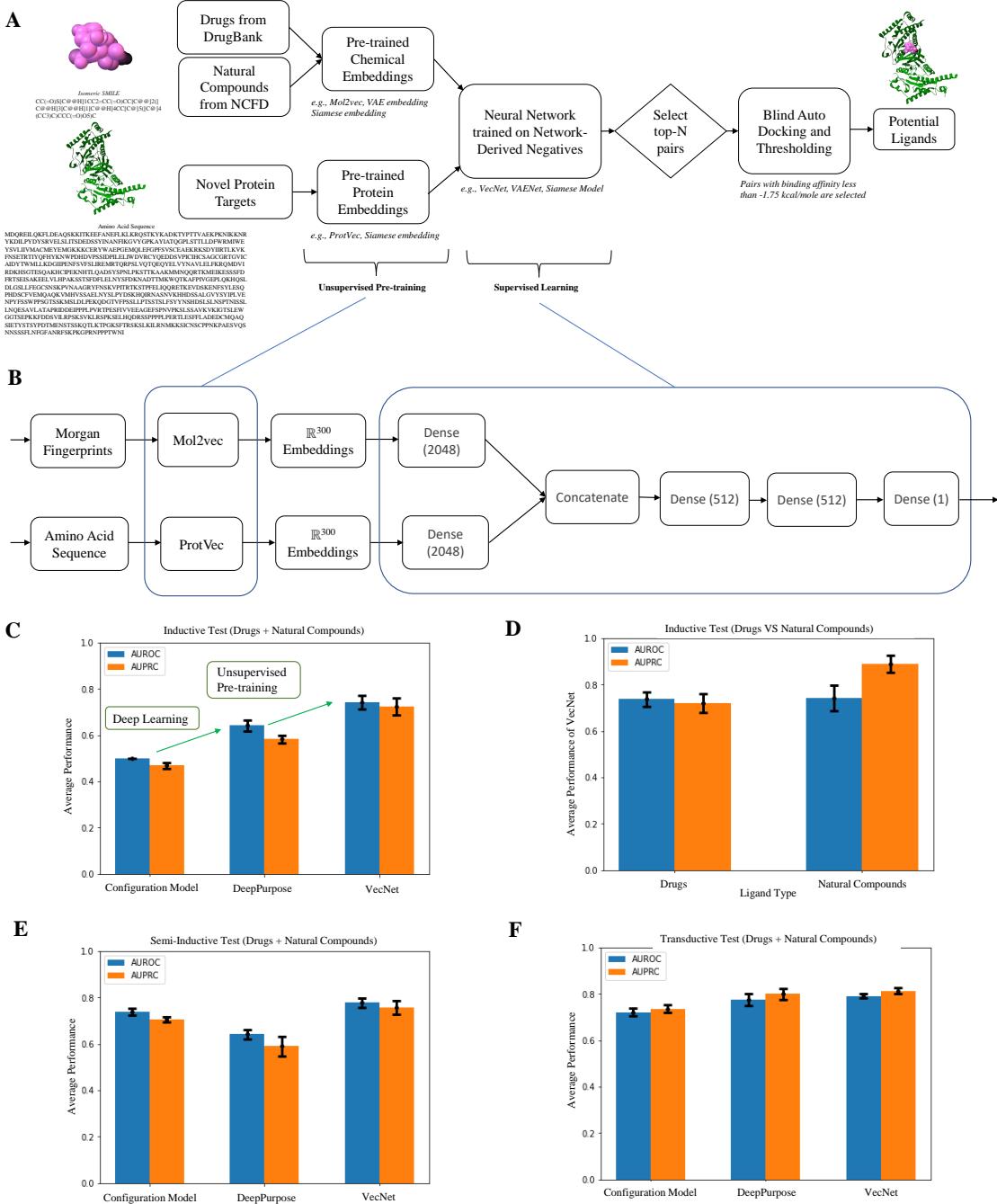
Figure 4: **AI-Bind pipeline: VecNet Performance and Validation. (A)** AI-Bind pipeline generates embeddings for ligands (drugs and natural compounds) and proteins using unsupervised pre-training. These embeddings are used to train the deep models. Top predictions are validated using docking simulations and are used as potential binders to test experimentally. **(B)** AI-Bind's VecNet architecture uses Mol2vec and ProtVec for generating the node embeddings. VecNet is trained in a 5-fold cross-validation set-up. Averaged prediction over the 5 folds is used as the final output of VecNet. **(C)-(F)** 5-fold cross-validation performance of VecNet, DeepPurpose, and Configuration Model. All the models perform similarly in case of predicting binding for unseen edges and unseen targets. The advantage of using deep learning and unsupervised pre-training is observed in the case of unseen nodes (inductive test). AI-Bind's VecNet is the best performing model across all the scenarios. Additionally, we observe similar performance of VecNet for both drugs and natural compounds.
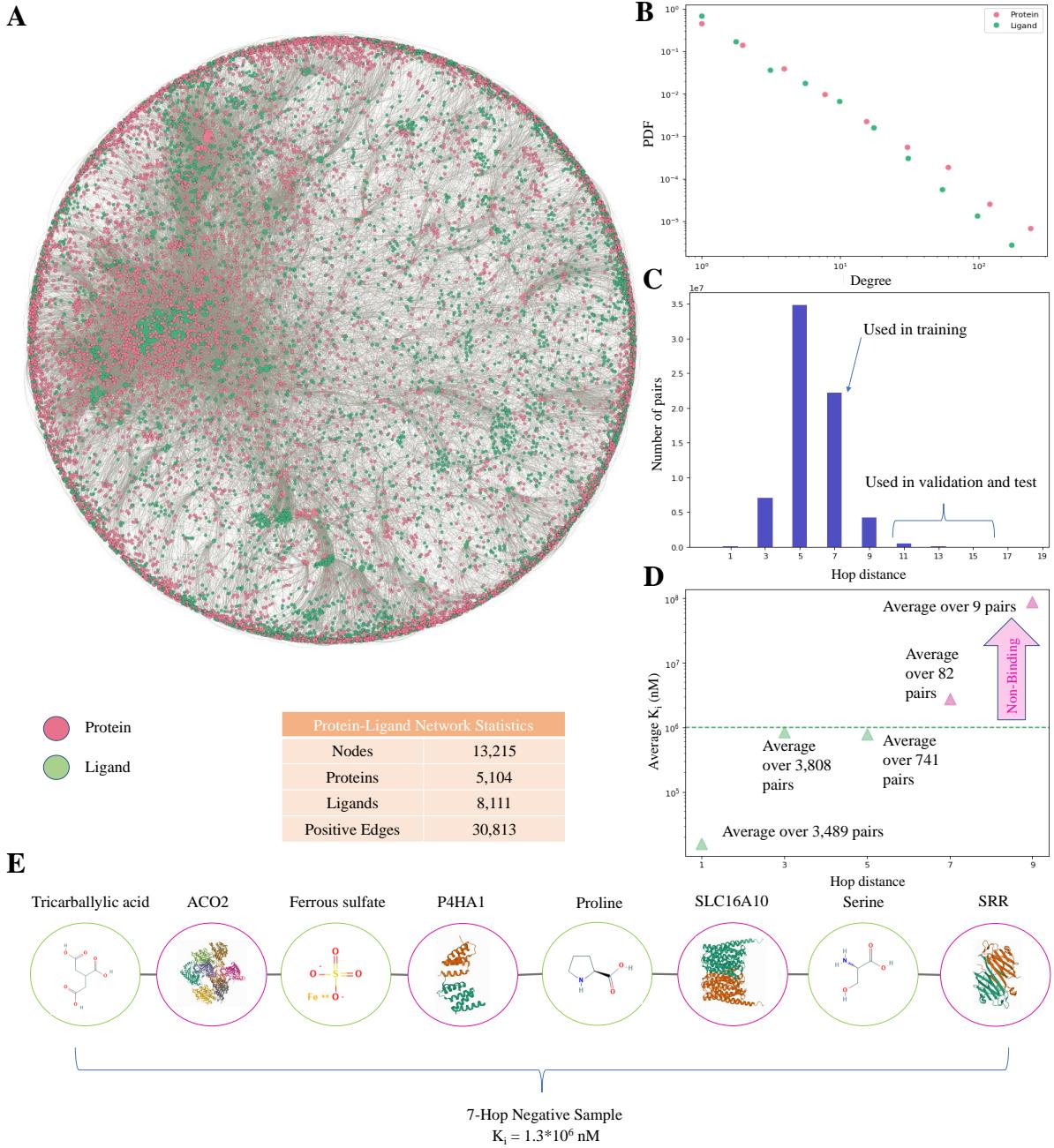
Figure 5: **Network-Derived Negatives. (A)** Protein-ligands bipartite network consisting of only binding (positive) annotations for drugs and natural compounds. **(B)** Degree distributions of ligands and proteins are fat-tailed in nature. **(C)** Shortest-path length distribution capturing all possible protein-ligand pairs. We use protein-ligand pairs with shortest path distance of 7 for training, while absolute negatives obtained from BindingDB and pairs with shortest path distances $\geq 11$ are used for validation and test. **(D)** Average experimental kinetic constant as a function of the shortest path distance. Higher path distance corresponds to higher $K_i$ in BindingDB. Beyond 7 hops, the expected constant exceeds the binding threshold of $10^6 nM$. **(E)** An example of a protein-ligand pair which is 7 hops apart and is used as a negative sample in the AI-Bind training set.
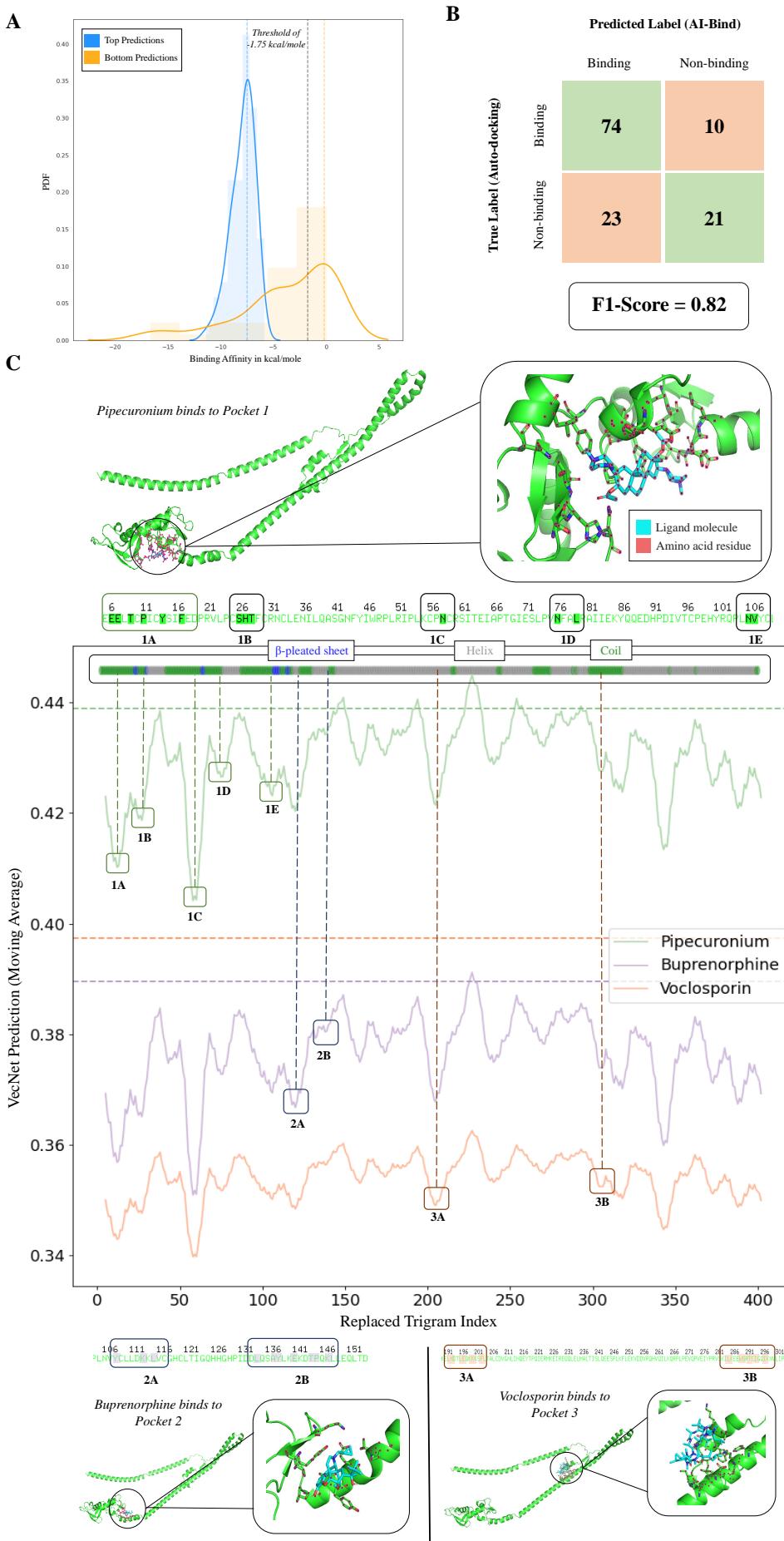
**A**

**B**

Predicted Label (AI-Bind)

|  |  | Binding | Non-binding |
|---|---|---|---|
| True Label (Auto-docking) | Binding | 74 | 10 |
|  | Non-binding | 23 | 21 |

**F1-Score = 0.82**

**C**

*Pipecuronium binds to Pocket 1*

Ligand molecule
Amino acid residue

Pipecuronium
Buprenorphine
Voclosporin

*Buprenorphine binds to Pocket 2*

*Voclosporin binds to Pocket 3*

25

Figure 6: **Validating and interpreting AI-Bind predictions. (A)** Distribution of binding affinities for top and bottom 100 predictions made by AI-Bind's VecNet over viral and human proteins associated with COVID-19. We ran docking on top 84 predictions and bottom 44 predictions. We observe that the top binding predictions of AI-Bind show lower binding energies (better binding) compared to the bottom predictions. Considering the binding threshold of $-1.75$ kcal/mole, 88% of the top predicted pairs by AI-Bind are inline with the docking simulations. **(B)** We construct the confusion matrix for the top and bottom predictions from AI-Bind. We obtain the true labels using the threshold of $-1.75$ kcal/mole on the binding affinities from docking. We observe that AI-Bind predictions produce excellent F1-Score, offering predictions significantly better than random selection. **(C)** Binding probability profile for human protein TRIM59. Multiple valleys in the profile directly map to the amino acid residues to which the ligands bind and are indicative of the active binding sites on the amino acid sequence. We identify the valleys on the binding probability profiles for three ligands Pipecuronium, Buprenorphine and Voclosporin, which bind at different pockets on TRIM59. Valleys for these pockets have been mapped back to the amino acid sequence (valleys 1A, 1B, 1C and 1D for pocket 1, valleys 2A and 2B for pocket 2, and valleys 3A and 3B for pocket 3). Furthermore, we highlight the secondary structure of TRIM59 from the amino acid sequence. Valleys containing the $\beta$-pleated sheets and the coils are more prone to binding compared to the one with the $\alpha$-helices [80–83]. Combining the binding probability profile and the secondary protein structure allows us to identify active binding sites, guiding the design of an optimal search grid for docking simulations.

# SUPPLEMENTARY INFORMATION

## Table of Contents

## S1 Emergence of topological shortcuts

Decision rules learned by many machine learning (ML) models tend to perform well on benchmark datasets, but fail to generalize well when given never-before-seen data. Instead of learning generalizable patterns from features observed during training, these models leverage shortcuts in the data to maximize transductive performance, i.e., the performance on seen data [25]. In this section, we investigate how the properties of the network data used in training can drive ML models to learn topological shortcuts, rather than taking into account node features that would allow better generalizability to unseen data. We assess the emergence of topological shortcuts by null models (configuration model) achieving good transductive test performance.

### BindingDB data observations

First, we investigated the statistical properties of the training database used by DeepPurpose [9], a modeling pipeline offering state-of-the-art neural architectures to predict protein-ligand binding. The training data is based on all the records in BindingDB [102] characterized by a kinetic disassociation constant $K_d$. The distribution of the number of annotations per protein $P(k_p)$ is well fitted by a power law distribution using [38]

$$P(k_p) \sim k^{-\gamma_p}, \tag{10}$$

with $\gamma_p = 2.84$, $k_p^{min} = 1$, and $k_p^{max} = 1{,}426$ (Figure 1A in the main text). We observe similar results for the ligands, with $\gamma_l = 2.94$, $k_l^{min} = 1$, and $k_l^{max} = 1{,}161$.

From the original annotations in BindingDB, a binary classification dataset is derived using a threshold of $30nM$ for $K_d$ [9]. Protein-ligand pairs with $K_d < 30nM$ are binding or positive pairs, and everything else is labeled as non-binding or negative. Overall, we observe that 16% of the records are labeled as positive, a characteristic of the database that we summarize as $p_{bind} = 0.16$, the probability to observe a binding annotation independently from the identity of the protein and the ligand.

Additionally, we find that the number of annotations $k$ and the average disassociation constant $\langle K_d \rangle$ per degree $k$ are not independent but show a negative rank correlation. In particular, for proteins we find $r_{Spearman}(k_p, \langle K_d \rangle) = -0.47$, and for ligands $r_{Spearman}(k_l, \langle K_d \rangle) = -0.29$ (see Figs. 2A, B). Alongside this negative correlation, we observe higher variance for $\langle K_d \rangle$ values for the low degree nodes compared to higher degree nodes (Figs. 2A, B). We find that it is easier to capture the properties of the $K_d$ distribution in BindingDB by modeling it as a log-normal

(Figure S2A, B). Thus, we select the $K_d$ values in the log-space for exploring the emergence of topological shortcuts. Consistently, we observe that in the log $K_d$ space the variance is larger for the lower degree nodes compared to the hubs with $r_{Spearman}(k, \sigma_{\langle \log K_d \rangle}) = -0.71$, where $\sigma_{\langle \log K_d \rangle}$ is the standard deviation of $\langle \log K_d \rangle$ (Figure S2C). This means that the variance is progressively decreasing for higher degree nodes and hubs. Indeed, the hubs are less in number and show similar kinetic features leading to less variance in both $\langle K_d \rangle$ and $\langle \log K_d \rangle$. The relation between degree and kinetic constant makes the link prediction task easier for the hubs compared to the low degree nodes, using only degree information. Since most links in the network are associated with the hubs, the configuration model is able to achieve excellent transductive test performance.

## Toy model set-up

In order to test our hypotheses regarding the creation of topological shortcuts, we simulate synthetic network data that we call *toy models* (Figure S3). We create a duplex of unipartite networks with features inspired by the protein sample captured in BindingDB, as similar considerations extend to bipartite networks. Specifically, we vary the degree distribution $P(k)$ and $r_{Spearman}(k, \langle K_d \rangle)$ to explore when the output of the duplex configuration model $\{p_{ij}^{conditional}\}$ (Eq. 7) becomes highly variable and thus informative, creating the potential for topological shortcuts (see Figure 3A and Methods). In other words, the closer $\{p_{ij}^{conditional}\}$ gets to a Naive Bayes classifier, the less attractive it will be for machine learning models learning a classification task as the predictions would provide information with no discrimination power.

We generate random duplexes of unipartite networks with Poisson or power law degree distributions and different correlations $r_{Spearman}(k, \langle K_d \rangle)$ according to four different specifications:

- Poisson degree distribution and $r_{Spearman}(k, \langle K_d \rangle) \approx -0.47$;

- Poisson degree distribution and $r_{Spearman}(k, \langle K_d \rangle) \approx 0$;

- Power law degree distribution and $r_{Spearman}(k, \langle K_d \rangle) \approx -0.47$;

- Power law degree distribution and $r_{Spearman}(k, \langle K_d \rangle) \approx 0$.

We generate random unipartite toy networks inspired by the topological and kinetic features of the protein training data used in DeepPurpose. We fix the size of the network to $N = 1,507$ and we create randomized networks using the same degree sequence as in BindingDB. For the Poisson case, the link density is constrained by the average number of annotations in the original network. The weight $K_d^{(i,j)}$ assigned to edge $(i,j)$ represents a kinetic constant, and it is derived

as the geometric mean of the contribution $K_d^i$ from node $i$ and the contribution $K_d^j$ from node $j$, namely,

$$K_d^{(i,j)} = \sqrt{K_d^i K_d^j}. \tag{11}$$

We explore multiple scenarios to validate our hypothesis on the emergence of topological short-cuts: in presence of a complex correlated relation between $k$ and $\langle K_d \rangle$ as observed in real-world scenarios, affecting both average values and fluctuations, only power law networks will lead to topological shortcuts. To simplify the modeling of different correlation structures, we use the log-space of kinetic constants, and explore three different sampling strategies: (a) sampling without any variance in the $\log K_d^i$ values contributed by node $i$, (b) sampling with variance in the $\log K_d^i$ values, independent from the degree $k$ of node $i$ and equal to logarithmic variance of BindingDB annotations, (c) sampling with variance in the $\log K_d^i$ values, decreasing as a function of $k$, as observed in the BindingDB data.

According to the sampling strategy, each node contributes to Eq. 11 for all the associated edges with a different extent of variability. In particular, in (a) each node is assigned to a single $\log K_d^i$ for all its edges, sampled according its degree $k$. For a fixed degree $k$, $\log K_d^i$ is sampled from a normal distribution with mean $m = \mu + r_{Spearman} * \sigma * ECDF(k)$ and standard deviation $s = \sqrt{((1 - r_{Spearman}^2) * \sigma^2)}$, where $\mu$ and $\sigma$ are the mean and standard deviation of the $\log K_d$ values in BindingDB. In (b) we follow a similar approach to (a), but instead of sampling a single value, we assign to each node a sample of 5,000 i.i.d. $\log K_d$ instances. Thereafter, when assigning $\log K_d^i$ to each link, we sample uniformly from the generated list of values. In (c), the scenario observed in real data, we first calculate the mean $\langle \log K_d \rangle$ and the standard deviation $\sigma_{\langle \log K_d \rangle}$ for all unique $k$ values. Then, for each link associated with node $i$ we sample a $\log K_d^i$ value from a normal distribution with mean and standard deviation equal to the parameters corresponding to the degree of node $i$. This methodology enforces the same type of complex correlated scenario observed in the original data.

We select as threshold for $K_d^{(i,j)}$ the value for which 16% of the annotations become positive (binding), enforcing the constraint on the observed $p_{bind} = 0.16$. Based on this threshold, we generate the duplex layers with positive and negative edges and calculate the multilink degree sequences, input to the configuration model (see Methods). For the uncorrelated version, we fix the topology while shuffling the $K_d^{(i,j)}$ values at random, which removes any correlation between $k$ and $\langle K_d \rangle$.

## Mathematical formalism for the uncorrelated scenario

When $\langle K_d \rangle$ and $k$ are independent, we can analytically derive the statistical behavior of positive degree $k^+$, negative degree $k^-$, and degree ratio $\rho$ (Eq. 1 in the main text). For each node, the probability of observing $k^+$ positive annotations out of $k$ links is binomial

$$P(k^+|k) = \binom{k}{k^+} p_{bind}^{k^+} (1 - p_{bind})^{(k-k^+)}, \tag{12}$$

where $p_{bind}$ encodes the percentage of positive records observed in the database.

The distribution of positive annotations $k^+$ for the whole database is then a compound distribution

$$P(k^+) = \int P(k)P(k^+|k)dk, \tag{13}$$

where $P(k)$ is the candidate probability distribution for the number of annotations $k$.

From the laws of total expectation and total variance we derive

$$\langle k^+ \rangle = p_{bind}\langle k \rangle, \tag{14}$$

$$\sigma^2(k^+) = p_{bind}(1 - p_{bind})\langle k \rangle + p_{bind}^2(\langle k^2 \rangle - \langle k \rangle^2), \tag{15}$$

where similar equations hold for $k^-$, with $(1 - p_{bind})$ replacing $p_{bind}$. When $P(k)$ is fat-tailed, $(\langle k^2 \rangle - \langle k \rangle^2)$ becomes dominant, and the random variable $k^+ \approx p_{bind}\langle k \rangle$. This formulation suggests that, even in the presence of fat-tailed $P(k)$, the lack of correlations between $\langle K_d \rangle$ and $k$ would determine a distribution of degree ratio $\rho$ well represented by the average

$$\langle \rho \rangle \approx \frac{\langle k^+ \rangle}{\langle k \rangle} = p_{bind}, \tag{16}$$

with noise determined by $p_{bind}$ and $P(k)$. As the duplex configuration model constrains the degree ratio sequence $\{\rho_i\}$, the variability of $\{p_{ij}^{conditional}\}$ drops significantly in absence of correlation between $\langle K_d \rangle$ and $k$, bringing the model closer to a Naive Bayes classifier (Figure S4).

We can clearly derive the behavior of $p_{ij}^{conditional}$ in the case of *uncorrelated networks*, i.e., networks with no degree correlation and an upper bound for the maximum degree equal to $\sqrt{\langle k \rangle N}$, where $N$ is the size of the unipartite network (Advanced Topics 7.B in [36]). In this

scenario the Lagrangian multipliers satisfy:

$$e^{-\lambda_i^{(1,0)}} = \frac{k_i^+}{\sqrt{\langle k^+ \rangle N}} \approx \frac{p_{bind} k_i}{\sqrt{p_{bind}\langle k \rangle N}}, \tag{17}$$

$$e^{-\lambda_i^{(0,1)}} = \frac{k_i^-}{\sqrt{\langle k^- \rangle N}} \approx \frac{(1-p_{bind}) k_i}{\sqrt{(1-p_{bind})\langle k \rangle N}}, \tag{18}$$

$$p_{ij}^{(1,0)} = e^{-(\lambda_i^{(1,0)}+\lambda_j^{(1,0)})}, \tag{19}$$

$$p_{ij}^{(0,1)} = e^{-(\lambda_i^{(0,1)}+\lambda_j^{(0,1)})}. \tag{20}$$

It follows that Eq. 7 in the main text for $p_{ij}^{conditional}$ in the transductive test becomes independent from the identity of node $i$ and $j$, as the product $k_i k_j$ simplifies, leading to $p_{ij}^{conditional} \approx p_{bind}$.

## Observations

The major driving factor is the emergence of topological shortcuts in the relation between $k$ and $\langle K_d \rangle$. The monotonicity of the relation between $k$ and $\langle K_d \rangle$ helps the configuration model to predict the link probabilities using the degree sequence as the $K_d$ values are directly associated with the link types after thresholding. When $k$ and $\langle K_d \rangle$ values are anti-correlated and $\langle K_d \rangle$ values have negligible fluctuations for a fixed $k$, degree becomes a strong predictor of $K_d^{(i,j)}$ and subsequently the link types. Hence we observe excellent transductive test performance of the configuration model, for any topology. But when we introduce variance over the $\langle K_d \rangle$ values, the monotonic relation between $k$ and the link types is disrupted. Hence it is difficult for the configuration model to predict the link type only using the degree information. This observation is consistent for networks with both power law and Poisson degree distributions. Yet, we observe that the variance of $\langle K_d \rangle$ is not uniform for different $k$ values in BindingDB. The hubs encounter less variance in $\langle K_d \rangle$ compared to the low degree nodes. Thus, the configuration model is able to predict the link types associated with the hubs. Since these hubs are associated with the majority of the links in the protein-ligand interaction network, making correct predictions using only the degree information of the hubs helps the configuration model achieve commendable transductive test performance. The performance drops for networks with Poisson degree distributions, where hubs are absent, despite enforcing the same type of correlation structure. When we remove the anti-correlation between $k$ and $\langle K_d \rangle$, irrespective of the variance of $\langle K_d \rangle$ values, the configuration model fails to predict the link types using only the degree information. In this scenario, the configuration model performs similar to a Naive Bayes classifier. Related observations are summarized in Table S1. Given the observed correlation structure in BindingDB (real world scenario), which affects both expected values and

fluctuations in the kinetic constants, topological shortcuts emerge in presence of power law.

## S2  Naturally occurring ligands

We extend the drug repurposing task to additional ligands which are not necessarily considered drugs but may nonetheless bind to protein targets. Specifically, we look into the *Natural Compounds in Food Database (NCFD)* (see Section S4.4), which contains food-borne natural compounds, some of which are potential protein binders. Although these ligands have known chemical structures, they lack adequate binding annotations for training ML models. Binding predictions for these ligands largely depend on comparing their chemical features to other ligands, for which more binding data is available. Figure S5 shows that the naturally occurring compounds in NCFD are larger in size and are more diverse in terms of atomic constituents compared to the drug molecules in DrugBank. This suggests that the binding prediction task on these natural compounds is challenging, which we tackle by maximizing the amount of training data for these natural compounds, and pre-training the chemical embeddings on large chemical libraries.

## S3  DeepPurpose false negative predictions due to annotation imbalance

A false negative prediction corresponds to a low binding probability output by the ML model for a protein-ligand pair which does, in fact, bind. In Figures S6A and S6B, we observe that DeepPurpose produces false negative predictions more often for ligands and proteins with low degree ratios. We notice the opposite for the false positives; nodes with high degree ratios contribute more to the false positive predictions in DeepPurpose predictions (see Figure 2C in the main text).

## S4  Databases

AI-Bind combines data from four databases: DrugBank, Drug Target Commons (DTC), BindingDB, and Natural Compounds in Food Database (NCFD).

## S4.1  DrugBank

DrugBank [16] consists of 7,307 drugs and 4,762 protein targets, which form 25,373 drug-target binding pairs. 167 of these drugs are found in NCFD, and we classify them as naturally occurring and food-borne. We consider all reported protein-ligand pairs from DrugBank as positive samples in our dataset, except 53 pairs which have kinetic constants $\geq 10^6 nM$ in

BindingDB. The protein sequences included in DrugBank are derived from a wide variety of organisms, including human and different viruses.

We observe that the annotation distribution of the proteins and the ligands in DrugBank is fat-tailed (see Figure S7A). This observation is similar to the annotation distributions in BindingDB. The fat-tailed nature of the degree distribution in the binding datasets is a result of the experimentation associated with studying protein-ligand binding. Some proteins and ligands are indeed studied more than others, and hence appear as hubs in such datasets.

## S4.2  Drug Target Commons

We use Drug Target Commons (DTC) [21] for obtaining binding information related to the natural compounds in NCFD. The intersection of NCFD and DTC contains 1,820 natural ligands and 466 associated proteins.

## S4.3  BindingDB

BindingDB [102] consists of protein-ligand pairs along with associated kinetic constants and physical conditions related to the reactions such as pH and temperature. We use BindingDB to extend the number of binding pairs in our training data, filter out the non-binding ones from DrugBank, and obtain absolute negative samples.

## S4.4  Natural Compounds in Food Database

Multiple existing databases contain information about the compounds present in different food items. As a part of the Foodome project at Center for Complex Network Research (CCNR), we curated external databases like FooDB [103], Dictionary of Food Compounds (DFC) [104], and KNApSAcK [105] to gather information about the compounds in food. Metabolomic experiments were performed to further enrich the database. NCFD contains 20,700 compounds found in different food items, among which $\approx$ 19,000 contain isomeric SMILES [35], a plain-text encoding of the chemical structures of each molecule[1]. AI-Bind uses SMILES as input to its ML models for learning useful chemical embeddings.

Figure S7 shows the detailed breakdown of the protein-ligand binding pairs obtained from different databases.

---

[1]NCFD data was accessed on 7.14.2021. As this database undergoes constant change, we have included a description of the dataset at the time of download in the SI.

## S5    7-hop threshold for network-derived negatives

We use shortest path distances to generate negative samples. We consider the node pairs which have shortest path distance $\geq 7$ in the network as non-binding. We derive this 7-hop threshold based on two observations. First, 7 hops is the minimum shortest path distance at which the average kinetic constant value is above the non-binding threshold of $10^6 nM$ (See Figure 5D in the main text). Second, 7 hops is small enough that the negative samples for a given node are not easily distinguishable from positive samples, making the learning task more complex, which helps to defeat shortcut learning in ML models. The latter observation is based on EigenSpokes [106] analysis, a network-based dimensionality reduction procedure inspired by Principal Component Analysis (PCA). Let $A$ be the square adjacency matrix of the protein-ligand network. Since $A$ is real symmetric, it is orthogonally diagonalizable. Let $e_1, \ldots, e_n$ be the eigenvectors of $A$ sorted by eigenvalue magnitude $|\lambda_1| \geq |\lambda_2| \geq \ldots \geq |\lambda_n|$. Given a node $i$, we write the row $a_i$ of $A$ in terms of the eigenbasis $a_i = u_{i1}e_1 + \ldots + u_{in}e_n$. Truncating after the first 5 eigenvectors (with highest eigenvalue magnitudes) gives a low-dimensional embedding $\bar{u}_i = (u_{i1}, u_{i2}, u_{i3}, u_{i4}, u_{i5})$ of each node. The choice of 5 dimensions gives a useful low-dimensional embedding, while still capturing the most significant degrees of variation.

Now, consider a fixed protein $i$. Then ligands $\{j_1, j_2, ...\}$ which bind to $i$ (1-hop) have high magnitude and variance in this 5-dimensional space. On the other hand, ligands $\{k_1, k_2, ...\}$ that are at a distance of 13 hops from $i$ have $\bar{u}_k$ very close to the origin (Table S2). When 13 hops is chosen as the threshold for negative samples, it would thus be trivial for ML models to distinguish nodes $\{j_1, j_2, ...\}$ 1 hop away apart from nodes $\{k_1, k_2, ...\}$ 13 hops away, resulting in shortcut learning. Indeed, the same low-degree ligands on the periphery of the network would become negative samples for all the proteins.

We observe a similar behavior for 11-hop and 9-hop thresholds. However, at 7 hops, we see significantly higher magnitude and variance in $\bar{u}_k$, indicating more diverse negative samples for each protein. In Figure S8A, we visualize $(u_3, u_4)$ for ligands, colored based on the hop-distances from the example protein BPT4. We see that at shortest path distances $\geq 7$, most nodes are very close to the origin. In Figure S8B, we show the mean of all $\|\bar{u}_j\|$ values averaged over all pairs $(i, j)$ of a given path length. Similarly to what we observed for BPT4, we observe a significant fall-off in magnitude as the shortest path length increases.

## S6    Novel deep learning models

We observe that neural networks exploit the topology of the protein-ligand bipartite network used in training to achieve good performance, and lack node-level generalizability when trained in an end-to-end fashion. AI-Bind circumvents these issues by training its ML models in two phases. First, AI-Bind learns the node features using unsupervised pre-training, and then it separately trains its classifiers in a supervised manner to predict binding. To show that AI-Bind is not specific to a certain neural network architecture, we experiment with 3 two-phase networks: VecNet, Siamese model, and VAENet. AI-Bind first trains a neural network in an unsupervised manner to embed the nodes into a low-dimensional latent space, learning generalizable node representations based on the node features alone (chemical structures of ligands and amino acid sequences of proteins). For example, one of the AI-Bind architectures, VecNet, uses unsupervised node representations from Mol2vec [107] and ProtVec [51], which are trained separately from each other and from the protein-ligand bipartite network used in training. Mol2vec and ProtVec are both based on Word2Vec [108], and are designed to create low-dimensional vector representations which retain contextual information for "words" in "sentences", where the "sentences" are formed by molecular sequence descriptions such as Morgan fingerprints [109] or protein sequences. In the second phase, these node representations are passed as input to a binding prediction network, which is trained in a supervised manner. In AI-Bind's VecNet, the binding prediction network uses fully-connected layers and ReLU non-linearities.

The Siamese model uses triplet similarity to find a representation for the node (protein and ligand) features based on their common bindings. The embeddings are then used as inputs to a multilayer perceptron, which learns bindings in a separate supervised training. The last of AI-Bind's three models, VAENet, uses a Variational Auto-Encoder [110] in order to learn unsupervised ligand representations.

### S6.1    VecNet

We use the pre-trained Mol2vec [50] and ProtVec [51] models for node representations. The pre-trained Mol2vec and ProtVec models create 300 and 100 dimensional embeddings for ligands and proteins, respectively. They are based on Word2Vec [98], and treat the Morgan fingerprints [99] and amino acid sequences as sentences in which substructure fingerprints (fragments) and trigrams are the words, respectively. They are trained in an unsupervised manner to create the representations independent of the binding information. Namely, they are trained to predict which words appear near each other in sentences.

Given a fingerprint $x^0$ and an amino acid sequence $x^1$, we encode them using Mol2vec and ProtVec, and then pass them through a simple decoder. We experimented with different neural network architectures with differing number of layers (up to 6 dense layers) and number of neurons per layer (selected from powers of 2 starting at 128 to 2048) and picked one that performed best in inductive tests. This architecture is shown in Figure S9A.

More formally, VecNet computes :

$$\bar{x}^0 = \texttt{mol2vec}(x^0) \in \mathbb{R}^{300}, \qquad \bar{x}^1 = \texttt{protvec}(x^1) \in \mathbb{R}^{100}$$
$$\tilde{x}^0 = \texttt{ReLU}(W^0\bar{x}^0) \in \mathbb{R}^{2048}, \qquad \tilde{x}^1 = \texttt{ReLU}(W^1\bar{x}^1) \in \mathbb{R}^{2048}$$

$$h^0 = \texttt{Concatenate}(\tilde{x}^0, \tilde{x}^1) \in \mathbb{R}^{4096}$$
$$h^1 = \texttt{ReLU}(W^2 h^0) \in \mathbb{R}^{512}$$
$$h^2 = \texttt{ReLU}(W^3 h^1) \in \mathbb{R}^{512}$$
$$\hat{y} = \sigma(W^4 h^2) \in [0, 1]$$

where $\sigma$ is the sigmoid function and $\sigma(x) = \frac{1}{1+e^{-x}}$.

**Prior use of Mol2Vec and ProtVec in binding prediction.** Mol2Vec has previously been used for binding prediction, but only for pre-specified proteins [107], where the ML model is trained on one protein at a time. No information is encoded regarding the protein except for its binding scores with other chemicals in the training data. In contrast, AI-Bind's VecNet attempts to generalize for different proteins, which we encode using ProtVec. Jaeger et al. [107] also propose PCM2vec, in which they predict properties of proteins by concatenating Mol2Vec and ProtVec vectors for the same protein read in as a molecule and amino acid sequence, respectively. However, they do not attempt to combine these vectors for different inputs corresponding to a protein-ligand pair.

## S6.2  Siamese model

The Siamese model uses one-shot learning to embed proteins and ligands into the same latent space [101]. For a given node, the Siamese model minimizes the Euclidean distances of that node from the nodes which bind to it, while maximizing the distances to the nodes which do not. This process is executed in triplets of the forms ⟨protein, non-binding ligand, binding ligand⟩. For the first kind, AI-Bind trains the network to maximize the Euclidean distance between the protein target and the non-binding ligand, while minimizing the distance

of the target from the binding ligand (Figure S10A). AI-Bind uses these embeddings, generated by the Siamese architecture, to train a separate model for the downstream classification task of predicting binding. We studied the inductive test performance by changing the number of convolutional layers and the number of embedding dimensions. The final Siamese model consists of 4 convolutional layers and creates 128-dimensional output vectors. The classification network concatenates the embeddings for a protein and a ligand, and then passes it through two fully connected layers, similar to VecNet, to predict the binding probabilities (Figure S9C).

## S6.3   VAENet

VAENet uses a Variational Auto-Encoder [100], an unsupervised learning technique, to embed ligands onto a latent space. Morgan fingerprints are directly fed into a convolutional neural network. The auto-encoder minimizes the loss of structural information while reconstructing the molecule back from the latent representation (Figure S10B). We generate 300-dimensional ligand embeddings using the auto-encoder, which is consistent with the dimensionality of the Mol2vec embeddings used in VecNet. The variational nature of this 300-dimensional space allows it to be continuous, allowing for better generalizability. We achieve this generalizability by using the re-parametarization trick from [111] to sample from the latent space, instead of directly connecting the latent space to the decoder. Having a generalizable continuous space allows us to map novel ligands into the latent space.

The downstream classification task is achieved by training a fully connected neural network on the concatenated embeddings generated from the Variational Auto-Encoder and ProtVec. The non-end-to-end nature of this architecture ensures that the learned molecular features are independent from the classification task, which has a tendency to exploit shortcuts related to the topology of the protein-ligand interaction network. We observe lower performance for VAENet compared to VecNet (Table S3) mainly for two reasons: (i) VAE has a smaller training dataset of $\sim 9$ million chemicals from ZINC, whereas Mol2vec uses 19.9 million chemicals in training. Thus, Mol2vec is better at generalizing to unknown ligand structures. (ii) VAE uses an auto-encoder to embed the ligand molecules, which is a dimensionality reduction approach. Mol2vec uses skipgrams to embed the molecular structures, which is better at capturing contextual information for different fragments in the molecule and provides a more intelligible representation of the ligand structures for the downstream classification task.

## S7   Additional deep learning model results

Table S3 contains the performances of AI-Bind's novel deep learning architectures, a Deep-Purpose baseline (Transformer-CNN), and the duplex network configuration model on the network-derived dataset. We also report the performances for models trained with randomized node features. This removes structural information about the proteins and ligands, helping us understand whether the deep learning models leverage structure to learn binding or take topological shortcuts. We observe that DeepPurpose's performance does not change if the inputs are randomly shuffled, which suggests that DeepPurpose learns the topology of the protein-ligand interaction network instead of the node features (see Table 3 in the main).

In AI-Bind, network-derived negatives and unsupervised pre-training allow the deep learning models to learn binding patterns using the chemical structures instead of the topology of the protein-ligand interaction network. Thus, we observe diminished performance while using random features to make predictions for unseen nodes (inductive test). In this case network-derived negatives remove the annotation imbalance from the training data and prohibit the ML models from taking topological shortcuts.

Figure S11 shows the training curves averaged over 5 data splits (85 : 15 split to create train and validation-test datasets) for AI-Bind's three novel models. We set the stopping criterion for training to maximize the inductive test performance (AUPRC) on the validation set. Figure S12 shows the F1-scores for the trained models relative to the classification threshold. We obtain the optimal threshold from this curve, which corresponds to the highest F1-score. This optimal threshold is used to obtain the binary labels from the predicted continuous outputs of the AI-Bind architectures. For AI-Bind's VecNet, we obtain an optimal threshold of 0.09 ($\pm$ 0.015) in the inductive test scenario. We observe a low optimal threshold as AI-Bind's VecNet predicts high binding probability ($p_{ij}^{VecNet}$) for a few protein-ligand pairs, but we have roughly the same number of positive and negative samples in the test data. We recommend to use $p_{ij}^{VecNet}$ values to select the top-N predictions, rather than using this optimal threshold to derive the binary labels for novel protein-ligand pairs absent in AI-Bind test data.

## S8   Comparison with MolTrans

We compare the performance of AI-Bind with the Molecular Interaction Transformer (MolTrans) [29], a state-of-the-art protein-ligand binding prediction model which uses a combination of sub-structural pattern mining algorithm, interaction modeling module, and an augmented transformer encoder to better learn the molecular structures. Innovative representation of the

molecules improves the transductive test performance upon DeepPurpose. MolTrans achieves transductive AUROC of $0.952 \pm 0.051$ and AUPRC of $0.872 \pm 0.131$ on the BindingDB data, while DeepPurpose achieves transductive AUROC of $0.775 \pm 0.0.25$ and AUPRC of $0.800 \pm 0.025$. However, MolTrans performs poorly in inductive tests, i.e., while predicting over novel proteins and ligands. We observe that AI-Bind's VecNet performs better than both DeepPurpose and MolTrans in transductive, semi-inductive, and inductive tests. VecNet's improved inductive performance validates that unsupervised pre-training improves the generalizability of the protein-ligand binding models. The results are summarized in Table S4.

## S9 Interpretability of AI-Bind: Identifying active binding sites

AI-Bind may be used to find active binding sites on the amino acid sequence. We plan to leverage this information to define an optimal search grid for docking simulations. Specifically, we use AI-Bind to identify which trigrams in the amino acid sequence are most significant in predicting binding, thus indicating potential binding locations. This is achieved via an ablation study [112], where each trigram in the input amino acid sequence is mutated, that is, replaced with *xxx*, which maps to the *unknown* vector $\langle unk \rangle$ in the ProtVec model. The *unknown* vector is a learned 100-dimensional vector to which all out-of-vocabulary entries, i.e., amino acid trigrams not present in the ProtVec training corpus, are mapped. This *unknown* vector is set to the mean of all the other learned vectors. We predict the probable binding locations by mutating each trigram in the amino acid sequence one at a time and observing the fluctuations in the AI-Bind predictions (Figure S13A). We then smooth the fluctuations using a moving average with a window size of 10 (to eliminate the auxiliary valleys) and obtain a *binding probability profile*.

The suggested binding sites correspond to the amino acid trigrams in the valleys of the binding probability profile. We validate that ligands bind at these valleys by visualizing the docking outputs (see Results) using PyMOL [78] and identifying the region around the ligand with a radius of 5Å, corresponding to the active binding sites (Figure 6C). These regions enclose the amino acid residues which form different bonds with the ligand molecule. Bond distances are measured between the centers of two atoms. Length of hydrogen bonds are typically between 2.3 and 3.9Å [113, 114]. London dispersion forces or Van der Waals interaction between non-polar chains have bond length between 3.8 to 4.2Å [115]. Thus, selecting a sphere with a radius of 5Å around the ligand encloses all possible bonds between the ligand and the protein. We identify the amino acid residues inside this sphere, map them to the regions on the amino acid sequence and compare the results with the valleys in the binding probability profile.

To test this method in a specific case, we identify the active binding sites on the human

protein TRIM59. For Pipecuronium, Buprenorphine and Voclosporin, three ligands binding to TRIM59 at three different pockets, we study the valleys in the binding probability profile which predict binding locations on the amino acid sequence (Figure 6C). More generally, considering a broad range of ligands, we predicted a total of four active binding sites on the protein TRIM59, three of which have been validated in the docking simulations (Figure S13B). We group the ligands binding to TRIM59 according to the different binding sites (Figure S14). From our analysis, a possible fourth active binding site emerges, based on a valley in the binding probability profile, but not associated with any ligand tested in the docking simulations. The shape of the binding probability profiles remains the same across different ligands, but the drop from the original VecNet prediction (depth of a valley or $\Delta p_{ij}^{VecNet}$) fluctuates for different ligands (Figure S14). We observe a moderate positive correlation between the depth of the valleys and binding affinities ($r_{Spearman}(\Delta p_{ij}^{VecNet}, \Delta G) = 0.13$). This indicates the depth of the valleys could be indicative of the binding strength and help in identifying the exact binding site on the protein.

Furthermore, we performed unsupervised hierarchical clustering on the binding probability profiles for different ligands of TRIM59. By clustering first the ligand structures using Tanimoto similarity [116], we find that the ligands binding to TRIM59 are diverse in structure, irrespective of the pockets they bind to (Figure S15A). Thus, we cannot identify the specific binding pocket using only ligand structure. When we compare this result to the clustering emerging from the binding probability profiles, we observe a grouping more correlated with the pocket labeling (Figure S15B).

## S10   Validation using gene phylogeny and bias in false predictions

As additional validation, we investigate if AI-Bind's VecNet predictions are biased towards certain protein structures. The inductive test sets contain a total of 4,583 proteins which are unseen during training in different splits of the 5-fold cross-validation set-up. On 3,162 of these proteins, AI-Bind's VecNet makes at least one false prediction, meaning that our model incorrectly labels at least one ligand as a binder which is not (false positive), or labels a ligand as non-binder which does, in fact, bind (false negative). Among these targets, we find that only 228 (5% of all the proteins) are indeed over-represented (proportions test [117]; $p_{BH-fdr}$[2]-value $\leq 0.05$), meaning that over half the predictions involving these proteins are false. To assess the nature of these false predictions, we test their bias for false positives or false negatives. We find that 168 proteins are biased towards false positive predictions, whereas 68 are biased towards

---

[2]$p$-value (Benjamini Hochberg - False Discovery Rate corrected)

false negatives (proportions test; $p_{BH-fdr}$-value $\leq 0.05$).

To understand whether these biases are intrinsic to the evolutionary origins of certain proteins and if AI-Bind's biases are associated with certain protein domains, we perform a phylogenetic analysis. We use MUSCLE [118], a tool for multiple protein sequence alignment, to understand the similarity between these 228 over-represented protein sequences. We observe only weak similarities between these over-represented proteins. We also reconstruct their phylogenetic tree using the neighbor joining tree method [119] on their amino acid sequences and visualize the results using `treeio` and `ggtree` R packages [120, 121]. The results suggest that the false predictions for AI-Bind's VecNet have no bias towards any particular protein structure (Figure S16).

## S11 Optimal representation of protein and ligand molecules

AI-Bind's VecNet achieves the highest inductive performance, i.e., the performance on never-before-seen proteins and ligands. VecNet uses pre-trained Mol2vec (300-dimensional) and ProtVec (100-dimensional) embeddings. These embeddings encode the structural information from the whole protein and ligand molecules [107, 122]. However, protein-ligand binding is influenced by specific molecular properties, hence we focused on the structural features that are believed to be important to binding in the literature, the so-called *engineered features* [123]. For ligands, we construct the features using the counts of different atoms in the molecule (B, Br, C, Cl, F, I, P, N, O, S), total count of atoms, count of heavy atoms, rings, hydrogen donors, hydrogen acceptors, chiral centers, molecular weight and solubility. For proteins, we use the count of each amino acid, total number of amino acids, and sum of their overall molecular weight. In this set-up, ligands and proteins are represented using 18- and 22-dimensional features, respectively, instead of the original 300 dimensions for Mol2vec and 100 dimensions in ProtVec. Leveraging these engineered features, we observe inductive performance proximal to VecNet (Table S5).

We further explore which dimensions of Mol2vec and ProtVec are the best in explaining the engineered features. We do so by learning matrix $E$ through algebraic decomposition, with $VE = F$, $V \in \mathbb{R}^{N_{ligands},300}$ for ligands, and $V \in \mathbb{R}^{N_{proteins},100}$ for proteins. Matrix $F$ encodes the engineered features: for ligands we have $F \in \mathbb{R}^{N_{ligands},18}$, while for proteins $F \in \mathbb{R}^{N_{proteins},22}$ [124]. We re-scale Mol2vec and ProtVec embeddings, as well as the engineered features, between $[0, 1]$ and perform non-negative matrix factorization to obtain $E$. The rows of $E$ explain the relevance of each Mol2vec or ProtVec dimension to the engineered features. While investigating the relation between engineered features and embeddings, we observed that 15 dimensions of the 300 for Mol2Vec showed the highest variance, suggesting that relevant information is embedded

in a smaller dimensional space compared to the standard dimension used in the literature. Similarly, for ProtVec we found a subset of 16 dimensions (Figure S17). On the same note, concatenating the engineered features with Mol2vec and ProtVec does not change the inductive performance of VecNet (Table S5). This experiment suggests that the engineered features do not add any extra information to the binding prediction task, i.e the two representations are highly correlated.

We further investigated the engineered features to understand which of them contribute most to protein-ligand binding as they have an intuitive and straightforward interpretation. SHAP [125] values show that count of carbon atoms, hydrogen acceptor count, number of chirals, count of fluorine atoms and count of oxygen atoms are the top 5 properties of a ligand that determine its binding to a protein. Presence of amino acids like Glutamic acid, Tryptophan, Asparagine, Methionine, and Threonine in a protein target, presence of aromatic rings (helps in $\pi$-stacking), presence of R groups, and N or C terminus of the protein molecules, drive protein-ligand binding (see Tables S6 and S7).

VecNet with engineered features achieves similar inductive test score as the original version. Yet, the predictions from VecNet using engineered features $\{p_{ij}^{VecNet-Engineered}\}$ show poor negative correlation with $\Delta G$ binding affinities obtained from docking simulations in the Results Section ($r_{Spearman} = -0.10$) when obtain the binary labels by thresholding using the median predicted probabilities, compared to the original VecNet with Mol2vec and ProtVec embeddings ($r_{Spearman} = -0.51$). We also observe a significant reduction in F1-score, from 0.82 to 0.64 (see Results).

Overall, when representing protein and ligand molecules in 2D, we find that only a small subset of the features drive protein-ligand binding and are able to explain intuitive properties of the molecules. Simple molecular descriptors like the presence of R groups in the amino acids, different atom counts, charge distribution in the ligand molecule represented by hydrogen acceptor, and donor counts have significant predictive power for protein-ligand binding. However, these features do not provide insight into the surface structure of the molecules or their rigidity. Indeed, presence of binding pockets on proteins and rotatability of bonds in ligand molecules significantly influence protein-ligand binding. Including these relevant aspects into the prediction task would reduce the number of false positives, often determined by the lack of 3D structures in the model. Adding 3D features of ligands and proteins (e.g., shape of the molecules, rotation of bonds in ligand, location of binding pockets etc.) will help AI-Bind to learn the detailed mechanism behind protein-ligand binding and make more accurate predictions.

## S12    Random Negative Sampling

Existing binding prediction models do not consider any balancing between the binding and the non-binding pairs. In DeepPurpose, the non-binding pair generation is done via selecting random pairs of proteins and ligands which do not appear as binding pairs in the training data. As a result, an imbalance is created between the positive and negative samples for certain nodes based on their degrees in the network, and the deep learning models learn from the network topology of the protein-ligand network instead of learning the binding patterns from the molecular structures. Researchers are aware of this imbalanced training caused by binding data-sets like Tox21 and have proposed an oversampling-based approach to resolve the issue [126]. This method, however, did not improve prediction accuracy and generalizability since the root cause of degree bias is not resolved via oversampling.

In this section, we propose different methods for generating the negative pairs in a balanced fashion from the protein-ligand bipartite network. As we use a batch size of 16 in AI-Bind training [127], our ML models observe 16 data instances corresponding to a protein-ligand pair, 15 of which are negative samples, and the remaining being the positive edge.

We generate random negatives for each positive edge $(t, d)$ representing the binding pair of target $(t)$ and drug $(d)$. This is achieved by randomly selecting drugs with no known binding information to $t$ and randomly selecting proteins with no known binding information to $d$. A list of 15 random negative edges is generated where 7(8) random negatives relate to the target of the positive pair and 8(7) random negatives relate to the drug of the positive pair. Since this method produces negative samples surmounting the number of positives, we use a smaller class weight for the negative samples during training.

In Figure 18, we explore the plausibility of using the network-derived negatives for training ML models instead of the random negative samples. We show that the non-binding (or negative) degrees of the nodes in random negative sampling are correlated with the binding (or positive) degrees. Thus, the random negative samples accommodate the same topological information on the protein-ligand network as the positives, providing no additional information on the negative annotations to training. This issue is resolved by creating the network-derived negative samples, which are based on the shortest path distances in the protein-ligand bipartite network.

Finally, we studied the inductive performance of VecNet on both random negatives and network-derived negatives in a 5-fold cross validation set-up. We observe lower inductive test performance on the random negatives (AUROC of $0.709 \pm 0.011$ and AUPRC of $0.566 \pm 0.013$) compared to the network-derived negatives (AUROC of $0.745 \pm 0,032$ and AUPRC of $0.729 \pm 0.038$).

## S13 Gold standard validation of binding probability profile

In this section we use gold standard protein-ligand binding data to validate the binding probability profiles predicted by AI-Bind. In other words, we validate our hypothesis that ligands bind to proteins at the valleys on the binding probability profile with high confidence gold standard experimental protein-ligand binding data [84]. This validation also shows higher propensity of the $\beta$-sheets and the coils regions to bind with the ligands.

First, we obtain the binding probability profiles generated by AI-Bind for two different ligand-protein pairs. We chose *E. Coli* protein Thymidylate Synthase. The ligands are SP-722 and SP-876. We obtain the experimentally validated secondary structure from the RCSB website, and overlay it over the binding probability profile. We then extract from the PDB file the primary binding sites of the ligand molecules. These binding locations (amino acid residues) are represented by the AC1 keyword in the PDB file. The site lists the amino acids that the ligands bind to, which are represented by red dots on the binding probability profiles (see Figure S19). In both cases, the binding sites lie in the valleys of the probability profile, and overlay on the the $\beta$-sheets and the coils regions. Figure S20 depicts similar observations on human proteins. We have also compared the binding sites predicted by AI-Bind with p2rank, another state-of-the-art site detection method [85].

We compare the binding sites predicted by both AI-Bind and p2rank to the gold standard experimental data. For determining the binding sites from the valleys of AI-Bind's binding probability profile, we fit a sine curve to the valleys and consider the region between the points of inflection of the sine curve as the AI-Bind predicted binding site. On the other hand, p2rank predicts the amino acid residues and the associated pockets as the binding locations. Thereafter, we compare the binding pockets predicted by AI-Bind and p2rank to the gold standard experimental data. We observe that the AI-Bind predicted binding sites cover 64.05% of all pockets on the 195 different proteins mentioned in the gold standard dataset, whereas p2rank is able to identify 53.57% of all of these pockets.

## References

[1] Hughes, J., Rees, S., Kalindjian, S. & Philpott, K. Principles of early drug discovery. *British Journal of Pharmacology* **162**, 1239–1249 (2011). URL https://doi.org/10.1111/j.1476-5381.2010.01127.x.

[2] Thafar, M., Raies, A. B., Albaradei, S., Essack, M. & Bajic, V. B. Comparison study of computational prediction tools for drug-target binding affinities. *Frontiers in Chemistry*

**7** (2019). URL https://doi.org/10.3389/fchem.2019.00782.

[3] U.S. Food & Drug Administration. The Drug Development Process. URL https://www.fda.gov/patients/learn-about-drug-and-device-approvals/drug-development-process.

[4] Hollingsworth, S. A. & Dror, R. O. Molecular dynamics simulation for all. *Neuron* **99**, 1129–1143 (2018). URL https://doi.org/10.1016/j.neuron.2018.08.011.

[5] Vivo, M. D., Masetti, M., Bottegoni, G. & Cavalli, A. Role of Molecular Dynamics and Related Methods in Drug Discovery. *Journal of Medicinal Chemistry* **59**, 4035–4061 (2016). URL https://pubs.acs.org/doi/full/10.1021/acs.jmedchem.5b01684.

[6] Meng, X.-Y., Zhang, H.-X., Mezei, M. & Cui, M. Molecular docking: A powerful approach for structure-based drug discovery. *Current Computer Aided-Drug Design* **7**, 146–157 (2011). URL https://doi.org/10.2174/157340911795677602.

[7] Chen, H., Engkvist, O., Wang, Y., Olivecrona, M. & Blaschke, T. The rise of deep learning in drug discovery. *Drug Discovery Today* **23**, 1241–1250 (2018). URL https://doi.org/10.1016/j.drudis.2018.01.039.

[8] Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021). URL https://doi.org/10.1038/s41586-021-03819-2.

[9] Huang, K. *et al.* DeepPurpose: a deep learning library for drug–target interaction prediction. *Bioinformatics* (2020). URL https://doi.org/10.1093/bioinformatics/btaa1005.

[10] Zhang, H. *et al.* DeepBindPoc: a deep learning method to rank ligand binding pockets using molecular vector representation. *PeerJ* **8**, e8864 (2020). URL https://doi.org/10.7717/peerj.8864.

[11] Huang, G., Liu, Z., van der Maaten, L. & Weinberger, K. Q. Densely connected convolutional networks (2018). 1608.06993.

[12] Verma, N. *et al.* SSnet: A deep learning approach for protein-ligand interaction prediction. *bioRxiv* (2019). URL https://doi.org/10.1101/2019.12.20.884841.

[13] Cui, Y., Dong, Q., Hong, D. & Wang, X. Predicting protein-ligand binding residues with deep convolutional neural networks. *BMC Bioinformatics* **20** (2019). URL https://doi.org/10.1186/s12859-019-2672-1.

[14] Zhao, J., Cao, Y. & Zhang, L. Exploring the computational methods for protein-ligand binding site prediction. *Computational and Structural Biotechnology Journal* **18**, 417–426 (2020). URL https://doi.org/10.1016/j.csbj.2020.02.008.

[15] Xia, C.-Q., Pan, X. & Shen, H.-B. Protein–ligand binding residue prediction enhancement through hybrid deep heterogeneous learning of sequence and structure data. *Bioinformatics* **36**, 3018–3027 (2020). URL https://doi.org/10.1093/bioinformatics/btaa110.

[16] Wishart, D. S. *et al.* DrugBank: a knowledgebase for drugs, drug actions and drug targets. *Nucleic Acids Research* **36**, D901–D906 (2007). URL https://doi.org/10.1093/nar/gkm958. https://academic.oup.com/nar/article-pdf/36/suppl\_1/D901/18782085/gkm958.pdf.

[17] Gilson, M. K. *et al.* BindingDB in 2015: A public database for medicinal chemistry, computational chemistry and systems pharmacology. *Nucleic Acids Research* **44**, D1045–D1053 (2015). URL https://doi.org/10.1093/nar/gkv1072. https://academic.oup.com/nar/article-pdf/44/D1/D1045/9482229/gkv1072.pdf.

[18] Richard, A. M. *et al.* The tox21 10k compound library: Collaborative chemistry advancing toxicology. *Chemical Research in Toxicology* **0**, null (0). URL https://doi.org/10.1021/acs.chemrestox.0c00264. PMID: 33140634, https://doi.org/10.1021/acs.chemrestox.0c00264.

[19] Davies, M. *et al.* ChEMBL web services: streamlining access to drug discovery data and utilities. *Nucleic Acids Research* **43**, W612–W620 (2015). URL https://doi.org/10.1093/nar/gkv352.

[20] Davis, M. I. *et al.* Comprehensive analysis of kinase inhibitor selectivity. *Nature Biotechnology* **29**, 1046–1051 (2011). URL https://doi.org/10.1038/nbt.1990.

[21] Tang, J. *et al.* Drug target commons: A community effort to build a consensus knowledge base for drug-target interactions. *Cell Chemical Biology* **25**, 224–229.e2 (2018). URL https://doi.org/10.1016/j.chembiol.2017.11.009.

[22] Hu, F., Jiang, J., Wang, D., Zhu, M. & Yin, P. Multi-PLI: interpretable multi-task deep learning model for unifying protein–ligand interaction datasets. *Journal of Cheminformatics* **13** (2021). URL https://doi.org/10.1186/s13321-021-00510-6.

[23] van Laarhoven, T., Nabuurs, S. B. & Marchiori, E. Gaussian interaction profile kernels for predicting drug–target interaction. *Bioinformatics* **27**, 3036–3043 (2011). URL `https://doi.org/10.1093/bioinformatics/btr500`.

[24] Öztürk, H., Ozkirimli, E. & Özgür, A. A comparative study of SMILES-based compound similarity functions for drug-target interaction prediction. *BMC Bioinformatics* **17** (2016). URL `https://doi.org/10.1186/s12859-016-0977-x`.

[25] Geirhos, R. *et al.* Shortcut learning in deep neural networks. *Nature Machine Intelligence* **2**, 665–673 (2020). URL `https://doi.org/10.1038/s42256-020-00257-z`.

[26] van Laarhoven, T. & Marchiori, E. Biases of drug–target interaction network data. In *Pattern Recognition in Bioinformatics*, 23–33 (Springer International Publishing, 2014). URL `https://doi.org/10.1007/978-3-319-09192-1_3`.

[27] Lee, A. A., Brenner, M. P. & Colwell, L. J. Predicting protein–ligand affinity with a random matrix framework. *Proceedings of the National Academy of Sciences* **113**, 13564–13569 (2016). URL `https://doi.org/10.1073/pnas.1611138113`.

[28] Wang, Z., Liang, L., Yin, Z. & Lin, J. Improving chemical similarity ensemble approach in target prediction. *Journal of Cheminformatics* **8** (2016). URL `https://doi.org/10.1186/s13321-016-0130-x`.

[29] Huang, K., Xiao, C., Glass, L. M. & Sun, J. MolTrans: Molecular interaction transformer for drug–target interaction prediction. *Bioinformatics* **37**, 830–836 (2020). URL `https://doi.org/10.1093/bioinformatics/btaa880`.

[30] Li, S. *et al.* MONN: A multi-objective neural network for predicting compound-protein interactions and affinities. *Cell Systems* **10**, 308–322.e11 (2020). URL `https://doi.org/10.1016/j.cels.2020.03.002`.

[31] Ye, Q. *et al.* A unified drug–target interaction prediction framework based on knowledge graph and recommendation system. *Nature Communications* **12** (2021). URL `https://doi.org/10.1038/s41467-021-27137-3`.

[32] Kalakoti, Y., Yadav, S. & Sundar, D. TransDTI: Transformer-based language models for estimating DTIs and building a drug recommendation workflow. *ACS Omega* **7**, 2706–2717 (2022). URL `https://doi.org/10.1021/acsomega.1c05203`.

[33] Lee, I. & Nam, H. Sequence-based prediction of protein binding regions and drug–target interactions. *Journal of Cheminformatics* **14** (2022). URL `https://doi.org/10.1186/s13321-022-00584-w`.

[34] Jiang, L. *et al.* Identifying drug–target interactions via heterogeneous graph attention networks combined with cross-modal similarities. *Briefings in Bioinformatics* **23** (2022). URL `https://doi.org/10.1093/bib/bbac016`.

[35] Weininger, D. SMILES, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of Chemical Information and Modeling* **28**, 31–36 (1988). URL `https://doi.org/10.1021/ci00057a005`.

[36] Barabási, A.-L. *Network Science* (Cambridge University Press, 2016).

[37] Barabási, A.-L. & Albert, R. Emergence of scaling in random networks. *Science* **286**, 509–512 (1999). URL `https://doi.org/10.1126/science.286.5439.509`.

[38] Alstott, J., Bullmore, E. & Plenz, D. powerlaw: A python package for analysis of heavy-tailed distributions. *PLOS ONE* **9**, e85777 (2014). URL `https://doi.org/10.1371/journal.pone.0085777`.

[39] Yang, J., Shen, C. & Huang, N. Predicting or pretending: Artificial intelligence for protein-ligand interactions lack of sufficiently large and unbiased datasets. *Frontiers in Pharmacology* **11** (2020). URL `https://doi.org/10.3389/fphar.2020.00069`.

[40] Newman, M. *Networks: An Introduction* (OUP Oxford, 2010).

[41] Menichetti, G., Remondini, D., Panzarasa, P., Mondragón, R. J. & Bianconi, G. Weighted multiplex networks. *PLoS ONE* **9**, e97857 (2014). URL `https://doi.org/10.1371/journal.pone.0097857`.

[42] Menichetti, G. & Remondini, D. Entropy of a network ensemble: definitions and applications to genomic data. *Theor Biol Forum* **107**, 77–87 (2014).

[43] Kim, S. *et al.* PubChem in 2021: new data content and improved web interfaces. *Nucleic Acids Research* **49**, D1388–D1395 (2020). URL `https://doi.org/10.1093/nar/gkaa971`.

[44] Pdb database. `www.rcsb.org/`.

[45] and Alex Bateman *et al.* UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Research* **49**, D480–D489 (2020). URL `https://doi.org/10.1093/nar/gkaa1100`.

[46] Stelzer, G. *et al.* The GeneCards suite: From gene data mining to disease genome sequence analyses. *Current Protocols* **54** (2016). URL `https://doi.org/10.1002/cpbi.5`.

[47] Irwin, J. J., Sterling, T., Mysinger, M. M., Bolstad, E. S. & Coleman, R. G. ZINC: A free tool to discover chemistry for biology. *Journal of Chemical Information and Modeling* **52**, 1757–1768 (2012). URL `https://doi.org/10.1021/ci3001277`.

[48] Gaulton, A. *et al.* The ChEMBL database in 2017. *Nucleic Acids Research* **45**, D945–D954 (2016). URL `https://doi.org/10.1093/nar/gkw1074`.

[49] Bairoch, A. The SWISS-PROT protein sequence data bank and its new supplement TREMBL. *Nucleic Acids Research* **24**, 21–25 (1996). URL `https://doi.org/10.1093/nar/24.1.21`.

[50] Jaeger, S., Fulle, S. & Turk, S. Mol2vec: Unsupervised machine learning approach with chemical intuition. *Journal of Chemical Information and Modeling* **58**, 27–35 (2018). URL `https://doi.org/10.1021/acs.jcim.7b00616`.

[51] Asgari, E. & Mofrad, M. R. K. Continuous distributed representation of biological sequences for deep proteomics and genomics. *PLOS ONE* **10**, e0141287 (2015). URL `https://doi.org/10.1371/journal.pone.0141287`.

[52] Haykin, S. *Neural networks: a comprehensive foundation* (Prentice Hall PTR, 1994).

[53] Trott, O. & Olson, A. J. AutoDock vina: Improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *Journal of Computational Chemistry* NA–NA (2009). URL `https://doi.org/10.1002/jcc.21334`.

[54] Gysi, D. M. *et al.* Network medicine framework for identifying drug-repurposing opportunities for COVID-19. *Proceedings of the National Academy of Sciences of the United States of America* **118**, e2025581118 (2021). URL `https://doi.org/10.1073/pnas.2025581118`.

[55] Gordon, D. E. *et al.* A SARS-CoV-2 protein interaction map reveals targets for drug repurposing. *Nature* **583**, 459–468 (2020). URL `https://doi.org/10.1038/s41586-020-2286-9`.

[56] Kruskal, W. H. & Wallis, W. A. Use of ranks in one-criterion variance analysis. *Journal of the American Statistical Association* **47**, 583–621 (1952). URL `https://doi.org/10.1080/01621459.1952.10483441`.

[57] scipy.stats.kruskal. {https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.kruskal.html}.

[58] Smith, R. D., Engdahl, A. L., Dunbar, J. B. & Carlson, H. A. Biophysical limits of protein–ligand binding. *Journal of Chemical Information and Modeling* **52**, 2098–2106 (2012). URL https://doi.org/10.1021/ci200612f.

[59] Jeon, S. *et al.* Identification of antiviral drug candidates against SARS-CoV-2 from FDA-approved drugs. *Antimicrobial Agents and Chemotherapy* **64** (2020). URL https://doi.org/10.1128/aac.00819-20.

[60] Cour, M., Ovize, M. & Argaud, L. Cyclosporine a: a valid candidate to treat COVID-19 patients with acute respiratory failure? *Critical Care* **24** (2020). URL https://doi.org/10.1186/s13054-020-03014-1.

[61] Hughes, J., Rees, S., Kalindjian, S. & Philpott, K. Principles of early drug discovery. *British Journal of Pharmacology* **162**, 1239–1249 (2011). URL https://doi.org/10.1111/j.1476-5381.2010.01127.x.

[62] Dey, S. K. *et al.* Suramin, penciclovir, and anidulafungin exhibit potential in the treatment of COVID-19 via binding to nsp12 of SARS-CoV-2. *Journal of Biomolecular Structure and Dynamics* 1–17 (2021). URL https://doi.org/10.1080/07391102.2021.2000498.

[63] Kondo, T., Watanabe, M. & Hatakeyama, S. TRIM59 interacts with ECSIT and negatively regulates NF-$\kappa$b and IRF-3/7-mediated signal pathways. *Biochemical and Biophysical Research Communications* **422**, 501–507 (2012). URL https://doi.org/10.1016/j.bbrc.2012.05.028.

[64] Li, S., Wang, L., Berman, M., Kong, Y.-Y. & Dorf, M. E. Mapping a dynamic innate immunity protein interaction network regulating type i interferon production. *Immunity* **35**, 426–440 (2011). URL https://doi.org/10.1016/j.immuni.2011.06.014.

[65] Dr. duke's phytochemical and ethnobotanical databases, u.s. department of agriculture. https://phytochem.nal.usda.gov/phytochem/search.

[66] Jeon, D., Son, M. & Choi, J. Effect of spironolactone on COVID-19 in patients with underlying liver cirrhosis: A nationwide case-control study in south korea. *Frontiers in Medicine* **8** (2021). URL https://doi.org/10.3389/fmed.2021.629176.

[67] Cadegiani, F. A., Goren, A. & Wambier, C. G. Spironolactone may provide protection from SARS-CoV-2: Targeting androgens, angiotensin converting enzyme 2 (ACE2), and

renin-angiotensin-aldosterone system (RAAS). *Medical Hypotheses* **143**, 110112 (2020). URL https://doi.org/10.1016/j.mehy.2020.110112.

[68] Cadegiani, F. A., Wambier, C. G. & Goren, A. Spironolactone: An anti-androgenic and anti-hypertensive drug that may provide protection against the novel coronavirus (SARS-CoV-2) induced acute respiratory distress syndrome (ARDS) in COVID-19. *Frontiers in Medicine* **7** (2020). URL https://doi.org/10.3389/fmed.2020.00453.

[69] Liaudet, L. & Szabo, C. Blocking mineralocorticoid receptor with spironolactone may have a wide range of therapeutic actions in severe COVID-19 disease. *Critical Care* **24** (2020). URL https://doi.org/10.1186/s13054-020-03055-6.

[70] Pawełczyk, A. & Zaprutko, L. Anti-COVID drugs: repurposing existing drugs or search for new complex entities, strategies and perspectives. *Future Medicinal Chemistry* **12**, 1743–1757 (2020). URL https://doi.org/10.4155/fmc-2020-0204.

[71] Carino, A. *et al.* Hijacking SARS-CoV-2/ACE2 receptor interaction by natural and semi-synthetic steroidal agents acting on functional pockets on the receptor binding domain. *Frontiers in Chemistry* **8** (2020). URL https://doi.org/10.3389/fchem.2020.572885.

[72] Kumar, A. *et al.* Identification of phytochemical inhibitors against main protease of COVID-19 using molecular modeling approaches. *Journal of Biomolecular Structure and Dynamics* **39**, 3760–3770 (2020). URL https://doi.org/10.1080/07391102.2020.1772112.

[73] Joh, E.-H., Gu, W. & Kim, D.-H. Echinocystic acid ameliorates lung inflammation in mice and alveolar macrophages by inhibiting the binding of LPS to TLR4 in NF-$\kappa$b and MAPK pathways. *Biochemical Pharmacology* **84**, 331–340 (2012). URL https://doi.org/10.1016/j.bcp.2012.04.020.

[74] Joh, E.-H., Jeong, J.-J. & Kim, D.-H. Inhibitory effect of echinocystic acid on 12-o-tetradecanoylphorbol-13-acetate-induced dermatitis in mice. *Archives of Pharmacal Research* **37**, 225–231 (2013). URL https://doi.org/10.1007/s12272-013-0092-8.

[75] Ryu, S. *et al.* Echinocystic acid isolated from eclipta prostrata suppresses lipopolysaccharide-induced iNOS, TNF-$\alpha$, and IL-6 expressions via NF-$\kappa$b inactivation in RAW 264.7 macrophages. *Planta Medica* **79**, 1031–1037 (2013). URL https://doi.org/10.1055/s-0032-1328767.

[76] Tong, X., Lin, S., Fujii, M. & Hou, D.-X. Echinocystic acid induces apoptosis in HL-60 cells through mitochondria-mediated death pathway. *Cancer Letters* **212**, 21–32 (2004). URL `https://doi.org/10.1016/j.canlet.2004.03.035`.

[77] ting Deng, Y., bo Kang, W., ning Zhao, J., Liu, G. & gao Zhao, M. Osteoprotective effect of echinocystic acid, a triterpone component from eclipta prostrata, in ovariectomy-induced osteoporotic rats. *PLOS ONE* **10**, e0136572 (2015). URL `https://doi.org/10.1371/journal.pone.0136572`.

[78] Schrödinger, L. & DeLano, W. Pymol. URL `http://www.pymol.org/pymol`.

[79] Moffat, L. & Jones, D. T. Increasing the accuracy of single sequence prediction methods using a deep semi-supervised learning framework. *Bioinformatics* **37**, 3744–3751 (2021). URL `https://doi.org/10.1093/bioinformatics/btab491`.

[80] Kutchukian, P. S., Yang, J. S., Verdine, G. L. & Shakhnovich, E. I. All-atom model for stabilization of $\alpha$-helical structure in peptides by hydrocarbon staples. *Journal of the American Chemical Society* **131**, 4622–4627 (2009). URL `https://doi.org/10.1021/ja805037p`.

[81] Fujiwara, K., Toda, H. & Ikeguchi, M. Dependence of alpha-helical and beta-sheet amino acid propensities on the overall protein fold type. *BMC Structural Biology* **12**, 18 (2012). URL `https://doi.org/10.1186/1472-6807-12-18`.

[82] Cheng, P.-N., Pham, J. D. & Nowick, J. S. The supramolecular chemistry of $\beta$-sheets. *Journal of the American Chemical Society* **135**, 5477–5492 (2013). URL `https://doi.org/10.1021/ja3088407`.

[83] Remaut, H. & Waksman, G. Protein–protein interaction through $\beta$-strand addition. *Trends in Biochemical Sciences* **31**, 436–444 (2006). URL `https://doi.org/10.1016/j.tibs.2006.06.007`.

[84] Cheng, T., Li, X., Li, Y., Liu, Z. & Wang, R. Comparative assessment of scoring functions on a diverse test set. *Journal of Chemical Information and Modeling* **49**, 1079–1093 (2009). URL `https://doi.org/10.1021/ci9000053`.

[85] Krivák, R. & Hoksza, D. P2rank: machine learning based tool for rapid and accurate prediction of ligand binding sites from protein structure. *Journal of Cheminformatics* **10** (2018). URL `https://doi.org/10.1186/s13321-018-0285-8`.

[86] Heaven, W. D. The way we train ai is fundamentally flawed. *Technology Review* (2020). URL `https://www.technologyreview.com/2020/11/18/1012234/training-machine-learning-broken-real-world-heath-nlp-computer-vision/`.

[87] Guney, E., Menche, J., Vidal, M. & Barábasi, A.-L. Network-based in silico drug efficacy screening. *Nature Communications* **7** (2016). URL `https://doi.org/10.1038/ncomms10331`.

[88] do Valle, I. F. *et al.* Network medicine framework shows that proximity of polyphenol targets and disease proteins predicts therapeutic effects of polyphenols. *Nature Food* **2**, 143–155 (2021). URL `https://doi.org/10.1038/s43016-021-00243-7`.

[89] Ferreira de Freitas, R. & Schapira, M. A systematic analysis of atomic protein-ligand interactions in the PDB. *Medchemcomm* **8**, 1970–1981 (2017).

[90] Satorras, V. G., Hoogeboom, E., Fuchs, F. B., Posner, I. & Welling, M. E (n) equivariant normalizing flows for molecule generation in 3d. *arXiv preprint arXiv:2105.09016* (2021).

[91] Fuchs, F. B., Worrall, D. E., Fischer, V. & Welling, M. Se (3)-transformers: 3d roto-translation equivariant attention networks. *arXiv preprint arXiv:2006.10503* (2020).

[92] Jumper, J. *et al.* Highly accurate protein structure prediction with alphafold. *Nature* **596**, 583–589 (2021).

[93] Stärk, H., Ganea, O.-E., Pattanaik, L., Barzilay, R. & Jaakkola, T. Equibind: Geometric deep learning for drug binding structure prediction (2022). URL `https://arxiv.org/abs/2202.05146`.

[94] Finan, C. *et al.* The druggable genome and support for target identification and validation in drug development. *Science Translational Medicine* **9**, eaag1166 (2017). URL `https://doi.org/10.1126/scitranslmed.aag1166`.

[95] Gorgulla, C. *et al.* An open-source drug discovery platform enables ultra-large virtual screens. *Nature* **580**, 663–668 (2020). URL `https://doi.org/10.1038/s41586-020-2117-z`.

[96] Heller, S. R., McNaught, A., Pletnev, I., Stein, S. & Tchekhovskoi, D. InChI, the IUPAC international chemical identifier. *Journal of Cheminformatics* **7** (2015). URL `https://doi.org/10.1186/s13321-015-0068-4`.

[97] Menichetti, G., Bianconi, G., Castellani, G., Giampieri, E. & Remondini, D. Multiscale characterization of ageing and cancer progression by a novel network entropy measure. *Mol. BioSyst.* **11**, 1824–1831 (2015). URL `http://dx.doi.org/10.1039/C5MB00143A`.

[98] Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S. & Dean, J. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, 3111–3119 (2013).

[99] Rogers, D. & Hahn, M. Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling* **50**, 742–754 (2010). URL `https://doi.org/10.1021/ci100050t`.

[100] Doersch, C. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908* (2016). URL `https://arxiv.org/abs/1606.05908`.

[101] Koch, G., Zemel, R. & Salakhutdinov, R. Siamese neural networks for one-shot image recognition. In *ICML deep learning workshop*, vol. 2 (Lille, 2015).

[102] Liu, T., Lin, Y., Wen, X., Jorissen, R. N. & Gilson, M. K. BindingDB: a web-accessible database of experimentally determined protein-ligand binding affinities. *Nucleic Acids Research* **35**, D198–D201 (2007). URL `https://doi.org/10.1093/nar/gkl999`.

[103] Foodb database. `www.foodb.ca`.

[104] Yannai, S. (ed.) *Dictionary of Food Compounds* (CRC Press, New York, NY, USA, 2012), 2nd edn.

[105] Afendi, F. M. *et al.* KNApSAcK family databases: Integrated metabolite–plant species databases for multifaceted plant research. *Plant and Cell Physiology* **53**, e1–e1 (2011). URL `https://doi.org/10.1093/pcp/pcr165`.

[106] Prakash, B. A., Sridharan, A., Seshadri, M., Machiraju, S. & Faloutsos, C. EigenSpokes: Surprising patterns and scalable community chipping in large graphs. In *Advances in Knowledge Discovery and Data Mining*, 435–448 (Springer Berlin Heidelberg, 2010). URL `https://doi.org/10.1007/978-3-642-13672-6_42`.

[107] Jaeger, S., Fulle, S. & Turk, S. Mol2vec: unsupervised machine learning approach with chemical intuition. *Journal of chemical information and modeling* **58**, 27–35 (2018).

[108] Mikolov, T., Chen, K., Corrado, G. & Dean, J. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).

[109] Rogers, D. & Hahn, M. Extended-connectivity fingerprints. *Journal of chemical information and modeling* **50**, 742–754 (2010).

[110] Kingma, D. P. & Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).

[111] Kingma, D. P., Salimans, T. & Welling, M. Variational dropout and the local reparameterization trick. In Cortes, C., Lawrence, N. D., Lee, D. D., Sugiyama, M. & Garnett, R. (eds.) *Advances in Neural Information Processing Systems 28*, 2575–2583 (Curran Associates, Inc., 2015). URL `http://papers.nips.cc/paper/5666-variational-dropout-and-the-local-reparameterization-trick.pdf`.

[112] Meyes, R., Lu, M., de Puiseau, C. W. & Meisen, T. Ablation studies in artificial neural networks (2019). `1901.08644`.

[113] Harris, T. K. & Mildvan, A. S. High-precision measurement of hydrogen bond lengths in proteins by nuclear magnetic resonance methods. *Proteins: Structure, Function, and Genetics* **35**, 275–282 (1999). URL `https://doi.org/10.1002/(sici)1097-0134(19990515)35:3<275::aid-prot1>3.0.co;2-v`.

[114] Laskowski, R. A., Moss, D. S. & Thornton, J. M. Main-chain bond lengths and bond angles in protein structures. *Journal of Molecular Biology* **231**, 1049–1067 (1993). URL `https://doi.org/10.1006/jmbi.1993.1351`.

[115] Roth, C., Neal, B. & Lenhoff, A. Van der waals interactions involving proteins. *Biophysical Journal* **70**, 977–987 (1996). URL `https://doi.org/10.1016/s0006-3495(96)79641-8`.

[116] Bajusz, D., Rácz, A. & Héberger, K. Why is tanimoto index an appropriate choice for fingerprint-based similarity calculations? *Journal of Cheminformatics* **7** (2015). URL `https://doi.org/10.1186/s13321-015-0069-3`.

[117] Newcombe, R. G. Two-sided confidence intervals for the single proportion: comparison of seven methods. *Statistics in Medicine* **17**, 857–872 (1998). URL `https://doi.org/10.1002/(sici)1097-0258(19980430)17:8<857::aid-sim777>3.0.co;2-e`.

[118] Papadopoulos, J. S. & Agarwala, R. COBALT: constraint-based alignment tool for multiple protein sequences. *Bioinformatics* **23**, 1073–1079 (2007). URL `https://doi.org/10.1093/bioinformatics/btm076`.

[119] Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high through-put. *Nucleic Acids Research* **32**, 1792–1797 (2004). URL `https://doi.org/10.1093/nar/gkh340`.

[120] Yu, G. Using ggtree to visualize data on tree-like structures. *Current Protocols in Bioinformatics* **69** (2020). URL `https://doi.org/10.1002/cpbi.96`.

[121] Wang, L.-G. *et al.* Treeio: An r package for phylogenetic tree input and output with richly annotated and associated data. *Molecular Biology and Evolution* **37**, 599–603 (2019). URL `https://doi.org/10.1093/molbev/msz240`.

[122] Asgari, E. & Mofrad, M. R. Continuous distributed representation of biological sequences for deep proteomics and genomics. *PloS one* **10**, e0141287 (2015).

[123] Rohrer, S. G. & Baumann, K. Maximum unbiased validation (MUV) data sets for virtual screening based on PubChem bioactivity data. *Journal of Chemical Information and Modeling* **49**, 169–184 (2009). URL `https://doi.org/10.1021/ci8002649`.

[124] Henderson, K. *et al.* Rolx: structural role extraction & mining in large graphs. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, 1231–1239 (2012).

[125] Lundberg, S. M. & Lee, S.-I. A unified approach to interpreting model predictions. In Guyon, I. *et al.* (eds.) *Advances in Neural Information Processing Systems 30*, 4765–4774 (Curran Associates, Inc., 2017).

[126] Idakwo, G. *et al.* Structure–activity relationship-based chemical classification of highly imbalanced tox21 datasets. *Journal of Cheminformatics* **12** (2020). URL `https://doi.org/10.1186/s13321-020-00468-x`.

[127] Bengio, Y. Practical recommendations for gradient-based training of deep architectures (2012). URL `https://arxiv.org/abs/1206.5533`.

Table SI. 1: Transductive test performance of a duplex configuration model on unipartite layers with varying annotation distribution $P(k)$ and correlation $r_{Spearman}(k, \langle K_d \rangle)$. The network has $N = 1,507$ nodes, same as the number of unique proteins in the BindingDB training data. The constrained features are consistent with the protein sample in BindingDB, e.g., for the power law network we use the degree sequence derived from the protein network in BindingDB data, while the Poisson network has the same average degree $\langle k \rangle = 47$ of the power law network. To achieve $r_{Spearman}(k, \langle K_d \rangle) \approx 0$, we shuffle the edges of the original network, which removes the negative correlation between $k$ and $\langle K_d \rangle$.

| Variance of $\langle K_d \rangle$ | Annotation distribution | $r_{Spearman}(k, \langle K_d \rangle)$ | $p_{bind}$ | AUROC | AUPRC |
|---|---|---|---|---|---|
| Negligible variance | Power law | $\approx -0.47$ | 0.16 | 0.95 | 0.89 |
| $k$-independent variance | Power law | $\approx -0.47$ | 0.16 | 0.64 | 0.26 |
| $k$-dependent variance | Power law | $\approx -0.47$ | 0.16 | 0.84 | 0.54 |
| Negligible variance | Poisson | $\approx -0.47$ | 0.16 | 0.94 | 0.87 |
| $k$-independent variance | Poisson | $\approx -0.47$ | 0.16 | 0.68 | 0.29 |
| $k$-dependent variance | Poisson | $\approx -0.47$ | 0.16 | 0.69 | 0.27 |
| Negligible variance | Power law | $\approx 0$ | 0.16 | 0.54 | 0.17 |
| $k$-independent variance | Power law | $\approx 0$ | 0.16 | 0.51 | 0.17 |
| $k$-dependent variance | Power law | $\approx 0$ | 0.16 | 0.49 | 0.14 |
| Negligible variance | Poisson | $\approx 0$ | 0.16 | 0.50 | 0.16 |
| $k$-independent variance | Poisson | $\approx 0$ | 0.16 | 0.49 | 0.16 |
| $k$-dependent variance | Poisson | $\approx 0$ | 0.16 | 0.49 | 0.16 |

Table SI. 2: As path length to a fixed protein $i$ increases, the mean and variance of the length of the low-dimensional embedding of the ligand $\|\bar{u}_j\|$ decrease.

| Path Length $i$ to $j$ | 1 | 3 | 5 | 7 | 9 | 11 | 13 |
|---|---|---|---|---|---|---|---|
| Mean $\|\bar{u}_j\|$ | 0.045 | 0.035 | 0.014 | 0.004 | 0.001 | 0.0001 | $5 \cdot 10^{-7}$ |
| Std. dev. $\|\bar{u}_j\|$ | 0.075 | 0.035 | 0.025 | 0.014 | 0.008 | 0.001 | $4 \cdot 10^{-6}$ |

Table SI. 3: **Results across different models.** We summarize all performances on the network-derived negative samples. We perform 5-fold cross-validation, reporting AUROC and AUPRC averaged over the 5 runs with random initialization and data split. Results are reported separately on 3 different train-validation-test splits with different data held out in the validation and testing sets: (1) **Unseen edges (Transducitve test)** - test sets contain unseen edges in the train network, (2) **Unseen targets (Semi-inductive test)** - test sets contains pairs with proteins that do not appear in train or validation set, (3) **Unseen nodes (Inductive test)** - nodes in test set pairs are completely disjoint from train set. *Random Input Tests:* We train and test AI-Bind's VecNet replacing node features with random entries drawn from a uniform distribution in the range $U([-1, 1]^d)$. We run two tests (1) Unif. - All node features are replaced by vectors drawn from a uniform distribution $U([-1, 1]^d)$, (2) Unif.Targ. - Only the target node features are replaced by vectors from $U([-1, 1]^d)$; drug features remaining the same. Note that the transductive (unseen edges) performance is reported based on the models optimized for unseen nodes, except for the case of the Random Inputs, where we report performance based on models optimized for unseen targets.

| Model | Test Data Division | | | | | |
| | Transd. | | Semi-induc. | | Induc. | |
| | AUROC | AUPRC | AUROC | AUPRC | AUROC | AUPRC |
| Configuration | $.738 \pm .014$ | $.739 \pm .017$ | $.754 \pm .021$ | $.691 \pm .038$ | $.500 \pm .000$ | $.469 \pm .014$ |
| VecNet | $.794 \pm .008$ | $.817 \pm .013$ | $.779 \pm .025$ | $.752 \pm .039$ | $.745 \pm .032$ | $.729 \pm .038$ |
| Siamese | $.664 \pm .027$ | $.637 \pm .003$ | $.666 \pm .031$ | $.621 \pm .032$ | $.639 \pm .026$ | $.583 \pm .025$ |
| VAENet | $.777 \pm .010$ | $.701 \pm .048$ | $.756 \pm .022$ | $.710 \pm .030$ | $.740 \pm .024$ | $.701 \pm .048$ |
| DeepPurpose | $.775 \pm .025$ | $.800 \pm .025$ | $.642 \pm .022$ | $.591 \pm .042$ | $.642 \pm .025$ | $.583 \pm .016$ |
| *Random Inputs* | | | | | | |
| VecNet Unif. | $.668 \pm .012$ | $.702 \pm .015$ | $.539 \pm .019$ | $.541 \pm .013$ | $.466 \pm .054$ | $.456 \pm .041$ |
| VecNet - Targ. | $.704 \pm .008$ | $.725 \pm .009$ | $.575 \pm .032$ | $.556 \pm .025$ | $.558 \pm .009$ | $.501 \pm .021$ |

Table SI. 4: **Comparing AI-Bind with MolTrans.** We compare transductive, semi-inductive, and inductive performances of MolTrans with AI-Bind's VecNet. MolTrans uses a combination of sub-structural pattern mining algorithm, interaction modeling module and an augmented transformer encoder to better learn the molecular structures. VecNet performs better compared to MolTrans in semi-inductive and inductive tests. This analysis validates that unsupervised pre-training improves the generalizability of the protein-ligand binding models. We have trained and tested MolTrans on both BindingDB (used in the original paper), and network-derived negatives (AI-Bind data).

| Model | Test Data Division | | | | | |
| | Transd. | | Semi-induc. | | Induc. | |
| | AUROC | AUPRC | AUROC | AUPRC | AUROC | AUPRC |
|---|---|---|---|---|---|---|
| MolTrans[a] | $.952 \pm .051$ | $.872 \pm .131$ | $.653 \pm .041$ | $.415 \pm .095$ | $.575 \pm .059$ | $.430 \pm .098$ |
| MolTrans[b] | $.860 \pm .074$ | $.805 \pm .092$ | $.641 \pm .015$ | $.489 \pm .051$ | $.619 \pm .021$ | $.480 \pm .028$ |
| DeepPurpose | $.775 \pm .025$ | $.800 \pm .025$ | $.642 \pm .022$ | $.591 \pm .042$ | $.642 \pm .025$ | $.583 \pm .016$ |
| VecNet | $.794 \pm .008$ | $.817 \pm .013$ | $.779 \pm .025$ | $.752 \pm .039$ | $.745 \pm .032$ | $.729 \pm .038$ |

[a]BindingDB data

[b]Network-derived Negatives

Table SI. 5: **Optimal feature selection:** We observe that AI-Bind's VecNet shows similar performances in inductive tests, when Mol2vec and ProtVec are replaced by simple engineered features encoding certain properties of protein and ligand molecules. Furthermore, we observe that only 15 dimensions of Mol2vec and 16 dimensions of ProtVec embeddings encode these molecular properties driving the binding task. Using these feature subsets of Mol2vec and ProtVec helps VecNet achieving similar inductive performance. Concatenating the engineered features with Mol2vec and ProtVec does not improve inductive performance. This suggests that the information encoded in the engineered features strongly correlates with Mol2vec and ProtVec embeddings.
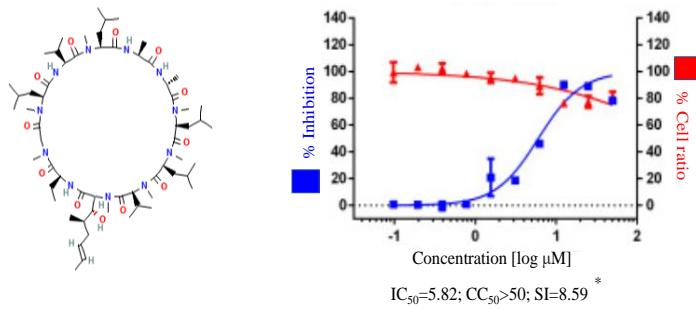
| Performance | Original | Engineered features | Mol2vec and ProtVec dimensions explaining Engineered Features | Concatenated: Mol2vec/ProtVec + Engineered features |
|---|---|---|---|---|
| AUROC | $0.75 \pm 0.032$ | $0.73 \pm 0.032$ | $0.72 \pm 0.066$ | $0.75 \pm 0.033$ |
| AUPRC | $0.73 \pm 0.038$ | $0.74 \pm 0.033$ | $0.72 \pm 0.057$ | $0.73 \pm 0.042$ |

Table SI. 6: **Engineered feature importance for ligands:** We tabulate the engineered features in descending order of average absolute SHAP importance over AI-Bind data. A higher SHAP value represents more relevance of the molecular property in predicting protein-ligand binding.

| Feature | Average SHAP Importance |
|---|---|
| Count of Carbon Atom | 0.012546 |
| Hydrogen Acceptor Count | 0.012362 |
| Number of Chirals | 0.008750 |
| Count of Flourine Atoms | 0.006527 |
| Count of Oxygen Atoms | 0.006184 |
| Hydrogen Donor Count | 0.005647 |
| Number of Atoms | 0.004165 |
| Count of Heavy Atoms | 0.003468 |
| Solubility in Water ($\log p$) | 0.003202 |
| Molecular Weight | 0.003161 |
| Count of Nitrogen Atoms | 0.002007 |
| Count of Chlorine Atoms | 0.001720 |
| Count of Sulphur Atoms | 0.001483 |
| Number of Rings | 0.000525 |
| Count of Phosphorus Atoms | 0.000191 |
| Count of Iodine Atoms | 0.000130 |
| Count of Bromine Atoms | 0.000083 |
| Count of Boron Atoms | 0.000070 |

Table SI. 7: **Engineered feature importance for proteins:** We tabulate the engineered features in descending order of average absolute SHAP importance over AI-Bind data. A higher SHAP value represents more relevance of the molecular property in predicting protein-ligand binding.

| Feature | Average SHAP Importance |
| --- | --- |
| Count of Glutamic acid (E) | 0.036747 |
| Count of Tryptophan (W) | 0.033210 |
| Count of Asparagine (N) | 0.024770 |
| Count of Methionine (M) | 0.022734 |
| Count of Threonine (T) | 0.021194 |
| Count of Glycine (G) | 0.020832 |
| Count of Arginine (R) | 0.019599 |
| Count of Phenylalanine (F) | 0.017040 |
| Count of Cysteine (C) | 0.016428 |
| Count of Isoleucine (I) | 0.016215 |
| Count of Alanine (A) | 0.015732 |
| Count of Histidine (H) | 0.014813 |
| Count of Leucine (L) | 0.014026 |
| Count of Tyrosine (Y) | 0.013844 |
| Count of Proline (P) | 0.013014 |
| Count of Valine (V) | 0.011152 |
| Count of Serine (S) | 0.010930 |
| Count of Lysine (K) | 0.008689 |
| Count of Aspartic acid (D) | 0.008303 |
| Total amino acid count | 0.003381 |
| Count of Glutamine (Q) | 0.002957 |
| Molecular Weight | 0.002088 |

**A** *Anidulafungin*

IC$_{50}$=4.64; CC$_{50}$>50; SI=10.78 *

**B** *Cyclosporin*

IC$_{50}$=5.82; CC$_{50}$>50; SI=8.59 *

Figure SI. 1: **Experimental evidence to validate AI-Bind predictions. (A)-(B)** Anidulafungin and Cyclosporin, two FDA approved anti-fungal agents predicted by AI-Bind, show potential antiviral activities against SARS-CoV-2, with $IC_{50}$ values $4.64\mu M$ and $5.82\mu M$, respectively.

Figure SI. 2: **Disassociation constant $K_d$ and its relation with the number of annotations/records $k$ in BindingDB. (A)-(B)** Density distribution and cumulative distribution of $\log K_d$ in BindingDB training data. With threshold $30nM$, we obtain an average binding probability of $p_{bind} = 0.16$. **(C)-(D)** Each node is characterized by the number of annotations $k$, and the average $\langle K_d \rangle$ over its records. We select the $K_d$ values in the log-space for creating the toy model. We do not observe significant correlation between $k$ and $\langle \log K_d \rangle$ opposed to the anti-correlation observed in the linear space, but $k$ and the variance of $\langle \log K_d \rangle$ values are highly anti-correlated with are $r_{Spearman}(k, \sigma_{\langle \log K_d \rangle}) = -0.71$. This observation implies that the lower degree nodes have higher fluctuations in the associated $\langle \log K_d \rangle$ values compared to the higher degree nodes and hubs.

**A**

*Benchmark BindingDB data*

Protein degree sequence
{$d_1$, $d_2$, ..., $d_N$}

PDF

Annotations (k)

*Configuration Model*

Power law degree
distribution with same
degree sequence

PDF

Annotations (k)

*ER Graph*

Poisson degree
distribution with
N=1,507
<$k_p$>=47

PDF

Annotations (k)

**B**

Anti-correlated
k and (<$K_d$>)

*Negligible variance*

Each node is assigned to a
single log $K_d^i$ for all its
edges, sampled according
its degree k.

*k-independent variance*

Assign to each node a
sample of 5,000 i.i.d. log $K_d$
instances. Sample uniformly
from the generated list of
values when assigning log
$K_d^i$ to each link.

*k-dependent variance*

Calculate the mean log<$K_d$>
and the standard deviation
$\sigma_{logKd}$ for all unique k values.
For each link associated with
node i, sample a log $K_d^i$ value
from a normal distribution
with mean and standard
deviation equal to the
parameters corresponding to
the degree of node i.

Get $K_d^i$ values
for each link
associated with
node i

$K_d^i$

$K_d^{(i,j)} = \sqrt{K_d^i K_d^j}$

$K_d^j$

**C**

CDF

$P_{bind}$

log(Threshold)

log $K_d^{(i,j)}$

$K_d^{(i,j)}$ < Threshold

Positive Edge

$K_d^{(i,j)} \geq$ Threshold

Negative Edge

*Unipartite Version*

Figure SI. 3: **Experimental set-up for studying the emergence of topological short-cuts.** **(A)** We generate random unipartite networks inspired by the topological and kinetic features of the protein sample in BindingDB. In particular, we fix the size of the network to $N = 1,507$ and use the same degree distribution as in BindingDB, while for the Poisson case the link density is constrained by the average number of annotations in the power law network. **(B)** We explore three different strategies of sampling the kinetic constants: (a) sampling without any variance in the $\log K_d^i$ values contributed by node $i$ to its links, (b) sampling with variance in the $\log K_d^i$ values, independent from the degree $k$ of node $i$ and equal to logarithmic variance of $\log K_d$ in the BindingDB protein sample, and (c) sampling with variance in the $\log K_d^i$ values, decreasing as a function of $k$, as observed in the BindingDB data. According to the sampling strategy, each node contributes to all its edges with a different extent of variability. In particular, in (a) each node is assigned to a single $\log K_d^i$ for all its edges, sampled according its degree $k$. In (b) we follow a similar approach to (a), but instead of sampling a single value, we assign to each node a sample of 5,000 i.i.d. $\log K_d$ instances. Thereafter, when assigning $\log K_d^i$ to each link associated with node $i$, we sample uniformly from the generated list of values. In (c), the scenario observed in BindingDB data, we first calculate the mean $\langle \log K_d \rangle$ and the standard deviation $\sigma_{\langle \log K_d \rangle}$ for all unique $k$ values. Then, for each link associated with node $i$ we sample a $\log K_d^i$ value from a normal distribution with mean and standard deviation equal to the parameters corresponding to the degree of node $i$. The final disassociation constant $K_d^{(i,j)}$ assigned to edge $(i,j)$ is the geometric average of the contribution $K_d^i$ from node $i$ and the contribution $K_d^j$ from node $j$. In the uncorrelated scenario, we randomly shuffle the $K_d$ values associated with the links, which removes the anti-correlation between $k$ and $\langle K_d \rangle$. **(C)** We select as threshold for $K_d^{(i,j)}$ the value for which a fixed percentage of the annotations become positive or binding, enforcing the constraint on the observed $p_{bind} = 0.16$. Based on this threshold, we generate the duplex layers with positive and negative edges and calculate the multilink degree sequences, input to the configuration model (see Methods).

Figure SI. 4: **Emergence of topological shortcuts in scale-free networks.** In absence of variance in the $\langle K_d \rangle$ values, the relation between $k$ and $\langle K_d \rangle$ is monotonic and the configuration model is able to predict the link types using only the degree information of the nodes. When variance is introduced, the monotonicity of the relation between $k$ and $\langle K_d \rangle$ is disrupted. Thus, the configuration model is unable to learn the link types only using the degree information. The scenario with varying variance resembles the data in BindingDB. Less fluctuations for the hubs makes the link classification task easier for the hubs. Since majority of the links in the protein-ligand interaction network are associated with the hubs, we observe the configuration model achieving excellent transductive test performance.

Figure SI. 5: **Naturally occurring compounds are structurally more complex than drugs. (A)** Prevalence of different atoms in DrugBank and natural ligands present in NCFD. Natural ligands show more diversity in terms of the constituent atoms. **(B)** The distribution of the radius across the ligand molecules in DrugBank and NCFD, and **(C)** The distribution of the number of atoms across the ligand molecules in DrugBank and NCFD.

Figure SI. 6: **Annotation bias in top 100 false negative predictions made by DeepPurpose. (A)-(B)** Degree ratio distribution of the nodes involved in the false negative predictions is shown compared to all the nodes in the BindingDB data. The false negative predictions correspond to proteins and ligands with low degree ratios. **(C)-(D)** DeepPurpose predicts lower binding probabilities for the nodes with lower degree ratios.

Figure SI. 7: **Network property of the DrugBank DTI and Venn diagram of positive binding samples across different databases. (A)** Annotation distribution of the proteins and the drugs in DrugBank are fat-tailed. The nature of the annotation distribution is similar to our observations in BindingDB. **(B)** AI-Bind training data combines protein-ligand binding data from three databases: DrugBank, BindingDB, and Drug Target Commons (DTC). Majority of the binding examples are taken from DrugBank. BindingDB and DTC are used to obtain additional protein-ligand pairs, especially to maximize the binding information involving naturally occurring ligands.

Figure SI. 8: **EigenSpokes Analysis.** **(A)** Network-based dimension reduction of nodes in the full protein-ligand network. Node $i$ is represented by the vector $\bar{u}_i = (u_1, u_2, u_3, u_4, u_5) \in \mathbb{R}^5$. Here we visualize $(u_3, u_4)$ for only the ligands. Coloring is based on the hop-distances from an example target BPT4: Green = 1 hop, Blue = 3 hops, Yellow = 5 hops, Orange = 7 hops, Red $\geq 9$ hops. We see that at $> 7$ hop, most nodes are very close to the origin. **(B)** Mean of all reduced vector magnitudes $\|\bar{u}_j\|$ averaged over all pairs $(i, j)$ of a given path length. We see a significant decrease in magnitude as the shortest path length increases.

Figure SI. 9: **Deep architectures of VecNet, VAENet, and Siamese model. (A)** Vec-Net uses Mol2vec and ProtVec as the unsupervised pre-trained models for ligand and protein embeddings respectively. The dense layers act as decoders, and are trained using the network-derived dataset. **(B)** VAENet architecture is similar to VecNet, where Mol2vec embeddings are replaced with embeddings obtained from a variational auto-encoder. This auto-encoder is trained on $\approx 9.5$ million compounds from the ZINC database. **(C)** Siamese model embeds both proteins and ligands onto the same latent space. Siamese ConvNet blocks minimize the triplet loss between the proteins binding to the same ligand. We follow a similar approach for generating the ligand embeddings.

Figure SI. 10: **Logical flow of the Siamese Model and Variational Auto-Encoder**. **(A)** We minimize the embedded Euclidean distances between the proteins which bind to the same ligand, and maximize the distance between the non-binding ones. Similar logic is applied for creating the ligand embeddings. **(B)** Variational auto-encoder minimizes the reconstruction loss for the ligands to create a latent space embeddings. We generate Morgan fingerprints from the isomeric SMILES and feed that to the auto-encoder. The auto-encoder generates latent space representations by minimizing reconstruction loss on the fingerprints.

Figure SI. 11: **Training curves for three AI-Bind architectures.** We plot the training curves for **(A)** VecNet, **(B)** VAENet, and **(C)** Siamese model over 30 epochs. The AUROC and the AUPRC are separately shown for the transductive (unseen edges) and inductive (unseen nodes) test scenarios.

Figure SI. 12: **F1-Score and Optimal Threshold.** We plot the F1-scores for the trained **(A)** VecNet, **(B)** VAENet, and **(C)** Siamese models relative to the classification threshold in the inductive test scenario. The threshold value corresponding to the highest F1-score is considered as the optimal threshold, and is used to obtain the binary labels from the predicted binding probabilities. For VecNet, we obtain an optimal threshold of 0.09.

Figure SI. 13: **Methodology for interpreting AI-Bind predictions.** **(A)** We replace each amino acid trigram in the amino acid sequence with *xxx*, which maps to the out-of-vocabulary ⟨*unk*⟩ vector in ProtVec, and observe the fluctuations in AI-Bind prediction. Removal of some trigrams affect AI-Bind prediction more than others. These trigrams are indicative of the binding location(s) on the protein. **(B)** We identify three active binding sites on the protein TRIM59 from the docking simulations and map them to AI-Bind's binding probability profile. We also identify a possible pocket location from the binding probability profile and visualize that on the 3D protein structure.

Figure SI. 14: **Binding probability profiles for different active binding sites on TRIM59. (A)-(C)** We group the ligands based on the binding pockets on TRIM59 and plot the binding probability profiles, highlighting the binding locations on the amino acid sequence. We observe similar shape of the binding probability profiles for different ligands, but the deviation from the original AI-Bind prediction varies across the ligands, which conveys the dependency of the binding probability profile on the ligand structure.

Figure SI. 15: **Hierarchical clustering on binding probability profiles.** **(A)** We plot the heatmap of the Tanimoto similarities between the ligands binding to TRIM59. We do not observe significant grouping of the ligands solely based on their molecular structures. **(B)** We cluster the ligands based on he similarities of their binding probability profiles. We observe that multiple ligands binding to the same pocket are clustered together in the clustermap. Thus, the binding probability profiles generated by AI-Bind are not only specific to a protein, but carry information about the ligand structures.

Figure SI. 16: **Phylogenetic tree of genes enriched towards prediction bias.** We compare proteins associated with the false predictions (both false positives and false negatives) made by AI-Bind's VecNet to uncover structural similarities. AI-Bind does not show any bias towards certain protein structures in the false predictions, and can be used for binding prediction involving protein structures emerging from different organisms.
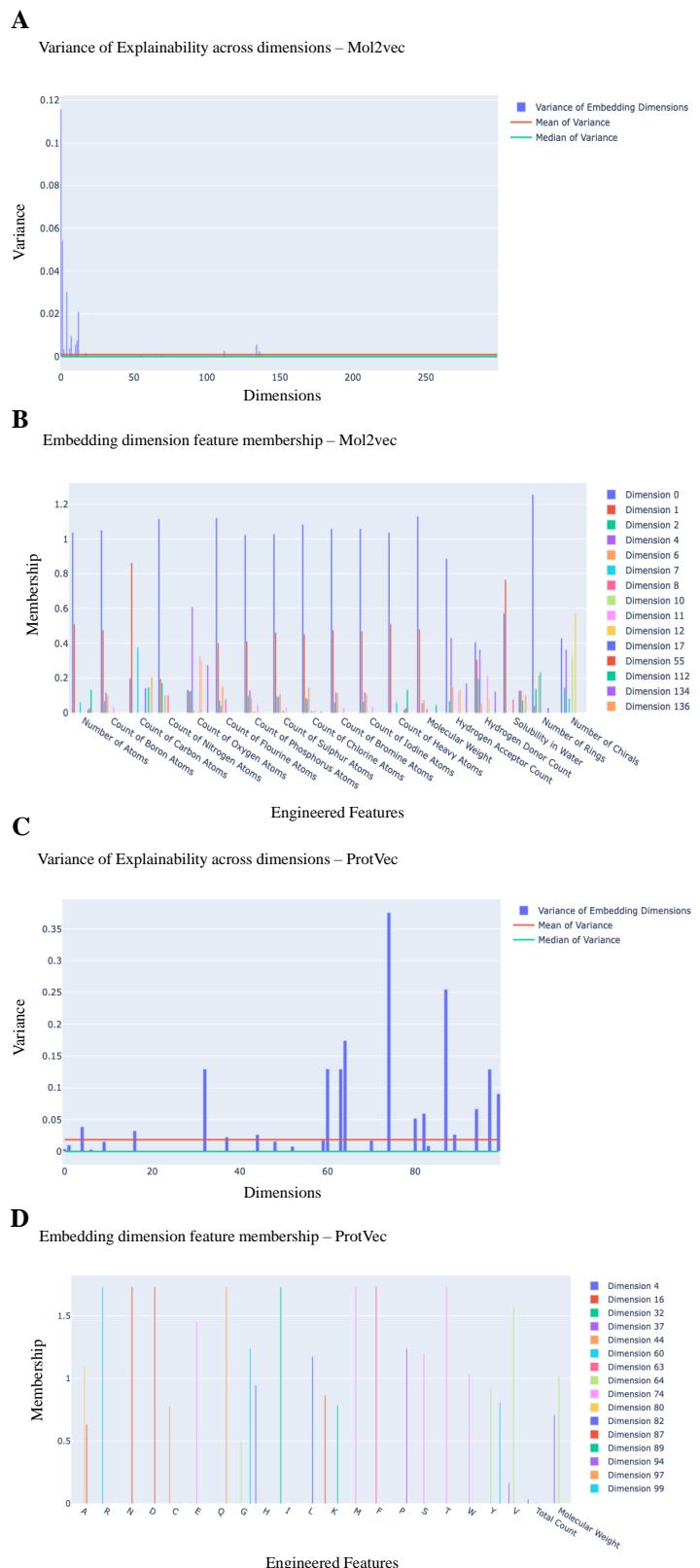
Figure SI. 17: **Dimensions of Mol2vec and ProtVec contributing to protein-ligand binding.** **(A)-(B)** Only 15 Mol2vec dimensions show high variability when explaining the engineered features representing ligand molecules. **(C)-(D)** We find similar results for 16 ProtVec dimensions.
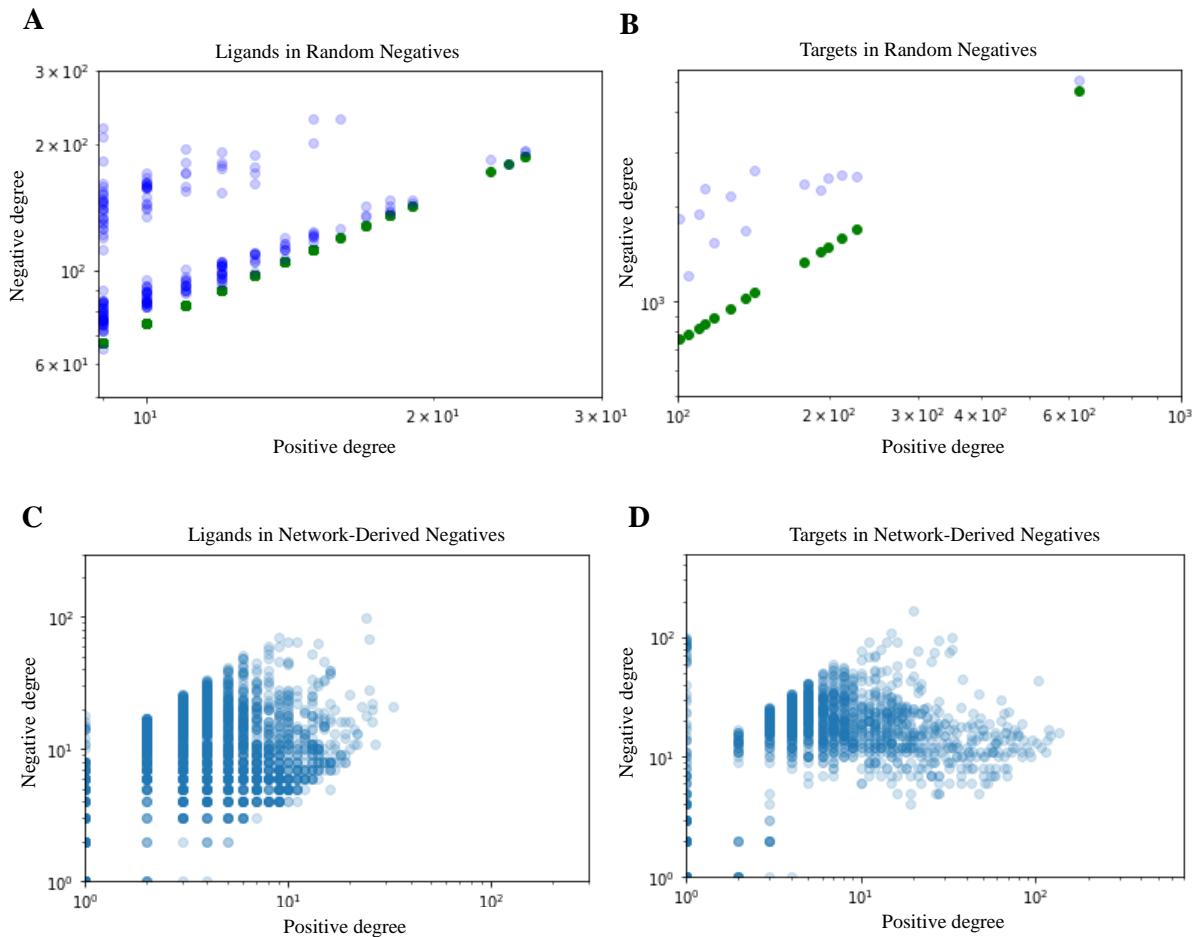
Figure SI. 18: **Random vs. network-derived negatives**. **(A)-(B)** In random negative sampling, both the ligand and the protein on a positive edge has the lower bound of negative degree equal to 7.5 times its positive degree. Higher positive degree nodes have lower probabilities of being present in a random negative sample, as they are present in many positive edges and are discarded more often from getting included in a negative sample. Thus, the negative degree diminishes as the positive degree increases. **(C)-(D)** We observe less correlation between positive and negative degrees for the network-derived negatives. This helps in removing the annotation imbalance we observe in the existing protein-ligand databases.
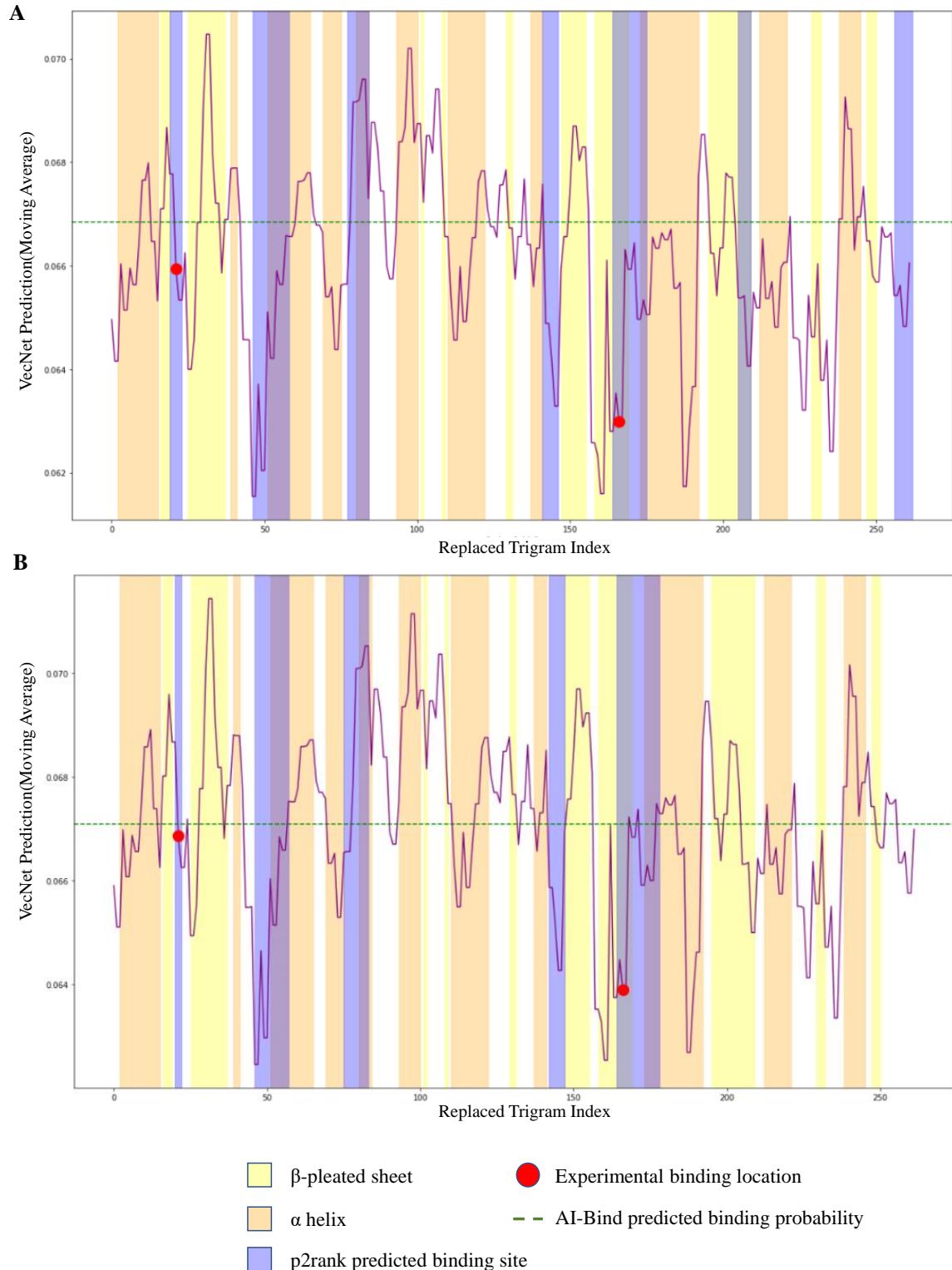
**Figure SI. 19: Experimental binding sites lie on the valleys of the binding probability profile.** We plot the binding probability profile for the *E. Coli* protein Thymidylate Synthase, and the ligands **(A)** SP-722 and **(B)** SP-876. We observe that the experimentally obtained binding sites are in the valleys of the binding probability profile, and overlay on the the $\beta$-sheets and the coils regions. These binding locations also overlap with the binding locations predicted by p2rank, a state-of-the-art binding site prediction algorithm.
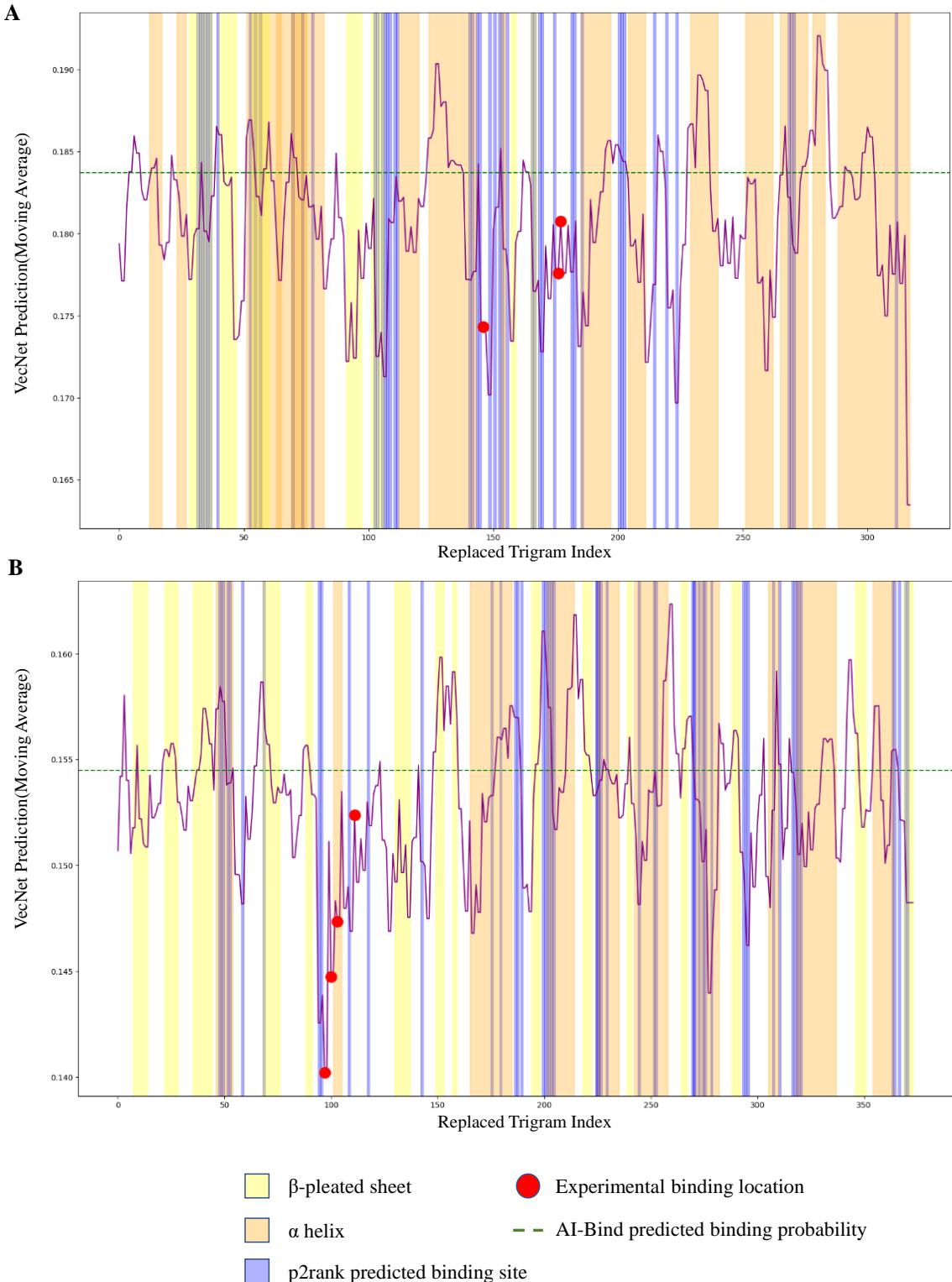
Figure SI. 20: **Experimental binding sites lie on the valleys of the binding probability profile for human proteins TAO3 Kinase and Human Alcohol Dehydrogenase.** We plot the binding probability profile for the human protein and ligand pairs **(A)** Human TAO3 Kinase and ADP **(B)** Human Alcohol Dehydrogenase and Nicotinamide Adenine Dinucleotide. We observe that the experimentally obtained binding sites are in the valleys of the binding probability profile, and often overlay on the the $\beta$-sheets and the coils regions. These binding locations also overlap with the binding locations predicted by p2rank, a state-of-the-art binding site prediction algorithm.