

# PLAS-20k: Extended Dataset of Protein-Ligand Affinities from MD Simulations for Machine Learning Applications

Divya B. Korlepara<sup>1,3,†</sup>, C. S. Vasavi<sup>1,†</sup>, Rakesh Srivastava<sup>2</sup>, Pradeep Kumar Pal<sup>2</sup>, Saalim H. Raza<sup>1</sup>, Vishal Kumar<sup>2</sup>, Shivam Pandit<sup>1</sup>, Aathira G Nair<sup>1</sup>, Sanjana Pandey<sup>1</sup>, Shubham Sharma<sup>1</sup>, Shruti Jeurkar<sup>2</sup>, Kavita Thakran<sup>1</sup>, Reena Jaglan<sup>1</sup>, Shivangi Verma<sup>1</sup>, Indhu Ramachandran<sup>1</sup>, Prathit Chatterjee<sup>1</sup>, Divya Nayar<sup>4,\*</sup>, and U. Deva Priyakumar<sup>2,\*</sup>

<sup>1</sup>IHub-Data, International Institute of Information Technology, Hyderabad, 500032, India.

<sup>2</sup>Centre for Computational Natural Sciences and Bioinformatics, International Institute of Information Technology, Hyderabad, 500032, India.

<sup>3</sup>Division of Physics, School of Advanced Sciences, Vellore Institute of Technology, Chennai, 600127, India.

<sup>4</sup>Department of Materials Science and Engineering, Indian Institute of Technology Delhi, Hauz Khas, New Delhi, 110016, India.

<sup>†</sup>These authors contributed equally to this work: Divya B. Korlepara, C.S. Vasavi.

<sup>\*</sup>corresponding author(s): Divya Nayar, U. Deva Priyakumar (divyanayar@mse.iitd.ac.in, deva@iiit.ac.in)

## ABSTRACT

Computing binding affinities is of great importance in drug discovery pipeline and its prediction using advanced machine learning methods still remains a major challenge as the existing datasets and models do not consider the dynamic features of protein-ligand interactions. To this end, we have developed PLAS-20k dataset, an extension of previously developed PLAS-5k, with 97,500 independent simulations on a total of 19,500 different protein-ligand complexes. Our results show good correlation with the available experimental values, performing better than docking scores. This holds true even for a subset of ligands that follows Lipinski's rule, and for diverse clusters of complex structures, thereby highlighting the importance of PLAS-20k dataset in developing new ML models. Along with this, our dataset is also beneficial in classifying strong and weak binders compared to docking. Further, OnionNet model has been retrained on PLAS-20k dataset and is provided as a baseline for the prediction of binding affinities. We believe that large-scale MD-based datasets along with trajectories will form new synergy, paving the way for accelerating drug discovery.

## Background & Summary

High-throughput screening plays a crucial role in the drug discovery process. However, this approach to identifying lead molecules is time-consuming and labour-intensive. On the other hand, computational methods offer a promising solution by significantly reducing the cost, time, and resources required for physical experiments in screening potential hit molecules. High-throughput docking and molecular dynamics (MD) simulations provide an appealing virtual screening approach to expedite the discovery of biologically active hit compounds<sup>1</sup>. Despite the advantages of these methods, certain limitations and drawbacks still exist in docking. These include a restricted sampling of both protein and ligand conformation during pose prediction and the use of approximated scoring functions that often yield docking scores with poor correlation to experimental binding affinities<sup>2</sup>. On the other hand, MD simulations offer several benefits for investigating the structural and dynamical properties of a Protein-Ligand (PL) system and accurately predicting binding affinities. However, screening of umpteen molecules consumes prohibitively expensive computational resources rendering the prediction of binding affinity (MD based) on a large scale infeasible<sup>3</sup>.

In recent years, machine learning (ML) has emerged as a powerful tool to accelerate various aspects of drug development<sup>4</sup>. ML has already shown to be successful in the hunt for antibiotics<sup>5</sup>, drug re-purposing for emerging diseases<sup>6,7</sup>, virtual screening<sup>8,9</sup>, bio-molecular interactions, prediction of binding site and protein folding<sup>10-14</sup>. Notably, enormous ML models have been developed to predict PL binding affinity<sup>15</sup>. These data-driven approaches have been successful in attaining a high level of accuracy by learning the binding modes directly from rapidly growing experimental three-dimensional (3D) PL structural data deposited in Protein Data Bank (PDB)<sup>16,17</sup>. Numerous attempts have been made to enhance the performance of machine learning (ML) models through different types of encoding, topology, spectral sequence, and atom pairs. These

37 approaches have predominantly relied on feature engineering from static 3D structures<sup>18</sup>. However, this static picture of PL  
38 interactions often lacks dynamic features. Incorporating dynamic properties can provide crucial insights into bio-molecular  
39 processes such as protein folding, conformational changes, and ligand binding. In addition, considering dynamic features can  
40 help address fundamental questions related to binding affinity and specificity<sup>19,20</sup>. The greatest strength of MD simulations lies  
41 in their ability to reveal dynamic effects of the bio-molecules that go beyond the experimentally determined structures available  
42 in PDB<sup>21,22</sup>. Furthermore, MD simulations capture the interactions and energy exchanges between the protein, ligand (solute),  
43 and solvent (water, buffer ions) to dictate the binding event through both long-range and short-range interactions<sup>23–26</sup>. While  
44 existing ML models have shown promise in predicting binding affinity, they often rely on training datasets composed of only a  
45 few hundred static binding poses of PL complexes. With the continuous growth in the number of ligands and proteins, there is  
46 an increasing demand for massive and dynamic data to improve the ML model's accuracy in predicting binding affinities.

47 By integrating MD simulations with ML techniques, researchers can leverage the dynamic nature of biomolecular systems  
48 and incorporate a broader range of data, leading to more accurate and reliable predictions of binding affinities. The combination  
49 of MD simulations and ML holds great potential for accelerating drug discovery efforts in an ever-expanding chemical space.  
50 To this end, in our previous work, we developed a MD-based dataset called PLAS-5k<sup>27</sup>. This dataset included binding affinities  
51 averaged over conformations of each of 5000 PL complexes, representing various classes of enzymes. In addition to the binding  
52 affinities, the dataset also included energy components contributing to the binding free energy.

53 When attempting to accurately predict PL interactions through ML models, a labyrinth of interactions needs to be  
54 accounted for. In continuation of our previous dataset, the current work focuses on expanding heterogeneous proteins and a  
55 large spectrum of ligand types, including small organic molecules and peptides. The extended dataset, encompasses 19,500 PL  
56 structures, providing protein-ligand affinities and non-covalent interaction components, along with accompanying trajectories  
57 suitable for machine learning applications.

58 The creation of the PLAS dataset was primarily motivated by the need for high-quality datasets that can support the  
59 development of advanced algorithms and drive significant advancements in drug development. The PLAS-20k dataset  
60 comprises a diverse collection of protein-ligand (PL) complexes, providing a valuable resource for researchers in the field.  
61 To assess the performance of calculated binding affinities, we conducted comparisons by calculating correlation coefficients  
62 between experimentally determined values and the affinities obtained through molecular mechanics/Poisson-Boltzmann surface  
63 area (MM-PBSA) and docking methods. This evaluation allowed us to validate the accuracy and reliability of the computational  
64 approaches employed. Based on the experimental binding affinities within the PLAS-20k dataset, we categorized the complexes  
65 into strong binders (SB) and weak binders (WB). This classification helps to differentiate between PL complexes with high and  
66 low affinities, providing valuable insights into the range of binding strengths within the dataset. Furthermore, we assessed the  
67 ligand's adherence to Lipinski's Rule of 5, which offers insights into their drug-like properties. As a baseline for comparison,  
68 we retrained the OnionNet framework using our dataset. The availability of large datasets is often considered essential for  
69 successful deep learning applications. Thus, we believe that the PLAS-20k dataset will serve as a catalyst for the development  
70 of data-driven methods in various drug design tasks, including hit identification, lead optimization, and de novo molecular  
71 design. By providing a comprehensive and diverse dataset, the PLAS-20k dataset empowers researchers to more effectively  
72 explore and apply data-driven approaches, leading to advancements in drug discovery and design processes. The dataset's  
73 availability will drive further innovation and contribute to significant progress in the field of drug development.

## 74 **Methods**

### 75 **Data Curation**

76 In this article, we have chosen a set of 14,500 complexes from the Protein Data Bank (PDB)<sup>17</sup>, expanding upon our previous  
77 PLAS-5k<sup>27</sup> dataset. The selection criteria for these complexes focused on proteins that are complex with small molecules  
78 (ligands) or peptides.

### 79 **Dataset Preparation**

80 In this study, we followed the preprocessing and calculation protocol similar to our previous work<sup>27</sup>. A brief account of the  
81 methods is given here. The initial structure of the complexes was taken from PDB<sup>17</sup>. Protein chains with missing residues were  
82 modelled as loop regions using UCSF Chimera<sup>2,28</sup>. Further, the protein chains were protonated at a physiological pH, 7.4 using  
83 H++ server<sup>29</sup>. The tleap program of ambertools<sup>30,31</sup> was used to build the input files of each complex system (protein-ligand,  
84 cofactors and crystal water molecules) files required for MD simulations. The crystal waters were modelled using a TIP3P  
85 force field<sup>32</sup>. The proteins were modelled using Amber ff14SB force field<sup>33</sup> in the all-atom model, and parameters of the ligand  
86 and cofactors were taken from General AMBER force field (GAFF2)<sup>34</sup> using antechamber program<sup>35</sup>. Each complex was  
87 solvated in an orthorhombic TIP3P water box with a 10 Å extension from the protein surface. More detailed information on the  
88 dataset preparation is discussed in our earlier work with 5k complexes<sup>27</sup> and the flowchart for data preparation is shown in  
89 Figure 1. The counter ions were added to maintain the charge neutrality of the system.

MD simulations were performed using OpenMM 7.2.0 program<sup>36</sup>. The simulation protocol involved several steps as described below. To initiate the simulations, we performed a minimization process using the L-BFGS minimizer with a harmonic potential applied to the atoms of the protein backbone. The force constant for this potential was set to 10 kcal/mol/Å<sup>2</sup>. The minimization consisted of 1000 steps, and after every 10 steps, the restraint force on the backbone atoms was reduced by half. Subsequently, an additional 1000 steps of minimization were conducted after removing the harmonic potential entirely.

During the simulation, a time step of 2 fs was used, and constraints were applied to the bonds involving hydrogen atoms. We implemented a Langevin thermostat with a friction coefficient of 5 ps<sup>-1</sup> to maintain the temperature. The system was gradually heated from an initial temperature of 50 K to the target temperature of 300 K, increasing by 1 K every 100 steps (200 fs). The backbone atoms of the protein were restrained using harmonic potentials during this heating process. Once the target temperature was reached, the simulations were performed for 1 ns in the NVT ensemble.

In the next step, the systems were equilibrated in NPT ensemble at 300 K and 1 atm using a Langevin thermostat and Monte Carlo barostat for 2 ns. Finally, a production run of 4 ns in NPT ensemble is performed and the trajectory is saved every 100 ps for post-processing analysis. The final coordinates of the systems were subjected to minimization for 4000 steps. The coordinates at every 1000 steps were saved and used as the initial structures to start the four more independent simulations.

MD trajectories from five independent simulations were used to calculate the binding affinity using MMPBSA (Molecular-Mechanics Poisson Boltzmann Surface Area) method. Here we used a single trajectory approach to estimate the contribution of the complex, ligand, and receptors separately. We considered two explicit water molecules near the active site. The binding affinity is calculated as follows:

$$\Delta G_{MM-PBSA} = \Delta E_{MM} + \Delta G_{Sol} \quad (1)$$

Electrostatic interaction energy  $\Delta E_{ele}$ , and Van der Waals interaction energy  $\Delta E_{vdw}$  contributes to  $\Delta E_{MM}$  (equation (2)) and  $\Delta G_{Sol}$ , is defined as sum of polar  $\Delta G_{pol}$ , and non-polar contributions  $\Delta G_{np}$  (equation (3))

$$\Delta E_{MM} = \Delta E_{ele} + \Delta E_{vdw} \quad (2)$$

$$\Delta G_{Sol} = \Delta G_{pol} + \Delta G_{np} \quad (3)$$

## Data Records

The PLAS-20k dataset is available publicly and can be accessed at (<https://healthcare.iit.ac.in/d4/plas20k/plas20k.html>). The list of PDB ids that are part of PLAS-20k is provided and can be downloaded from the website. The PDB id search icon in the database opens a specific 3D structure along with energy components (Van der Waals interaction energy, electrostatic energy, polar and non-polar solvation free energies in conjunction with binding affinity) from the MD trajectories using the MM-PBSA method. An example of HIV-1 protease complex (PDB id: 1hwx) is shown in Supplementary Figure S1. The binding affinity and energy components for all the complexes can be accessed through <https://figshare.com/s/05a562608b47d1682b8f> in csv format.

## Technical Validation

### Overall Structures of the Protein-Ligand Complexes

Though there are a lot of advances in predicting PL binding affinity through machine learning methods, the incorporation of receptor flexibility remains a major bottleneck. In the present work, we propose a novel dataset based on binding affinities of PL complexes retrieved from MD simulations. The binding affinities were calculated by considering the flexibility of both protein and ligand. The simulated complexes were validated by calculating the RMSD with respect to the experimental structure. The protein structures were superimposed to calculate RMSDs of protein and ligand. These calculations have been performed over 200 frames (40 from each simulation trajectory) and the corresponding distributions are shown in Supplementary Figure S2. The long tails of RMSD distributions of protein and ligand are evident due to the flexibility of the complex during the simulations.

### Comparison of experimental vs computed binding affinities

Experimentally, the binding affinity of a protein-ligand complex is expressed in terms of dissociation constant ( $K_d$ ) or inhibition constant ( $K_i$ ). This experimentally determined binding equilibrium constant is related to binding free energy as,

$$\Delta G_{expt} = -k_B T \ln K_i = -k_B T \ln(1/K_d) \quad (4)$$

131 In this work, for a comparison study, we selected a subset of 6842 complexes of the PLAS-20k dataset, whose experimental  
132 binding affinities are available. To assess the performance of our dataset, the Pearson correlation coefficient ( $R_p$ ) and Spearman  
133 rank correlation coefficient ( $R_s$ ) were calculated. Both these correlation coefficients showed that, studies based on MM-PBSA  
134 have superior performance with ( $R_p$ ) of 0.50 and ( $R_s$ ) of 0.56 compared to docking studies whose ( $R_p$ )&( $R_s$ ) are 0.39 and 0.41  
135 respectively. The corresponding plots are shown in Figure 2. The results highlight the importance of considering both protein  
136 and ligand flexibility. We expect that ML-based scoring functions developed using the PLAS-20k dataset could be more reliable  
137 than classical scoring functions. The distribution of the calculated binding affinity is shown in Supplementary Figure S3.

## 138 Classification of Binders

139 Drug discovery is the process by which lead molecules are identified by screening chemical space based on binding affinity. The  
140 existing ML models or scoring functions were formulated based on several assumptions but they still have certain limitations.  
141 Mostly, researchers are interested in identifying only strong binders (SB), and one of the major reasons for neglecting weak  
142 binding molecules in drug discovery is because of its cross reactivity<sup>37,38</sup>. However, these weak binders (WB) are also equally  
143 important as they play a key role in fragment-based drug design<sup>39</sup> and they serve as a foundation towards the development of  
144 more potent and selective drug candidates with improved therapeutic efficacy.

145 In our dataset, 4343 PL complexes with experimental  $K_{i/d}$  fall into SB and WB categories. This subset is used to classify  
146 SB and WB based on experimental vs MMPBSA and experimental vs docking binding affinities. For experimental binding  
147 affinities, the strong and weak binders were classified with a predefined cut-off value of -8.18 kcal/mol. The corresponding  
148 MMPBSA and docking cut-offs are -38.70 kcal/mol and -6.35 kcal/mol respectively. A brief discussion of the binding affinity  
149 cutoff values is given in detail in Supplementary Information.

150 The classification based on MMPBSA and Docking is shown in Figure 3 and the qualitative performance was evaluated  
151 using the metrics given in Tables 1-2. In Figure 3, the diagonal elements of the confusion matrix represent the number of  
152 correct predictions, while the off-diagonal elements represent incorrect predictions. Based on the evaluation metrics, given in  
153 Tables 1-2 and correlation coefficients (Supplementary Figure S4) it can be observed that MMPBSA classification is performing  
154 better compared to docking scores. Also, the confusion matrix revealed that the majority of SB (true positives) and WB (true  
155 negatives) were correctly identified with respect to MMPBSA, indicating the dataset is good enough to distinguish SB and WB.  
156 The definitions of the evaluation metrics are provided in SI.

## 157 Performance of Diverse Protein Sequences

158 The central goal of any machine learning (ML) model is to get the best model, and its performance depends on training data.  
159 More diverse the training data, one can expect a better model. We have collected a humongous number of complex structures  
160 for this dataset preparation. Our dataset covers 1856 protein families which are of functional significance and a pie chart of the  
161 highly populated family is shown in supplementary Figure S5. Proteins with sequence similarity of  $\leq 40\%$  are grouped and the  
162 correlation coefficients are shown in Supplementary Figure S6. The results highlight the importance of the PLAS-20k dataset  
163 as it shows a good correlation for a diverse set of proteins.

## 164 Performance Based on Ligand Structural Properties

165 In the field of drug discovery, prediction of bio-active molecules are based on several rules such as Lipinski,<sup>40</sup> MDDR-like  
166 rule,<sup>41</sup> Veber rule,<sup>42</sup> and Ghose filter<sup>43</sup>. The physicochemical properties like molecular weight and hydrogen bonding capacity  
167 are important to design drug-like molecules. For a comparison study, we chose a set of ligands with drug-like properties  
168 (Molecular weight  $\leq 500$ , number of hydrogen bond donors  $\leq 5$ , number of hydrogen bond acceptors  $\leq 10$ ) and evaluated the  
169 performance of those complexes based on docking and MMPBSA calculations.

170 As seen in Figure 4, MMPBSA calculations showed good correlation with ( $R_p$ ) of 0.55 and ( $R_s$ ) of 0.57 compared to  
171 docking with ( $R_p$ ),( $R_s$ ) 0.41 and 0.43 respectively. Also, for each of the individual components of drug-like properties,  
172 MMPBSA showed a good correlation compared to docking and the results are shown in Supplementary Figure S7-S9. Further,  
173 as seen in Supplementary Figure S10 our dataset holds diverse ligands highlighting a few molecular descriptors, as they play an  
174 important role in drug discovery.

## 175 Components of the Binding Free Energies

176 Binding free energy is the most important initial indicator of drug potency and remains a major challenge in predicting affinities.  
177 In this work, we have provided binding energies for 19,500 PL complexes along with energy components ( $\Delta E_{ele}$ ,  $\Delta E_{vdw}$ ,  
178 and  $\Delta G_{sol}$ ). This PLAS-20k dataset could be helpful in training ML models for predicting the binding affinities and energy  
179 components. The knowledge of these components can help in lead optimization. The distribution of the energy components is  
180 shown in Supplementary Figure S11. Moreover, the availability of dynamic binding poses from the PLAS-20k dataset can help  
181 in building ML models that can screen lead compounds in a more efficient manner compared to existing methods.

## Machine Learning Baseline

PLAS-20k data was also trained and tested using a deep Convolutional Neural Network (CNN) based model, OnionNet. As ML and deep learning methods have begun to make significant contributions in predicting the binding affinity of a PL complex. The OnionNet model extracts various features from the 3D molecular structure of each PL complex and corresponding binding affinities as input, it then predicts the binding affinity of unknown complexes using deep CNN. The model trained on PLAS-20k data gave an  $R_p$  of 0.91 with an RMSE of 8.15 kcal/mol as shown in Figure 5. This further shows that the PLAS-20k dataset can be used effectively for training various ML and deep learning models.

## Code availability

There is no in-house code used for ML model. We used OnionNet<sup>44</sup> <http://github.com/zhenglz/onionnet/> ML model to train on PLAS-20k dataset.

## References

1. Shim, H., Kim, H., Allen, J. E. & Wulff, H. Pose classification using three-dimensional atomic structure-based neural networks applied to ion channel–ligand docking. *J. Chem. Inf. Model.* **62**, 2301–2315 (2022).
2. Gilson, M. K. & Zhou, H.-X. Calculation of protein-ligand binding affinities. *Annu. review biophysics biomolecular structure* **36**, 21–42 (2007).
3. Osaki, K., Ekimoto, T., Yamane, T. & Ikeguchi, M. 3d-rism-ai: A machine learning approach to predict protein–ligand binding affinity using 3d-rism. *The J. Phys. Chem. B* **126**, 6148–6158 (2022).
4. Karthikeyan, A. & Priyakumar, U. D. Artificial intelligence: machine learning for chemical sciences. *J. Chem. Sci.* **134**, 1–20 (2022).
5. Stokes, J. M. *et al.* A deep learning approach to antibiotic discovery. *Cell* **180**, 688–702 (2020).
6. Morselli Gysi, D. *et al.* Network medicine framework for identifying drug-repurposing opportunities for covid-19. *Proc. Natl. Acad. Sci.* **118**, e2025581118 (2021).
7. Choudhury, C., Murugan, N. A. & Priyakumar, U. D. Structure-based drug repurposing: Traditional and advanced ai/ml-aided methods. *Drug Discov. Today* (2022).
8. Goel, M., Aggarwal, R., Sridharan, B., Pal, P. K. & Priyakumar, U. D. Efficient and enhanced sampling of drug-like chemical space for virtual screening and molecular design using modern machine learning methods. *Wiley Interdiscip. Rev. Comput. Mol. Sci.* **13**, e1637 (2023).
9. Mehta, S., Goel, M. & Priyakumar, U. D. Mo-memes: A method for accelerating virtual screening using multi-objective bayesian optimization. *Front. Medicine* **9** (2022).
10. Chelur, V. R. & Priyakumar, U. D. Birds-binding residue detection from protein sequences using deep resnets. *J. Chem. Inf. Model.* **62**, 1809–1818 (2022).
11. Aggarwal, R., Gupta, A., Chelur, V., Jawahar, C. & Priyakumar, U. D. Deep-pocket: ligand binding site detection and segmentation using 3d convolutional neural networks. *J. Chem. Inf. Model.* **62**, 5069–5079 (2021).
12. Huang, K., Xiao, C., Glass, L. M., Zitnik, M. & Sun, J. Skipggnn: predicting molecular interactions with skip-graph networks. *Sci. reports* **10**, 1–16 (2020).
13. Gainza, P. *et al.* Deciphering interaction fingerprints from protein molecular surfaces using geometric deep learning. *Nat. Methods* **17**, 184–192 (2020).
14. Žitnik, M. *et al.* Gene prioritization by compressive data fusion and chaining. *PLoS computational biology* **11**, e1004552 (2015).
15. Ashtawy, H. M. *Data-Driven and Task-Specific Scoring Functions for Predicting Ligand Binding Poses and Affinity and for Screening Enrichment* (Michigan State University, 2017).
16. Avery, C., Patterson, J., Grear, T., Frater, T. & Jacobs, D. J. Protein function analysis through machine learning. *Biomolecules* **12**, 1246 (2022).
17. Berman, H. M. *et al.* The protein data bank. *Nucleic acids research* **28**, 235–242 (2000).
18. Yang, J., Shen, C. & Huang, N. Predicting or pretending: artificial intelligence for protein-ligand interactions lack of sufficiently large and unbiased datasets. *Front. pharmacology* **11**, 69 (2020).



- 228 **19.** Sinha, S., Tam, B. & Wang, S. M. Applications of molecular dynamics simulation in protein study. *Membranes* **12**, 844  
229 (2022).
- 230 **20.** Du, X. *et al.* Insights into protein–ligand interactions: mechanisms, models, and methods. *Int. journal molecular sciences*  
231 **17**, 144 (2016).
- 232 **21.** Childers, M. C. & Daggett, V. Insights from molecular dynamics simulations for computational protein design. *Mol.*  
233 *systems design & engineering* **2**, 9–33 (2017).
- 234 **22.** Kanakala, G. C., Aggarwal, R., Nayar, D. & Priyakumar, U. D. Latent biases in machine learning models for predicting  
235 binding affinities using popular data sets. *ACS Omega* (2023).
- 236 **23.** Defelipe, L. A. *et al.* Solvents to fragments to drugs: Md applications in drug design. *Molecules* **23**, 3269 (2018).
- 237 **24.** Seo, M.-H., Park, J., Kim, E., Hohng, S. & Kim, H.-S. Protein conformational dynamics dictate the binding affinity for a  
238 ligand. *Nat. communications* **5**, 1–7 (2014).
- 239 **25.** Bronowska, A. K. Thermodynamics of ligand-protein interactions: implications for molecular design. In *Thermodynamics-*  
240 *Interaction Studies-Solids, Liquids and Gases* (IntechOpen, 2011).
- 241 **26.** Wang, C., Greene, D., Xiao, L., Qi, R. & Luo, R. Recent developments and applications of the mmpbsa method. *Front.*  
242 *molecular biosciences* **4**, 87 (2018).
- 243 **27.** Korlepara, D. B. *et al.* Plas-5k: Dataset of protein-ligand affinities from molecular dynamics for machine learning  
244 applications. *Sci. data* **9**, 1–10 (2022).
- 245 **28.** Pettersen, E. F. *et al.* Ucsf chimera—a visualization system for exploratory research and analysis. *J. computational*  
246 *chemistry* **25**, 1605–1612 (2004).
- 247 **29.** Gordon, J. C. *et al.* H++: a server for estimating p k as and adding missing hydrogens to macromolecules. *Nucleic acids*  
248 *research* **33**, W368–W371 (2005).
- 249 **30.** Case, D. A. *et al.* The amber biomolecular simulation programs. *J. computational chemistry* **26**, 1668–1688 (2005).
- 250 **31.** Salomon-Ferrer, R., Case, D. A. & Walker, R. C. An overview of the amber biomolecular simulation package. *Wiley*  
251 *Interdiscip. Rev. Comput. Mol. Sci.* **3**, 198–210 (2013).
- 252 **32.** Jorgensen, W. L., Chandrasekhar, J., Madura, J. D., Impey, R. W. & Klein, M. L. Comparison of simple potential functions  
253 for simulating liquid water. *The J. chemical physics* **79**, 926–935 (1983).
- 254 **33.** Maier, J. A., Martinez, C., Kasavajhala, L., Koushik a nd Wickstrom, Hauser, K. E. & Simmerling, C. ff14sb: improving  
255 the accuracy of protein side chain and backbone parameters from ff99sb. *J. chemical theory computation* **11**, 3696–3713  
256 (2015).
- 257 **34.** Wang, J., Wolf, R. M., Caldwell, J. W., Kollman, P. A. & Case, D. A. Development and testing of a general amber force  
258 field. *J. computational chemistry* **25**, 1157–1174 (2004).
- 259 **35.** Wang, J., Wang, W., Kollman, P. A. & Case, D. A. Automatic atom type and bond type perception in molecular mechanical  
260 calculations. *J. molecular graphics modelling* **25**, 247–260 (2006).
- 261 **36.** Eastman, P. *et al.* Openmm 7: Rapid development of high performance algorithms for molecular dynamics. *PLoS*  
262 *computational biology* **13**, e1005659 (2017).
- 263 **37.** Wang, J. *et al.* Weak-binding molecules are not drugs?—toward a systematic strategy for finding effective weak-binding  
264 drugs. *Briefings Bioinforma.* **18**, 321–332 (2017).
- 265 **38.** Buratto, R., Mammoli, D., Canet, E. & Bodenhausen, G. Ligand–protein affinity studies using long-lived states of  
266 fluorine-19 nuclei. *J. medicinal chemistry* **59**, 1960–1966 (2016).
- 267 **39.** Ohlson, S. Designing transient binding drugs: a new concept for drug discovery. *Drug Discov. Today* **13**, 433–439 (2008).
- 268 **40.** Ivanović, V., Rančić, M., Arsić, B. & Pavlović, A. Lipinski’s rule of five, famous extensions and famous exceptions. *Pop.*  
269 *Sci. Article* **3**, 171–177 (2020).
- 270 **41.** Oprea, T. I. Property distribution of drug-related chemical databases. *J. computer-aided molecular design* **14**, 251–264  
271 (2000).
- 272 **42.** Veber, D. F. *et al.* Molecular properties that influence the oral bioavailability of drug candidates. *J. medicinal chemistry* **45**,  
273 2615–2623 (2002).

- 274 **43.** Ghose, A. K., Viswanadhan, V. N. & Wendoloski, J. J. A knowledge-based approach in designing combinatorial or  
275 medicinal chemistry libraries for drug discovery. 1. a qualitative and quantitative characterization of known drug databases.  
276 *J. combinatorial chemistry* **1**, 55–68 (1999).
- 277 **44.** Zheng, L., Fan, J. & Mu, Y. Onionnet: a multiple-layer intermolecular-contact-based convolutional neural network for  
278 protein–ligand binding affinity prediction. *ACS omega* **4**, 15956–15965 (2019).

## 279 Acknowledgements

280 We thank Akash Ranjan for the website development. We thank IHub-Data for its support. The authors thank IIT Delhi and  
281 IIIT Hyderabad HPC facilities for computational resources. DN acknowledges financial support by INSPIRE faculty research  
282 grant (DST/INSPIRE/04/2018/000455) provided by the Department of Science and Technology, India. UDP thanks DST-SERB  
283 (CRG/2021/008036) and Kohli Center on Intelligent Systems, IIIT Hyderabad for support.

## 284 Author contributions statement

285 UDP conceived the study, D.B.K. and S.H.R wrote the codes and analyzed the data. C.S.V. and D.B.K. contributed to the  
286 writing of the manuscript. S.H.R trained ML model. D.B.K., C.S.V., and S.P. performed docking studies. D.B.K., C.S.V., R.S.,  
287 P.K.P., S.H.R., V.K., S.P., S.S., S.J., S.P., K.T., R.J., S.V., A.G.N., contributed to the preparation of dataset and simulation. D.N.,  
288 and UDP supervised the project. Indhu Ramachandran for coordinating this project and P.C contributed in checking data.

## 289 Competing interests

290 The authors declare no competing interests.

## 291 Figures & Tables

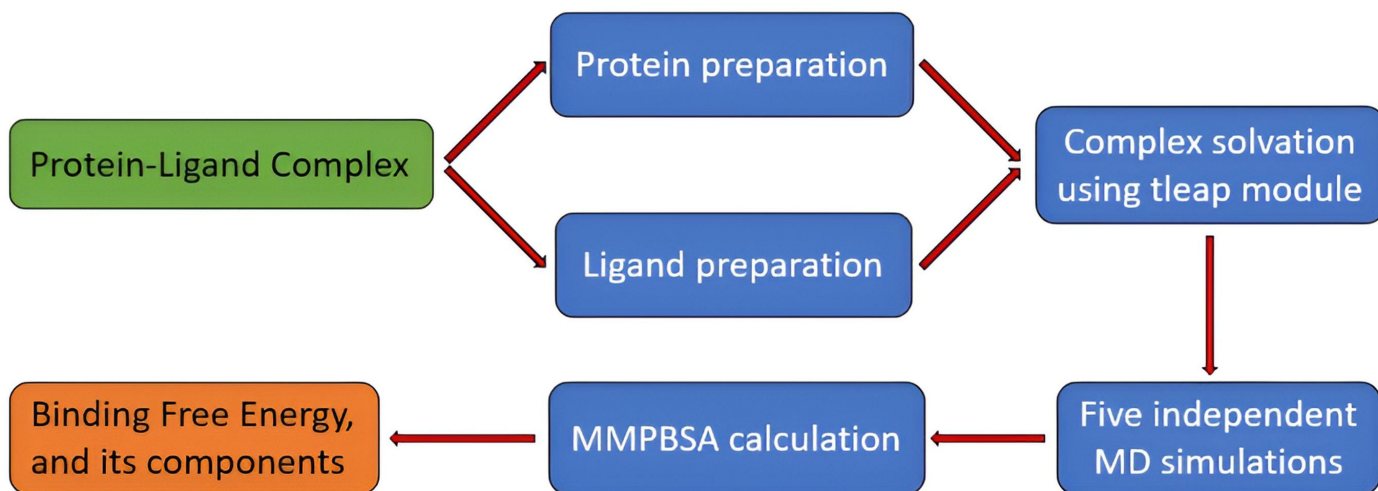
Exp. vs MMPBSA	Precision	Recall	f1-score	support
Strong Binders	0.86	0.78	0.82	2579
Weak Binders	0.72	0.81	0.76	1764
Accuracy			0.79	4343
Macro Average	0.79	0.79	0.79	4343
Weighted Average	0.80	0.79	0.79	4343

**Table 1.** Performance metrics from confusion matrix to evaluate the classification models performance in distinguishing strong and weak binders based on MMPBSA calculations.

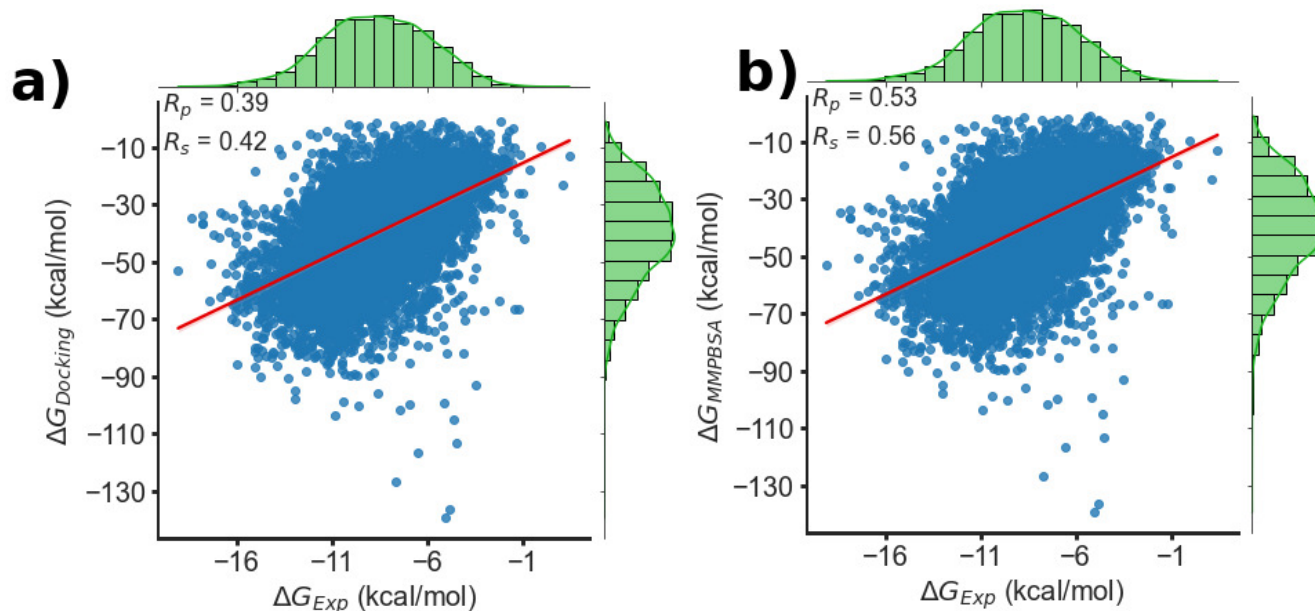
Exp. vs Docking	Precision	Recall	f1-score	support
Strong Binders	0.79	0.74	0.76	2579
Weak Binders	0.65	0.72	0.68	1764
Accuracy			0.73	4343
Macro Average	0.72	0.73	0.72	4343
Weighted Average	0.74	0.73	0.73	4343

**Table 2.** Performance metrics from confusion matrix to evaluate the classification models performance in distinguishing strong and weak binders based on docking simulations.

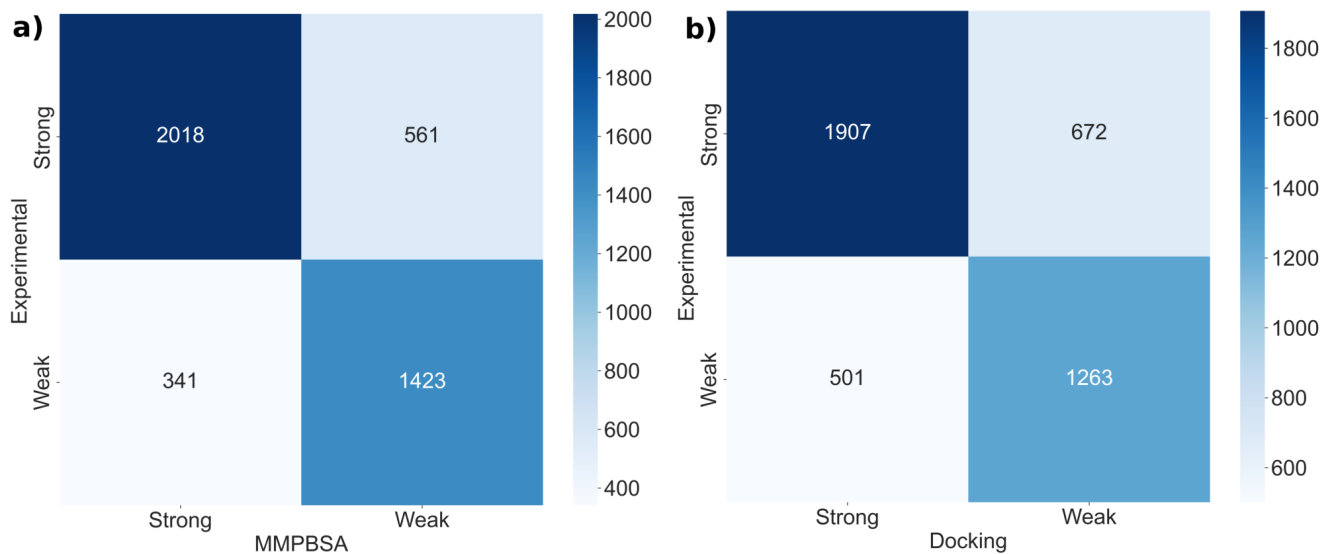




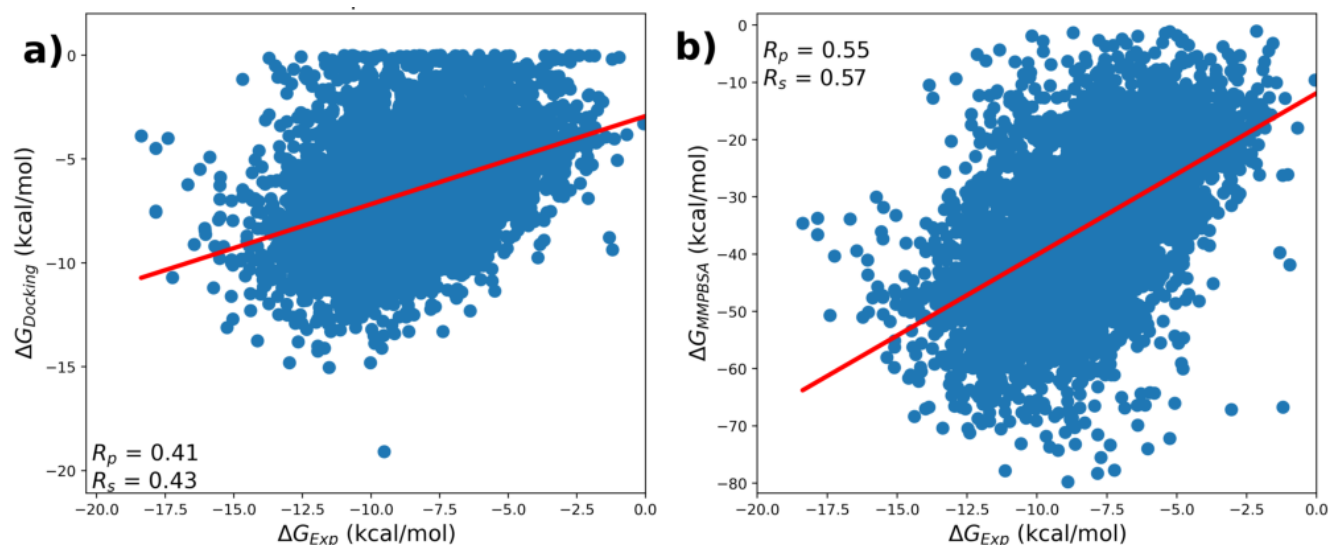
**Figure 1.** Protocol for input preparation and simulations. A similar approach to our earlier work has been followed.<sup>27</sup>



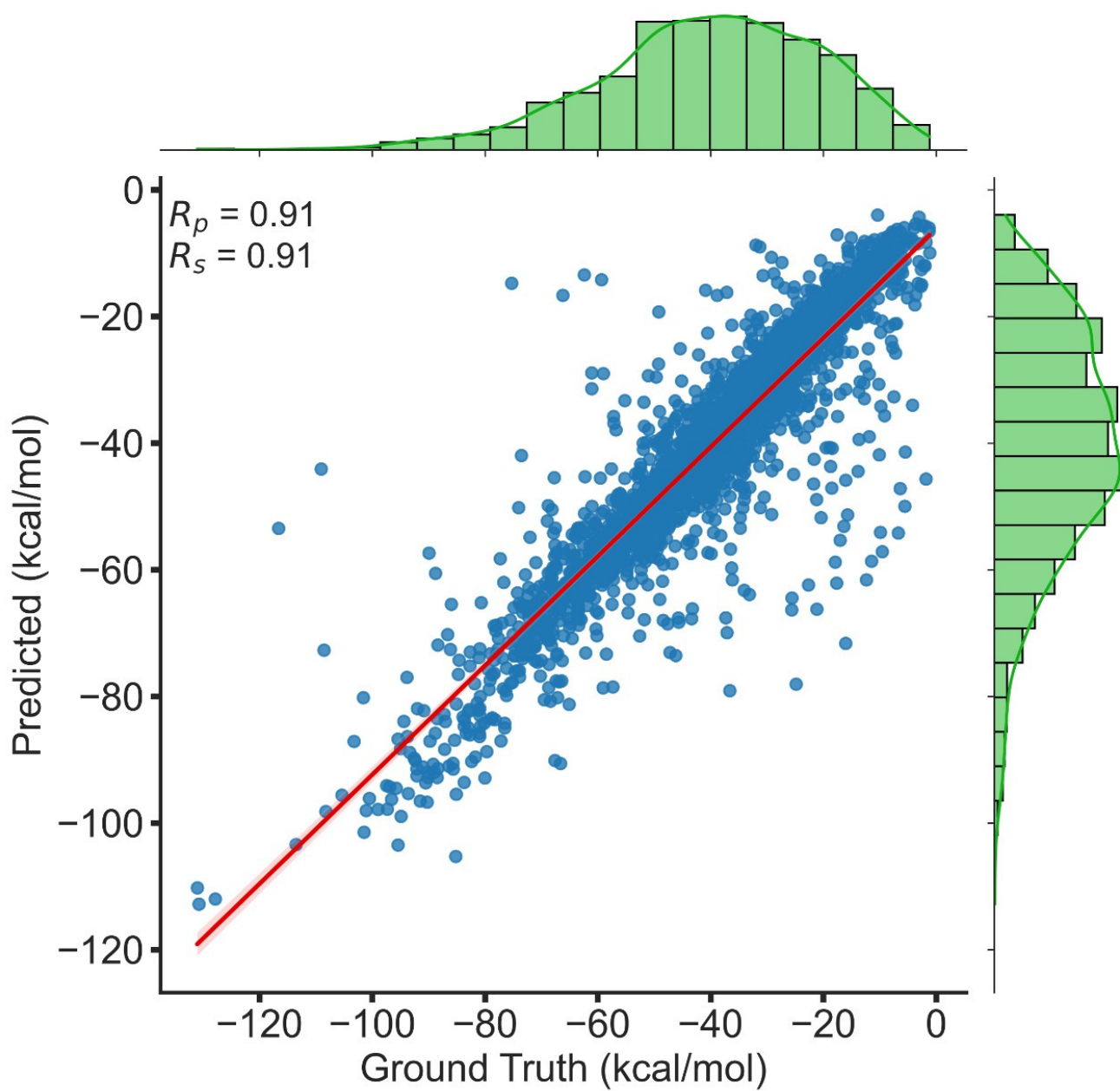
**Figure 2.** Correlation plots between the experimental and calculated binding affinities for a subset with 6842 (includes 2000 data points from PLAS-5k dataset<sup>27</sup>) pdbids. The calculated binding affinities are calculated (a) using Auto-dock Vina, and (b) using MM-PBSA.



**Figure 3.** Confusion matrix to distinguish strong and weak binders (a) Experimental vs MMPBSA, (b) Experimental vs Docking.



**Figure 4.** Correlation plots for a set of PDB ids from PLAS-20k (which follows Lipinski rule of five - Molecular weight, number of donors and number of acceptors of the ligand) for which experimental binding affinities are known - (a) Experimental vs Docking, (b) Experimental vs MM-PBSA



**Figure 5.** Pearson correlation coefficient of OnionNet trained on PLAS-20k dataset.