# Advanced Machine Learning Approaches for Predicting Protein-Ligand Binding Affinity Using 3D Structural and Sequence Information

**Faraz Sarmeili**

**July 2024**

## Abstract

The accurate prediction of binding affinities between proteins and drugs is crucial in drug discovery and development. Traditional protein-ligand docking methods, however, face significant limitations, such as the reliance on high-quality, rigid protein structures, and often fail to account for protein flexibility and conformational changes upon ligand binding. This proposal aims to address these challenges by developing a novel AI-based model that predicts the binding affinity between proteins and drugs using their 3D structures and sequence information. Leveraging recent advances in machine learning and deep learning, the proposed model will predict flexible, all-atom structures of protein-ligand complexes, providing a more accurate and comprehensive approach to binding affinity prediction. This research holds the potential to significantly advance the field of drug discovery, leading to more efficient development pipelines and the discovery of novel therapeutic compounds.

# 1. Introduction

Protein-ligand interactions play a critical role in various biological processes, and understanding these interactions is essential for developing effective drugs. Traditional protein-ligand docking methods, such as AutoDock Vina and Gold, have been widely used for predicting the binding affinity of potential drug candidates. However, these methods are limited by their dependence on high-quality, rigid protein structures. In reality, proteins are dynamic molecules that undergo conformational changes upon ligand binding, and these changes can significantly impact binding affinity. As a result, traditional docking methods often fail to accurately predict binding affinities, particularly for proteins with unknown structures or those exhibiting significant flexibility [1, 2].

In recent years, machine learning and deep learning approaches have shown promise in improving protein-ligand docking accuracy. For instance, Bryant et al. (2024) developed the Umol AI system, which predicts flexible all-atom structures of protein-ligand complexes directly from sequence information. Despite these advancements, challenges remain in accurately predicting binding affinities for proteins with unknown structures or those exhibiting significant flexibility. The proposed research aims to build upon these foundational studies by developing a model that leverages both protein sequence and 3D structural data to predict binding affinities more accurately and comprehensively [3].

The development of a robust model for predicting protein-ligand binding affinities has significant implications for drug discovery and development. Accurate predictions can streamline the drug discovery process, reducing the need for extensive experimental testing and accelerating the identification of promising drug candidates. By addressing the limitations of traditional docking methods and leveraging advances in machine learning, this research has the potential to transform the field of drug discovery, making it more efficient and effective.

## 2. Literature Review

Recent studies have explored various machine learning approaches to improve protein-ligand docking accuracy. Bryant et al. (2024) developed the Umol AI system, which predicts flexible all-atom structures of protein-ligand complexes from sequence information, demonstrating significant improvements over traditional methods. Umol leverages deep learning techniques to predict fully

flexible, all-atom structures of protein-ligand complexes, addressing the limitations of traditional docking methods that treat proteins as rigid or partially rigid entities [3].

In addition to Umol, Harren et al. (2024) provided a comprehensive review of modern machine-learning techniques for binding affinity estimation. The review highlights the potential of deep learning models in improving the accuracy of binding affinity predictions, particularly for proteins with unknown structures or those exhibiting significant flexibility. The authors discuss various approaches, including convolutional neural networks (CNNs), graph neural networks (GNNs), and attention-based models, and their applications in predicting protein-ligand interactions [4].

Other notable contributions include the development of DiffDock by Corso et al. (2022), which uses diffusion models for molecular docking, and RoseTTAFold All-Atom by Baek et al. (2021), which incorporates 3D structural information for enhanced prediction accuracy. DiffDock leverages diffusion processes and deep learning to improve the accuracy of docking predictions, while RoseTTAFold All-Atom integrates structural information to predict interactions between proteins and ligands more accurately. These studies underscore the potential of integrating advanced machine learning techniques with structural bioinformatics to improve binding affinity predictions [5, 6].

While these advancements represent significant progress in the field of protein-ligand docking, challenges remain in predicting accurate binding affinities for proteins with unknown structures or those exhibiting significant flexibility. Traditional docking methods often fail to account for the dynamic nature of proteins, leading to inaccurate predictions. Moreover, many existing machine learning models rely on high-quality structural data, limiting their applicability to proteins with unknown or poorly resolved structures [4].

The proposed research aims to address these limitations by developing a machine learning model that leverages both protein sequence and 3D structural data to predict binding affinities more accurately and comprehensively. By incorporating flexibility and conformational changes into the model, the proposed approach aims to overcome the challenges faced by traditional docking methods and existing machine learning models.

## 3. Research Question

How can a machine learning model be developed to accurately predict the binding affinity of protein-ligand complexes using their 3D structures and sequence information, accounting for protein flexibility and conformational changes?

## 4. Methodology

### 4.1. Data Collection

The dataset for this research will be compiled from existing protein-ligand complexes available in databases such as PDBbind and BindingDB. These databases contain a wealth of experimental data on protein-ligand interactions, including high-quality affinity measurements and diverse structural features. To ensure the robustness and generalizability of the model, the dataset will include a wide range of protein-ligand complexes, covering various protein families, ligand types, and binding modes.

In addition to the existing data, molecular dynamics simulations will be used to generate additional data, capturing a range of protein conformations. These simulations will provide valuable insights into the dynamic nature of protein-ligand interactions, enabling the model to account for flexibility and conformational changes. The generated data will be integrated with the experimental data to create a comprehensive dataset for training and evaluation.

### 4.2. Model Architecture

The proposed model will extend the EvoFormer architecture from AlphaFold2, integrating elements from Umol and other state-of-the-art methods. The EvoFormer architecture is a powerful framework for protein structure prediction, leveraging multiple sequence alignments (MSAs) and structural templates to predict protein structures accurately. In this research, the EvoFormer architecture will be adapted to predict protein-ligand complexes, incorporating flexibility and conformational changes [7].

The model will consist of multiple blocks processing both sequence and structural data, incorporating attention mechanisms to capture interactions between protein residues and ligand atoms. The attention mechanisms will enable the model to focus on relevant regions of the protein

and ligand, capturing the key interactions that determine binding affinity. The model will also include modules for predicting the flexible, all-atom structure of protein-ligand complexes, leveraging the powerful capabilities of deep learning to capture the complex nature of protein-ligand interactions.

## 4.3. Training and Evaluation

The model will be trained using supervised learning techniques, with loss functions designed to optimize both the accuracy of structural predictions and binding affinity estimations. The training process will involve multiple stages, including pre-training on large datasets of protein sequences and structures, followed by fine-tuning on the dataset of protein-ligand complexes. During training, the model will learn to predict the 3D structures of protein-ligand complexes from sequence information, as well as the corresponding binding affinities.

To ensure the robustness and generalizability of the model, cross-validation will be employed during training. The dataset will be divided into multiple folds, with each fold used for training and validation in turn. This approach will help to mitigate overfitting and ensure that the model performs well on unseen data.

Performance will be evaluated using a range of metrics, including Root Mean Square Deviation (RMSD) for structural accuracy and Pearson correlation for binding affinity predictions. RMSD will measure the accuracy of the predicted 3D structures, while Pearson correlation will assess the relationship between the predicted and experimental binding affinities. Additional metrics, such as precision and recall, will also be used to evaluate the model's performance in different scenarios.

## 5. Results and Discussion

The results of the proposed research are expected to demonstrate significant improvements in the accuracy of binding affinity predictions for protein-ligand complexes. By leveraging both sequence and structural data, and incorporating flexibility and conformational changes, the proposed model aims to overcome the limitations of traditional docking methods and existing machine learning models.

The anticipated outcomes include:

1. **Improved Prediction Accuracy:** The model is expected to achieve higher accuracy in predicting the 3D structures of protein-ligand complexes, as measured by RMSD. This will demonstrate the model's ability to capture the dynamic nature of protein-ligand interactions and account for flexibility and conformational changes.

2. **Enhanced Binding Affinity Predictions:** The model is expected to show strong correlations between predicted and experimental binding affinities, as measured by Pearson correlation. This will validate the model's ability to accurately predict binding affinities from sequence and structural data.

3. **Robustness and Generalizability:** The use of cross-validation and comprehensive evaluation metrics will ensure the robustness and generalizability of the model. The model is expected to perform well on unseen data, demonstrating its applicability to a wide range of protein-ligand complexes.

4. **Insights into Protein-Ligand Interactions:** The research will provide valuable insights into the key factors that determine binding affinity, including the roles of flexibility and conformational changes. These insights can inform the development of more effective therapeutic agents and guide future research in the field.

## 6. Future Directions

The proposed research represents a significant step forward in the field of protein-ligand docking and binding affinity prediction. However, there are several avenues for future research that can build upon the findings of this study:

1. **Integration with Experimental Data:** Future research can explore the integration of additional experimental data, such as cryo-electron microscopy (cryo-EM) and nuclear magnetic resonance (NMR) data, to further enhance the accuracy of binding affinity predictions.

2. **Expansion to Other Biomolecular Interactions:** The methodologies developed in this research can be extended to other types of biomolecular interactions, such as protein-

protein and protein-DNA interactions. This can broaden the applicability of the model and contribute to a deeper understanding of biomolecular interactions.

3. **Development of Generative Models:** Building on the predictive capabilities of the proposed model, future research can explore the development of generative models for designing novel ligands with desired binding affinities. These models can leverage the learned representations of protein-ligand interactions to generate new compounds with high binding affinity.

4. **Exploration of Transfer Learning:** Transfer learning techniques can be investigated to improve the model's performance on specific protein families or ligand types. By fine-tuning the model on domain-specific data, researchers can enhance its predictive accuracy for targeted applications.

5. **Collaborative Efforts and Open Science:** Encouraging collaborative efforts and open science practices can accelerate progress in the field. Sharing data, models, and findings with the scientific community can foster innovation and lead to new breakthroughs in drug discovery and development.

## Conclusion

The proposed research aims to develop a novel AI-based model for predicting the binding affinity between proteins and drugs using their 3D structures and sequence information. By leveraging advances in machine learning and deep learning, the proposed model will predict flexible, all-atom structures of protein-ligand complexes, providing a more accurate and comprehensive approach to binding affinity prediction. The research holds the potential to significantly advance the field of drug discovery, leading to more efficient development pipelines and the discovery of novel therapeutic compounds.

The anticipated outcomes include improved prediction accuracy, enhanced binding affinity predictions, robustness, and generalizability, and valuable insights into protein-ligand interactions. The research represents a significant step forward in the field of protein-ligand docking and binding affinity prediction and paves the way for future advancements in drug discovery and development.

# References

1.      Eberhardt, J., et al., *AutoDock Vina 1.2. 0: New docking methods, expanded force field, and python bindings.* Journal of chemical information and modeling, 2021. **61**(8): p. 3891-3898.

2.      David, L., et al., *A toolkit for covalent docking with GOLD: from automated ligand preparation with KNIME to bound protein–ligand complexes.* Bioinformatics Advances, 2022. **2**(1): p. vbac090.

3.      Bryant, P., et al., *Structure prediction of protein-ligand complexes from sequence information with Umol.* Nature Communications, 2024. **15**(1): p. 4536.

4.      Harren, T., et al., *Modern machine-learning for binding affinity estimation of protein–ligand complexes: Progress, opportunities, and challenges.* Wiley Interdisciplinary Reviews: Computational Molecular Science, 2024. **14**(3): p. e1716.

5.      Corso, G., et al., *Diffdock: Diffusion steps, twists, and turns for molecular docking.* arXiv preprint arXiv:2210.01776, 2022.

6.      Baek, M., et al., *Accurate prediction of protein structures and interactions using a three-track neural network.* Science, 2021. **373**(6557): p. 871-876.

7.      Jumper, J., et al., *Highly accurate protein structure prediction with AlphaFold.* nature, 2021. **596**(7873): p. 583-589.