

---

# GraphPrint: Extracting Features from 3D Protein Structure for Drug Target Affinity Prediction

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Accurate drug target affinity prediction can improve drug candidate selection, accel-  
2 erate the drug discovery process, and reduce drug production costs. Previous work  
3 focused on traditional fingerprints or used features extracted based on the amino  
4 acid sequence in the protein, ignoring its 3D structure which affects its binding  
5 affinity. In this work, we propose GraphPrint: a framework for incorporating 3D  
6 protein structure features for drug target affinity prediction. We generate graph  
7 representations for protein 3D structures using amino acid residue location coordi-  
8 nates and combine them with drug graph representation and traditional features  
9 to jointly learn drug target affinity. Our model achieves a mean square error of  
10 0.1378 and a concordance index of 0.8929 on the KIBA dataset and improves  
11 over using traditional protein features alone. Our ablation study shows that the 3D  
12 protein structure-based features provide information complementary to traditional  
13 features.

## 14 1 Introduction

15 Predicting drug and target affinity (DTA) can improve drug candidate selection, accelerate drug  
16 discovery duration, and reduce drug production costs, repurposing existing drugs (1). Different  
17 techniques such as molecular coupling (2; 3; 4), similarity-based methods, (6; 7), network-based  
18 methods (8), and deep learning-based methods. Deep learning methods have been used to perform  
19 DTA, and have gained popularity in recent times as they can have higher accuracy and lower prediction  
20 time. (9; 10; 11; 12; 17)

21 Existing deep learning-based methods use amino acid residue sequences in protein to extract features  
22 directly or learn feature extractors. (9; 10; 11; 12; 17). The amino acid sequence can be of primary,  
23 secondary, tertiary, and sometimes quaternary structure, as shown in figure 1. The primary structure  
24 consists of a sequence of amino acids, and the secondary structure represents the alpha-helical or  
25 beta-sheet structure. The tertiary structure consists of folding the chain itself, leading to a 3D structure.  
26 Proteins with more than one peptide sequence can have additional folding among each other, leading  
27 to a quaternary structure. These higher structures affect the bindings and docking sites of these  
28 proteins. Incorporating the 3D structure of a protein can help improve performance on tasks such as  
29 drug target affinity prediction. With the advent of accurate 3D protein structure prediction (5), it is  
30 now possible to predict the 3D structure of proteins based on amino acid sequence. With the growing  
31 size of the 3D protein structure database, incorporating 3D features to learn drug-related features can  
32 be a new direction to explore in the field of drug discovery.

33 DTA has been framed as a binary classification and regression problem. In binary classification,  
34 binary labels 0 and 1 are predicted. In a regression problem, drug-target affinity is quantified as a

35 regression task, which can be used to rank drug targets in order of their binding affinity. This can be  
36 important for selecting only a limited number of drug candidates for further investigation in the lab.  
37 In this work, we propose a graph neural network-based architecture to learn features based on the 3D  
38 structure of proteins. More specifically, our contributions are:

- 39 1. We propose GraphPrint, a framework to integrate 3D protein structural representation to  
40 learn features using graph neural networks.
- 41 2. We perform an ablation study to show that the 3D structure of protein provides complemen-  
42 tary information to traditional handcrafted features.
- 43 3. We share a curated version of the KIBA dataset along with its 3D protein structure to help  
44 future works.

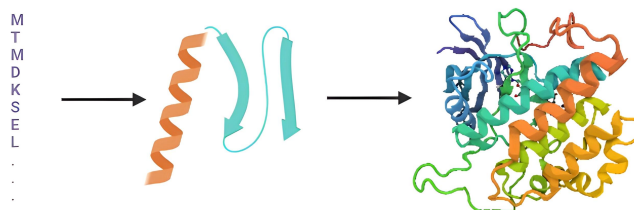


Figure 1: Protein structure can be visualized at several levels. The primary structure involves the structure of amino acids in proteins. The secondary structure involves the formation of helix and sheet structures. The tertiary structure involves the folding of the amino acid chain into a 3D space. Proteins with more than one peptide chain can have a quaternary structure, which involves further folding of chains over each other.

## 45 2 Related works

46 Different techniques such as molecular coupling (2; 3; 4), similarity-based methods, (6; 7), network-  
47 based methods (8), and deep learning-based methods. Deep learning methods have increasingly been  
48 used for drug target affinity prediction. Table 1 shows some of the previous state-of-the-art models  
49 and their backbone architecture. DeepDTA(9) was one of the first deep learning-based architectures  
50 for DTA tasks. The authors used 1D convolution to embed the drug’s SMILES representation and  
51 protein sequence separately, followed by concatenation and classifier training.

52 To improve feature extraction, graph neural network-based architectures have been proposed.  
53 GraphDTA(10) is one such earlier work, to use graph neural networks to learn a drug structure  
54 representation, combined with CNN for protein sequence. iEdgeDTA () treats protein sequence as a  
55 1d graph and uses graph convolutional layers to extract protein features, combined with edge-graph  
56 convolutional for drug embeddings. BiCompDTA use normalized compression distance and Smith-  
57 Waterman measures to generate protein embeddings, that is later used to train deep learning network.  
58 Table 1 summarizes the architecture and input data representations used in these methods.

59 Most previous works treat protein as a sequence of amino acids to either extract hand-crafted features  
60 or learn embeddings using neural network-based architectures. This work ignores the secondary,  
61 tertiary, or quaternary structure of proteins. These higher structures affect the bindings and docking  
62 sites of these proteins. Incorporating the 3D structure of a protein can help improve performance on  
63 tasks such as drug target affinity prediction. With recent works like AlphaFold (5), it is now possible  
64 to predict 3D structures of proteins using their aminoacid sequence.

## 65 3 Methods

66 In this section, we describe a framework for a graph representation of protein 3D structure and our  
67 architecture, followed by the data set details.

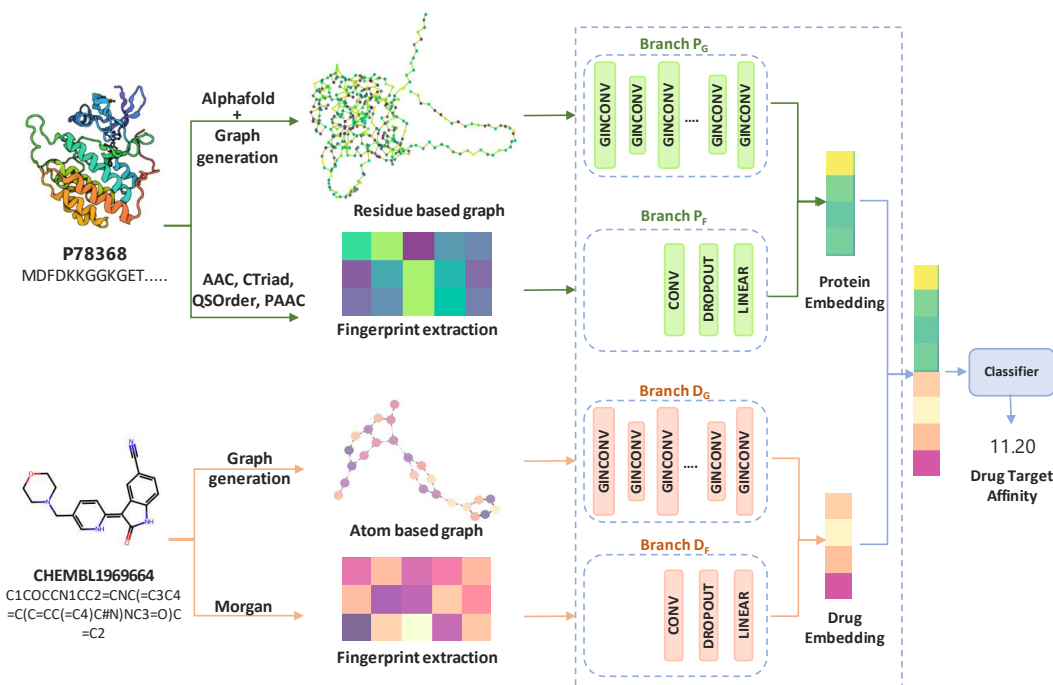


Figure 2: Diagram showing pipeline for feature extraction for protein using AlphaFold (5) and architecture for GraphPrint, with multihead architecture containing graph isomorphic convolution layers (GINCONV) and 1D convolutional blocks, followed by concatenation of features into a multilayer perceptron as a classifier.

Table 1: Previous work for drug target affinity prediction, along with their architectures and inputs used. Keywords: NCD= Normalised compression distance.

Model	Year	Protein		Drug	
		Input	Backbone	Input	Backbone
DeepDTA	2018	AA sequence	CNN	SMILES	CNN
GraphDTA	2020	AA sequence	CNN	SMILES	GNN
iEdgeDTA	2023	AA sequence	1D-GCN	SMILES	Edge-GCN
BiComp-DTA	2023	AA sequence	NCD features	SMILES	CNN

### 68 3.1 Architecture

69 We explore a multimodality approach by using a multi-head network to learn protein/drug embeddings  
 70 directly from their structure as well as their extracted fingerprints. This allows the model to leverage  
 71 the complementary information in these representations.

72 We create a multihead neural network consisting of four branches. We use a graph convolution-based  
 73 network for graph representations of proteins/drugs. For fingerprints, we use 1D convolutional blocks  
 74 to learn embeddings. Figure 2 shows a pictorial representation of the architecture used.

75 Mathematically, for any drug molecule  $D_i$  and protein sequence  $P_i$ , we learn a model  $F_c$ , as shown  
 76 in equation 1, where  $F_{P_G}$ ,  $F_{P_F}$ ,  $F_{D_G}$ ,  $F_{D_F}$  represent branches  $P_G$ ,  $P_F$ ,  $D_G$ ,  $D_F$  of the architecture,  
 77 respectively. The output from three branches is concatenated, before passing through the classifier  
 78  $F_c$ , generating a final prediction of the affinity score  $y_i$ .

$$y_i = F_c([\text{concat}(F_{P_G}(P_i), F_{P_F}(P_i), F_{D_G}(D_i), F_{D_F}(D_i))]) \quad (1)$$

79 Branch  $P_G$  and  $D_G$  consist of 5 graph blocks in series, followed by a global pool average and a  
 80 linear block. This helps to extract graph-level feature embeddings from the drug molecule. We add

81 bottlenecks in the architecture to improve the information transfer enforcing model to learn efficient  
82 information on embedding information on embedding and lower parameter size. The branches  
83  $P_F$  and  $D_F$  consist of a 1D convolutional layer and a linear layer. Outputs from all branches  
84 are concatenated and passed through fully connected layers to generate the final affinity score.

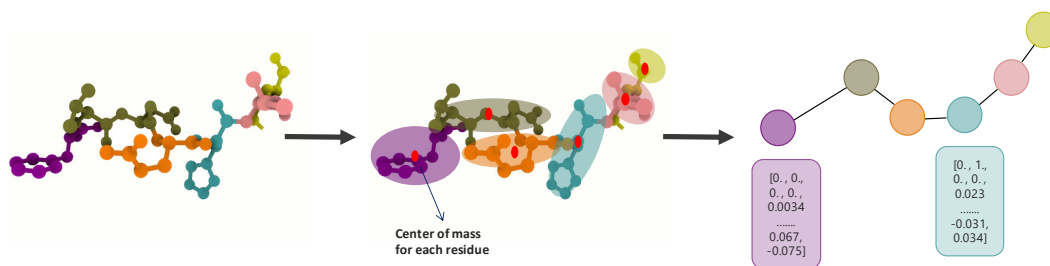


Figure 3: Protein graph representation: We calculate the center of mass of the amino acid and use this as a center of mass for the amino acid residue. Each residue represents a node, with node features containing position and amino acid properties.

### 85 3.1.1 Graph embeddings

86 To generate protein graph representations, we use Alphafold () to generate the 3D protein structure  
87 of each protein target. Using the 3D structure, we define each amino acid residue as a node. We  
88 use the residue’s coordinates of the center of mass to represent amino acids. We add additional  
89 features, such as amino acid encoding, molecular weight, polarity, solubility, and pKa values, to each  
90 node feature. Figure 3 shows the creation of graph representation for protein structure. For drugs,  
91 we convert the SMILES string into a molecular graph representation using the open-source library  
92 RDKit(21). Atoms are represented as nodes, and the bonds between them are represented as edges.  
93 For node embeddings, we use one-hot encodings of atom type, number of neighboring atoms, number  
94 of neighboring hydrogens (explicit and implicit), and implicit valence and aromaticity of molecule.

### 95 3.1.2 Fingerprint extraction

96 We extract traditional hand-designed features to augment the graph representation. We choose AAC,  
97 Conjoint triad fingerprint, and quasi-sequence fingerprint, as they only require the protein’s amino  
98 acid sequence to acquire. AAC encodes  $k$ -mers of amino acids into an 8,420-length bit vector. It  
99 can contain sequence neighborhood (local) information. The conjoint triad fingerprint utilizes a  
100 hand-made 7-letter alphabet to encode a continuous frequency distribution of three amino acids.  
101 This transforms the protein into a homogeneous vector space of a 343-length vector. The quasi-  
102 sequence fingerprint contains a residue pair correlation within its 100-length vector. All three of these  
103 representations are concatenated before being passed to the network. We use the open-source library  
104 iFeature (20) to calculate protein representations.

105 For drug molecules, we extract Morgan and Daylight fingerprints using the drug SMILES sequence,  
106 generating output of size 1024 and 2048 respectively. Morgan’s fingerprint encodes circular radius-2  
107 substructures of a molecule with a 1024-length bit vector. This includes partially disambiguated atom  
108 identifiers. The daylight fingerprint encodes path-based substructures into a 2048-length bit vector.  
109 We use the Therapeutics Data Commons (TDC) open-source library(22) to calculate these features.

## 110 3.2 Datasets

111 In this work, we use the KIBA dataset for evaluation (19). The KIBA dataset contains an affinity  
112 score that combines  $K_d$ , an inhibitor constant ( $K_i$ ), and the half maximum inhibition concentration  
113 ( $IC_{50}$ ). For comparison with previous works, we use a similar dataset generated by filtering all the  
114 drug-protein combinations with fewer than 10 interactions. The final dataset contains a total of 2,111  
115 drugs and 229 proteins. The KIBA score ranges from 0.0 to 17.2 and a larger score represents a  
116 weaker binding affinity.

### 117 3.3 Evaluation

118 After model training is complete, we freeze the architecture and measure the performance on the  
119 test set. We will use the following metrics: Mean Squared Error (MSE) as defined as  $MSE =$   
120  $\frac{1}{n} \sum_{i=1}^n (P_i - Y_i)^2$ , where  $n$  is the number of data points,  $P$  are the predicted affinity values, and  $Y$   
121 are the expected affinity values. Lower MSE is better. Concordance Index (CI): (9) compares the predicted  
122 order of the binding affinity values corresponding to drug-target interactions with ground truth. CI  
123 values greater than 0.8 indicate a strong model. It is defined as  $CI = \frac{1}{Z} \sum_{y_i > y_j} h(p_i - p_j)$ , where  
124  $h(x) = \begin{cases} 1 & x > 0 \\ 0.5 & x = 0 \\ 0 & x < 0 \end{cases}$ .  $r_m^2$  metric: measures external prediction performance  
125 of Quantitative structure-activity relationship (QSAR) models. This metric is defined as  $r_m^2 =$   
126  $r^2 * \left(1 - \sqrt{r^2 - r_0^2}\right)$ , where  $r^2$  is the squared correlation coefficient with intercept, and  $r_0^2$  is  
127 without intercept. A value above 0.5 is good. Spearman correlation measures if two variables are  
128 monotonically related, as defined by  $Spearman = 1 - \frac{6 \sum (P_i - Y_i)^2}{n(n^2 - 1)}$ . Pearson Correlation: a measure  
129 of the linear correlation between two variables, as defined by:  $Pearson = \frac{cov(P, Y)}{\sigma_P \sigma_Y}$ .

130 We explore the ablation of branches to understand their contributions to the learning of our architecture.  
131 We remove one layer at a time and report the metrics of the modified architecture. We also explore  
132 errors contributed by individual drugs and proteins to the overall error, by plotting the sum of MAE  
133 error per drug or protein.

## 134 4 Results and Discussion

135 In this study, we evaluate our trained model on the test subset of the KIBA dataset and compare it to  
136 previously reported baselines. Table 2 shows a comparison of the performance metrics of our models  
137 with DeepDTA, GraphDTA, iEdgeDTA, and BiCompDTA. Our model achieved 0.3713 RMSE, 0.1378  
138 MSE, 0.8929 CI, Spearman correlation of 0.8852, and Pearson correlation of 0.8920. Our model  
139 results are competitive against state-of-the-art models.

Table 2: Performance metrics our architecture with comparison to previous work on KIBA test dataset.  
Keywords: RMSE=root mean square, MSE= mean square error, CI=conformance index

Model	RMSE	MSE	CI	Spearman	Pearson	Epochs
DeepDTA	—	0.194	0.863	—	—	-
GraphDTA	—	0.139	0.891	—	—	1000
iEdgeDTA	—	0.139	0.89	—	—	-
BiComp-DTA	—	0.167	0.891	—	—	-
Ours ( $P_G, P_F, D_G$ )	0.3926	0.1542	0.8790	0.8632	0.8748	300
Ours ( $P_G, P_F, D_G, D_F$ )	0.3713	<b>0.1378</b>	<b>0.8929</b>	0.8852	0.892	300

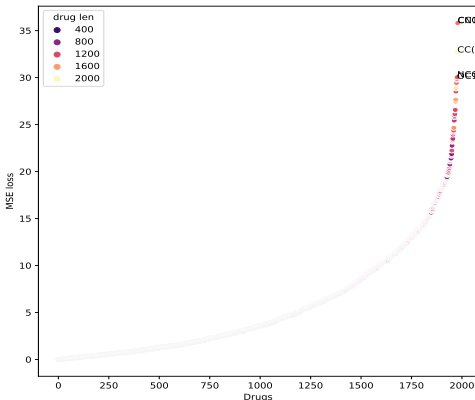
140 We look into the MAE error contribution of individual proteins and drugs to the overall model  
141 MAE error. Figure 4a plots the sum of errors for each protein and drug molecule. As shown in the  
142 figure, only a small number of protein and drug molecules are responsible for the majority of errors.  
143 Figure 5a, 5b, 5c show scatterplot for the sum of error contribution per drug vs a number of drugs  
144 atom counts, aromatic atoms and bonds respectively. We can see that there is a linear relation with  
145 increasing atom count making drug affinity prediction harder. The same goes for the number of  
146 aromatic and number of bonds.

### 147 4.1 Ablation study

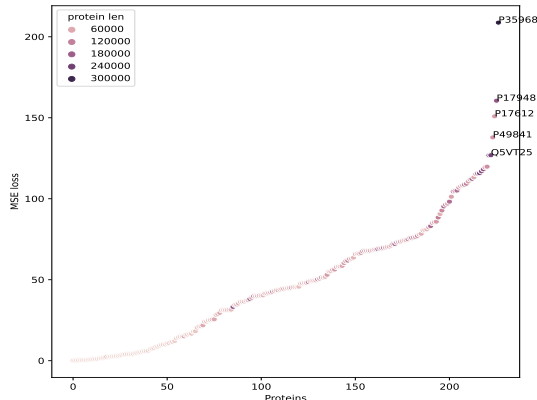
148 We perform architecture ablations by removing one or more branches of the network. At each  
149 stage, the model receives one of the two embeddings for both protein and drug. Table 3 shows  
150 ablation results for our architecture. All the ablations have lower performance compared to our main  
151 architecture. Removing branch  $P_G$  causes a dip in the concordance index and an increase in the  
152 MSE score. This supports our hypothesis that 3D structure provides complementary information  
153 to traditional hand-crafted fingerprints. Further, replacing simple architecture with a bottleneck  
154 architecture causes an increase of 1.3% in the concordance index.

Table 3: Ablations results on different branches of our architecture on KIBA Dataset. Removing branch  $P_G$  causes a dip in the concordance index, suggesting the complementary nature of 3D structure to traditional hand-crafted fingerprints.

Removed Branch	CI	RMSE	MSE	Spearman	Pearson
-	<b>0.8929</b>	<b>0.3713</b>	0.1378	0.8852	<b>0.892</b>
$P_G$	0.8911	0.3740	0.1399	0.8824	0.8899
$D_G$	0.8902	0.3716	<b>0.1381</b>	<b>0.8924</b>	0.8825
$P_F$	0.8820	0.3857	0.1488	0.8838	0.8746
$D_F$	0.8790	0.3926	0.1542	0.8632	0.8748
$D_F, P_F$	0.8783	0.3938	0.1551	0.8626	0.8736
$D_F, P_G$	0.8699	0.4104	0.168	0.8549	0.868



(a) Drug ID vs mse contribution



(b) Protein ID vs mse contribution

Figure 4: Error breakdown based on drug ID, protein ID, aromatic compounds in drugs, bonds inside drug. A small amount of drugs and proteins contribute the most amount of error.

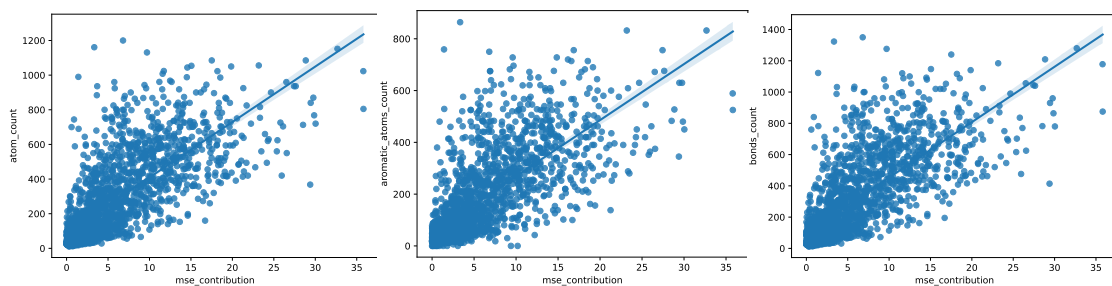
## 155 4.2 Limitations and Future work

156 There are a few limitations to improve on. Since the generation of 3D structures is a computationally  
 157 expensive process, we present our results only on the KIBA dataset. More evaluation on multiple  
 158 datasets is required for a more comprehensive evaluation. In this work, we focus on integrating 3D  
 159 protein structure representation, we do not perform extensively on graph neural architecture. Using  
 160 more recent architecture with attention can provide further boost performance. In the future, this  
 161 work can be further extended and directly implement an explainability layer, pointing to amino acid  
 162 regions interacting with drugs. It would be worthwhile to quantify error-contributing correlations in  
 163 protein structures, as understanding the reasoning behind this might also lead to the development of  
 164 better model architectures.

## 165 5 Conclusion

166 In conclusion, our GraphPrint framework is a novel approach for drug target affinity prediction that  
 167 leverages 3D protein structure features in addition to traditional protein and drug representations. We  
 168 have demonstrated the potential of our model in improving drug candidate selection, accelerating  
 169 drug discovery, and reducing production costs. By incorporating protein 3D structure information,  
 170 our model achieved a mean square error of 0.1378 and a concordance index of 0.8929 on the KIBA  
 171 dataset, surpassing the performance of models that rely solely on traditional protein features.

172 With the advent of accurate 3D protein structure prediction and increasing 3D protein databases,  
 173 incorporating 3D protein structure for drug target affinity and other approaches can be a new direction  
 174 to explore in the field of drug discovery and multimodal learning. This research not only enhances our  
 175 understanding of drug-target interactions but also holds promise for more efficient and cost-effective  
 176 drug development processes in the future. Further exploration and refinement of our model may pave  
 177 the way for even more accurate predictions in this critical area of pharmaceutical research.



(a) Drug atom count vs mse contribution

(b) Drug aromatic atoms count vs mse contribution

(c) Drug bonds count vs mse contribution

Figure 5: Scatter plots showing the mse error contribution vs. different parameters. There is a linear relation between the number of atoms, aromatic atoms, and bonds to the error contribution of the respective drug molecule.

## 178 6 Code and dataset availability

179 To honor the blinded peer review, code and curated dataset of 3D protein structure will be released  
180 after the peer review process. code link:

## 181 7 Acknowledgement

182 The author extends their gratitude to XXXXX XXXXXXXX for his essential role in data processing  
183 and drug fingerprint integration for the pipeline.

## 184 References

- 185 [1] Paul, Steven M., et al. "How to improve R&D productivity: the pharmaceutical industry's grand  
186 challenge." *Nature reviews Drug discovery* 9.3 (2010): 203-214.
- 187 [2] Alonso, Hernan, Andrey A. Bliznyuk, and Jill E. Gready. "Combining docking and molecular  
188 dynamic simulations in drug design." *Medicinal research reviews* 26.5 (2006): 531-568.
- 189 [3] Kontoyianni, Maria. "Docking and virtual screening in drug discovery." *Proteomics for drug  
190 discovery: Methods and protocols* (2017): 255-266.
- 191 [4] Mousavian, Zaynab, and Ali Masoudi-Nejad. "Drug–target interaction prediction via chemoge-  
192 nomic space: learning-based methods." *Expert opinion on drug metabolism & toxicology* 10.9  
193 (2014): 1273-1287.
- 194 [5] "Highly accurate protein structure prediction with AlphaFold | Nature." Accessed: Oct. 05,  
195 2023. [Online]. Available: <https://www.nature.com/articles/s41586-021-03819-2>
- 196 [6] Pahikkala, Tapio, et al. "Toward more realistic drug–target interaction predictions." *Briefings in  
197 bioinformatics* 16.2 (2015): 325-337.
- 198 [7] He, Tong, et al. "SimBoost: a read-across approach for predicting drug–target binding affinities  
199 using gradient boosting machines." *Journal of cheminformatics* 9.1 (2017): 1-14.
- 200 [8] Yang, Xixi, et al. "Modality-dta: Multimodality fusion strategy for drug-target affinity predic-  
201 tion." *IEEE/ACM Transactions on Computational Biology and Bioinformatics* (2022).
- 202 [9] Öztürk, Hakime, Arzucan Özgür, and Elif Ozkirimli. "DeepDTA: deep drug–target binding  
203 affinity prediction." *Bioinformatics* 34.17 (2018): i821-i829.
- 204 [10] Nguyen, Thin, et al. "GraphDTA: predicting drug–target binding affinity with graph neural  
205 networks." *Bioinformatics* 37.8 (2021): 1140-1147.

- 206 [11] Zhijian, Lyu, et al. "GDGRU-DTA: Predicting Drug-Target Binding Affinity Based on GNN  
207 and Double GRU." arXiv preprint arXiv:2204.11857 (2022).
- 208 [12] Thafar, Maha A., et al. "Affinity2Vec: drug-target binding affinity prediction through represen-  
209 tation learning, graph mining, and machine learning." Scientific reports 12.1 (2022): 4751.
- 210 [13] Wan, Fangping, et al. "DeepCPI: a deep learning-based framework for large-scale in silico drug  
211 screening." Genomics, proteomics & bioinformatics 17.5 (2019): 478-495.
- 212 [14] Ma, Dong, Shuang Li, and Zhihua Chen. "Drug-target binding affinity prediction method based  
213 on a deep graph neural network." Mathematical Biosciences and Engineering 20.1 (2023):  
214 269-282.
- 215 [15] Pei, Qizhi, et al. "SMT-DTA: Improving Drug-Target Affinity Prediction with Semi-supervised  
216 Multi-task Training." arXiv preprint arXiv:2206.09818 (2022).
- 217 [16] Gu, Yuliang, et al. "Protein–ligand binding affinity prediction with edge awareness and super-  
218 vised attention." Iscience 26.1 (2023).
- 219 [17] Suviriyapaisal, Natchanon, and Duangdao Wichadakul. "iEdgeDTA: integrated edge information  
220 and 1D graph convolutional neural networks for binding affinity prediction." (2023).
- 221 [18] Lin, Xuan. "DeepGS: Deep representation learning of graphs and sequences for drug-target  
222 binding affinity prediction." arXiv preprint arXiv:2003.13902 (2020).
- 223 [19] Tang, Jing, et al. "Making sense of large-scale kinase inhibitor bioactivity data sets: a compar-  
224 ative and integrative analysis." Journal of Chemical Information and Modeling 54.3 (2014):  
225 735-743.
- 226 [20] Chen, Zhen, et al. "iFeature: a python package and web server for features extraction and  
227 selection from protein and peptide sequences." Bioinformatics 34.14 (2018): 2499-2502.
- 228 [21] Bento, A. Patrícia, et al. "An open source chemical structure curation pipeline using RDKit."  
229 Journal of Cheminformatics 12 (2020): 1-16.
- 230 [22] Huang, Kexin, et al. "Therapeutics data commons: Machine learning datasets and tasks for drug  
231 discovery and development." arXiv preprint arXiv:2102.09548 (2021).
- 232 [23] Zhe Quan, Xuan Lin, Zhi-Jie Wang, Yan Liu, Fan Wang, and Kenli Li, 'A system for learning  
233 atoms based on long short-term memory recurrent neural networks', in 2018 IEEE International  
234 Conference on Bioinformatics and Biomedicine, pp. 728–733, (2018).
- 235 [24] Ehsaneddin Asgari and Mohammad RK Mofrad, 'Continuous distributed representation of  
236 biological sequences for deep proteomics and genomics', PloS one, 10(11), e0141287, (2015).