# METADTA: META-LEARNING-BASED DRUG-TARGET BINDING AFFINITY PREDICTION

**Eunjoo Lee, Jiho Yoo, Huisun Lee**
Standigm Inc.
Seoul, Korea, 06234
{eunjoo.lee,jiho.yoo,huisun.lee}@standigm.com

**Seunghoon Hong**
School of Computing, KAIST
Daejeon, Korea, 34141
seunghoon.hong@kaist.ac.kr

## ABSTRACT

We propose a meta-learning-based model for drug-target binding affinity prediction (MetaDTA), for which no information on the protein structures or binding sites is available. We formulate our method based on the Attentive Neural Processes (ANPs) (Kim et al., 2019), where the binding affinities for each target protein are modeled as a regression function of the compounds. Known drug-target binding affinity pairs are used as support set data to determine the regression function. We designed few-shot prediction experiments with a small number of support set data, similar to the typical situations in actual drug discovery processes. Experimental results showed that the proposed method outperforms the sequence-based baseline models with the same amount of limited data.

## 1 INTRODUCTION

In the early stage of drug discovery, accurate prediction of the binding behavior between the target protein and drugs is crucial for discovering the candidate molecule with decent potency and selectivity (Hughes et al., 2011). The drug-target binding affinity (DTA) prediction is a regression problem that aims to predict experimentally measured binding affinity values, which is helpful for ranking and optimizing the compounds. Although it is usually more difficult than the drug-target interaction (DTI) prediction problem, which is a binary classification of the active/inactive compounds, the DTA prediction has been actively tackled based on the recent advances in deep learning (Ragoza et al., 2017; Stepniewska-Dziubinska et al., 2018; Jiménez et al., 2018; Zhang et al., 2019; Jones et al., 2021; Abbasi et al., 2020; Öztürk et al., 2018; Abbasi et al., 2020; Nguyen et al., 2020b; 2021).

DTA prediction becomes challenging when the exact structural information of the target protein and the binding location are not precisely determined. This often occurs, especially when the drug discovery program aims to conquer disease by controlling a novel target protein that has not been studied much. For those cases, structure-free binding affinity prediction methods are required. Recently, the DTA prediction methods based on the protein sequence were proposed (Öztürk et al., 2018; Huang et al., 2021) which bypasses the use of three-dimensional structural information of the target proteins. Those methods successfully predict binding affinity without using the structural information. However, the performance is degraded when the number of training data samples for the target protein is limited (Huang et al., 2021).

In the hit identification and lead optimization procedure, designing, synthesizing, and measuring the potency of molecules are time-consuming and expensive processes. These limit the number of molecules to be tested for a target protein, reducing the size of the dataset. In
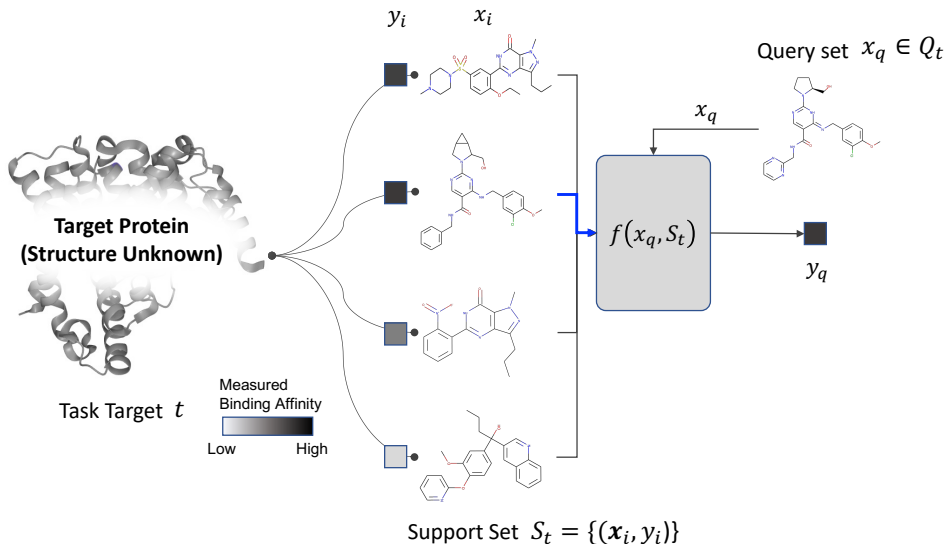
Figure 1: Schematic of ligand-based meta-learning protein drug-target binding affinity prediction

such a small data set, the performance of the deep learning approaches mentioned above is reduced, so developing a model that can predict well even in a few-shot manner is challenging. We can think of this as a few-shot learning problem where a binding affinity prediction function is determined with a small number of given data (Figure 1). Instead of using the intrinsic protein features such as the sequences, we use the known (compound, binding affinity) pairs for a target protein as the features. Based on this setting, we propose a meta-learning model for drug-target binding affinity prediction (MetaDTA), where the regression functions that map between the given compounds to corresponding binding affinities are modeled with Attentive Neural Processes (ANPs) (Kim et al., 2019). Each regression function for each task represents the binding affinities for a specific target protein. By stochastically splitting the compounds into support and query sets, we generate various episodes for a task and use them for training the model.

ANPs have many advantages to be applied for binding affinity predictions. First, ANPs are few-shot learning models specialized for regression problems. Next, the model architecture, consisting of a cross-attention-based deterministic path and a global latent path, represents the prediction strategies well. The deterministic path can serve as a ligand similarity-based prediction (Altae-Tran et al., 2017), and the latent path can be viewed as a target protein feature extraction. Finally, ANPs do not underfit the given context data, which fits well for the binding affinity prediction problem where the context data is usually the most reliable experimental result.

Numerical experiments on the binding affinity dataset were performed to compare the performance with existing sequence-based prediction models (Öztürk et al., 2018; Huang et al., 2021). The advantages of the MetaDTA are summarized as follows:

- The accuracy of MetaDTA is better than the sequence-based models, especially when a small number ($\leq 100$) of data pairs are available. (§4.2)

- MetaDTA can accurately predict the binding affinity even when support set data is used only in test time. (§4.1.2)

- The accuracy of MetaDTA is gradually improved when the training data contains more diverse proteins, even though the model does not directly exploit structural information. (§4.3)

## 2 RELATED WORKS

### 2.1 DRUG-TARGET BINDING AFFINITY PREDICTION

Machine learning has been studied as a promising approach for drug-target binding affinity prediction. Structure-based methods use the experimentally measured 3D structure of the protein-ligand binding complexes (Ragoza et al., 2017; Stepniewska-Dziubinska et al., 2018; Jiménez et al., 2018; Zhang et al., 2019; Jones et al., 2021). However, obtaining such structural information in high resolution is still a challenging task, and virtual generation of the combined structure from docking simulations often produces inaccurate results; hence the prediction accuracy degrades in these cases.

Structure-free methods try to predict the affinity without using the exact 3D binding structures. Instead, the models use simple but abundant and diverse protein features, such as amino acid sequences (Öztürk et al., 2018; Abbasi et al., 2020; Nguyen et al., 2020b; 2021; Huang et al., 2021), the motifs (Öztürk et al., 2019), Protein Sequence Composition (PSC)(Feng et al., 2018), or Position-Specific Scoring Matrix (PSSM) features(Mousavian et al., 2016). These methods achieve somewhat promising results. However, learning proper protein embeddings and interactions is difficult, especially when only limited data for the target protein is available.

Most models try to learn the embedding of drugs and target proteins, as well as the interaction between them (Lim et al. (2021) and the references therein). In this case, the model for predicting the binding affinity $y$ between drug $d$ and target $t$ from drug information $\boldsymbol{x}_d$ and target information $\boldsymbol{x}_t$ is formulated as

$$y \quad = \quad f_\theta\big(g_\phi(\boldsymbol{x}_d), h_\psi(\boldsymbol{x}_t)\big) \tag{1}$$

where $g_\phi$ represents the embedding function for drugs, $h_\psi$ represents the embedding function for targets, and $f_\theta$ represents the interaction function between the drugs and targets. Although there exists millions of binding affinity pair data, the learned model does not generalize well for the protein targets with a small number of samples (Huang et al., 2021).

### 2.2 FEW-SHOT LEARNING FOR DTI PREDICTION

Meta-learning, which is often referred to as 'learning to learn,' aims to improve the ability to deal with a new task (Koch et al., 2015; Vinyals et al., 2016; Finn et al., 2017; Garnelo et al., 2018a; Hospedales et al., 2021). Few-shot learning is a representative problem of meta-learning, which makes predictions for the target task from only a few available examples.

Altae-Tran et al. (2017) proposed a few-shot learning method for DTI prediction. The binary interaction label was predicted as the weighted sum of context labels, where the weights were computed from the similarity between the learned embeddings of the drugs. Few-shot learning showed adequate performance for various molecular property prediction tasks, including DTI prediction (Nguyen et al., 2020a) and drug response prediction (Ma et al., 2021). Recently, a benchmark dataset (Stanley et al., 2021) was proposed for the few-shot DTI prediction problem. Existing few-shot learning methods focused on the DTI problems, and to the best of our knowledge, few-shot learning methods specialized for the DTA prediction problem have not been proposed.

## 3 METHOD

### 3.1 FEW-SHOT LEARNING FOR BINDING AFFINITY PREDICTION

We consider a problem of binding affinity prediction as a few-shot learning problem. To this end, we consider a task $\mathcal{T}_t$ as predicting the binding affinities for a specific target protein $t$. A support set $\mathcal{S}_t = \{(\boldsymbol{x}_i, y_i)\}$ for the task is the set of ligand $\boldsymbol{x}_i$ and binding affinity $y_i$ pairs, which represents the target protein. A query set $\mathcal{Q}_t$ for the task is the set of ligands whose binding affinities sholud be predicted by the model. Then, the problem of predict the binding affinity $y$ of a query ligand $\boldsymbol{x}_q \in \mathcal{Q}_t$ is formulated as

$$y_q \quad = \quad f_\theta(\boldsymbol{x}_q, \mathcal{S}_t). \tag{2}$$
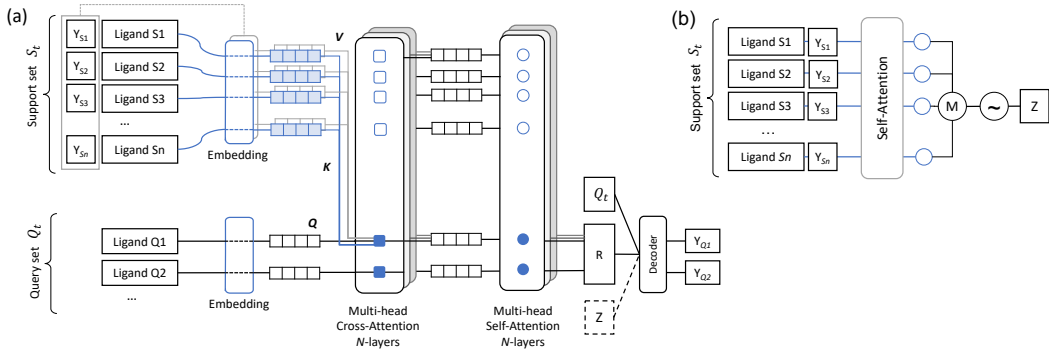
Figure 2: (a) MetaDTA model architecture. (b) Latent path prior (Kim et al., 2019)

Note that this formulation does not directly use the target protein information. Instead of extracting information from the target protein, the target-specific binding affinity function $f_\theta$ is determined from available data pairs. Various and abundant episodes for the task $\mathcal{T}_t$ are generated by splitting available training data pairs into support and query sets during training. The number of possible episode combinations is much greater than the number of training data pairs, so the meta-learning formation has the potential to overcome the over-fitting for a limited amount of training data pairs which is conventional in the methods following (Eq. 1).

## 3.2 METADTA MODEL

We employ the ANPs to model the binding affinity prediction functions. As in the Neural Processes (NPs) (Garnelo et al., 2018b), ANPs consist of deterministic and latent paths. The deterministic part models the conditional distributions conditioned on the support set $\mathcal{S}_t$ to model the query set $\mathcal{Q}_t$, which is formulated as

$$p(y_q|\boldsymbol{x}_q, \mathcal{S}_t) := p(y_q|\boldsymbol{x}_q, \boldsymbol{r}_q), \tag{3}$$

where a query-specific representation $\boldsymbol{r}_q := a_W(\mathcal{S}_t, \boldsymbol{x}_q)$ is modeled with a cross-attention mechanism between the query and the items in the support set. We employ the multi-head, scaled dot-product attention (Vaswani et al., 2017),

$$\text{MultiHead}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) \quad := \quad \text{concat}(\text{head}_1, \text{head}_2, \text{head}_N)\boldsymbol{W} \tag{4}$$

$$\text{head}_h \quad := \quad \text{DotProduct}(\boldsymbol{Q}\boldsymbol{W}_h^Q, \boldsymbol{K}\boldsymbol{W}_h^K, \boldsymbol{V}\boldsymbol{W}_h^V) \tag{5}$$

$$\text{DotProduct}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) \quad := \quad \text{softmax}\left(\frac{\boldsymbol{Q}\boldsymbol{K}^\top}{\sqrt{d_k}}\right)\boldsymbol{V}. \tag{6}$$

In the cross-attention module, $\boldsymbol{Q}$ is the feature of the query set ligands, and $\boldsymbol{K}$ is the feature of the support set ligands. $\boldsymbol{V}$ is the embedded vector of the support set binding affinities.

The differences between our model and ANPs are in this cross-attention module. We used the embedded vector of $y_i$ as $\boldsymbol{V}$ to focus on the binding affinities, while ANPs used the embedded vector of concatenated $(\boldsymbol{x}_i, y_i)$. Based on the intuition of the ligand-based prediction (Altae-Tran et al., 2017), the attention which represents the relationship between the support and the query ligands is used as a weight distribution of support binding affinities. We constructed deeper models by adding cross-attention modules or additional self-attention modules after cross-attention modules (Figure 2).

The latent path uses a global latent variable $\boldsymbol{z}$ to model different realizations of regression functions, which is written by

$$p(y_q|\boldsymbol{x}_q, \mathcal{S}_t) := \int p_\rho(y_q|\boldsymbol{x}_q, \boldsymbol{z})q_\xi(\boldsymbol{z}|\mathcal{S}_t)d\boldsymbol{z}, \tag{7}$$

where the prior function $q_\xi$ models the latent variable conditioned on the support set. In the binding affinity prediction problem, the latent path can be considered to extract the representation of target protein from the given support set $\mathcal{S}_t$ and use it to predict the binding affinities. Considering both paths lead to the following conditional distribution:

$$p(y_q|\boldsymbol{x}_q, \mathcal{S}_t) := \int p_\rho(y_q|\boldsymbol{x}_q, a_W(\mathcal{S}_t, \boldsymbol{x}_q), \boldsymbol{z})q_\xi(\boldsymbol{z}|\mathcal{S}_t)d\boldsymbol{z}. \tag{8}$$

Here the functions $a_W(\mathcal{S}_t, \boldsymbol{x}_q)$ and $q_\xi(\mathcal{S}_t)$ forms the encoder, and the likelihood $p_\rho(y_q|\boldsymbol{x}_q, \boldsymbol{r}_q, \boldsymbol{z})$ corresponds to the decoder.

The training of the model is done by maximizing the following ELBO

$$\begin{aligned}
\log p(y_{\mathcal{Q}_t}|\boldsymbol{x}_{\mathcal{Q}_t}, \mathcal{S}_t) \quad \geq \quad & \mathbb{E}\left[\log p(y_{\mathcal{Q}_t}|\boldsymbol{x}_{\mathcal{Q}_t}, a_W(\mathcal{S}_t, \boldsymbol{x}_{\mathcal{Q}_t}), \boldsymbol{z})\right] \\
& - D_{\mathrm{KL}}\left[q_\xi(\boldsymbol{z}|\mathcal{Q}_t \cup \mathcal{S}_t)||q_\xi(\boldsymbol{z}|\mathcal{S}_t)\right],
\end{aligned} \tag{9}$$

by using stochastically splitted query and support sets.

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETTING

The primary objective of our work is to develop a binding affinity prediction method for the practical drug-discovery scenario, where limited binding affinity measurements are only available for the target protein without three-dimensional protein structure information. To show the behavior of the proposed model, we designed two sets of experiments to answer the following questions:

- Does MetaDTA predict binding affinity well with small test samples?
- Does MetaDTA generalize well for the drugs/targets not contained in the training data?
- How much data is required for good generalization performance?

The first set of experiments tests the model performances for a limited number of few-shot test data (4.2), and the second one tests for a limited amount of training target diversities (4.3).

### 4.1.1 DATASET

We used the BindingDB (Gilson et al., 2016) for the experiments, since it is one of the largest databases of experimentally measured binding affinities. We used the measurements labeled with the half-maximal inhibitory concentration ($IC_{50}$), representing the compound concentration required for 50% inhibition. The BindingDB contains 1,390,284 inhibition affinities with 5,960 proteins, for 689,301 small molecules. Since $IC_{50}$ values distribute over several orders of magnitude, we used $pIC_{50} = -\log_{10}(IC_{50})$ values for the prediction targets.

The outlier detection using interquartile (IQR) method was applied to filter out extraordinarily large or small values. The multiple measurements for single molecule are averaged out, but if the standard deviation of the duplicated values is more than 0.5, the corresponding drug-target pair is discarded. Only the targets with the binding affinities with more than 20 ligands were used. As a result of the process, 1,040,153 inhibition affinities with 2,968 proteins and 662,438 small molecules have remained in the dataset.

For the ligand representation, we used the extended-connectivity fingerprints (ECFPs) (Rogers & Hahn, 2010) with 2,048 dimensions and a diameter of 6. The simplified molecular-input line-entry system (SMILES) string was used for the baseline models as a ligand representation. MetaDTA did not use the explicit protein features, but the support set data pairs were used to characterize the target protein. The $pIC_{50}$ values were converted to 32-dimensional vectors using a Gaussian kernel to increase model sensitivity towards affinity values. The baseline methods used protein amino acid sequences as protein features.
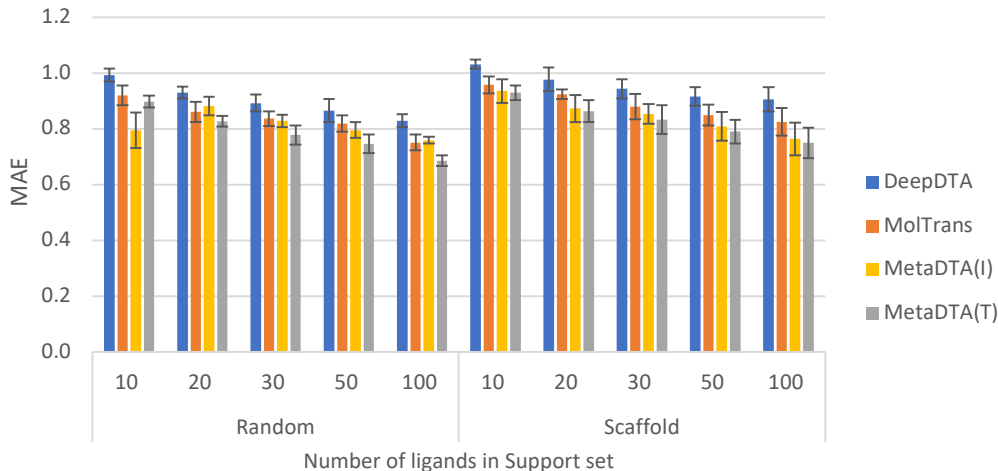
Figure 3: Mean Absolute Error (MAE) of prediction models with respect to the number of ligands in few-shot data by *Random* and *Scaffold* split. MetaDTA(T) shows the best performance in prediction with few-shot ligands data.

### 4.1.2 BASELINES AND METRICS

We compared our MetaDTA model with two baseline models, DeepDTA(Öztürk et al., 2018) and MolTrans(Huang et al., 2021), which do not rely on the three-dimensional structure of target proteins. Both methods use the protein sequences as protein features and SMILES strings as compound features. DeepDTA uses convolutional neural networks (CNNs), and MolTrans uses the Transformers to extract information of the target protein and compound features. MolTrans was originally proposed for DTI prediction methods, so we modified the final activation of MolTrans to a linear function.

For MetaDTA, we can use the few-shot data for test tasks in two different ways. The first is the inference mode (**MetaDTA(I)**), which does not use the few-shot data pairs in training, but only uses them in the test tasks as support sets. The other is the training mode (**MetaDTA(T)**), in which the few-shot data were used to train the model and predict the test task as support sets. The baseline methods used the few-shot data to train the model with other train data.

The metrics for model evaluation of regression problems are Root Mean Squared Error (RMSE), Mean Absolute Error(MAE), Concordance Index (CI), and $r_m^2$, and details are described in A.1. Averaged results of 5-fold cross-validation were provided as a model performance.

### 4.2 BINDING AFFINITY PREDICTION WITH LIMITED DATA

In this experiment, we evaluated model performances when the number of available data for test targets was limited. A limited number of data points were selected as the few-shot data for test targets. The maximum sizes of the few-shot data for a target were 10, 20, 30, 50, and 100. The models were evaluated for two different settings. The overall ligands were not explicitly split into the training and test types in **Random** setting. In **Scaffold** setting, the ligands were split into training and test types based on the Murcko scaffold (Bemis & Murcko, 1996). The similar ligands were either entirely in the training or test data in the latter case. *Scaffold* setting is much more difficult because the model should predict with a limited number of previously unseen ligands.

In general, the measured MAE of the models was improved as the size of few-shot data increased (Figure 3). Except for some cases, MetaDTA(I) outperforms the baseline models,

Table 1: Metrics of prediction models with 50 ligands of few-shot data split by random choice for test tasks. MetaDTA(T) model has the best performance in binding affinity prediction for test proteins.

|  | RMSE | MAE | $r_m^2$ | CI |
|---|---|---|---|---|
| **DeepDTA** | 1.109±0.039 | 0.866±0.042 | 0.427±0.014 | 0.746±0.014 |
| **MolTrans** | 1.070±0.038 | 0.819±0.029 | 0.456±0.012 | 0.764±0.009 |
| **MetaDTA(I)** | 1.070±0.041 | 0.795±0.029 | 0.460±0.036 | 0.759±0.009 |
| **MetaDTA(T)** | **1.009±0.026** | **0.747±0.033** | **0.521±0.036** | **0.779±0.007** |

Table 2: Metrics of prediction models with 50 ligands of limited few-shot data split by Murcko scaffold. MetaDTA(T) and MetaDTA(I) show better performance than sequence-based methods.

|  | RMSE | MAE | $r_m^2$ | CI |
|---|---|---|---|---|
| **DeepDTA** | 1.171±0.036 | 0.916±0.033 | 0.400±0.079 | 0.720±0.021 |
| **MolTrans** | 1.113±0.031 | 0.849±0.037 | 0.432±0.071 | 0.746±0.018 |
| **MetaDTA(I)** | 1.092±0.062 | 0.809±0.051 | 0.462±0.098 | 0.764±0.012 |
| **MetaDTA(T)** | **1.066±0.052** | **0.790±0.042** | **0.491±0.089** | **0.769±0.006** |

even though it does not update the model parameters using the few-shot data of the test tasks. MetaDTA(T) consistently showed the best performance for all cases. In *Scaffold* setting, the prediction errors increased than that of *Random* setting. The tendency is the same for the other metrics (Table 1, 2). The performance of $r_m^2 > 0.5$ is generally believed to be predictive, and MetaDTA(T) achieved this level with over 50 few-shot ligands in *Random* setting and 100 ligands in *Scaffold* setting (Appendix. Table 5).

### 4.3 PROTEIN CLUSTER SPLIT

Although the above experiments were conducted for the target proteins not contained in the training data, there is still a possibility of bias due to the structural similarity of the target proteins, which could result from very similar affinity patterns for ligands. For an appropriate train/test split of target proteins, Ragoza et al. (2017) exploited protein sequence similarity and binding site similarity. However, sequence similarity does not provide a complete split because some proteins have similar structures, but their sequences are different. Clustering based on binding site similarity excludes the protein targets whose binding sites are not determined, resulting in a decrease in the size of the dataset. To obtain an appropriate protein split for a larger number of protein targets, we exploited AlphaFold (Jumper et al., 2021) predicted protein structures to cluster the protein domains.

All available structure models for human proteins (23,391 structures as of Feb 10, 2022) were downloaded from the AlphaFold Protein Structure Database's FTP site (Varadi et al., 2022). We divided a single protein into domains by identifying long disordered regions and only used 1,778 proteins that have a single domain and exist in bindingDB. The application calculates the number of inter-residue contacts using a distance cutoff of 4 Å for every residue. The starting position of a domain boundary is determined as the first residue of a 10-residue window when all the consecutive residues in the window have inter-residue contacts and the ending position as the first residue of a 5-residue window with no inter-residue contact under the condition of the defined starting position. Domains were extracted from each AlphaFold structure model using a set of the starting and end positions (i.e., domain boundaries). We performed an all-against-all structure comparison of 1,778 protein domains. Structural similarities between two domains were measured using TM-align (Zhang & Skolnick, 2005) and normalized by the length of the smaller domain. The similarity score, TM-score (Zhang & Skolnick, 2004) has the value in (0, 1], where 1 indicates a perfect match. The proteins are grouped by hierarchical clustering into 32 different classes with respect to the distances (1 − TM-score).
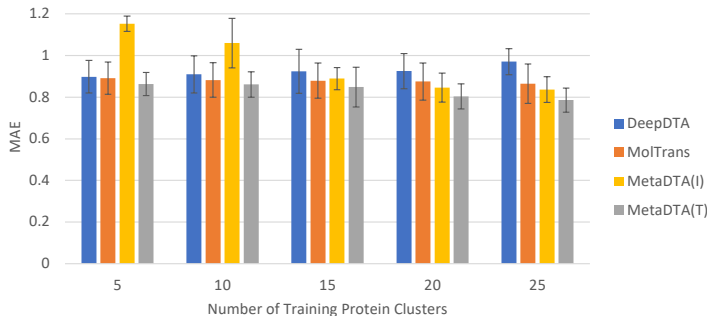
Figure 4: Mean Absolute Error (MAE) of prediction models with respect to the number of training protein clusters. MetaDTA gets better prediction performance while the available target clusters increase.

Table 3: Metrics of model prediction trained with 25 protein clusters.

|  | RMSE | MAE | $r_m^2$ | CI |
|---|---|---|---|---|
| **DeepDTA** | 1.220±0.083 | 0.971±0.063 | 0.262±0.054 | 0.641±0.081 |
| **MolTrans** | 1.122±0.112 | 0.865±0.095 | 0.326±0.047 | 0.665±0.066 |
| **MetaDTA(I)** | 1.107±0.103 | 0.837±0.063 | 0.408±0.071 | 0.704±0.059 |
| **MetaDTA(T)** | **1.058±0.095** | **0.786±0.058** | **0.442±0.049** | **0.720±0.062** |

We applied 5-fold cross-validation by using the 32 protein clusters and varied the number of protein clusters from 5 to 25 for each training fold. We limited the number of ligands in the few-shot data to 50. MetaDTA models showed steady improvements in test accuracy as the number of training protein groups increased (Figure 4). For the MetaDTA, the training dataset's diversity helps to improve the model performance for the unseen, structurally dissimilar protein targets. Although MetaDTA(I) performances were relatively poor at small training protein clusters, they became better than the baseline models with sufficiently diverse training data. The model's performance is generally expected to improve with more diverse training data. However, such improvement was not evident in the baseline models. For sufficiently diverse protein data given (25 clusters), MetaDTA(T) showed the best performance overall metrics (Table 3).

## 5    CONCLUSION

We proposed a meta-learning model for drug-target binding affinity prediction (MetaDTA). The model predicts the binding affinity for query compounds by applying attention mechanisms to the support compounds and known affinity data. Numerical experiments showed that the proposed method achieves higher prediction accuracy than the baseline models without additional information about the target protein. In the experiments using the AlphaFold-predicted structure-based protein clustering, the proposed method showed gradual performance improvement for more diverse protein information. In contrast, the other methods did not show such improvements.

## REFERENCES

Karim Abbasi, Parvin Razzaghi, Antti Poso, Massoud Amanlou, Jahan B Ghasemi, and Ali Masoudi-Nejad. DeepCDA: deep cross-domain compound–protein affinity prediction through LSTM and convolutional neural networks. *Bioinformatics*, 36(17):4633–4642, 05 2020.

Han Altae-Tran, Bharath Ramsundar, Aneesh S. Pappu, and Vijay Pande. Low data drug discovery with one-shot learning. *ACS Central Science*, 3(4):283–293, 2017.

Guy W Bemis and Mark A Murcko. The properties of known drugs. 1. molecular frameworks. *Journal of medicinal chemistry*, 39(15):2887–2893, 1996.

Qingyuan Feng, Evgenia V. Dueva, Artem Cherkasov, and Martin Ester. PADME: A deep learning-based framework for drug-target interaction prediction. *arXiv preprint*, arXiv:1807.09741, 2018.

Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, pp. 1126–1135, 2017.

Marta Garnelo, Dan Rosenbaum, Chris J. Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo J. Rezende, and S. M.Ali Eslami. Conditional Neural Processes. In *Proceedings of the 35th International Conference of Machine Learning (ICML)*, volume 4, pp. 2738–2747, 2018a.

Marta Garnelo, Jonathan Schwarz, Dan Rosenbaum, Fabio Viola, Danilo J. Rezende, S. M. Ali Eslami, and Yee Whye Teh. Neural processes. In *ICML Workshop on Theoretical Foundations and Applications of Deep Generative Models*, 2018b.

Michael K Gilson, Tiqing Liu, Michael Baitaluk, George Nicola, Linda Hwang, and Jenny Chong. BindingDB in 2015: A public database for medicinal chemistry, computational chemistry and systems pharmacology. *Nucleic Acids Research*, 44(D1):D1045–D1053, 2016.

Mithat Gönen and Glenn Heller. Concordance probability and discriminatory power in proportional hazards regression. *Biometrika*, 92(4):965–970, 2005.

Timothy M. Hospedales, Antreas Antoniou, Paul Micaelli, and Amos J. Storkey. Meta-learning in neural networks: A survey. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, preprint, 2021.

Kexin Huang, Cao Xiao, Lucas M. Glass, and Jimeng Sun. Moltrans: Molecular interaction transformer for drug-target interaction prediction. *Bioinformatics*, 37(6):830–836, 2021.

JP Hughes, S Rees, SB Kalindjian, and KL Philpott. Principles of early drug discovery. *British Journal of Pharmacology*, 162(6):1239–1249, 2011.

José Jiménez, Miha Škalič, Gerard Martínez-Rosell, and Gianni De Fabritiis. $K_{DEEP}$: Protein–Ligand Absolute Binding Affinity Prediction via 3D-Convolutional Neural Networks. *Journal of Chemical Information and Modeling*, 58(2):287–296, 2018.

Derek Jones, Hyojin Kim, Xiaohua Zhang, Adam Zemla, Garrett Stevenson, W. F. Drew Bennett, Daniel Kirshner, Sergio E. Wong, Felice C. Lightstone, and Jonathan E. Allen. Improved Protein–Ligand Binding Affinity Prediction with Structure-Based Deep Fusion Inference. *Journal of Chemical Information and Modeling*, 61(4):1583–1592, 2021.

John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A A Kohl, Andrew J Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, 2021.

Hyunjik Kim, Andriy Mnih, Jonathan Schwarz, Marta Garnelo, Ali Eslami, Dan Rosenbaum, Oriol Vinyals, and Yee Whye Teh. Attentive neural processes. In *International Conference on Learning Representations (ICLR)*, 2019.

Gregory Koch, Richard Zemel, and Ruslan Salakhutdinov. Siamese neural networks for one-shot image recognition. In *Deep Learning Workshop, International Conference of Machine Learning (ICML)*, 2015.

Sangsoo Lim, Yijingxiu Lu, Chang Yun Cho, Inyoung Sung, Jungwoo Kim, Youngkuk Kim, Sungjoon Park, and Sun Kim. A review on compound-protein interaction prediction methods: Data, format, representation and model. *Computational and Structural Biotechnology Journal*, 19:1541–1556, 2021.

Jianzhu Ma, Samson H Fong, Yunan Luo, Christopher J Bakkenist, John Paul Shen, Soufiane Mourragui, Lodewyk F A Wessels, Marc Hafner, Roded Sharan, Jian Peng, and Trey Ideker. Few-shot learning creates predictive models of drug response that translate from high-throughput screens to individual patients. *Nature Cancer*, 2(2):233–234, 2021.

Zaynab Mousavian, Sahand Khakabimamaghani, Kaveh Kavousi, and Ali Masoudi-Nejad. Drug–target interaction prediction from PSSM based evolutionary information. *Journal of Pharmacological and Toxicological Methods*, 78:42–51, 2016.

Cuong Q. Nguyen, Constantine Kreatsoulas, and Kim M. Branson. Meta-Learning GNN Initializations for Low-Resource Molecular Property Prediction. *arXiv preprint*, arXiv:2003.05996, 2020a.

Thin Nguyen, Hang Le, Thomas P Quinn, Tri Nguyen, Thuc Duy Le, and Svetha Venkatesh. GraphDTA: predicting drug–target binding affinity with graph neural networks. *Bioinformatics*, 37(8):1140–1147, 10 2020b.

Tri Minh Nguyen, Thin Nguyen, Thao Minh Le, and Truyen Tran. GEFA: early fusion approach in drug-target affinity prediction. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 2021.

Hakime Öztürk, Arzucan Özgür, and Elif Ozkirimli. DeepDTA: Deep Drug–Target Binding Affinity Prediction. *Bioinformatics*, 34(17):i821–i829, 2018.

Hakime Öztürk, Elif Ozkirimli, and Arzucan Özgür. WideDTA: prediction of drug-target binding affinity. *arXiv preprint*, arXiv:1902.04166, 2019.

Partha Pratim Roy, Somnath Paul, Indrani Mitra, and Kunal Roy. On Two Novel Parameters for Validation of Predictive QSAR Models. *Molecules*, 14(5), 2009.

Matthew Ragoza, Joshua Hochuli, Elisa Idrobo, Jocelyn Sunseri, and David Ryan Koes. Protein–ligand scoring with convolutional neural networks. *Journal of Chemical Information and Modeling*, 57(4):942–957, 2017.

David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling*, 50(5):742–754, 2010.

Megan Stanley, John F Bronskill, Krzysztof Maziarz, Hubert Misztela, Jessica Lanini, Marwin Segler, Nadine Schneider, and Marc Brockschmidt. FS-mol: A few-shot learning dataset of molecules. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.

Marta M Stepniewska-Dziubinska, Piotr Zielenkiewicz, and Pawel Siedlecki. Development and evaluation of a deep learning model for protein-ligand binding affinity prediction. *Bioinformatics*, 34(21):3666–3674, 2018.

Mihaly Varadi, Stephen Anyango, Mandar Deshpande, Sreenath Nair, Cindy Natassia, Galabina Yordanova, David Yuan, Oana Stroe, Gemma Wood, Agata Laydon, Augustin Žídek, Tim Green, Kathryn Tunyasuvunakool, Stig Petersen, John Jumper, Ellen Clancy, Richard Green, Ankur Vora, Mira Lutfi, Michael Figurnov, Andrew Cowie, Nicole Hobbs, Pushmeet Kohli, Gerard Kleywegt, Ewan Birney, Demis Hassabis, and Sameer Velankar. AlphaFold Protein Structure Database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Research*, 50(D1):D439–D444, 2022.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NuerIPS)*, pp. 5998–6008, 2017.

Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, volume 29, 2016.

Haiping Zhang, Linbu Liao, Konda Mani Saravanan, Peng Yin, and Yanjie Wei. DeepBindRG: a deep learning based method for estimating effective protein–ligand affinity. *PeerJ*, 7(e7362), 2019.

Yang Zhang and Jeffrey Skolnick. Scoring function for automated assessment of protein structure template quality. *Proteins*, 57(4):702–710, 2004.

Yang Zhang and Jeffrey Skolnick. TM-align: a protein structure alignment algorithm based on the TM-score. *Nucleic Acids Research*, 33(7):2302–2309, 2005.

## A    APPENDIX

### A.1    EVALUATION METRICS

We evaluate the performance of the proposed model using most common evaluation indexes Root Mean Squared Error (RMSE), Mean Absolute Error (MAE) and Concordance Index (CI) (Gönen & Heller, 2005) for regression problems. Equation for CI is described in 10 and 11. The CI index measures whether the predicted values of two random binding affinities are predicted in the same order as the true values. $r_m^2$ index also be used to evaluate the predictive performance of QSAR models and model with $r_m^2 > 0.5$ for the test set was considered acceptable (Pratim Roy et al., 2009). The metric for $r_m^2$ is described in Equation 12 where $r^2$ and $r_0^2$ are the squared correlation coefficients with and without intercept, respectively.

$$CI \quad = \quad \frac{1}{Z} \sum_{y_i > y_j} h\left(p_i - p_j\right) \tag{10}$$

$$h\left(x\right) \quad = \quad \begin{cases} 1 & , x > 0 \\ 0.5 & , x = 0 \\ 0 & , x < 0 \end{cases} \tag{11}$$

$$r_m^2 \quad = \quad r^2 \cdot \left(1 - \sqrt{r^2 - r_0^2}\right) \tag{12}$$

### A.2    ABLATION STUDY

We present the results of the ablation study for proposed model in Table 4. We evaluated the performance of cross-attention module (C.A) , self-attention module (S.A) and latent encoder path of MetaDTA model on benchmark dataset with support set length 100. In this dataset, cross-attention alone has strong predictive performance. The use of latent path and self-attention achieves similar performance with cross-attention and model configuration would be better to optimized in differenct dataset and different condition.

Table 4: Comparison of the prediction performance between the different configuration of the models. C.A and S.A means cross-attention and self-attention module. The use of cross-attention shows best performance in three metrics and both self-attention and latent path improves the prediction accuracy in CI score. The model performance along to the configuration depends on dataset and hyper parameters.

|  | RMSE | MAE | $r_m^2$ | CI |
|---|---|---|---|---|
| ANPs (original) | 0.973±0.022 | 0.711±0.018 | 0.538±0.033 | 0.790±0.033 |
| **C.A+S.A+Latent** | 0.979±0.063 | 0.694±0.026 | 0.553±0.038 | **0.795±0.014** |
| **C.A+Latent** | 0.961±0.021 | 0.701±0.016 | 0.557±0.032 | 0.793±0.009 |
| **C.A+S.A** | 1.022±0.069 | 0.733±0.029 | 0.534±0.037 | 0.782±0.012 |
| **C.A** | **0.952±0.015** | **0.694±0.011** | **0.564±0.044** | 0.794±0.008 |

## A.3 $r_m^2$ INDEXES FOR FEW-SHOT SCENARIO LIGANDS EXPERIMENT

Table 5: $r_m^2$ indexes for few-shot scenario ligands split experiments

|  |  | DeepDTA | MolTrans | MetaDTA(I) | MetaDTA(T) |
|---|---|---|---|---|---|
| Random | 10 | 0.299±0.039 | 0.363±0.029 | 0.468±0.067 | 0.374±0.033 |
|  | 20 | 0.360±0.022 | 0.410±0.012 | 0.383±0.029 | 0.433±0.035 |
|  | 30 | 0.404±0.011 | 0.431±0.020 | 0.435±0.035 | 0.477±0.033 |
|  | 50 | 0.427±0.014 | 0.456±0.012 | 0.460±0.036 | **0.521±0.036** |
|  | 100 | 0.473±0.032 | 0.488±0.018 | **0.505±0.034** | **0.557±0.022** |
| Scaffold | 10 | 0.301±0.061 | 0.351±0.076 | 0.364±0.050 | 0.366±0.067 |
|  | 20 | 0.352±0.085 | 0.352±0.018 | 0.408±0.073 | 0.412±0.071 |
|  | 30 | 0.378±0.068 | 0.416±0.086 | 0.434±0.074 | 0.461±0.077 |
|  | 50 | 0.400±0.079 | 0.433±0.071 | 0.461±0.098 | 0.491±0.089 |
|  | 100 | 0.403±0.069 | 0.440±0.057 | **0.508±0.088** | **0.516±0.074** |