OXFORD

Systems biology

# DeepREAL: a deep learning powered multi-scale modeling framework for predicting out-of-distribution ligand-induced GPCR activity

Tian Cai [1], Kyra Alyssa Abbu[2], Yang Liu[2] and Lei Xie[1,2,3,*]

[1]Ph.D. Program in Computer Science, The Graduate Center, The City University of New York, New York, NY 10016, USA, [2]Department of Computer Science, Hunter College, The City University of New York, New York, NY 10065, USA and [3]Helen and Robert Appel Alzheimer's Disease Research Institute, Feil Family Brain & Mind Research Institute, Weill Cornell Medicine, Cornell University, New York, NY 10021, USA

*To whom correspondence should be addressed.
Associate Editor: Pier Luigi Martelli

## Abstract

**Motivation:** Drug discovery has witnessed intensive exploration of predictive modeling of drug–target physical interactions over two decades. However, a critical knowledge gap needs to be filled for correlating drug–target interactions with clinical outcomes: predicting genome-wide receptor activities or function selectivity, especially agonist versus antagonist, induced by novel chemicals. Two major obstacles compound the difficulty on this task: known data of receptor activity is far too scarce to train a robust model in light of genome-scale applications, and real-world applications need to deploy a model on data from various shifted distributions.

**Results:** To address these challenges, we have developed an end-to-end deep learning framework, DeepREAL, for multi-scale modeling of genome-wide ligand-induced receptor activities. DeepREAL utilizes self-supervised learning on tens of millions of protein sequences and pre-trained binary interaction classification to solve the data distribution shift and data scarcity problems. Extensive benchmark studies on G-protein coupled receptors (GPCRs), which simulate real-world scenarios, demonstrate that DeepREAL achieves state-of-the-art performances in out-of-distribution settings. DeepREAL can be extended to other gene families beyond GPCRs.

**Availability and implementation:** All data used are downloaded from Pfam (Mistry *et al.*, 2020), GLASS (Chan *et al.*, 2015) and IUPHAR/BPS and the data from reference (Sakamuru *et al.*, 2021). Readers are directed to their official website for original data. Code is available on GitHub https://github.com/XieResearchGroup/DeepREAL.

**Contact:** lei.xie@hunter.cuny.edu

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

## 1 Introduction

Over the past two decades, drug discovery has been dominated by target-based high-throughput compound screening. Unfortunately, this 'one-drug–one-gene' approach has been costly and had a low success rate due to our limited understanding of molecular and cellular mechanisms of drug actions (DiMasi *et al.*, 2016; Wong *et al.*, 2019). Drugs from the target-based screening often interact with unexpected off-targets, leading to serious side effects (Lin *et al.*, 2019; Lynch III *et al.*, 2017). Furthermore, a polypharmacology approach is often needed to achieve desired therapeutic efficacy and overcome drug resistance for complex diseases (Xie *et al.*, 2012). To predict drug phenotypic response at the organismal level, it is necessary to not only elucidate genome-scale drug–target

interactions (DTIs) but also reveal how DTIs collectively modulate a biological system.

The drug mode of action is a multi-scale process that starts with drug binding to its targets, principally proteins. Then the drug can act as an antagonist or an agonist to block or enhance downstream biological processes, respectively. Therefore, it is critically important to model the change of receptor activities or functional selectivity upon the drug binding for understanding how the drug modulates pathophysiological functions. The information on the receptor activity following the ligand binding will fill in a critical knowledge gap in correlating DTIs to clinical outcomes. Although a great deal of efforts have been devoted to predict genome-wide DTIs using deep learning (Cai *et al.*, 2021; Karimi *et al.*, 2019; Wan and Zeng, 2016), few large-scale experimental and computational studies have

been able to specify the ligand-induced receptor activity, i.e. the functional selectivity of the ligand as an antagonist or an agonist (Sakamuru et al., 2021).

In this research, we aim to predict not only whether any pairs of proteins and chemicals interact with each other or not but also the receptor activity upon the binding, especially, making reliable predictions for understudied 'dark' proteins that do not have any ligand annotations (Oprea, 2019) and novel chemicals whose structures are different from those in the training data. To our knowledge, only a recent work has used chemical features to train an independent machine learning model for each individual Opioid receptor for predicting their receptor activity (Sakamuru et al., 2021). Unfortunately, labeled data for the receptor activity are scarce. Only a limited number of receptors have sufficient function selectivity data to train a robust machine learning model. Thus, the one-protein–one-model approach cannot be extended to majority of proteins that have few or no labeled data (Sakamuru et al., 2021). An early work applied a neural network model to predict multiple interaction types for annotated proteins (Wang and Zeng, 2013). However, this work neither included antagonist/agonist as prediction tasks nor was tested for dark proteins. It is a challenging task to predict the function for dark proteins in general using machine learning. Conventional machine learning methods assume that the distribution of unseen data and training data is identically and independently distributed (IID). This assumption may not hold for the dark proteins that are dissimilar from those in the training data. In other words, many dark proteins are out-of-distribution (OOD) in terms of the training samples. Similarly, unseen novel chemicals whose structures are different from those in the training set are also OOD cases. To address the data scarcity and OOD challenges, we have developed an artificial intelligence (AI)-powered multi-scale modeling framework, DeepREAL, to simulate the multi-scale drug actions and predict the ligand-induced receptor activity for dark proteins and novel chemicals. We first apply self-supervised learning to train a protein sequence model for a universal protein sequence embedding on a genome scale. This allows us to detect subtle relationships between dark proteins and ligand-annotated proteins as demonstrated in other studies (Cai et al., 2021; Rao et al., 2019; Rives et al., 2021). We then train a binary classification deep learning model to predict whether a chemical binds to a protein and extract a latent presentation of DTIs. Because there is a large amount of binary interaction data, it is possible to train a robust deep learning model. Finally, we integrate chemical embedding model, sequence embedding model and DTI latent representation model to train an end-to-end deep learning model for predicting the ligand-induced receptor activity using limited data. In the rigorous benchmark studies on GPCRs, which simulate real-world applications, DeepREAL significantly improves the generalization ability in the OOD setting compared with the state-of-the-art methods (Cai et al., 2021; Sakamuru et al., 2021; Wang and Zeng, 2013).

The contributions of DeepREAL can be summarized in twofolds:

1. DeepREAL aims to address an unsolved but important challenging problem for drug discovery: robustly predicting genome-wide ligand-induced receptor activities or function selectivity under various data distribution shifts.

2. DeepREAL is based on a new multi-stage deep transfer learning architecture that combines binary DTI pre-training and embedding with a three-way receptor activity fine-tuning to address OOD challenges using sparse receptor activity data.

# 2 Materials and methods

## 2.1 Data
Four datasets were used in this study. Pfam, v33.1 (Mistry et al., 2020) was used to pre-train protein descriptors. GPCR–ligand binding binary data were obtained from GLASS, v2019.2 (Chan et al., 2015). Agonist/antagonist data were downloaded from the International Union of Basic and Clinical Pharmacology/British

Pharmacological Society (IUPHAR/BPS) Guide to Pharmacology, v2020.5. Additional Opioid receptor activity data were from the study by Sakamuru et al. (2021). The protein descriptor pre-training exactly followed DISAE (Cai et al., 2021). In brief, DISAE built up a distilled triplet sequence dictionary for the whole Pfam proteins based on multiple sequence alignments (MSA). Every input protein was mapped to its distilled triplets representation according to the protein dictionary, as illustrated in Figure 1. Chemical-protein pairs with the receptor activity annotation was treated as positive in the binary DTI setting and combined with GLASS for the binary classification pre-training. In terms of pre-training, only Stage 1 protein descriptor pre-training was self-supervised as described in the study by Cai et al. (2021). Stage 2 uses CLASS data for supervised pre-training. IUPHAR/BPS combined with Sakamuru et al. (2021) Opioid data were used in the final Stage 3 three-way classification. Detailed data statistics is found in Table 1.

## 2.2 State-of-the-art baselines
We compared DeepREAL with Random Forest (RF) models for three Opiod receptors (Sakamuru et al., 2021) that used PubChem fingerprints (ftp://ftp.ncbi.nlm.nih.-gov/pubchem/specifications/pubchem_fingerprints.txt) (Bolton et al., 2008) as features. To our knowledge, the RF/protein baseline was the first and only work for the ligand-induced receptor activity prediction. Keeping other hyper-parameters the same as those in the study by Sakamuru et al. (2021), the RF depth was tuned to find the best performance model for each Opioid receptor. An example performance curve is shown in Supplementary Figure S6. For each experiment, one Random Forest is trained for each Opioid receptor. An average RF test performance was calculated by weighting the sample size of each Opioid receptor. When evaluating the variance of model performance, different random seeds were used.

Another baseline model is similar to restricted Boltzmann machines from an earlier work (Wang and Zeng, 2013) which is designed to predict DTI types. We built a multi-task deep learning model that consisted of two layer vanilla MLP (Goodfellow et al., 2016) for every single target, i.e. one Opioid receptor, with the same number of hidden units and the same definitions of visible units (Goodfellow et al., 2016) by optimizing the average cross entropy loss of the model for each target. The constructed multi-task MLP for a multidimensional DTI network was associated with the same parameters. The input feature was also PubChem fingerprints (Bolton et al., 2008).

## 2.3 DeepREAL framework
### 2.3.1 Architecture
DeepREAL has a novel three-stage framework. There are four major modules in DeepREAL model: protein sequence embedding, chemical structure descriptor, binary interaction learner and multi-class receptor activity classifier as shown in Supplementary Figure S1. Under this framework, the state-of-the-art model DISAE (Cai et al., 2021) was employed as the backbone for learning DTI embeddings, which includes ALBERT- (Lan et al., 2019) based protein descriptor, and attentive pooling- (Santos et al., 2016) based binary interaction learner. Different from DISAE that uses neuro-fingerprint for the chemical representation, the chemical descriptor in DeepREAL is state-of-the-art unpretrained graph neural network GIN (Xu et al., 2018).

The unique component of multi-class receptor activity classifier includes two sub-modules. A three-way interaction learner uses the same architecture as the binary interaction learner. After concatenating all related embeddings, the concatenated tensor goes through a ResNET (He et al., 2015) layer and MLP (Hastie et al., 2019) transformation to generate the final logit vector used in cross entropy loss calculation (Hu et al., 2019).

### 2.3.2 Information flow of DeepREAL
The knowledge transfer across stages is realized by sharing weights on the first three modules in DeepREAL architecture, i.e. protein descriptor, chemical descriptor and binary interaction learner. Protein
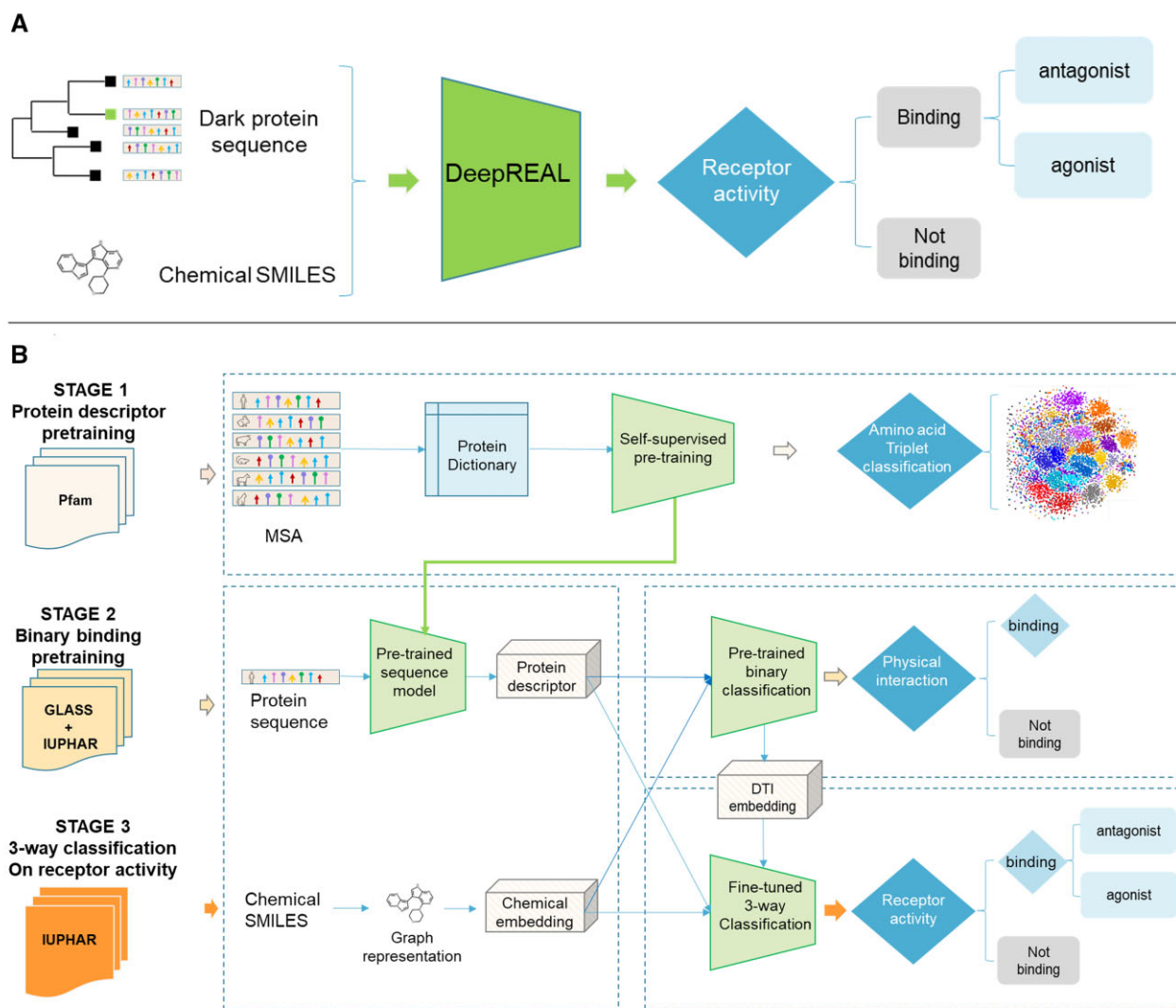
**Fig. 1.** Illustration of DeepREAL. (**A**) Given a chemical and a protein sequence as inputs, DeepREAL will predict not only if the chemical is the ligand of the protein but also the ligand-induced receptor activity. (**B**) DeepREAL is an end-to-end deep learning model trained using three stages of pre-training and fine-tuning. See text for details

**Table 1.** Training, validation and testing data used in this study

| | Unique protein | Unique chemical | Agonist | Antagonist | Not-binding | Binding | (OOD TEST imbalance ratio control) not-binding:agonist:antagonist |
|---|---|---|---|---|---|---|---|
| IUPHAR | 450 | 13 126 | 14 412 | 14 488 | 144 500 | 28 900 | 5:1:1 |
| Opioid receptors related | 3 | 2483 | 2920 | 2996 | 29 580 | 5916 | |
| GLASS | 689 | 181 114 | — | — | 70 089 | 270 545 | — |

descriptor first goes through Stage 1 sequence pre-training in a self-supervised fashion. The pre-trained sequence embeddings are then transferred to Stage 2 binary pre-training. Together with initialized chemical descriptor, a binary interaction learner learns to predict whether or not a protein and a chemical would interact in a supervised learning manner. The learned weights of these three modules are all transferred to Stage 3. In Stage 3, the three modules are first duplicated: one copy has frozen weights whereas the other copy updates its weights for $n$ epochs with multi-class-learner on DeepREAL receptor activity information in a supervised learning manner, where $n$ is a hyper-parameter as shown in Supplementary Table S1. In our experiments, we find that a small $n$ such as 50 would help to improve model generalization performance when the training data size was smaller. This phenomenon is due to the fact

that a complete model with a large number of trainable parameters is capable of memorizing a small training set, resulting in over-fitting and poor generalization. More frozen weights would limit the over-fitting and put more pressures on the multi-class-learner to learn a robust representation.

As illustrated in Supplementary Figure S1, the protein embedding vector, chemical embedding vector and binary interaction embedding vector that is the output of binary pre-trained module are fine-tuned via a three-way receptor activity learner that also learns a three-way receptor activity embedding. Seven embedding vectors, which include the protein embedding, the chemical embedding and the binary interaction embedding, both fine-tuned and frozen after the pre-training, along with the three-way receptor activity embedding, are concatenated and fed into a ResNET (He *et al.*,

2015) followed by a MLP (Hastie *et al.*, 2019) to make the final three-way classification.

The three-stage model is designed to train sequentially and separately. Only optimized weights are transferred. The Stage 1 optimization procedure has been described in Cai *et al.* (2021). Stages 2 and 3 optimizations are both driven by a cross-entropy loss in a stochastic manner using Adam (Kingma and Ba, 2014).

### 2.3.3 Pre-training implementation and module frozen strategy

A key element of success in multi-stage pre-training is to transfer knowledge. A major challenge in the three-stage pipeline is to prevent the previously learned knowledge from being lost during the weight update in the subsequent stage. DISAE has reported the benefits of a frozen mechanism. This strategy is adopted in DeepREAL Stages 2 and 3 as well. In Stage 2, following the experience of DISAE, part of the transformer (Vaswani *et al.*, 2017) layers is frozen. In Stage 3, the binary pre-trained modules are duplicated to have one copy always frozen and the other copy fine-tuned for only *n* epochs. Without tuning, *n* is empirically set to 50 in the Opioid receptor focused experiments, while on the complete DeepREAL receptor dataset involving 450 proteins, *n* is set as infinity until the model converges.

### 2.3.4 Data splitting for training and testing

In terms of data splitting, IID setting splits the data randomly as conventional cross-validations, except for the Opioid-context experiments where all three Opioid proteins are ensured to appear in both training and testing datasets. The OOD data split is carried out using a spectral clustering algorithm (Luxburg, 2007) based on pair-wise chemical similarity measured by Tanimoto coefficient and sequence similarity measured by sequence identity. The similarity distributions could be found in Supplementary Figure S3. In our experiments, the Stage 2 binary training is always carried out with the same data. The pairwise scores in Supplementary Figure S3 are measured for each pair of a chemical from training and a chemical from test as well as a protein from training and a protein from test. For more than 95% chemicals in the test set, less than 2% chemicals in the training set have Tanimoto coefficient larger than 0.6.

Because we studied several OOD and IID scenarios, in each scenario the number of proteins in the testing set is different.

1. IUPHAR OOD-protein-distribution-shift. The split is made upon protein similarity. 49 out of 450 proteins in the test set. Proteins in the training and testing set have no overlaps. As shown in Supplementary Figure S3, majority of proteins in the testing test are not similar to those in the training set with the sequence identity less than 10%.
2. IUPHAR IID setting. Data are randomly split. 298 out of 450 proteins are in the test set. 246 out of the 298 proteins in the test set are also in the training set, but there are no overlapped protein-chemical pairs between training and testing set.
3. Opioid OOD-chemical-distribution-shift. The split is made upon pair-wise chemical similarity between chemicals in the training set and chemicals in the test set as shown in Supplementary Figure S3, where only around 0.6% of chemicals in the testing set are similar to those in the training set with Tanimoto coefficient larger than 0.6. All three Opioid proteins are in the test set.
4. Opioid IID. Data are randomly split. All three Opioid proteins in the test set.

### 2.4 Ensemble model for novel receptor activity prediction

We build an ensemble of three DeepREAL models independently trained with different random seeds. The ensemble model is used to perform predictions on novel relations. Top predictions are selected by filtering out predictions agreed by all the three models in the ensemble.

## 3 Results and discussion

### 3.1 Overview of methods

Given a chemical structure and the sequence of a receptor protein, DeepREAL will predict whether the chemical is an agonist or an antagonist if it binds to the receptor, or not bind to it at all (Fig. 1A). As an end-to-end learning framework, the DeepReal is a three-way classifier: not-binding/agonist/antagonist. Intuitively, DeepREAL leverages large datasets to hierarchically inform predictions on the receptor activity whose labeled data are scarce along a three-stage pre-training-fine-tuning pipeline as illustrated in Figure 1B. In Stage 1, protein descriptor was pre-trained using Pfam (Mistry *et al.*, 2020) data. In Stage 2, a binary DTI classifier was then pre-trained using GLASS (Chan *et al.*, 2015) and IUPHAR binary data (Armstrong *et al.*, 2019). Finally, in Stage 3, three-way classification on the receptor activity was fine-tuned using the outputs of Stages 1 and 2 as inputs with IUPHAR antagonist/agonist data (Armstrong *et al.*, 2019).

The Stage 1 self-supervised sequence embedding was based on DISAE (Cai *et al.*, 2021). DISAE distilled the protein sequence into an ordered list of triplets by excluding evolutionarily unimportant positions from a multiple sequence alignment. Then long range residue interactions were learned via the self-attention in a transformer module. A self-supervised masked language modeling approach was used to train sequence embeddings. By pre-training protein sequences on whole Pfam in Stage 1, DeepREAL equipped itself with genome-scale protein representations that captured novel relationships between proteins beyond sequence homology as demonstrated by several studies (Cai *et al.*, 2021; Rao *et al.*, 2019; Rives *et al.*, 2021). The second stage was a binary DTI pre-training which predicts binding/not-binding. By pre-training on a large scale of binary DTI data in Stage 2, DeepREAL builds knowledge of chemical–protein interactions which is the initial step in the ligand binding event and generates DTI embeddings. Finally, in Stage 3, information learned from sequence embeddings and DTI embeddings were transferred into predicting receptor activities using a small amount of data. This hierarchy design maintained knowledge learned from heterogeneous resources and enhanced model robustness when facing shifted data distribution during the deployment. The model was trained in an end-to-end fashion without feature engineering. The embedding from the pre-training is not fixed, but can be fine-tuned by the subsequent training stage. More details of DeepREAL design and implementation could be found in Methods section. Detailed model architecture is in Supplementary Figure S1 and Supplementary Table S1.

It notes that DeepREAL is an extension of DISAE (Cai *et al.*, 2021) but with several major new contributions. DISAE is a general-purpose protein language model and has been applied to predict chemical-protein interactions (Cai *et al.*, 2021), while DeepREAL is a framework designed to tackle a different task that has not been explored: predicting out-of-distribution ligand-induced receptor activity (agonist versus antagonist). This task cannot be solved by the original DISAE architecture (Cai *et al.*, 2021). In the DeepREAL framework, DISAE was mainly used as Stage 1 pre-training. In addition, DeepREAL included two more components beyond DISAE: Stage 2 DTI interaction embedding and Stage 3 three-way classification.

As shown in Table 1, 689 unique human GPCRs was used for the Stage 2 DTI pre-training. These GPCRs consist of six Pfam families: PF00001, PF00002, PF00003, PF052496, PF01534 and PF02101. Among them, 450 GPCRs have known labeled receptor activity data, and was used for the Stage 3 fine-tuning. Among 180 000 ligands of GPCRs, only 3303 ligands have known agonist/antagonist activities. Moreover, majority GPCRs have less than 100 ligands that are labeled with receptor activities, as shown in Supplementary Figure S2. Only 3 Opioid receptors (P35372—Mu Opioid receptor, P41145—Kappa Opioid receptor, P41143—Delta Opioid receptor) have more than 300 chemicals with known receptor activities. Thus, the labeled receptor activity data are not large enough to train a robust machine learning model on the basis of a single protein for most GPCRs.

To evaluate DeepREAL performance in light of real-world applications for dark proteins and novel chemicals, both data preprocessing and controlled experiment are designed to simulate various scenarios of data distribution shifts and to answer the following questions.

Q1: Is the pre-training helpful to improve the performance of receptor activity prediction using a small amount of data?

Q2: When DeepREAL is applied to unseen *dark proteins* that have low sequence similarity to those in the training data, what is the OOD generalization performance of DeepREAL?

Q3: When DeepREAL is used to predict unseen novel *chemicals* that are significantly different from those in the training data, what is the OOD generalization performance of DeepREAL?

Q4: When the test set label (agonist/antagonist/not-binding) distribution is close to reality and imbalanced compared to the training data, what is the generalization performance of DeepREAL?

Q5: How does DeepREAL perform compared to the state-of-the-art baseline models in both OOD and IID settings for predicting Opioid receptor activity?

We used three metrics, AUC–ROC, MCC and Cohen's kappa to evaluate the performance of various models under different settings.

## 3.2 Pre-training enables DeepREAL to generalize genome-scale receptor activity predictions using a relatively small dataset

Pre-training has been demonstrated to be effective in several recent works (Karimi *et al.*, 2019; Wan and Zeng, 2016) for predicting protein–ligand interactions. DeepREAL used Stages 1 and 2 as pre-trainings for learning knowledge in the protein sequence space and binary interaction space, respectively. To answer Q1, the same model architecture is trained on the same IID and OOD settings using four procedures: (i) from total scratch without any pre-training, i.e. Stage 3 only, (ii) going through Stage 1 whole Pfam pre-training but not the Stage 2 binary DTI classification pre-training, which is equivalent to the DISAE model (Cai *et al.*, 2021), (iii) going through only Stage 2 but not Stage 1 and (iv) complete three-stage pre-training/fine-tuning as DeepREAL. As shown in Figures 2 and 3 on the evaluation cross three classes (no-binding, agonist, antagonist), the model without any pre-training (i.e. only Stage 3) has the worst performance. Stages 1 or 2 both boosts performance and the complete three-stage pipeline yields the best performance. From the by-class evaluation for antagonist or agonist as shown in Figures 4

and 5, the precision and recall of DeepREAL is mostly higher than other variants in IID, protein OOD and chemical OOD settings. Furthermore, the training curves of DeepREAL in Figures 4 and 5 converges faster than other variants in most cases. The advantage of pre-training is particularly apparent in chemical distribution shift OOD in the cross-class and the by-class evaluation as shown in Figures 3–5. The chemical OOD is a more challenging OOD setting than other settings, where both chemical structure distribution and label ratio balance shift (more details in the following section). DISAE and the only-Stage 3 model have lower Cohen's kappa, ROC–AUC, MCC than DeepREAL and the Stage2+Stage3 model. The latter two models have relative close performance, suggesting that DTI pre-training plays a more important role than the sequence pre-training in the current training procedure. It may be because the whole-Pfam information learned at Stage 1 is more difficult to transfer to Stage 3, as supported by the observation shown in Figures 4 and 5. It will be interesting to use other advanced training procedures such as prompting (Gao *et al.*, 2020) or design different architectures [e.g. using skip connections (Dosovitskiy *et al.*, 2020; He *et al.*, 2015), etc.].

## 3.3 DeepREAL is robust in various shifted distribution scenarios

Q2, Q3, Q4 are three typical shifted distribution scenarios in real-world applications, i.e. the OOD generalization challenge. DeepREAL proves robust in each of the settings. It makes DeepREAL applicable to explore dark chemical genomics space.

Q2 focuses on the distribution shift coming from proteins. It is a dominant challenge when applying DeepREAL to a genome-scale given majority of proteins are dark without any receptor activity data. In this setting, 450 proteins and their associated interaction data are split into an OOD train/test sets such that the sequence similarities between proteins in the testing set and those in the training set are less than 10% (Supplementary Fig. S3). As shown in Figure 2, although the performance drops compared with the easier IID setting, the ROC-AUC score is still at 0.766, while existing state-of-the-art RF-based one-protein-one-model (RF/protein) approach (Sakamuru *et al.*, 2021) and multi-task neural network model (Wang and Zeng, 2013) are totally unable to make reliable predictions in the protein OOD setting.

In a similar fashion, by splitting the data based on the chemical similarity measured by Tanimoto coefficient, DeepREAL is evaluated in the setting of chemical distribution shift to answer Q3. Only Opioid receptors are used in the evaluation because only Opioid receptors have sufficient large numbers of labeled chemicals to generate OOD training/testing dataset, as shown in Supplementary Figure S2. In addition, we would like to reduce the
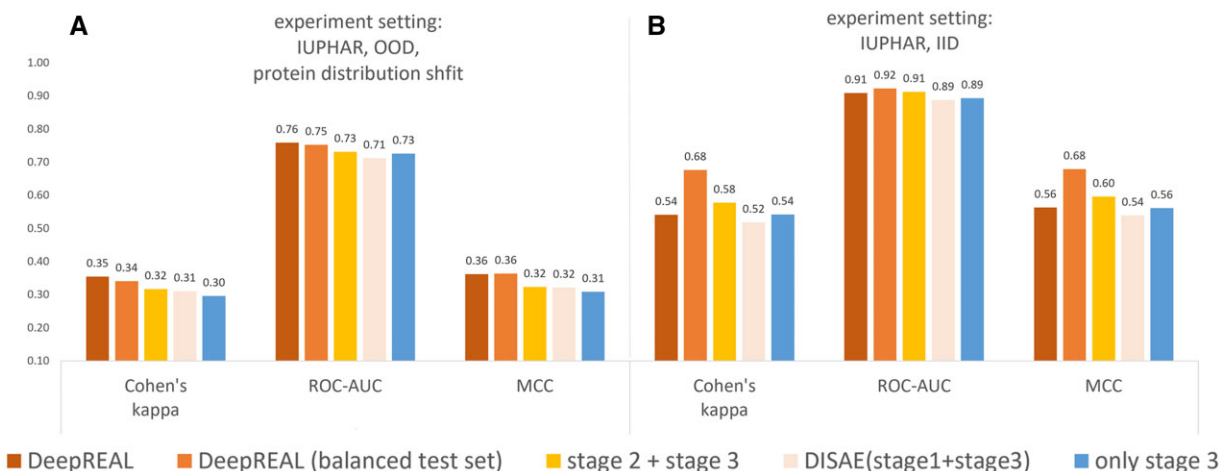


**Fig. 2.** Performance comparison of DeepREAL with its variants in (**A**) protein OOD and (**B**) protein IID settings. The performance is evaluated by multiple gene families in the IUPHAR database
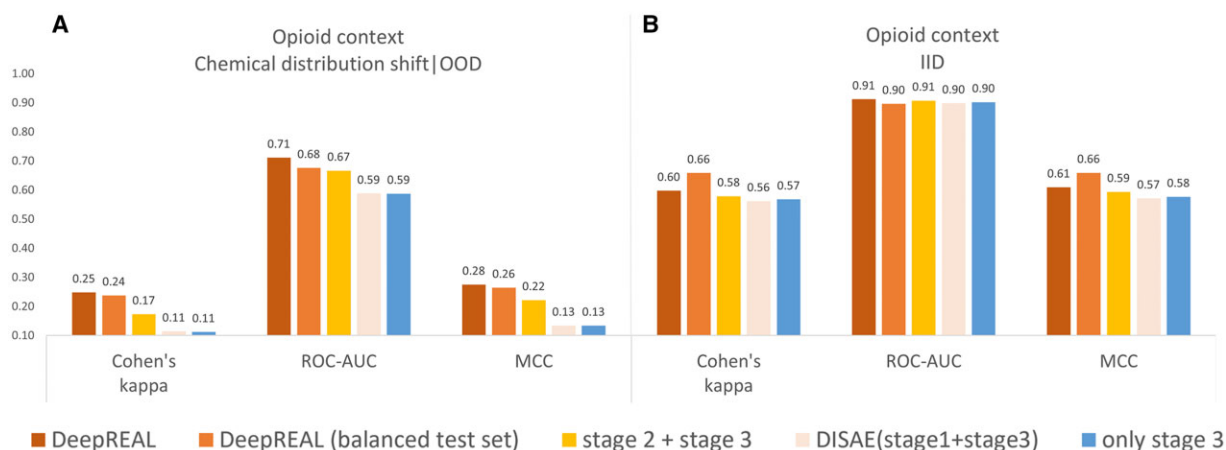
**Fig. 3.** Performance comparison of DeepREAL with its variants in (**A**) chemical OOD and (**B**) chemical IID settings. The performance is evaluated only by Opioid receptors

impact of the OOD from receptors. The advantage of DeepREAL is apparent over other configurations including DISAE, as shown in Figures 3 and 6A in the chemical OOD setting. It notes that only around 0.6% of chemicals in the testing set are similar to those in the training set with Tanimoto coefficient larger than 0.6, as shown in Supplementary Figure S3.

Although it is expected that a machine learning model performs the best when positive and negative data are balanced, the unseen binding/not-binding cases are imbalanced in reality, which has an estimated ratio of 1:5 (Lim *et al.*, 2016). The ratio of 1:5 is based on the estimated value in the published work (Lim *et al.*, 2016) when considering the chemical genomics space (millions of chemicals paired with thousands of proteins) as a whole, which is the same scenario as this manuscript. The ratio is lower than the observations from many compound screenings because a large number of potential off-targets are not taken into account in the existing target-based screening. It should be noted that the purpose here is to compare the use of imbalanced test data with the use of balanced ones, which is a common practice in most existing studies. Hence, to answer Q4 about label distribution shift, for all experiments the number of not-binding samples in the test set is about five times as large as that of the agonist/antagonist data while training data are balanced for each class. For a comparison, a balanced test set is also evaluated. In general, DeepREAL evaluated by the balanced test set in the IID setting, which represents conventional cross-validations, outperforms that evaluated by the imbalanced data. However, in both protein and chemical OOD settings that simulates a real application, DeepREAL evaluated by the imbalanced data performs the best, as shown in Figures 2 and 3. These observations suggest that the cross-validation in an IID setting is often over-optimistic and DeepREAL is more robust in a realistic application. To see if different imbalanced ratio will affect the result, we performed additional experiments with a ratio of 10:1. As shown in Supplementary Figure S4, the change of ratio will not change the results significantly.

### 3.4 DeepREAL significantly outperforms state-of-the-art models

To compare DeepREAL with the leading machine learning model (RF/protein) (Sakamuru *et al.*, 2021) that can only predict Opioid receptor activities as well as an earlier multi-task neural network model (Wang and Zeng, 2013). Only Opioid receptors are used in the comparison due to two reasons. First, the baseline models can be only trained using chemicals as input. Second, only Opioid receptors have sufficient large numbers of labeled chemicals for training the baseline model (Supplementary Fig. S2). If we include other proteins, the baseline model may have significant disadvantages.

Opioid receptor dataset is split in two different ways for IID and OOD experiments as described in the previous section. In both IID and OOD settings, DeepREAL significantly outperforms the

baselines in terms of precision and recall, as shown in Figure 6. Furthermore, the performance drop of DeepREAL from the IID setting to the OOD setting is less significant than that of the baseline. To prove the statistical significance of DeepREAL performance against the RF/protein baseline, the same training is repeated for five times under Opioid context with different random seeds. As shown in Supplementary Figure S5, the *P*-value of the hypothesis that the two models have the same average ROC-AUC is close to 0.0.

### 3.5 Application of DeepREAL to cocaine interacting proteins

We performed a screening for G-protein coupled receptors (GPCRs) that interact with cocaine and its analogs using trained DeepREAL model. Cocaine target, cocaine analogs and top ranked predictions could be found in Supplementary Tables S2–S4. 14 cocaine interacting GPCRs were collected from Fant *et al.* (2019). We collected 18 cocaine analogs and made predictions on them paired with the 14 targets. Among 14 proteins that we tested, cocaine or cocaine analogue was predicted as an agonist for glutamate metabotropic receptor 2 (GRM2) and 5-hydroxytryptamine receptor subtype 6 (5-HT6). As supporting evidences, Yang *et al.* (2017) have showed that GRM2 deletion decreases sensitivity to cocaine reward in rats. 5-HT6 antagonist blocks cocaine-induced DA release and cocaine self-administration, suggesting cocaine probably is a agonist for 5-HT6 (Valentini *et al.*, 2013). Our model also predicted cocaine's antagonist activity against 5-HT2C, delta-Opioid receptor and Cannabinoid Receptor 2 (CNR2). Injection of the 5-HT2C receptor agonist reduces cocaine self-administration in rats, suggesting cocaine is a potential antagonist for 5-HT2C receptor (Fletcher *et al.*, 2004). Similarly, dual kappa-delta Opioid receptor agonist blocks cocaine reward behavior intimating cocaine's antagonist role for delta Opioid receptor (Váradi *et al.*, 2015). Research shows CNR2 agonist dose-dependently inhibits cocaine self-administration, thus indicating cocaine negatively regulates CNR2's activity (Xi *et al.*, 2011). Overall, our predictions are largely consistent with existing experimental evidences.

## 4 Conclusion

This article proposed a deep learning framework DeepREAL that expands the traditional DTI task to predicting ligand-induced receptor activities of dark proteins and novel chemicals under various OOD settings. DeepREAL has several unique features. First, unlike the existing method that requires training one model for one protein and applying the trained model on the same protein, DeepREAL needs only to train only one model to make predictions on any proteins with improved accuracy. Second, DeepREAL has improved generalization power when facing all major types of data distribution shifts during deployment, making it robust in real-world
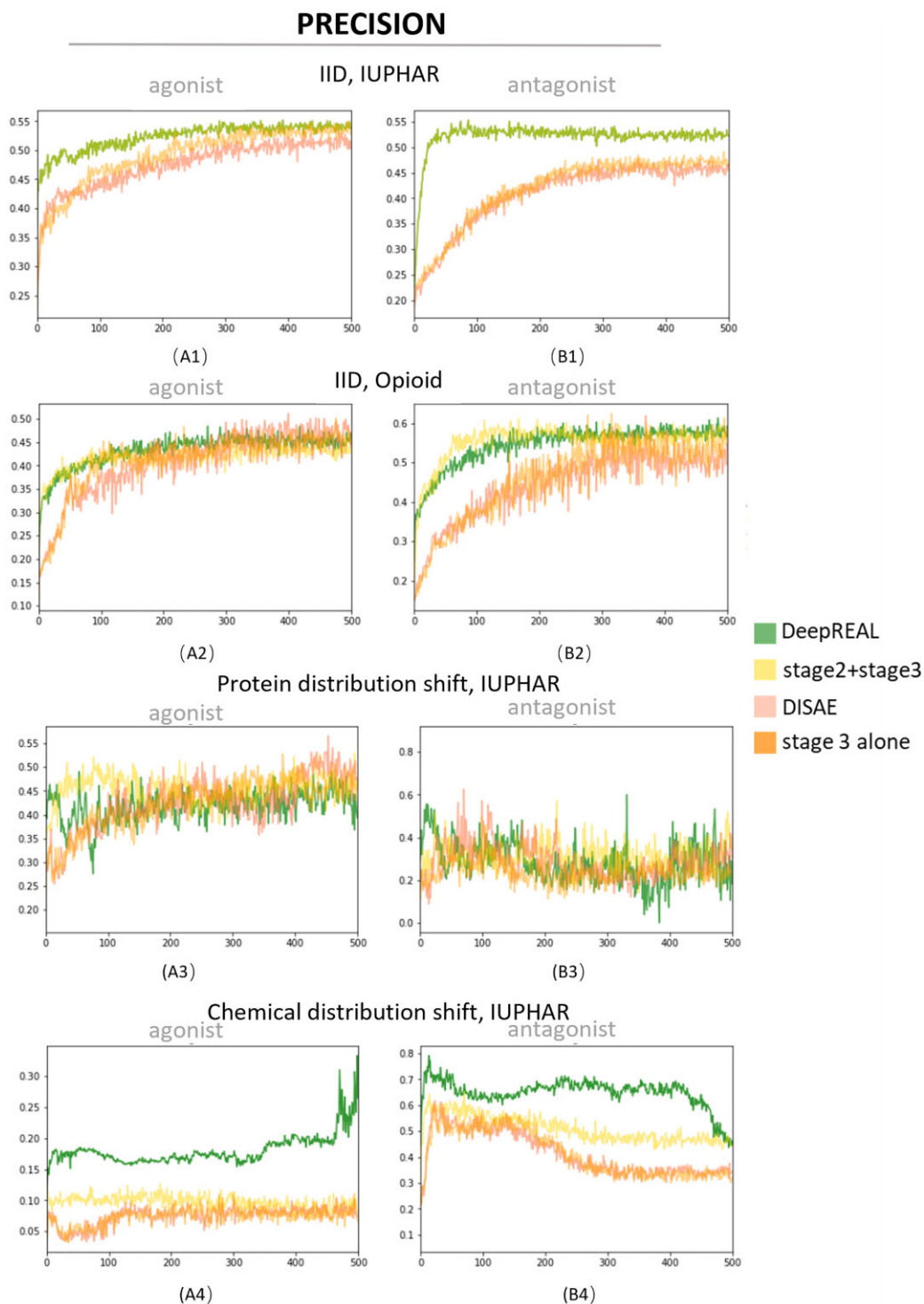
**Fig. 4.** Training curves of DeepREAL and its variants when measured by the precision for predicting agonists or antagonists. The *x*-axis is number of training epochs. The *y*-axis is precision

applications. Finally, by utilizing large unlabeled sequence data and rich binary bioassay data, DeepREAL models receptor activities on a multi-scale to alleviate data scarcity problem. Together, DeepREAL significantly outperforms existing algorithms for predicting ligand-induced receptor activities. The novelty of DeepREAL lies in the prediction of receptor activities for dark proteins (Stage 3)

using pre-trained protein sequence embedding (Stage 1) and binary DTI embedding (Stage 2). The incorporation of Stages 1 and 2 pre-training is motivated to achieve the OOD generalization in Stage 3. Additionally, the excellent performance of Stage 3 is not solely relying on the pre-training of Stages 1 and 2. The end-to-end model architecture as illustrated in Supplementary Figure S1 is designed to

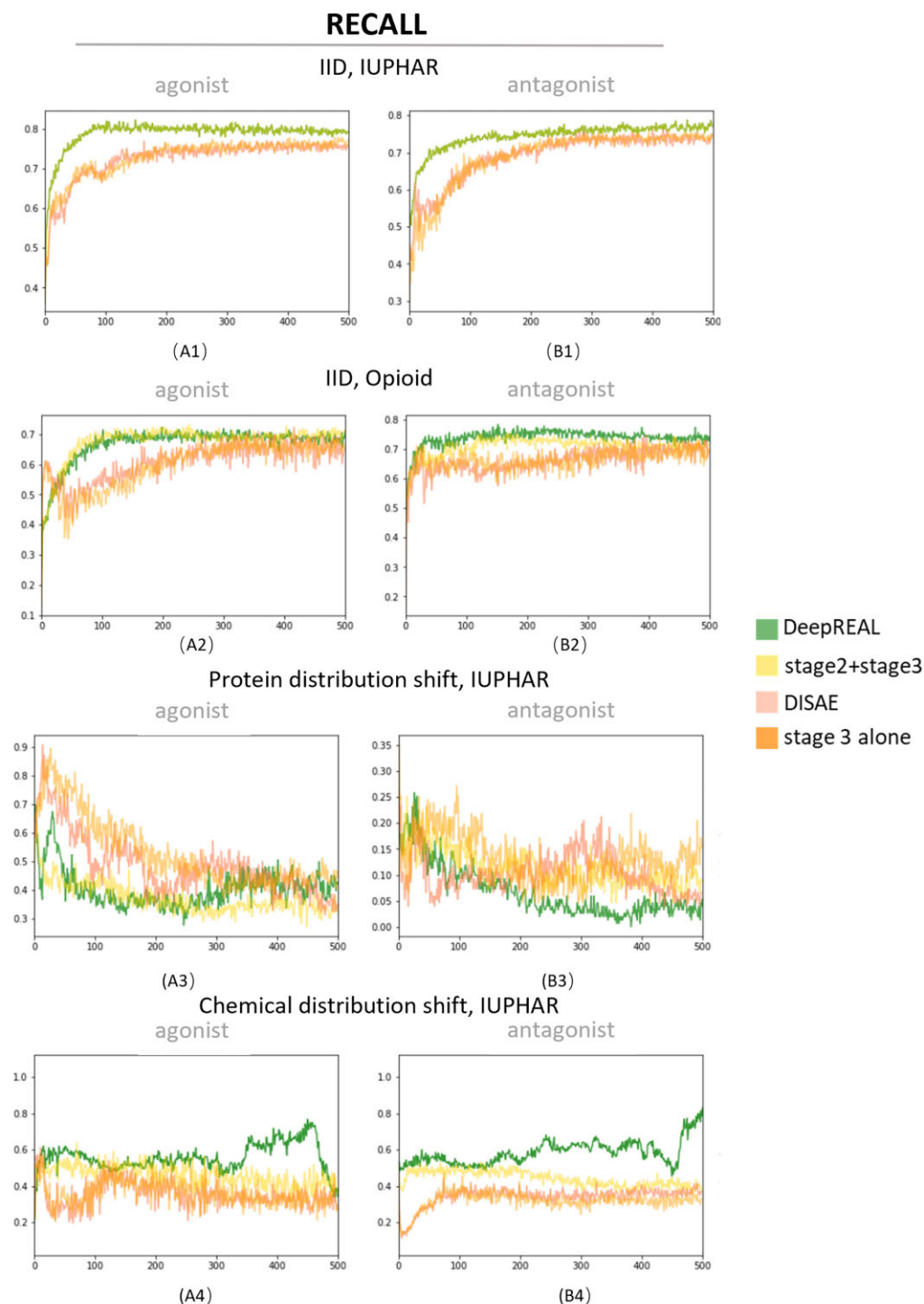# RECALL



**Fig. 5.** Training curves of DeepREAL and its variants when measured by the recall for predicting agonists or antagonists. The *x*-axis is number of training epochs. The *y*-axis is recall

ensure that knowledge transferred over stages will not be lost and get well utilized. Although DeepREAL was only tested using GPCRs, especially, Opioid receptors due to limited labeled data, it can be extended to other gene families when the ligand-induced receptor activity data are available.

The performance of DeepREAL can be further improved along several directions. For example, unsupervised pre-training of chemical space could improve DeepREAL's ability to detect novel chemicals (Hu *et al.*, 2019; Liu *et al.*, 2021). The sequence embedding method DISAE used in this study still has room for

improvement. Incorporating structure information into the protein sequence embedding could help the downstream prediction tasks for ligand binding and receptor activity. In addition, it may not perform well for small families similar to AlphaFold2 (Jumper *et al.*, 2021). It remains an open question to reliably prediction the structure and function of dark proteins from a small family. It will also interesting to test other state-of-the-art sequence embedding methods such as ESM (Rives *et al.*, 2021), ProtBERT (Elnaggar *et al.*, 2021) and TAPE (Rao *et al.*, 2019). We only predict two classes of receptor activity: agonist versus
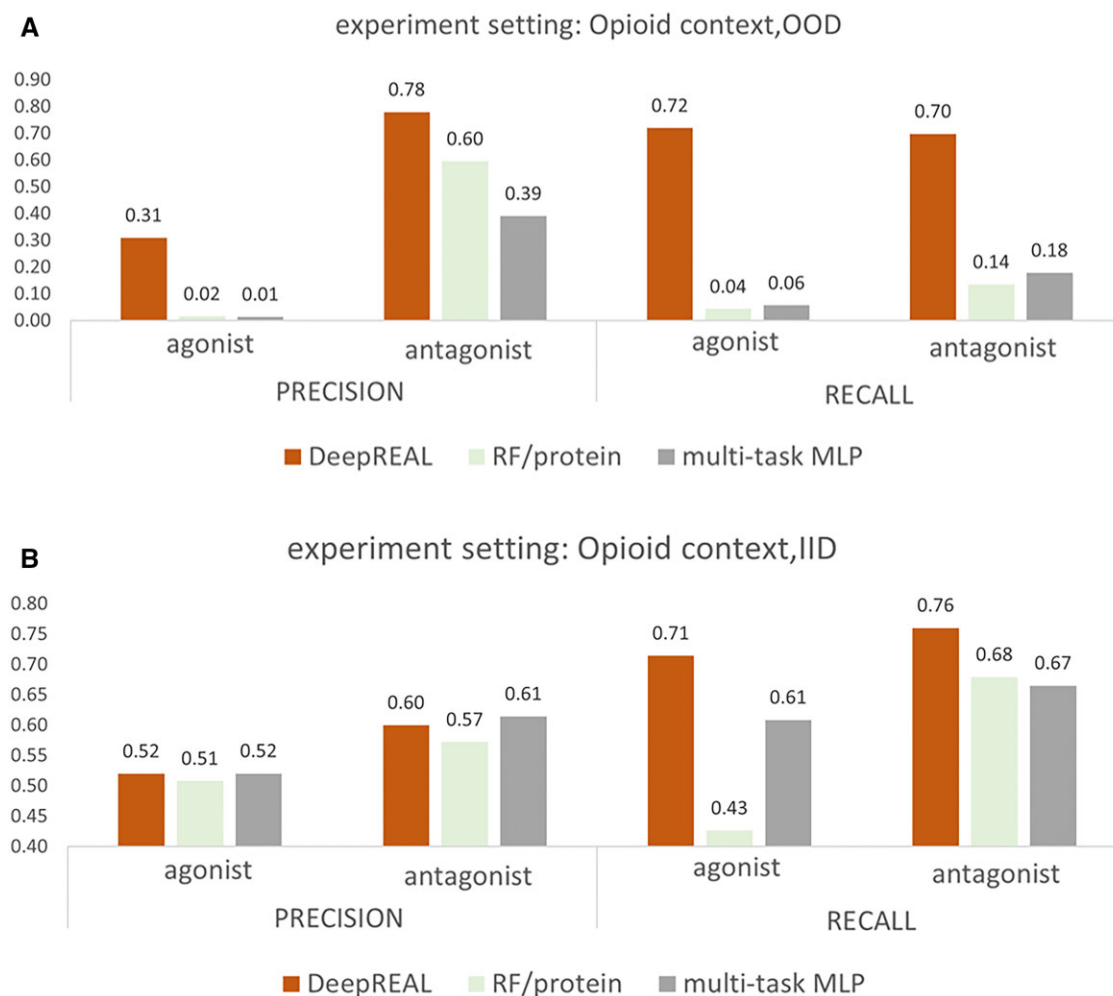
**Fig. 6.** Performance comparison of DeepREAL with the Random Forest model (RF/protein) and multi-task MLP in (**A**) chemical OOD and, (**B**) chemical IID settings. The performance is evaluated by precision and recall for agonist or antagonist predictions

antagonist. In fact, the receptor activity is more complex than two mutually exclusive classes. There are other subtle activity classes such as partial agonist. A multi-class model could be a more suitable choice and subject to future studies. In practice, detecting if an unseen case is OOD is an important but challenging problem. Few methods have been developed for protein or chemical data for the OOD detection. It is another direction for future works. Furthermore, there are more scenarios of distribution shift worth study such as compounding protein and chemical distribution shifts with various label distribution shifts for stress testing. In addition to the imbalanced ratio of binding/non-binding cases, the ratio of agonist/antagonist could vary a lot for different proteins and there is no generally known trend which one is more prevalent. This question remains unanswered, and will be addressed in the future.

## Author contributions

L.X. conceived and planned the experiments. T.C. developed and implemented the algorithm. T.C., K.A.A. and Y.L. carried out the experiments. L.X. and T.C. contributed to the interpretation of the results. T.C., K.A.A., Y.L. and L.X. wrote the manuscript. All authors provided critical feedback and helped shape the research, analysis and manuscript.[AQ: Please note that Author contributions section has been set per journal style.]

## Data availability

All data used are downloaded from Pfam (Mistry *et al.*, 2020), GLASS (Chan *et al.*, 2015) and IUPHAR/BPS and the state-of-the-art paper (Sakamuru *et al.*, 2021). Readers are directed to their official website for original data. Code is available on github https://github.com/XieResearchGroup/DeepREAL.

## References

Armstrong,J. *et al.*; NC-IUPHAR. (2019) The IUPHAR/BPS Guide to PHARMACOLOGY in 2020: extending immunopharmacology content and introducing the IUPHAR/MMV Guide to MALARIA PHARMACOLOGY. *Nucleic Acids Res.*, **48**, D1006–D1021.

Bolton,E.E. *et al.* (2008) PubChem: integrated platform of small molecules and biological activities. In: *Annual Reports in Computational Chemistry*, **Vol. 4**. Elsevier, pp. 217–241.

Cai,T. *et al.* (2021) MSA-regularized protein sequence transformer toward predicting genome-wide chemical-protein interactions: application to GPCRome deorphanization. *J. Chem. Inf. Model.*, **61**, 1570–1582.

Chan,W.K.B. *et al.* (2015) GLASS: a comprehensive database for experimentally validated GPCR-ligand associations. *Bioinformatics*, **31**, 3035–3042.

DiMasi,J.A. *et al.* (2016) Innovation in the pharmaceutical industry: new estimates of R&D costs. *J. Health Econ.*, **47**, 20–33.

Dosovitskiy,A. *et al.* (2020) An image is worth 16×16 words: transformers for image recognition at scale. *arXiv preprint* arXiv:2010.11929.

dos Santos,C. *et al.* (2016) Attentive pooling networks. *arXiv preprint arXiv: 1602.03609.*

Elnaggar,A. *et al.* (2021) ProtTrans: towards cracking the language of life's code through self-supervised learning. *bioRxiv*, pages 2020–07.

Fant,A.D. *et al.* (2019) Toward reducing hERG affinities for DAT inhibitors with a combined machine learning and molecular modeling approach. *Biophys. J.*, **116**, 562a.

Fletcher,P.J. *et al.* (2004) Injection of the 5-HT2C receptor agonist Ro60-0175 into the ventral tegmental area reduces cocaine-induced locomotor activity and cocaine self-administration. *Neuropsychopharmacology*, **29**, 308–318.

Gao,T. *et al.* (2020) Making pre-trained language models better few-shot learners. *arXiv preprint arXiv:2012.15723.*

Goodfellow,I. *et al.* (2016) *Deep Learning.* MIT Press, Cambridge, MA, USA, pp. 237–238.

Hastie,T. *et al.* (2019) *Statistical Learning with Sparsity: The Lasso and Generalizations.* Chapman and Hall/CRC, Boca Raton, FL, USA.

He,K. *et al.* (2015) Deep residual learning for image recognition. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA. pp. 770–778.

Hu,W. *et al.* (2019) Strategies for pre-training graph neural networks. *arXiv preprint arXiv:1905.12265.*

Jumper,J. *et al.* (2021) Highly accurate protein structure prediction with alphafold. *Nature*, **596**, 583–511.

Karimi,M. *et al.* (2019) DeepAffinity: interpretable deep learning of compound–protein affinity through unified recurrent and convolutional neural networks. *Bioinformatics*, **35**, 3329–3338.

Kingma,D.P. and Ba,J. (2014) Adam: a method for stochastic optimization. *arXiv preprint arXiv:1412.6980.*

Lan,Z. *et al.* (2019) ALBERT: A Lite BERT for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942.*

Lim,H. *et al.* (2016) Large-scale off-target identification using fast and accurate dual regularized one-class collaborative filtering and its application to drug repurposing. *PLoS Comput. Biol.*, **12**, e1005135.

Lin,A. *et al.* (2019) Off-target toxicity is a common mechanism of action of cancer drugs undergoing clinical trials. *Sci. Transl. Med.*, **11**, eaaw8412.

Liu,Y. *et al.* (2021) COVID-19 multi-targeted drug repurposing using few-shot learning. *Front. Bioinf.*, **1**, 18.

Luxburg,U.V. (2007) A Tutorial on Spectral Clustering. *Stat. Comput.*,**17**, 395–416.

Lynch,I.I.I. *et al.* (2017) Potential functional and pathological side effects related to off-target pharmacological activity. *J. Pharmacol. Toxicol. Methods*, **87**, 108–126.

Mistry,J. *et al.* (2020) Pfam: the protein families database in 2021. *Nucleic Acids Res.*, **49**, D412–D419.

Oprea,T.I. (2019) Exploring the dark genome: implications for precision medicine. *Mamm. Genome*, **30**, 192–200.

Rao,R. *et al.* (2019) Evaluating protein transfer learning with tape. *Adv. Neural Inf. Process. Syst.*, **32**, 9689–9701.

Rives,A. *et al.* (2021) Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proc. Natl. Acad. Sci. USA*, **118**, e2016239118.

Sakamuru,S. *et al.* (2021) Predictive models to identify small molecule activators and inhibitors of opioid receptors. *J. Chem. Inf. Model.*, **61**, 2675–2685. [CrossRef][*10.1021/acs.jcim.1c00439*]

Valentini,V. *et al.* (2013) Evidence for a role of a dopamine/5-HT6 receptor interaction in cocaine reinforcement. *Neuropharmacology*, **65**, 58–64.

Váradi,A. *et al.* (2015) Synthesis and characterization of a dual kappa-delta opioid receptor agonist analgesic blocking cocaine reward behavior. *ACS Chem. Neurosci.*, **6**, 1813–1824.

Vaswani,A. *et al.* (2017) Attention is all you need. In: *31th Advances in Neural Information Processing Systems*, Long Beach, CA, USA. pp. 5998–6008.

Wan,F. and Zeng,J.M. (2016) Deep learning with feature embedding for compound–protein interaction prediction. *bioRxiv*, 086033.

Wang,Y. and Zeng,J. (2013) Predicting drug–target interactions using restricted Boltzmann machines. *Bioinformatics*, **29**, i126–i134.

Wong,C.H. *et al.* (2019) Estimation of clinical trial success rates and related parameters. *Biostatistics*, **20**, 273–286.

Xi,Z. *et al.* (2011) Brain cannabinoid CB2 receptors modulate cocaine's actions in mice. *Nat. Neurosci.*, **14**, 1160–1166.

Xie,L. *et al.* (2012) Novel computational approaches to polypharmacology as a means to define responses to individual drugs. *Annu. Rev. Pharmacol. Toxicol.*, **52**, 361–379.

Xu,K. *et al.* (2018) How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826.*

Yang,H. *et al.* (2017) Deletion of type 2 metabotropic glutamate receptor decreases sensitivity to cocaine reward in rats. *Cell Rep.*, **20**, 319–332.