

Cross-modality and self-supervised protein embedding for compound–protein affinity and contact prediction

Yuning You¹ and Yang Shen ^{1,2,*}

¹Department of Electrical and Computer Engineering, Texas A&M University, College Station, TX 77843, USA and ²Department of Computer Science and Engineering, Texas A&M University, College Station, TX 77843, USA

*To whom correspondence should be addressed.

Abstract

Motivation: Computational methods for compound–protein affinity and contact (CPAC) prediction aim at facilitating rational drug discovery by simultaneous prediction of the strength and the pattern of compound–protein interactions. Although the desired outputs are highly structure-dependent, the lack of protein structures often makes structure-free methods rely on protein sequence inputs alone. The scarcity of compound–protein pairs with affinity and contact labels further limits the accuracy and the generalizability of CPAC models.

Results: To overcome the aforementioned challenges of structure naivety and labeled-data scarcity, we introduce cross-modality and self-supervised learning, respectively, for structure-aware and task-relevant protein embedding. Specifically, protein data are available in both modalities of 1D amino-acid sequences and predicted 2D contact maps that are separately embedded with recurrent and graph neural networks, respectively, as well as jointly embedded with two cross-modality schemes. Furthermore, both protein modalities are pre-trained under various self-supervised learning strategies, by leveraging massive amount of unlabeled protein data. Our results indicate that individual protein modalities differ in their strengths of predicting affinities or contacts. Proper cross-modality protein embedding combined with self-supervised learning improves model generalizability when predicting both affinities and contacts for unseen proteins.

Availability and implementation: Data and source codes are available at <https://github.com/Shen-Lab/CPAC>.

Contact: yshen@tamu.edu

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

1 Introduction

Most FDA-approved drug–target pairs are between small-molecule compounds and proteins (Santos *et al.*, 2017). Considering the enormous chemical space that is estimated to contain 10^{60} ‘drug-like’ compounds (Bohacek *et al.*, 1996), it is desirable to virtually screen compounds with high throughput and high accuracy, based on their computationally predicted properties as well as interactions with proteins (off)targets. Thanks to quickly growing data, modeling techniques and computing power, many machine-learning and deep-learning methods emerge for predicting compound–protein interactions, in particular, the structure-free ones addressing the often unavailability of protein structures (Gao *et al.*, 2018; Jiang *et al.*, 2020; Karimi *et al.*, 2019, 2021; Li *et al.*, 2020; Öztürk *et al.*, 2018; Tsubaki *et al.*, 2019).

Recent progress in structure-free methods includes increasing resolution of what they predict: from binary interactions (Gao *et al.*, 2018; Tsubaki *et al.*, 2019) to continuous affinity or activity values (Karimi *et al.*, 2019; Öztürk *et al.*, 2018). The progress also includes increasing explainability about how they predict such interactions: intermolecular atom–residue non-bonded contacts underlying compound–protein affinities are additionally predicted, often by introducing (Gao *et al.*, 2018; Karimi *et al.*, 2019), regularizing (Karimi *et al.*, 2021) and supervising (Karimi *et al.*, 2021; Li *et al.*, 2020) various attention mechanisms. We refer to such an

explainable affinity prediction problem as compound–protein affinity and contact (CPAC) prediction.

Despite the aforementioned progress, two challenges present major barriers to the accuracy and the generalizability. (i) **Lack of structure awareness.** While being generally applicable by assuming no co-crystal, docked or even unbound structures as protein inputs, structure-free methods rely on 1D amino-acid sequences (Li *et al.*, 2020; Öztürk *et al.*, 2018) and sequence-predicted 1D structural property sequences (Karimi *et al.*, 2019), thus lack the awareness of 3D structures that are critical to what they predict (affinity and contact labels). (ii) **Scarcity of labeled data.** Compared to the daunting size of compound–protein pairs, only a tiny fraction are labeled with affinity measurements and even less are labeled with non-bonded atomic contacts from co-crystal structures. This challenge for supervised models is known as ‘supervision starvation’.

To address the aforementioned challenges, we make two major contributions accordingly. First, to address structure naivety, without demanding co-crystal, compound-docked or even unbound protein 3D structures, we consider protein data as available in both modalities of 1D sequences and sequence-predicted 2D graphs (contact maps). Recent revolution in protein structure prediction (Baek *et al.*, 2021; Jumper *et al.*, 2021) is making the structure modality increasingly available. We introduce various neural network architectures to separately or jointly embed protein modalities and introduce **cross-modality learning** to inject structure-awareness into

resulting protein embeddings. Two cross-modality strategies, concatenation and cross-interactions, are introduced to encode the modalities independently and dependently. Second, to address supervision starvation, without demanding more labeled data, we leverage massive unlabeled protein data and introduce various **self-supervised learning** strategies to pre-train protein embedding. Specifically, we use masked language models (Devlin *et al.*, 2018) for pre-training protein sequence embedding and graph completion and graph contrastive learning (You *et al.*, 2021) for pre-training protein contact-map embeddings.

In cross-modality learning, we ask whether individual modalities could excel in predicting either affinities or contacts as well as whether and how their individual strengths could be combined for better accuracy and generalizability. Our results indicate that the 1D and 2D modalities of protein data do not dominate each other in CPAC prediction for proteins seen in the training set; however, they tend to generalize better for unseen proteins in affinity prediction and contact prediction, respectively. We thus provide a conjecture for such observations, which is verified numerically. To integrate knowledge from 1D and 2D protein modalities, two cross-modality schemes are proposed, with empirical demonstration that they achieve the state-of-the-art (SOTA) performance.

In self-supervised learning, we ask how to design self-supervised strategies, within and across individual protein modalities, in order to improve model accuracy and generalizability. We leverage rich unlabeled protein data and adopt self-supervised techniques for sequences and graphs so as to pre-train protein embeddings. Consistent with aforementioned results without pre-training, self-supervised pre-trainings of individual protein modalities differ in their strengths of predicting affinity or contacts. We further explore self-supervision on top of cross-modality learning, ask which pre-training scheme is beneficial in what circumstances of CPAC prediction, and provide conjectures to underlying reasons.

The rest of the manuscript is organized as follows. In Section 2, we will start with our curated, labeled and unlabeled data, to supervise model training and pre-train protein embedding, respectively. After introducing a backbone model for CPAC prediction and our modifications, we will introduce our methods of cross-modality learning and multi-modal self-supervised learning. In Section 3, we will first examine performances from single- and multi-modal learning without pre-training. We will then examine self-supervised pre-training within and across modalities.

2 Materials and methods

2.1 Data

Labeled dataset. We evaluate CPAC prediction methods through performing training and inference on a CPAC benchmark set (Karimi *et al.*, 2021; You and Shen, 2020) as follows.

(i) *Data source:* The diverse dataset contains 4446 pairs between 1287 proteins and 3672 compounds that are collected from PDBbind (Liu *et al.*, 2015) and BindingDB (Liu *et al.*, 2007) together with their affinity labels. In addition, their contact labels are gathered from the corresponding co-crystal structures deposited in the PDBsum database (Laskowski *et al.*, 2018) using LigPlot. Histograms of protein and compound lengths, measured in the number of protein residues and that of compound atoms, are shown in Supplementary Appendix SA, Fig. S1.

(ii) *Protein and compound graphs:* No 3D structures of proteins or compounds are used. Instead, RaptorX-Contact (Xu, 2019) is used to predict contact maps of proteins from sequences, where evolutionary information from multiple sequence alignment and structural information from its labels are additionally included. Only binary contact maps are used without 3D structural information, thus called 2D graphs. RDKit (Landrum *et al.*, 2006) is used to convert 1D SMILES into 2D chemical structures for compounds, after sanitization.

(iii) *Dataset split:* The labeled dataset is split into subsets of various challenging levels in generalizability: 795 pairs involving unseen proteins (proteins not present in the training set), 521 pairs

Table 1. Statistics of the dataset splits for affinity and contact prediction

		3672 compounds	
		3100	572
1287 proteins	1228	Training set: 2334 pairs Seen both test set: 591 pairs	Unseen-compound test set: 521 pairs
	59	Unseen-protein test set: 795 pairs	Unseen-both test set: 205 pairs

involving unseen compounds, and 205 for unseen both; whereas the rest is randomly split into training (2334) including validation and the default test (591) sets (Karimi *et al.*, 2021). Statistics of the dataset split is presented in Table 1.

Unlabeled datasets. We pre-train protein embeddings using two unlabeled datasets of different scales. Both are from Pfam-A, a database of protein domain sequences (Mistry *et al.*, 2021): (i) The *smaller set* with ground-truth structure information consists of 60 137 sequences from Pfam-A with PDB entries (Berman *et al.*, 2000), from which we extract contact maps from their PDB structures (two residues are deemed in contact if their C_β , or C_α for glycines, are within 8Å). (2) The *larger set* not necessarily with ground-truth structure information is Pfam-A RP15 which consists of 12 798 671 sequences with 15% Representative Proteomes co-membership (Chen *et al.*, 2011) threshold applied. Histograms of protein lengths are shown in Supplementary Appendix SA, Fig. S2.

2.2 Model backbone

The backbone of a CPAC prediction model is a system that is given a compound-protein pair as inputs and simultaneously predicts intermolecular affinity and atom-residue contacts as outputs. Here, we adopt the SOTA CPAC model, DeepAffinity+ (Karimi *et al.*, 2021), as our models' backbone.

Mathematically, given a compound-protein pair $(X_{\text{comp}}, X_{\text{prot}}) \in \mathbb{X}_{\text{comp}} \times \mathbb{X}_{\text{prot}}$ consisting of N_{comp} atoms in each compound and N_{prot} residues in each protein (padding is applied to ensure fixed sizes for all compounds or proteins), a CPAC model $f_{\text{CPAC}} : \mathbb{X}_{\text{comp}} \times \mathbb{X}_{\text{prot}} \rightarrow \mathbb{R}_{\geq 0} \times [0, 1]^{N_{\text{comp}} \times N_{\text{prot}}}$ aims at predicting both the compound-protein affinity z_{aff} and the intermolecular atom-residue contacts Z_{cont} . It includes the following three major components as shown in Figure 1.

- Neural-network encoders:** $f_{\text{comp}} : \mathbb{X}_{\text{comp}} \rightarrow \mathbb{R}^{N_{\text{comp}} \times D}$ and $f_{\text{prot}} : \mathbb{X}_{\text{prot}} \rightarrow \mathbb{R}^{N_{\text{prot}} \times D}$ that separately extract embeddings H_{comp} for the compound X_{comp} and H_{prot} for the protein X_{prot} where D is the hidden dimension. In DeepAffinity+, the compounds are available in 2D chemical graphs and proteins are only available in 1D amino-acid sequences. Accordingly, DeepAffinity+ used graph neural networks (GNNs) such as graph convolutional network (GCN) and GIN (Kipf and Welling, 2016; Veličković *et al.*, 2017) to encode 2D chemical graphs of compounds and hierarchical recurrent neural network (HRNN; El Hhihi and Bengio, 1996) to encode 1D amino-acid sequences of proteins.
- Contact module:** $f_{\text{cont}} : \mathbb{R}^{N_{\text{comp}} \times D} \times \mathbb{R}^{N_{\text{prot}} \times D} \rightarrow [0, 1]^{N_{\text{comp}} \times N_{\text{prot}} \times L \times D}$ takes molecular embeddings from the encoders H_{comp} and H_{prot} as inputs, employs a joint attention mechanism (Karimi *et al.*, 2019, 2021) to output the atom-residue interaction matrix Z_{cont} and jointly embeds the compound-protein pair into H_{cp} , where L is the hidden length determined by N_{comp} and N_{prot} .
- Affinity module:** $f_{\text{aff}} : \mathbb{R}^{L \times D} \rightarrow \mathbb{R}$ predicts the affinity z_{aff} given the joint embedding H_{cp} . It consists of 1D convolutional, pooling layers and multi-layer perceptron (MLP). Note that the contact-predicting interaction module feeds the affinity module,

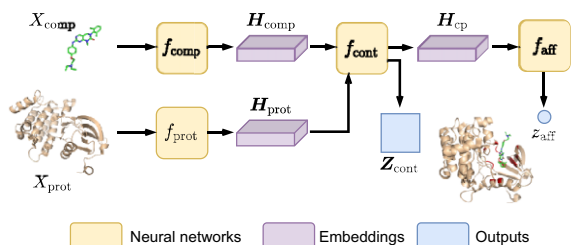


Fig. 1. Illustration of the backbone model f_{CPAC} for CPAC prediction

making affinity prediction intrinsically interpretable by the underlying contacts.

After the CPAC model f_{CPAC} forwardly generates the outputs $(z_{\text{aff}}, Z_{\text{cont}})$, true labels $(y_{\text{aff}}, Y_{\text{cont}})$ are provided to calculate the loss, l_{CPAC} , which consists of affinity loss l_{aff} , intermolecular atom-residue contact loss l_{cont} and three structure-aware sparsity regularization losses l_{group} , l_{fused} and l_{L1} as described in Karimi *et al.* (2021):

$$l_{\text{CPAC}} = l_{\text{aff}} + \lambda_{\text{cont}} l_{\text{cont}} + \lambda_{\text{group}} l_{\text{group}} + \lambda_{\text{fused}} l_{\text{fused}} + \lambda_{\text{L1}} l_{\text{L1}}. \quad (1)$$

The model is trained end to end while the training loss is minimized. More details for the pipeline can be found in Karimi *et al.* (2021).

2.3 Single-modality protein embeddings

In the conventional structure-free CPAC pipeline, compounds are represented as 2D chemical graphs since 1D SMILES strings have limited descriptive power and known worse performance in many tasks (Karimi *et al.*, 2019, 2021; Li *et al.*, 2020), whereas proteins are usually represented as 1D amino-acid sequences without exploration of other modalities. We delve into this under-explored area, proposing to utilize multi-modality protein data for CPAC prediction.

1D sequences. We follow DeepAffinity+ (Karimi *et al.*, 2021) as described in Section 2.2 and use HRNN to encode protein sequences. One change we made is replacing the hierarchical joint attention with naïve joint attention in the interaction module expressed as:

$$\begin{aligned} Z_{\text{cont}} &= Z'_{\text{cont}} / \text{sum}(Z'_{\text{cont}}), \\ z'_{\text{cont},i,j} &= (b_{\text{comp},i} W_{\text{comp,attn}})^T (b_{\text{prot},j} W_{\text{prot,attn}}), \end{aligned} \quad (2)$$

where $z_{i,j} = Z[i,j]$, $b_i = H[i,:]$, $i = 1, \dots, N_{\text{comp}}$, $j = 1, \dots, N_{\text{prot}}$; $W_{\text{comp,attn}}$ and $W_{\text{prot,attn}}$ are two learnable attention matrices.

2D contact maps. We propose to adopt the 2D modality of proteins as additional inputs and model them as graphs with the following reasons. (i) Graphs are more structure-aware compared to 1D sequences, potentially resulting in better generalizability. (ii) Graphs are concise yet informative (focusing on pairwise residue interactions) compared to the data structure of 3D coordinates (which are also harder to predict than contact maps; Cao and Shen, 2020). (iii) The recent surge of models for graph learning (Kipf and Welling, 2016; Veličković *et al.*, 2017) provides advanced tools to facilitate graph representation learning.

As unbound or ligand-bound structure data are not readily available for many proteins, we use sequence-predicted 2D contact maps (Xu, 2019) and can also use AlphaFold2 (Jumper *et al.*, 2021). Thereby, we additionally represent a protein input X_{prot} as a graph $\mathcal{G}_{\text{prot}} = \{\mathcal{V}_{\text{prot}}, \mathcal{E}_{\text{prot}}\}$ where vertices stand for residues and edges exist between residues predicted to be in contact. The graphs are associated with feature matrix $F_{\text{prot}} \in \mathbb{R}^{N_{\text{prot}} \times D}$ (embedded amino-acid types of residues) and the adjacency matrix $A_{\text{prot}} \in \{0,1\}^{N_{\text{prot}} \times N_{\text{prot}}}$ (binary contact map). We employ an expressive GNN model, graph attention network (GAT; Veličković *et al.*, 2017) with K layers as the protein encoder f_{prot} to extract graph embeddings, with the formulation of each layer's forward propagation as:

$$\begin{aligned} H_{\text{prot}}^{(k)} &= \text{MLP}(\tilde{S}^{(k-1)} H_{\text{prot}}^{(k-1)}), \\ \tilde{S}^{(k-1)} &= (D^{(k-1)})^{-1} (S^{(k-1)} \odot A_{\text{prot}}), \\ S^{(k-1)} &= \exp(H_{\text{prot}}^{(k-1)} W^{(k-1)} (H_{\text{prot}}^{(k-1)})^T), \end{aligned} \quad (3)$$

where $H_{\text{prot}} = H_{\text{prot}}^{(K)}$, $H_{\text{prot}}^{(0)} = F_{\text{prot}}$, the normalization matrix $D^{(k-1)} = \text{diag}((S^{(k-1)} \odot A_{\text{prot}}) J_{N_{\text{prot}},1})$, \odot is the element-wise multiplication, $J_{N_{\text{prot}},1}$ is an all-ones matrix with size $N_{\text{prot}} \times 1$ and $W^{(k-1)}$ is a learnable weight matrix. Comparison with the simplest GNN model, GCN is conducted in Supplementary Appendix SB to demonstrate the necessity of adopting the more expressive GAT.

2.4 Cross-modality protein embeddings

To integrate the knowledge from both 1D and 2D protein modalities, we introduce two cross-modality protein embedding schemes as follows.

Cross-modality concatenation. A simple integration model is to concatenate the extracted embeddings of the 1D and 2D modalities encoded by HRNN and GAT, respectively, as shown in Figure 2a. Indeed, concatenation is commonly used in previous work (Hamilton *et al.*, 2017; Xu *et al.*, 2018) to preserve information from different sources. The concatenated output is fed to an MLP for the final protein embedding H_{prot} .

Cross-modality cross-interaction. Although the aforementioned concatenation strategy preserves the information of individual modalities, the encoding processes for the two modalities are isolated. In other words, the two types of embeddings from different modalities were independently encoded and then mixed through concatenation. However, the different modalities of proteins are intrinsically correlated with each other and could be coupled in a properly designed representation-learning process. Therefore, we introduce a cross-interaction module to facilitate the encoder to learn protein embeddings from correlated data (1D and 2D modalities), as shown in Figure 2b. Specifically, given the outputs of encoders $H'_{\text{prot,seq}}$ and $H'_{\text{prot,graph}}$, we calculate sequence and graph cross-modality outputs $H_{\text{prot,seq}}$ and $H_{\text{prot,graph}}$, respectively:

$$\begin{aligned} b_{\text{prot,seq},n} &= (\text{sigmoid}((b'_{\text{prot,graph},n})^T b'_{\text{prot,seq},n}) + 1) b'_{\text{prot,seq},n}, \\ b_{\text{prot,graph},n} &= (\text{sigmoid}((b'_{\text{prot,seq},n})^T b'_{\text{prot,graph},n}) + 1) b'_{\text{prot,graph},n}, \end{aligned} \quad (4)$$

where $b_n = H[n,:]$, $H'_{\text{prot,graph}} = H'_{\text{prot,graph}} W_{\text{cross,graph}}$, $H'_{\text{prot,seq}} = H'_{\text{prot,seq}} W_{\text{cross,seq}}$; $W_{\text{cross,seq}}$ and $W_{\text{cross,graph}}$ are learnable weights.

Instead of independently extracting knowledge from protein modalities (1D sequences and 2D contact maps), the cross-interaction module enforces a learned relationship between the encoded embeddings of the two protein modalities, which is expected to better capture the information from the correlated modalities and to benefit the affinity and contact prediction. Again, $H_{\text{prot,seq}}$ and $H_{\text{prot,graph}}$ (now with information from each other) are concatenated and fed to an MLP for the final protein embedding H_{prot} .

The idea of cross-interaction was previously introduced in Tan and Bansal (2019) and modified here as follows. (i) We do not normalize cross-interaction along residues (sequence length is 1000 here) since it would significantly change the scale of the residue embeddings. (ii) We restrict the cross-interaction for each residue in the range of $[0, 1]$ with sigmoid function to represent the cross-modality ‘interaction strength’.

2.5 Multi-modality self-supervised pre-training

On top of the aforementioned cross-modality learning models, we further propose self-supervised pre-training for the following two reasons. (i) The paired and labeled data curated for CPAC (Karimi *et al.*, 2021) are limited (4446 compound-protein pairs in total), while there are more than billions of unpaired and unlabeled data available (here we make use of protein domain sequences as described in Section 2.1). Exploiting such abundant unlabeled data would generate context-relevant embeddings for downstream, as previously explored under unsupervised learning in CPAC

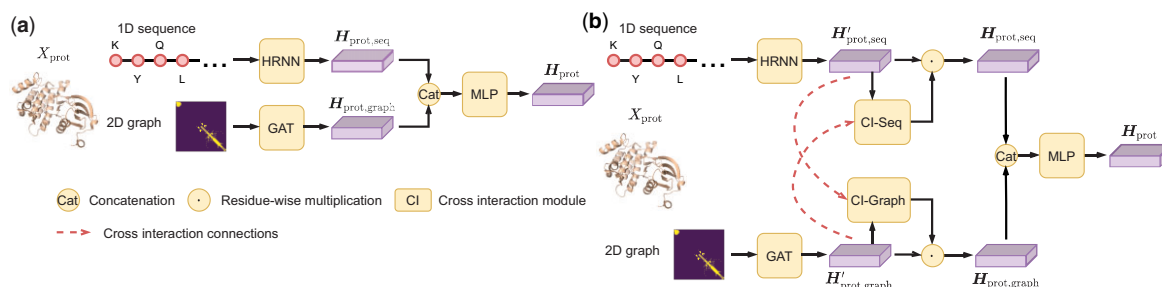


Fig. 2. Cross-modality encoders for proteins (f_{prot} in Fig. 1) to capture and integrate knowledge across data modalities. (a) Naïve concatenation preserves information from different sources and (b) cross-interaction additionally introduces information flows between modalities

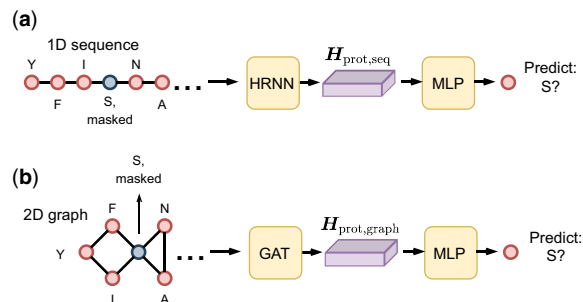


Fig. 3. Self-supervised tasks for pre-training cross-modality encoders (Fig. 2) in CPAC. (a) MLM takes the randomly masked amino-acid sequences as inputs, predicting the masked residues with network outputs and (b) graph completion (GraphComp) with inputting masked-residues contact maps, makes prediction for the masked tokens

prediction (Karimi *et al.*, 2019). (ii) Compared to conventional unsupervised learning, recently emerging self-supervised learning on both sequences (Devlin *et al.*, 2018) and graphs (You *et al.*, 2020a, b, 2021, 2022) further exploits the benefit from unlabeled data.

For reasons above, we introduce the following pre-training strategies, as illustrated in Figure 3. In addition, graph contrastive learning GraphCL (You *et al.*, 2021) is also applied.

Masked language modeling for sequences. We adopt masked language modeling (MLM) for the 1D sequence encoder HRNN, which is well-known as the dominant pre-training strategy in natural language processing (Devlin *et al.*, 2018). MLM takes the randomly masked amino-acid sequences as inputs and tries to predict the masked residues (we use residue types for a proper self-supervising ‘curriculum’) with network outputs, as illustrated in Figure 3a. The mathematical formulation of MLM optimization is expressed as:

$$\min_{\{\text{HRNN}, \text{MLP}\}} \mathcal{L}_{\text{CE}}(\text{MLP}(\text{HRNN}(\bar{F}_{\text{prot}})), Y_{\text{mask}}), \quad (5)$$

$$\text{s.t.} \quad \bar{F}_{\text{prot}}, Y_{\text{mask}} = \text{mask}(F_{\text{prot}}),$$

where $\mathcal{L}_{\text{CE}}(\cdot)$ is cross-entropy loss, \bar{F}_{prot} is the masked feature matrix, Y_{mask} is the masked residues and $\text{mask}(\cdot)$ is the masking function.

MLM reconstructs and enforces the missing knowledge through utilizing the sequential relation (where the information flow is specified by sequential inputs), which aligns with the 1D-modality model exploiting protein sequence information. We thus hypothesize that MLM pre-training provides performance gains in the tasks where the 1D-modality model has performed well, i.e. affinity prediction, which is supported by experimental results in Section 3.4.

Masked graph modeling (graph completion) for contact maps. Self-supervision on graph-structured data recently raises great interests with numerous self-supervised tasks proposed (You *et al.*, 2020a, b, 2021). We choose a simple and effective scheme, graph completion or GraphComp (You *et al.*, 2020b), to pre-train the 2D graph encoder GAT. GraphComp can be viewed as ‘the graph version of MLM’: it takes graphs with randomly masked residues as input and aims at making prediction for the masked tokens using the

structure-aware graph information, as illustrated in Figure 3b. GraphComp optimization is mathematically formulated as:

$$\min_{\{\text{GAT}, \text{MLP}\}} \mathcal{L}_{\text{CE}}(\text{MLP}(\text{GAT}(\bar{F}_{\text{prot}}, A_{\text{prot}})), Y_{\text{mask}}), \quad (6)$$

$$\text{s.t.} \quad \bar{F}_{\text{prot}}, Y_{\text{mask}} = \text{mask}(F_{\text{prot}}).$$

Joint self-supervised pre-training. Besides single-modality pre-training, we also propose joint pre-training for the cross-modality models that simultaneously perform MLM and GraphComp for self-supervision (since sequence and protein encoders share the amino-acid embedding layer, we cannot individually pre-train them and then load the checkpoints). Given benefits from single-modality pre-training, we expect more benefits can be achieved from multi-modality pre-training in both tasks of affinity prediction (where 1D modality models performed well) and contact prediction (where 2D modality models performed well). Results in Section 3.5 partly justified the added benefits.

Details about model training, including hyperparameters, are in Supplementary Appendix SH.

3 Results and discussion

We organize results and discussion as follows. Experiments on cross-modality protein embeddings are presented in Sections 3.1 and 3.2, with additional generalizability tests and case studies. Self-supervised pre-training experiments on top of cross-modality models are reported in Sections 3.4 and 3.5.

3.1 Individual modalities have strengths in different tasks

Without pre-training, Table 2 reports various models’ performances for affinity prediction and contact prediction for various test sets. Figure 4 further splits unseen molecules into proteins and compounds of different similarity bins compared to the training set.

In affinity prediction, 1D sequences or 2D graphs did not lead to significant difference for seen proteins. However, speaking of unseen proteins or even non-homologous proteins (sequence identity below 30%) where model generalizability is required, 1D sequences dominated over 2D graphs as inputs for affinity prediction (0.1 lower in RMSE).

One conjecture is that the information in graphs might be more difficult to learn compared to sequences (the training RMSE losses are 0.71 and 0.99 for 1D and 2D modalities, respectively). Moreover, affinity prediction for unseen-protein cases is not as challenging as intermolecular contact prediction to show the benefit of the 2D modality (shown next), as contact prediction often involves tens of thousands of values (rather than a single value) to fit for each compound–protein pair.

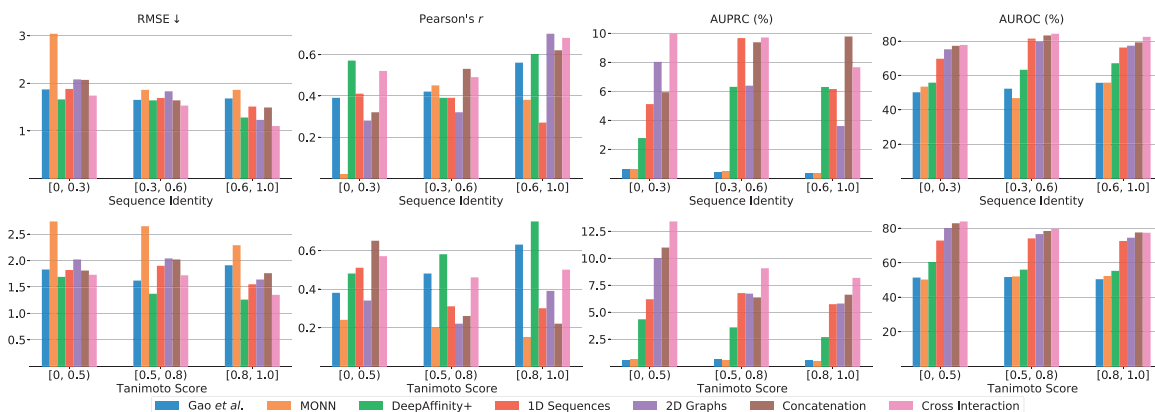
In contact prediction, encoding proteins as 1D sequences again performed better (+3.22% at AUPRC and +1.67% at AUROC) for seen proteins (the proteins in the training set). However, encoding 2D contact maps (graphs) significantly outperformed doing 1D protein sequences (+4.91% at AUPRC and +2.24% at AUROC) for unseen proteins (Table 2) and even more for non-homologous proteins (Figure 4). Using ‘true’ contact maps from (unbound) protein structures showed the same and improved AUROC.

We conjecture that sequential knowledge encoded in 1D amino-acid sequences is well captured especially for seen proteins after

Table 2 Comparison among competing methods and ours in compound–protein affinity prediction (measured by RMSE and Pearson’s r) and contact prediction (measured by AUPRC and AUROC)

Methods	Seen-protein sets		Unseen-protein sets	
	Seen-both	Unseen-compound	Unseen-protein	Unseen-both
Affinity prediction in RMSE (Pearson’s r in parentheses)				
Gao <i>et al.</i> *	1.87 (0.58)	1.75 (0.51)	1.72 (0.42)	1.79 (0.42)
MONN	1.44 (0.70)	1.28 (0.75)	1.67 (0.46)	1.75 (0.45)
DeepAffinity+*	1.49 (0.70)	1.34 (0.71)	1.57 (0.47)	1.61 (0.52)
1D sequences	1.57 (0.67)	1.38 (0.73)	1.63 (0.44)	1.79 (0.40)
Pred. 2D graphs	1.49 (0.68)	1.37 (0.70)	1.75 (0.43)	1.93 (0.34)
True 2D graphs	1.69 (0.59)	1.62 (0.58)	1.88 (0.33)	1.99 (0.25)
Concatenation	1.47 (0.68)	1.37 (0.71)	1.78 (0.47)	1.91 (0.40)
Cross-interaction	1.55 (0.65)	1.43 (0.68)	1.56 (0.50)	1.62 (0.53)
Contact prediction in AUPRC (AUROC in parentheses, %)				
Gao <i>et al.</i> *	0.60 (51.57)	0.57 (51.50)	0.48 (51.60)	0.48 (51.55)
MONN	0.98 (58.57)	0.99 (60.15)	0.99 (65.66)	0.98 (64.59)
DeepAffinity+*	19.74 (73.78)	19.98 (73.80)	4.77 (60.01)	4.11 (59.09)
1D sequences	20.51 (79.01)	20.80 (80.00)	6.54 (73.03)	6.36 (73.41)
Pred. 2D graphs	17.29 (77.34)	17.46 (78.70)	8.78 (77.94)	7.05 (76.59)
True 2D graphs	21.41 (84.60)	21.33 (85.17)	10.52 (84.08)	9.40 (84.29)
Concatenation	23.85 (80.90)	23.52 (81.64)	7.74 (80.59)	7.29 (78.95)
Cross-interaction	23.49 (81.30)	23.29 (82.07)	12.43 (80.64)	9.60 (79.78)

*denotes the cited performances. Boldfaced numbers are the best performances for given test sets. We note that, as intermolecular contacts only represent a minority (around 0.4%) of all compound–protein atom-residue pairs, AUPRC is a much more relevant measure than AUROC for assessing contact prediction.

**Fig. 4.** Generalizability test on various methods for predicting affinity (measured in RMSE and r) and contact (measured in AUPRC and AUROC)

training. The sequential dependency learned from the encoder could be accurate toward intermolecular contact prediction for close or even distant homologs of seen proteins. However, such dependency is less generalizable to unseen or non-homologous proteins. In contrast, the structural topology information encoded in protein 2D contact maps is more difficult for GNNs to capture even for seen proteins, leading to the worse contact predictions for seen proteins. But the information can generalize to unseen proteins well toward contact prediction. In particular, even when sequence similarity for non-homologous proteins (to training ones) is too low to be detectable using RNNS, binding-pocket (subgraph) similarity could still preserve and be detected in 2D contact maps using GNNs thus eventually leads to much better intermolecular contact prediction (Figure 4).

3.2 Cross-modality models combine the strengths

Fusing two modalities’ knowledge together, even by a simple concatenation strategy, could get the best of both modalities. Specifically, the cross-modality model by concatenation had better contact prediction than single-modality models (Table 2). It also had a boost in affinity prediction (better than the 2D single-modality model and slightly worse than the 1D single-modality model).

Enforcing a learned correlation between the 1D and 2D embeddings rather than independently learning two individual embeddings, the cross-modality model with cross-interaction further improved affinity prediction and actually had the best affinity accuracy among all methods for unseen proteins or unseen both. Moreover, it impressively achieved the best AUPRC for unseen proteins and unseen both. These results re-enforce our rationale that the learned correlation between embeddings from different modalities can better capture the data and better perform CPAC predictions.

Our models compare favorably to the SOTA models. They used similar backbone as DeepAffinity+ (Karimi *et al.*, 2021) and revised the joint attention mechanism as mentioned in Section 2.4; thus our 1D sequence-based single-modality model and DeepAffinity+, both using protein sequences, had similar performances in affinity prediction but ours improved contact prediction. Our cross-modality models further improved the performance compared with SOTAs including Gao *et al.*, (2018; after being converted from a binary predictor), MONN (Li *et al.*, 2020) and DeepAffinity+ (Karimi *et al.*, 2021), especially for unseen proteins (Table 2) and non-homologous proteins (Figure 4).

When the protein sequence encoder was changed from HRNN to a pre-trained transformer, no improvement was found (Supplementary Appendix SC).

3.3 Case studies for cross-modality models

All methods are compared in five case studies about compound–protein pairs (Karimi *et al.*, 2021). With detailed results included in [Supplementary Appendix SD](#), we conclude that one or both cross-modality models improved over DeepAffinity+ in AUPRC for four of the five cases. They performed on par with DeepAffinity+ in the precision of the predicted top-10 contacts. The case of LHL–LCK presented the most improvement in the precision of top 10 predicted contacts, from 0.4 to 0.6, as visualized in [Figure 5](#).

3.4 Single-modality pre-training further enhances individual modalities' strengths

We proceed to pre-train our cross-modality model (cross-interaction) in a single-modality setting. In other words, we pre-train the protein sequence and graph encoders using MLM and GraphComp, respectively. The results are detailed in [Table 3](#).

Different pre-training strategies showed different performances relative to no pre-training, depending on the task (affinity or contact prediction) and the test set (seen or unseen proteins/compounds). Consistent with our earlier observation of single-modality models without pre-training, pre-training the embedding of a single modality tended to enhance the strength of the corresponding modality. Specifically, sequence pre-training with MLM, especially with the smaller unlabeled protein dataset, improved upon what the 1D protein modality is good at—affinity prediction, for unseen proteins. MLM over the larger unlabeled set of protein sequences did not show much more benefits, possibly due to the fact that the smaller

unlabeled set and the labeled test sets are biased with protein of structures. Meanwhile, graph pre-training with GraphComp, over the smaller or the larger unlabeled protein dataset, improved upon what the 2D protein modality is good at—contact prediction, mainly for unseen both. Replacing GraphComp (You *et al.*, 2020b) with contrastive learning (GraphCL; You *et al.*, 2021) had similar performances ([Supplementary Appendix SE](#)).

We observe some trade-off between affinity and contact prediction while pre-training a single modality. Part of the reason could be that the two tasks compete with each other while their weighted losses are summed together. The question that remains is whether and how the pre-training strategies for individual modalities can be combined to further enhance model accuracy and generalizability, which is addressed next.

3.5 Multi-modal joint pre-training could further synergize 1D and 2D modalities

We further pre-train our cross-modality model in a multi-modal setting. In other words, we jointly pre-train both the sequence and the graph encoders that share layers. The results are reported in [Table 3](#) as before.

We found that jointly pre-training sequence and graph embeddings with the smaller unlabeled dataset did not change affinity prediction much for unseen proteins and improved contact prediction for the most challenging case of unseen both (+2.4% in AUPRC compared to no pre-training). Interestingly, doing so with the larger unlabeled dataset again improved contact prediction for the most challenging case of unseen both (+2.7% in AUPRC compared to no pre-training) and additionally did so for the unseen proteins (+2.1%

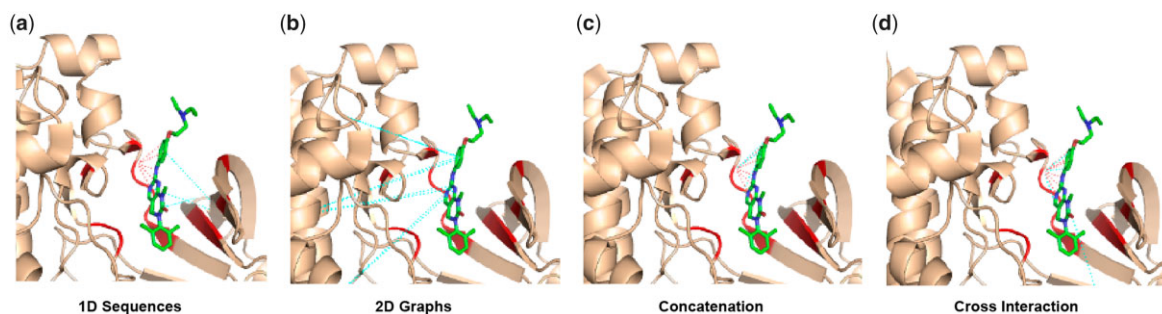


Fig. 5. Visualizing top-10 atom-residue contacts predicted by single- and cross-modal learning for the compound–protein pair of LHL–LCK. Compounds are shown in sticks, proteins in cartoons, and predicted contacts in dashed lines (darker/red for true positives and lighter/cyan for false positives) (A color version of this figure appears in the online version of this article.)

Table 3. Comparison among different pre-training settings (MLM and graph completion, with graph contrastive learning in [Supplementary Appendix SE](#)) based upon the cross-interaction model in compound–protein affinity and contact prediction

Cross-interaction	Seen-protein sets		Unseen-protein sets	
	Seen-both	Unseen-compound	Unseen-protein	Unseen-both
Affinity prediction in RMSE (Pearson's <i>r</i> in parentheses)				
Non pre-train	1.57 (0.66)	1.46 (0.68)	1.63 (0.49)	1.64 (0.54)
MLM-S	1.53 (0.64)	1.40 (0.68)	1.46 (0.56)	1.53 (0.58)
GraphComp-S	1.62 (0.59)	1.44 (0.66)	1.60 (0.43)	1.67 (0.47)
MLM+GraphComp-S	1.64 (0.58)	1.46 (0.65)	1.65 (0.39)	1.65 (0.50)
MLM-L	1.59 (0.62)	1.46 (0.65)	1.62 (0.47)	1.63 (0.57)
MLM+GraphComp-L	1.58 (0.62)	1.45 (0.66)	1.74 (0.33)	1.85 (0.32)
Contact prediction in AUPRC (AUROC in parentheses, %)				
Non pre-train	23.91 (79.48)	23.06 (80.60)	11.40 (77.73)	8.41 (76.42)
MLM-S	23.78 (80.34)	23.33 (81.09)	7.73 (77.44)	6.44 (76.42)
GraphComp-S	23.63 (79.71)	23.41 (81.31)	11.36 (76.67)	9.36 (76.00)
MLM+GraphComp-S	24.13 (82.09)	23.65 (82.70)	11.38 (78.75)	10.83 (78.63)
MLM-L	23.30 (80.40)	23.05 (81.18)	11.35 (81.01)	9.40 (79.46)
MLM+GraphComp-L	23.71 (81.21)	23.22 (82.33)	13.47 (82.00)	11.17 (80.10)

Boldfaced numbers are the best performances.

in AUPRC compared to no pre-training). Impressively, the joint pre-training strategies with predicted protein contact maps even outperformed non-pre-training with actual protein contact maps. In the end, the cross-modality model (cross-interaction) with joint sequence-graph pre-training over the larger set achieved the best contact prediction for both unseen proteins and unseen both. And doing that over the smaller set achieved best-balanced improvement in affinity and contact prediction, potentially suggesting the importance of data quality over data quantity.

We also tested additional pre-training for embedding 2D compound graphs on top of the cross-modality model with joint pre-training of protein data. To do so, we leveraged unlabeled compound data from STITCH. Further improvements, albeit moderate, were observed (Supplementary Appendix SG).

4 Conclusion

In this paper, we address two major challenges to advance explainable prediction of compound-protein affinity (or CPAC): the sequence-dominant yet structure-naïve models and the scarce labeled data. By introducing multi-modal and self-supervised learning for the first time to CPAC prediction, we address both challenges through fostering context- and task-relevant protein embedding. Specifically, to overcome structure naivety, we treat protein data as available in both modalities of 1D sequences and 2D graphs (predicted) and introduce cross-modality learning for sequence- and structure-aware protein embeddings. Empirical results indicated that individual modalities excel in different tasks and our approach of cross-modality learning could bring out the best of both modalities. Additionally, to overcome labeled-data scarcity, we design self-supervised learning strategies within and across modalities to pre-train cross-modal protein embedding. Empirical results indicated that cross-modal learning with joint pre-training can further improve model generalizability for unseen molecules and outperform the state-of-the-art. Meanwhile, there is still much to do for improving the synergy between both tasks of affinity and contact prediction.

Acknowledgements

A part of the computing resources is provided by Texas A&M University High Performance Research Computing (HPRC).

Funding

This paper was published as part of a special issue financially supported by ECCB2022. The study is in part supported by the National Science Foundation [CCF-1943008] and the National Institute of General Medical Sciences of the National Institutes of Health [R35GM124952].

Conflict of Interest: none declared.

References

Back, M. *et al.* (2021) Accurate prediction of protein structures and interactions using a three-track neural network. *Science*, **373**, 871–876.

Berman, H.M. *et al.* (2000) The protein data bank. *Nucleic Acids Res.*, **28**, 235–242.

Bohacek, R.S. *et al.* (1996) The art and practice of structure-based drug design: a molecular modeling perspective. *Med. Res. Rev.*, **16**, 3–50.

Cao, Y. and Shen, Y. (2020) Energy-based graph convolutional networks for scoring protein docking models. *Proteins*, **88**, 1091–1099.

Chen, C. *et al.* (2011) Representative proteomes: a stable, scalable and unbiased proteome set for sequence analysis and functional annotation. *PLoS One*, **6**, e18910.

Devlin, J. *et al.* (2018). Bert: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Vol. 1, Association for Computational Linguistics, Minneapolis, Minnesota. pp. 4171–4186. <https://doi.org/10.18653/v1/N19-1423>.

El Hihi, S. and Bengio, Y. (1996) Hierarchical recurrent neural networks for long-term dependencies. In Touretzky, D. *et al.* (eds) *Advances in Neural Information Processing Systems*, pp. 493–499.

Gao, K.Y. *et al.* (2018) Interpretable drug target prediction using deep neural representation. In *The Association for the Advancement of Artificial Intelligence, Stockholm, Sweden. IJCAI, 2018*, pp. 3371–3377. <https://dl.acm.org/doi/abs/10.5555/3304222.3304236>.

Hamilton, W.L. *et al.* (2017). Inductive representation learning on large graphs. *arXiv preprint arXiv:1706.02216*. <https://doi.org/10.5555/3294771.3294869>.

Jiang, M. *et al.* (2020) Drug-target affinity prediction using graph neural network and contact maps. *RSC Adv.*, **10**, 20701–20712.

Jumper, J. *et al.* (2021) Highly accurate protein structure prediction with alphafold. *Nature*, **596**, 583–589.

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

Karimi, M. *et al.* (2021) Explainable deep relational networks for predicting compound-protein affinities and contacts. *J. Chem. Inf. Model.*, **61**, 46–66.

Kipf, T.N. and Welling, M. (2016). Semi-supervised classification with graph convolutional networks. In *Proceedings of the 5th International Conference on Learning Representations (ICLR)*, Toulon, France. <https://doi.org/10.48550/arXiv.1609.02907>.

Landrum, G. *et al.* (2006). Rdkit: open-source cheminformatics. <https://doi.org/10.5281/zenodo.591637>.

Laskowski, R.A. *et al.* (2018) PDBsum: structural summaries of PDB entries. *Protein Sci.*, **27**, 129–134.

Li, S. *et al.* (2020) MONN: a multi-objective neural network for predicting compound-protein interactions and affinities. *Cell Systems*, **10**, 308–322.

Liu, T. *et al.* (2007) BindingDB: a web-accessible database of experimentally determined protein-ligand binding affinities. *Nucleic Acids Res.*, **35**, D198–D201.

Liu, Z. *et al.* (2015) PDB-wide collection of binding data: current status of the PDBbind database. *Bioinformatics*, **31**, 405–412.

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

Öztürk, H. *et al.* (2018) DeepDTA: deep drug-target binding affinity prediction. *Bioinformatics*, **34**, i821–i829.

Santos, R. *et al.* (2017) A comprehensive map of molecular drug targets. *Nat. Rev. Drug Discov.*, **16**, 19–34.

Tan, H. and Bansal, M. (2019). LXMERT: learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Association for Computational Linguistics, Hong Kong, China. pp. 5100–5111. <https://doi.org/10.18653/v1/D19-1514>.

Tsubaki, M. *et al.* (2019) Compound-protein interaction prediction with end-to-end learning of neural networks for graphs and sequences. *Bioinformatics*, **35**, 309–318.

Veličković, P. *et al.* (2017). Graph attention networks. In *Proceedings of the 6th International Conference on Learning Representations (ICLR)*, Vancouver, Canada. <https://doi.org/10.48550/arXiv.1710.10903>.

Xu, J. (2019) Distance-based protein folding powered by deep learning. *Proc. Natl. Acad. Sci. U S A*, **116**, 16856–16865.

Xu, K. *et al.* (2018). How powerful are graph neural networks? In *Proceedings of the 7th International Conference on Learning Representations (ICLR)*, New Orleans, Louisiana, USA. <https://doi.org/10.48550/arXiv.1810.00826>.

You, Y. and Shen, Y. (2020). Cross-modality protein embedding for compound-protein affinity and contact prediction. *arXiv preprint arXiv:2012.00651*. <https://doi.org/10.48550/arXiv.2012.00651>.

You, Y. *et al.* (2020a) Graph contrastive learning with augmentations. In Larochelle, H. *et al.* (eds) *Advances in neural information processing systems*. Vol. 33, pp. 5812–5823.

You, Y. *et al.* (2020b). When does self-supervision help graph convolutional networks? In *International Conference on Machine Learning*, PMLR, pp. 10871–10880. <https://proceedings.mlr.press/v119>.

You, Y. *et al.* (2021). Graph contrastive learning automated. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*. Proceedings of Machine Learning Research, Vol. 139, pp. 12121–12132. <https://doi.org/10.48550/arXiv.2106.07594>.

You, Y. *et al.* (2022). Bringing your own view: Graph contrastive learning without prefabricated data augmentations. In *Proceedings of the 15th ACM International Conference on Web Search and Data Mining (WSDM)*. Association for Computing Machinery, pp. 1300–1309. <https://doi.org/10.1145/3488560.3498416>.