

Geometric Graph Representation Learning on Protein Structure Prediction

Tian Xia
tianxia@auburn.edu
Auburn University
Auburn, AL, USA

Wei-Shinn Ku
weishinn@auburn.edu
Auburn University
Auburn, AL, USA

ABSTRACT

Determining a protein's 3D from its sequences is one of the most challenging problems in biology. Recently, geometric deep learning has achieved great success on non-Euclidean domains including social networks, chemistry, and computer graphics. Although it is natural to present protein structures as 3D graphs, existing research has rarely studied protein structures as graphs directly. The present research explores the geometry deep learning of three-dimensional graphs on protein structures and proposes a graph neural network architecture to address these challenges. The proposed Protein Geometric Graph Neural Network (PG-GNN) models both distance geometric graph representation and dihedral geometric graph representation by geometric graph convolutions. This research shed new light on protein 3D structure studies. We investigated the effectiveness of graph neural networks over five real datasets. Our results demonstrate the potential of GNNs for 3D structure prediction.

CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability.

KEYWORDS

graph neural networks; geometric deep learning; multi-attribute graph learning

ACM Reference Format:

Tian Xia and Wei-Shinn Ku. 2021. Geometric Graph Representation Learning on Protein Structure Prediction. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '21)*, August 14–18, 2021, Virtual Event, Singapore. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3447548.3467323>

1 INTRODUCTION

Prediction of a protein's structure from its amino acid sequence remains an open problem in the field of life science. The main practical problem confronting us is the challenge that comes from directly predicting protein structure from primary sequence. A common strategy used to study protein structure is to transform the direct

prediction of protein structure into several problems, including contact map prediction, secondary structure prediction, torsion angles prediction and others. Especially with the recent growth of convolutional neural networks (CNNs), several convolution neural networks were proposed to tackle the problem in this field, such as contact map prediction [28, 30], torsion angle prediction [13], and protein structure-property prediction [29].

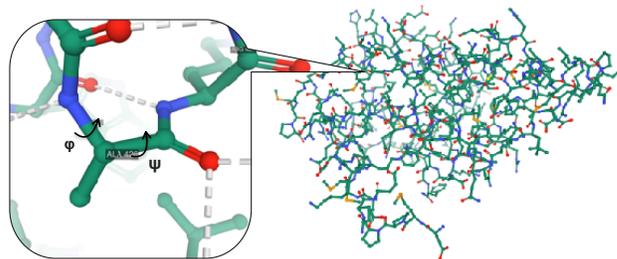


Figure 1: Graph representation of the protein structure.

Following this road, a number of works have been done using deep convolution neural networks to predict the contact map first, then recover the protein structure from contact maps. Several challenges remain to be tackled: **1) The foremost challenge is the fact that the protein contact/distance matrix cannot provide fully structural information needed for protein backbone structure modeling.** Contact maps are sparse matrices [23] that lack sufficient information to be treated by geometric representation because the only information that contact maps provide is the upper or lower bound of the contact threshold. Multiple conformations can be generated since alpha carbon can move freely in the space within the contact threshold. Even though several recent studies [1, 32] use the distance matrix instead of contact map which contains finer-grained information related to the alpha carbon pairwise distances than the contact map, they are still insufficient to model protein backbone geometry structure because of the absence of structural information of other backbone atoms in the protein. The free rotation of the chemical bonds around alpha carbon [25] cause the backbone atoms to rotate accordingly. Hence further constraint information is needed to model all atoms in the protein backbone structure. As illustrated in Fig. 1, the protein backbone structure composed of consecutive chains of coplanar units of $C_{\alpha} - CO - NH - C_{\alpha}$ with two primary degrees of freedom: dihedral angles named Phi (ϕ) and Psi (ψ). The rotation information on both sides of the C_{α} in addition to distance matrix could provide further geometric information of all the atoms in the protein backbone structure. Although both distance information and dihedral angle information

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
KDD '21, August 14–18, 2021, Virtual Event, Singapore

© 2021 Association for Computing Machinery.
ACM ISBN 978-1-4503-8332-5/21/08...\$15.00
<https://doi.org/10.1145/3447548.3467323>

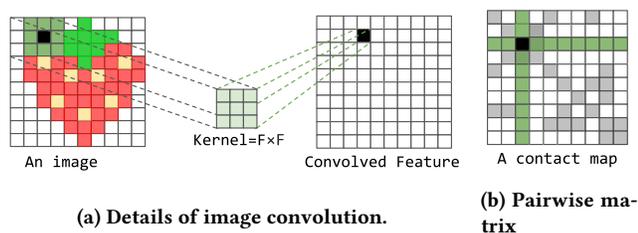


Figure 2: Illustration of image and pairwise distance matrix

are helpful to study the protein backbone structure, few studies have focused on modeling multiple attributes for protein backbone structural representations. 2) **Secondly, it remains challenging to capture the non-local relations in protein structures.** Existing works [28, 30] typically use image CNN-based methods which focus on local neighbor amino acids in the protein sequence but cannot consider those far away in the sequence. As shown in Fig. 2(a), the 2D image convolution operation only focuses on local information by multiplying the convolutional kernel over each pixel and its local neighbors. Although this operation is enough for image processing, one limitation of these methods is that our pairwise feature matrix is not image representation. Fig. 2(b) shows an example of the pairwise relation matrix, which represents the amino acid pairs in the feature representation space. The green cross in Fig. 2(b) shows the edges connected to the same nodes for the black square edge. Although all the pairwise features in the green cross could potentially influence the pairwise relation in the black square, only the local relations were considered in the image convolution process. Thus it remains challenging to model long-range relations in the protein sequences. 3) **The variable protein sequence lengths and the large size of the protein structures make the problem difficult to handle.** This problem is usually overcome by transforming the 2D feature maps via cropping [7, 28] or padding [11]. Although such operations generate fixed-sized inputs for the neural network, they could potentially cause the loss of relational details or undesired distortion of the contact/distance matrix. While there are few studies which model the chemical molecules as a graph, considering protein as a whole graph to study its structural properties remains an open problem in the area. The long length of protein sequences and the large sizes of the proteins have limited the approaches to model protein as one graph directly.

To rectify the above problems, we investigate the native structures of the protein and their common representations. Although the natural way to represent a protein structure is to model it as a 3D graph, the protein 3D graph structure has rarely been studied directly. Recent theoretical developments in graph neural networks inspired us to look at protein structure representation differently. Fig. 1 shows the structure of protein 2XSE. The backbone of the protein holds a protein structure together with residues of each amino acid. The C_{α} is the central atom in the backbone structure, which has two backbone angles (ϕ, ψ). Dihedral angles (ϕ, ψ) for alanine in protein 2XSE are marked in Fig. 1. We focus on overcoming the protein structure problem by modeling the protein structure as a 3D geometric graph and design a geometric graph convolutional network architecture based on this specific problem.

Our goal is to generate distance geometric graph representation and dihedral geometric graph representations together for protein structure modeling. Compared to positional 3D graph representations, i.e., Cartesian coordinates of nodes, our proposed 3D graph representation method gives a significant advantage because of its invariance to rotation and translation of the graph. To the best of our knowledge, this is the first work that can address all the above challenges. We summarize the main contributions as follows:

- We formally formulate the problem of protein backbone structure modeling as geometric 3D graph representations. We model the input graph into multi-attribute graphs in which the node represents the residues and the edge represents pairwise information between residues.
- A novel architecture is proposed for protein backbone 3D structure graph generation. Our proposed model could generate a protein graph with both geometric distance graph representations and geometric dihedral graph representations together.
- We propose the use of novel geometric graph convolution blocks for distance geometric graph representation generations. As the sizes of the proteins vary, our proposed approach can handle the sizes of protein graphs dynamically.
- Comprehensive experiments were conducted to validate the effectiveness of our proposed model in the generated 3D geometric protein graphs.

2 RELATED WORK

In this section, we will present and discuss three lines of research that are relevant to our work.

2.1 Protein structure prediction

Experiments for protein structure determination are time-consuming and expensive; thus, modeling the 3D structure of a protein remains one of the most important problems in bioinformatics [21]. Significant work has been done toward the construction of the protein 3D structure during recent decades. Critical Assessment of Protein Structure Prediction (CASP) [17] established benchmarks and assessed methods for protein structure prediction. Protein contact map and conformation prediction have proven to help the reconstruction of protein 3D structure [11, 24]. [13, 33] showed that inter-residue orientations in addition to residue distances can be used for protein structure prediction [3, 11].

2.2 Deep Learning for Protein Structure Prediction

Deep learning-based methods were already widely used to solve protein structure related problems, such as protein-protein interaction prediction, protein contact map prediction, protein secondary structure predictions, and protein dihedral angle predictions. Deep learning methods for inter-residue distance and contact prediction have considerably advanced protein structure prediction. [10] proposed an energy-based model using transformer architecture for protein conformation. [3] proposed end-to-end recurrent geometric network to predict 3D protein structure. [4] generated a protein 3D structure by using Generative Adversarial Networks. Currently, the most successful methods for residue-residue distance prediction

are CNN-based neural networks. [30] first applied CNN to predict protein contacts. [2], [12], [28] and [31] used dilated CNN method for protein contact map prediction.

2.3 Graph Neural Network

Graph neural network (GNN) attracts attention in a wide range of areas [35] including natural language processing, computer vision, traffic prediction, and so on. Graph convolutional network has shown practical utility in the field of chemistry. Inspired by GNN, several works treat molecules as graphs and achieved great progress. [22] used graph neural networks to learn energy function for small molecules. [9] proposed a three-dimensional graph convolution network to predict molecular properties. [14] uses graph convolutional networks to predict protein site-specific functions. [16] also utilizes GNN to tackle the molecule properties prediction problem.

3 PROBLEM FORMULATION

The protein backbone holds the protein together and generates the tertiary structure of the protein. This section introduces the protein backbone geometry problem: this is a sequence-to-structure task where we take protein amino acid sequences as input to predict the geometry of 3D protein backbone structure. To tackle this complex problem, we first formulate the protein structure into a graph representation as follows:

We consider each input protein as a graph $G(\mathcal{V}, \mathcal{E}, E, F)$, where \mathcal{V} is the set of L nodes in the graph representing amino acid residues and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ is the set of $L - 1$ edges in the protein sequence. $F \in \mathbb{R}^{L \times D}$ stands for the input node attribute matrix where $F_i \in \mathbb{R}^{1 \times D}$ refers to the node attribute of node i and D is the dimension of the node attribute vector. The input node attributes F include position-specific scoring matrix, predicted secondary structure, solvent accessibility, etc. $E \in \mathbb{R}^{L \times L \times K}$ is the edge attributes tensor, where $E_{i,j} \in \mathbb{R}^{1 \times K}$ refers to the edge attribute of edge $e_{i,j}$ and K is the dimension of edge attributes tensor. Likewise, the input edge attributes tensor E include co-evolution information, distance potential, and inter-residue coupling score.

Similarly, the target protein backbone geometry graph can be represented as $G(\mathcal{V}', \mathcal{E}', E', F')$, where \mathcal{V}' is the set of L nodes in the graph representing target amino acid residues and $\mathcal{E}' \subseteq \mathcal{V}' \times \mathcal{V}'$ is the set of M edges in the target graph. $E' \in \mathbb{R}^{K \times L \times L}$ denotes the target edge attribute matrix where $E'_{k,i,j}$ denotes the k -th feature of edge $e_{i,j}$. $F' \in \mathbb{R}^{H \times L}$ denotes the target node attribute matrix, where $F'_{k,i} \in \mathbb{R}^{1 \times H}$ is the node attributes of node i and H is the dimension of the node attributes. For example, the k -th node attribute phi of node i is the dihedral angle between the plane $C_{\beta}^{i-1}, N^i, C_{\alpha}^i$ and the plane $N^i, C_{\alpha}^i, C_{\beta}^i$ as shown in Fig. 11(a). In addition, psi is also a node feature which represents the dihedral angle between the plane $N^i, C_{\alpha}^i, C_{\beta}^i$ and the plane $C_{\alpha}^i, C_{\beta}^i, N^{i+1}$ as shown in Fig. 11(b). Without loss of generality, we assign the first node feature for torsion angle phi, namely F'_1 and the second node feature for torsion angle psi, namely F'_2 .

Our goal is to develop a graph translation model that can encode both the node and edge features extracted from input protein graph $G(\mathcal{V}, \mathcal{E}, E, F)$ and generate a graph-based geometric representation

for protein 3D structures $G(\mathcal{V}, \mathcal{E}', E', F')$. As the input graph and output graph have the same protein sequence composed of the same set of nodes, we have $\mathcal{V} = \mathcal{V}'$. Hence the output is graph-based representation for protein backbone geometry $G(\mathcal{V}, \mathcal{E}', E', F')$, where the pairwise distance information is represented by edge attributes E' and the dihedral angle information is represented by node attributes F' . The main advantage of such geometric representation is the invariant property of angle-geometric graph representation and distance-geometric graph representation. Hence they are invariant under rotation of the coordinate system and graph translation. As most of the research in this field is aimed at getting only distance representation or torsion representation separately, here we solve the above mentioned problems simultaneously by 3D geometric graph generation. Since the input graph node and edge attributes E, F and the target graph node and edge attributes E', F' are different, the learning from the multi-attribute graph input to the graph geometric representation output can be defined as learning a mapping: $G(\mathcal{V}, \mathcal{E}, E, F) \rightarrow G(\mathcal{V}, \mathcal{E}', E', F')$

4 METHODOLOGY

In this section, we propose Protein Geometric Graph Neural Network (PG-GNN) to model the geometric properties in terms of distance and angle representations and learn the geometric representations jointly from two separate translation paths. We first describe the overall architecture of the PG-GNN with the translation paths. We then describe in detail how the geometric representations are learned with our proposed edge translation path and node translation path collaboratively.

4.1 Model construction

Taking protein sequences as input to directly construct geometric 3D representation of the protein structures remains an open problem in the area [3]. In light of the above discussion, we need a framework that can dynamically handle different sized graph inputs and jointly generate output of both node and edge attributes together. With this aim in mind, in this paper we present a new framework composed of two translation paths to predict edge and node attributes separately. The key components of our framework are edge translation path and node translation path. The illustration of our proposed framework is shown in Fig. 3. For the edge translation path, a deep residual convolutional network is proposed that takes both node and edge features as the input and output information on the pairwise distance of all residue pairs in the protein. The objective of the edge translation path is to learn the mapping: $G(\mathcal{V}, \mathcal{E}, E, F) \rightarrow G(E')$. For the node translation path, we utilize a fast graph message-passing neural network which takes predicted pairwise distance potential and node features then outputs all node torsion angles (ϕ, ψ) in the protein graph. The objective of the node translation path is to learn the mapping: $G(\mathcal{V}, \mathcal{E}, E, F) \rightarrow G(F')$. The overall network is based on minimization with distance and orientation restraints derived from both edge translation path and node translation path outputs.

Although we can generate both node and edge attributes on separate paths based on the framework described above, the predicted attributes may not be consistent as they are generated from different paths.

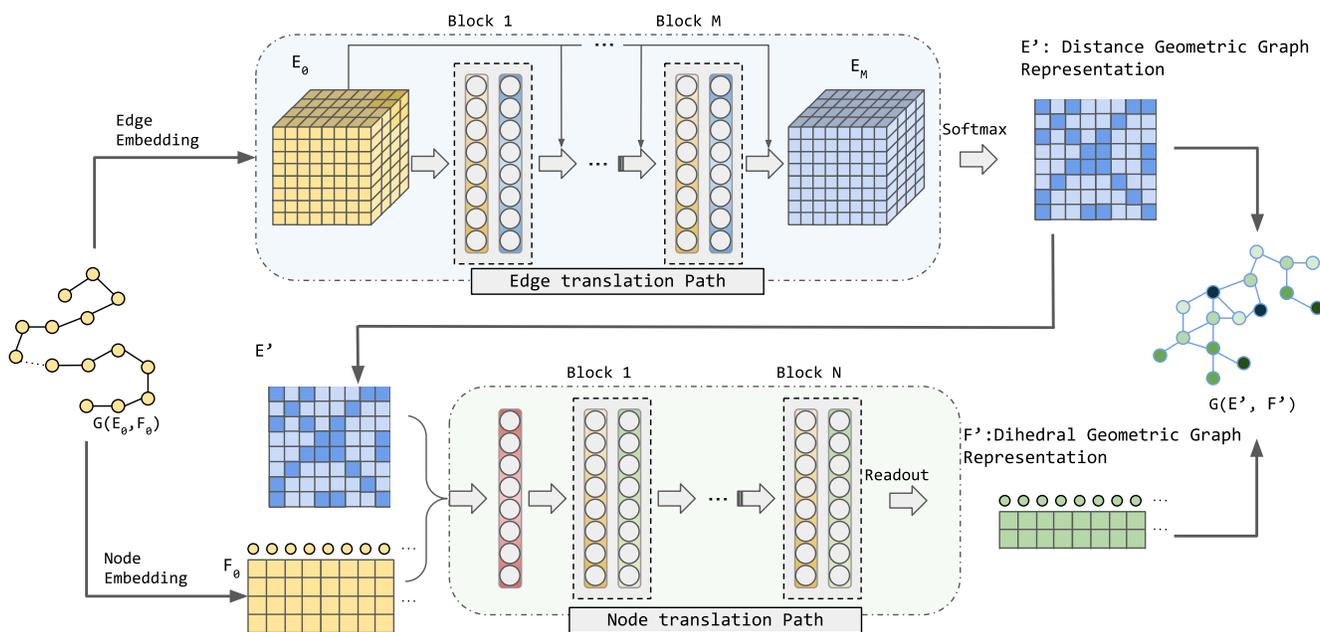


Figure 3: Overall architecture. The framework consists of two parts: the edge translation path is shown on top while node translation path is shown on the bottom. Given input protein sequences, we extract the features for both nodes and edges. Then the input can be denoted as $G(\mathcal{V}, \mathcal{E}, E, F)$ and fed into the framework as shown.

For the edge translation path, we use cross entropy loss; only the probability corresponding to ground-truth participates in calculation. The calculation can be written as $L(E, E') = -\frac{1}{L \times L} \sum_{i=1}^L \sum_{j=1}^L y_{ij} \times \log \hat{y}_{ij}$ where L is the number of residues for each protein, \hat{y}_{ij} is the predicted distance label between i -th residue and j -th residue in a protein.

For the node translation path, we use mean squared error as loss function. The equation is as follows:

$$\begin{aligned}
 L(F, F') &= L(F_1, F'_1) + L(F_2, F'_2) \\
 &= \frac{1}{L} \sum_{i=1}^L \left(\sin(\phi_i) - \sin(\hat{\phi}_i) \right)^2 + \frac{1}{L} \sum_{i=1}^L \left(\cos(\phi_i) - \cos(\hat{\phi}_i) \right)^2 \\
 &+ \frac{1}{L} \sum_{i=1}^L \left(\sin(\psi_i) - \sin(\hat{\psi}_i) \right)^2 + \frac{1}{L} \sum_{i=1}^L \left(\cos(\psi_i) - \cos(\hat{\psi}_i) \right)^2 \quad (1)
 \end{aligned}$$

where L is the number of residues for each protein.

For the overall training process, our PG-GNN is under the guidance of both edge translation path to learn the graph edge attributes and graph node translation path to learn the node attributes information. We set a parameter λ to balance the degree of the two models. Thus the overall loss function for the network is:

$$L = L(E, E') + \lambda \times L(F, F') \quad (2)$$

4.2 Edge Translation Path (ETP)

The aim for edge translation path is modeling all the interactions between edges and nodes to generate the geometric distance graph representation of the protein backbone structure. One challenging

problem for 3D protein structure modeling is that residues sequentially far apart might be in spatially close contact in protein's 3D structure.

Designing a translation path that can model different range interactions between residues and effects of all interactions connected to one residue is the key for our edge translation path. As the image convolution kernel only focuses on the source and its surrounding pixels, the long range patterns cannot be characterized by traditional image convolution kernels. Thus how to model both of the non-local and local pairwise relations in feature representation space becomes the crucial factor for protein structure modeling. Therefore, we propose a multi-branch convolution block for edge translation path to capture all the influences from different interaction types in proteins. Fig. 4 shows the operation branches in a single block, where the input graph attributes E_s updated to output graph edge attributes E_{s+1} . Each block composed of one identity mapping branch, one edge-to-edge convolution branch, and one two-dimensional (2D) convolution branch. The latter two branches used different hyper-parameters (filter type and sizes) separately. With 2D convolution operation focused on the local relations and edge-to-edge convolution operation focused on the long range contacts, our proposed edge translation block integrates both the local and non-local contacts features to generate each edge's attribute.

Edge-to-edge convolution layers We implement an edge-to-edge filter which can consider the edges that connect to one residue. As shown in Fig. 4, the protein contact map has its distinct features, unlike images. Thus, simply integrating the image convolution methods without modification will not work well for this problem. The edge-to-edge convolution is described as follows.

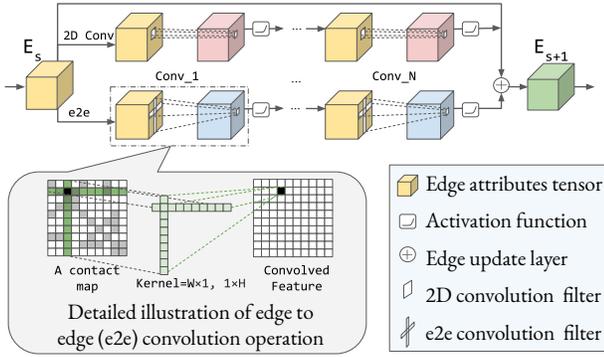


Figure 4: Details of protein graph edge convolution path for one protein in a single block

Formally, we denote $A \in \mathbb{R}^{L \times L}$ as an adjacency matrix where L is the length of each protein, $W^{i,j}$ as the shared weights for the edge $e_{i,j}$, and ρ as the non-linear differentiable function that computes the activations. The graph convolution operation over the edges $e_{i,j}$ for kernel W can be defined with $f_W^{i,j} = W^{i,j} \cdot A_{i,j}^{l-1,m}$, likewise the convolution operation $f_H^{j,i}$ for kernel H can be defined as $f_H^{j,i} = W^{i,j} \cdot A_{j,i}^{l-1,m}$. Then the edge-to-edge convolution is as follows:

$$A_{i,j}^{l,n} = \rho \left(\sum_{n=1}^L (f_W^{i,j} + f_H^{j,i}) \right) \quad (3)$$

For the edge $e_{i,j}$, the shared weights $W^{i,j}$ across kernel H and W make the size of the latent convolutional representation independent of the size of the input.

2D convolution layers We used 3×3 convolutional filters for the 2D convolution layers, followed by batch normalization and Exponential Linear Unit (ELU) activation function. Dilation convolution filters can also be used in this convolution branch. The illustration of the detailed convolution operation is shown in Fig. 2(a).

Edge updating layer Each edge is updated by integrating both of the convolution operations outputs and identity mapping of the inputs. All layer outputs and residual connections are added together and fed into the next block. The details are shown in Fig. 4 as \oplus operation. In this way, the edge will be generated by both related nodes in the input sequences and edges that connected to the same nodes.

4.3 Node Translation Path

For the node translation path, the model aims to generate the node attributes matrix F' based on learning the interactions between protein residues and the effects of the residue edges to the residues. As shown in Fig. 10, the torsion angles (ϕ, ψ) are influenced by both of the connected edges and nodes. Hence both node attributes and the output adjacent matrix of the edge translation path are used as the input for our node translation path. We use message passing to capture all the nodes and edges influences for the torsion attribute generation. The overall architecture of each block in node translation path contains two layers: message passing layer and node

update layers. The message passing layer learns all the influences from each pair of nodes and the node update layers aggregate all the influences and generate the new node attributes.

Message Passing on node layers: As shown in Fig. 3, the input of the proposed node translation path is sampled from both the node representation F and edge representation E . We take the output distance matrix from the edge translation path at each iteration to use as edge representation E_{vw} to feed into the node translation network [20]. The inputs are fed into the message-passing layer as Equation (4) to aggregate all incoming messages.

$$M_t(h_v^t, h_w^t, e_{vw}^t) = A_{e_{vw}} h_w^t \quad (4)$$

Node updating layers: After computing the message passing of all nodes, we update the hidden state by the Gated Recurrent Units (GRU) as Equation (5).

$$h_i^{t+1} = GRU(h_i^t, m_i^{t+1}) \quad (5)$$

Readout layer: We denote node attributes $F_i = (F_{1,i}, F_{2,i})$ for i -th residue in the protein sequence. We further represent the dihedral angles as $v_i = (v_{a,i}, v_{b,i}, v_{c,i}, v_{d,i})$, which denote the $\sin \phi_i, \cos \phi_i, \sin \psi_i, \cos \psi_i$ for training purposes. The readout layer is shown in Fig. 3 as R in the last layer of the node translation path. We have the readout function as follows:

$$\hat{v} = R(\{h_v^T | v \in G\}) \quad (6)$$

Then we derive the prediction of dihedral angles of $\hat{F}_i = (\arctan(\hat{v}_{a,i}/\hat{v}_{b,i}), \arctan(\hat{v}_{c,i}/\hat{v}_{d,i}))$ from the output of the model $\hat{v}_i = (\hat{v}_{a,i}, \hat{v}_{b,i}, \hat{v}_{c,i}, \hat{v}_{d,i})$.

5 EXPERIMENT

In this section, we present our experiments using the proposed methods on five real-world test datasets. The further details on the parameters and features used in the experiments are provided in the supplement material.

5.1 Dataset

The datasets that are used in the experiments are all real-world datasets. The proteins with a sequence length of more than 300 residues were excluded from the datasets. The contact matrix for the proteins in the test sets was calculated using the C_α^i atoms distance. We define the contact between two residues where the relative distances between C_α^i atoms of given residues are less than 8\AA .

Training/validation dataset: The PDB25 dataset were filtered as described in [30]. After the filtering process, 500 proteins were randomly sampled for the test set, and the remaining proteins were used for the training and validation of the models. The total number of proteins used in the training dataset is 4726 proteins and the validation dataset is 500 proteins. The validation dataset is used for hyperparameter tuning and to prevent model overfitting.

Testing dataset: The trained models were tested on the following datasets: the PDB25 dataset described above, the protein domains used in the CASP11 and CASP13 [17], 76 hard CAMEO benchmark proteins, and a test set that contains 400 membrane proteins (membrane dataset) [30].

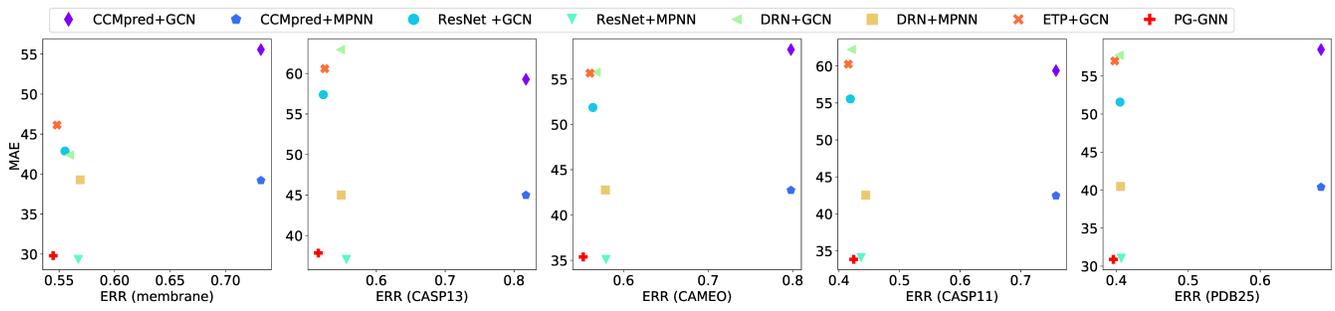


Figure 5: Overall performance comparison

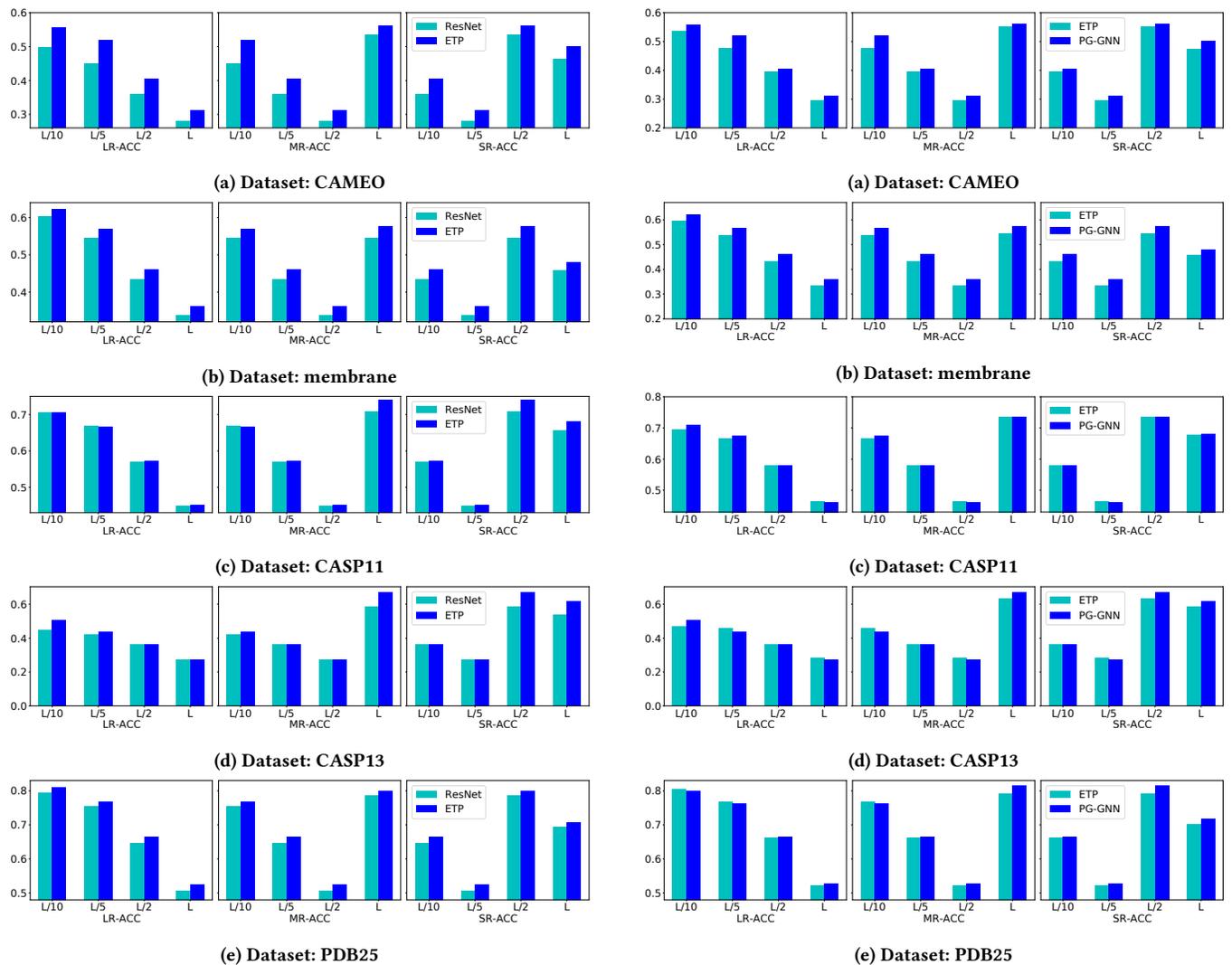


Figure 6: Performance comparison: proposed edge translation path (ETP) vs ResNet.

Figure 7: Performance comparison: Our proposed framework PG-GNN vs Edge translation path only (ETP)

5.2 Benchmarks

As we are the first to propose a multi-task framework for protein graph representation learning, our proposed method is compared with the two categories of methods: 1) the methods used for protein edge (residue-residue contact) prediction only or protein node (torsion angles (ϕ, ψ)) prediction only. We adapted the methods listed as baselines. The methods were implemented to reproduce the results provided by the authors. The goal of re-implementing the methods was to investigate the performance with the same data. The existing web-servers provided by authors only make the inference based on the pre-trained model which is trained with different datasets. Hence, our re-implementation ensures baselines to train the models with the same data and features. Each method represents state-of-the-art performance in either contact map prediction task, torsion angle prediction, or graph node prediction task.

- CCMpred: This method is proposed by [27]. The source code was directly used in the edge translation path of our framework as baseline.
- ResNet: This method was first adapted by [30] to be used for protein contact map prediction. We re-implemented the neural network and used the same parameters stated in [30] as our baseline for the edge translation path.
- Dilated Residual Networks (DRN): This method was used by [28], [1], for protein residue-residue contact prediction. We re-implemented the neural network and used the same parameters stated in Deepcon [1] as our baseline for edge translation path.
- Graph Convolution Network (GCN): We adapted plain GCN with dynamic k-NN described in [19] as our baseline for node translation path. We used $K=5, 10, 20$ for k-NN selection.
- ETP: This is our proposed method for edge translation path.
- Message passing neural network (MPNN): The method [20] is described in the node translation path in our proposed framework.

2) We adapted the above mentioned methods into our proposed framework as baselines, namely CCMpred + GCN, CCMpred + MPNN, ResNet + GCN, ResNet + MPNN, DRN + GCN, DRN + MPNN, and ETP + GCN. To further validate the effectiveness of the proposed node-edge joint convolution framework, we conduct comparison baseline models of node translation path and edge translation path trained separately as well.

5.3 Evaluation Metrics

5.3.1 Quantitative Evaluation. A set of metrics are used to measure the similarity between the generated graph and the real graphs in terms of node and edge attributes. To measure the generated edge attribute performance, we follow the CASP measurements [17]. We use the accuracy (ACC) and error (ERR) of the top L/k predicted contacts. L is the sequence length of nodes in each graph and $k = 10, 5, 2, 1$. We evaluate the non-local contacts where the sequence distance belongs to the three groups: short range (SR) [6,11], medium range (MR) [12,23], and long range (LR) [24,∞). To measure the generated node attributes performance, we use MAE (mean absolute error) between node attributes of generated and real target graphs.

5.3.2 Qualitative Evaluation. Qualitative evaluation are performed to compare the reconstructed protein structure with the X-ray crystallograph protein structures by the structure alignment of superimposed protein backbones [26]. Full atom structures were constructed by PyRosetta full-atom relaxation with our generated distance and torsion angle restraints [8][34]. The lowest-energy full atom model was selected for quality evaluation and visualization.

6 RESULT

6.1 Investigating Performance of Multi-Attribute Properties

The performances of the methods described in Section 5.2 are shown in Fig. 5. Fig. 5 is composed of five subplots, each illustrating the results of five different datasets. The color of the points refers to the specific method listed on the top of Fig. 5. The values on the x -axis show the average error rate of the models' performance for the edge translation path. Similarly, the values on the y -axis show the average MAE of the models' performance for the node translation path. Because lower is better for both the average error rate and average MAE, the best performance model is the one that is closest to the (0,0) point. Our proposed PG-GNN (red cross) shows the best performance across all five datasets. PG-GNN proves the ability to handle both the geometric node attributes and geometric edge attributes together. CCMPred was compared as a co-evolutionary analysis method for the edge translation path. All the deep neural network models used in edge translation path performs at least 47.1% better in predicting the edge representations than threading based method CCMPred. All of the proposed node translation models (MPNN) perform 25.2% better than the GCN models on average. We did an ablation study in the next section to further investigate the performances of the proposed edge-to-edge convolution branch and the node translation path.

6.2 Ablation Study

6.2.1 Evaluation of the proposed graph edge convolution block. We compared the robustness of the proposed edge translation path with the edge translation path without the edge-to-edge convolution shown in Fig. 6. The light blue bar shows the performance of the network without the edge-to-edge convolution branch. The dark blue bar shows the performance of our proposed framework. The results of all five datasets follow the same trend. We observe similar performance gains with the addition of the edge-to-edge convolution branch. The average ACCs with the proposed ETP for membrane, CASP13, CAMEO, CASP11 and PDB25 datasets are 45.5%, 48.4%, 44.8%, 57.6%, and 60.4%, while the average ACCs without the edge-to-edge convolution branch are 43.3%, 44.3%, 42.1%, 56.4%, and 59.3%, respectively. The proposed framework with the edge-to-edge convolution branch outperformed the model without the edge-to-edge convolution branch in all long-range, medium-range, and short-range metrics. The highest average ACC improvement outperforms the baseline by 36.8%. This proves that the proposed edge-to-edge convolution branch successfully helps the model's performance in learning edge-related representations.

6.2.2 Evaluation of the node translation path. Following the same approach used for the evaluation of the edge translation path, we

compared the performance of our node translation path (MPNN) with GCN for translating node-level features in our proposed framework. We use $k = 5, 10, 20$ for the training of GCN as the baseline. The results of average MAE for (ϕ, ψ) prediction compared with the ground truth are following the same trend. Table 1 shows the average MAE for ϕ prediction and Table 2 shows the average MAE for ψ prediction. The average MAE results show that our proposed node translation path successfully outperformed GCN models. Specifically, MPNN outperforms the comparison GCN methods by 20.6%, 22.5%, and 22.8% on node attribute ϕ on average. Comparing with GCN, we show similar performance gains by combining our proposed node translation path with edge translation path baselines in Fig. 5. This proves the superiority of the proposed node translation path in the interpretation of node-related representations.

Table 1: Performance comparison with different node translation path (ϕ)

Datasets	GCN(k=5)	GCN(k=10)	GCN(k=20)	MPNN
CASP11	32.350	34.107	33.343	25.696
CASP13	32.401	33.027	32.876	28.729
CAMEO	32.385	32.854	33.337	26.635
membrane	32.063	33.285	32.891	23.892
PDB25	32.143	33.027	34.642	24.114

6.2.3 Evaluation of the joint convolution framework. To evaluate the performance of the joint convolution framework, we compare our proposed joint convolution framework with a network that only contains our proposed edge translation path. The results are shown in Fig. 7. The light blue bar shows the performance of the network without the node translation path. The dark blue bar shows the performance of PG-GNN. The average ACC improvements of the joint convolution framework for dataset membrane, CASP13, CAMEO, CASP11, and PDB25 are 5.6%, 3.6%, 4.7%, 0.9%, and 1.2%. The performance improvements are consistent across all five test datasets. The joint convolution framework performs better in almost all metrics than the edge translation path alone consistently. Thus, the proposed PG-GNN can not only jointly predict the node and edge attributes, but also performs better than the edge translation path alone.

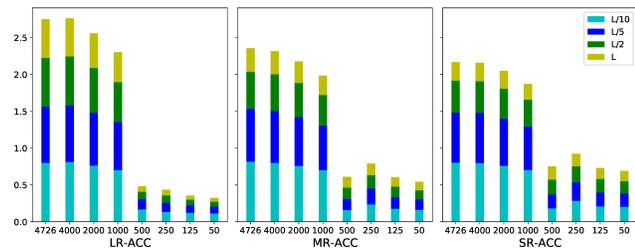


Figure 8: Sensitivity analysis of training dataset sizes

6.3 Sensitivity analysis on the effect of training dataset sizes

To further evaluate the effect of dataset sizes, we perform a sensitivity analysis over eight different sizes of the training datasets. The full training dataset contains 4725 proteins from which we randomly sampled 4000, 2000, 1000, 500, 250, 125, and 75 proteins to train the models. Fig. 10 shows the prediction results of the PDB25 test dataset. The yellow, green, dark blue and light blue bar segment represents the average ACC results of the top L/10, L/5, L/2, and L predicted contacts. We use a stacked bar graph with each bar represents the specific size of the training dataset and colored bar segment represent different ranges' metrics. The performance of the network continues to increase with the increase in the size of the training dataset. We noticed a clear trend where the ACC increases marginally from 500 to 1000 proteins compared to other samples. The average ACCs of PG-GNN using 1000 proteins training set increased by 240% comparing to 500 proteins training set. This trend from 500 to 1000 proteins fades when we go to larger datasets. The performances are consistent across all five test datasets.

6.4 Quality evaluation of generated protein structure

Fig. 9 shows the case study of predicted structure representation of five domains using the visualization program PyMol [26]. The predicted protein backbone traces were presented in the first row in yellow, where C_α atoms are balls and yellow lines between balls are backbone traces. The constructed full-atom model and the native structure are shown in the second row in cartoon representation. Our constructed models are light blue and their native structures are in pink. The constructed full atom models are consistent with the predicted backbone traces. The predicted structural models have very close topology to the native structures as shown in the figure.

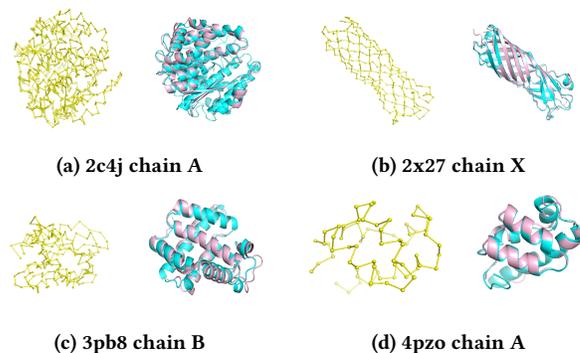


Figure 9: The generated protein backbone trace (yellow), the superimposition of the predicted full atom model (light blue) and its native structure (pink).

7 CONCLUSION

In this work, we revisit a long-existing challenge. Unlike treating the protein as a non-structural problem like an image, we propose

to encode protein topological information into the node and edge representation in graphs. The experiments conducted in this paper proved the robustness of graph neural networks. Using our proposed model, we can achieve comparable results with less training time and smaller training datasets. As we establish an open problem for future research, it would be interesting to see if the performance keeps increasing with more features used. Because protein structures and inter-residue relations are ubiquitous in the real world, we left the further improvement of this longstanding challenge to future work.

ACKNOWLEDGMENTS

This research has been funded in part by the U.S. National Science Foundation grants IIS-1618669 (III) and ACI-1642133 (CICI).

REFERENCES

- [1] Badri Adhikari. 2020. DEEPCON: protein contact prediction using dilated convolutional neural networks with dropout. *Bioinform.* 36, 2 (2020), 470–477. <https://doi.org/10.1093/bioinformatics/btz593>
- [2] Badri Adhikari. 2020. DEEPCON: protein contact prediction using dilated convolutional neural networks with dropout. *Bioinformatics (Oxford, England)* 36, 2 (2020), 470–477. <https://doi.org/10.1093/bioinformatics/btz593>
- [3] Mohammed AlQuraishi. 2019. End-to-End Differentiable Learning of Protein Structure. *Cell Systems* 8, 4 (2019), 292–301.e3. <https://doi.org/10.1016/j.cels.2019.03.006>
- [4] Namrata Anand, Raphael Eguchi, and Po Ssu Huang. 2019. Fully differentiable full-atom protein backbone generation. *Deep Generative Models for Highly Structured Data, DGS@ICLR 2019 Workshop* 11 (2019), 1–10.
- [5] Michael Beckstette, Robert Homann, Robert Giegerich, and Stefan Kurtz. 2006. Fast index based algorithms and software for matching position specific scoring matrices. *BMC Bioinformatics* 7 (2006), 1–25. <https://doi.org/10.1186/1471-2105-7-389>
- [6] Marcos R Betancourt and D Thirumalai. 1999. Pair potentials for protein folding: choice of reference states and sensitivity of predicted native states to variations in the interaction schemes. *Protein science* 8, 2 (1999), 361–369.
- [7] Wendy M Billings, Bryce Hedelius, Todd Millicam, David Wingate, and Dennis Della Corte. 2019. ProSPR: Democratized Implementation of AlphaFold Protein Distance Prediction Network. *BioRxiv* (2019), 830273.
- [8] Sidhartha Chaudhury, Sergey Lyskov, and Jeffrey J Gray. 2010. PyRosetta: a script-based interface for implementing molecular modeling algorithms using Rosetta. *Bioinformatics* 26, 5 (2010), 689–691.
- [9] Hyeoncheol Cho and Insung S. Choi. 2018. Three-Dimensionally Embedded Graph Convolutional Network (3DGCN) for Molecule Interpretation. *CoRR abs/1811.09794* (2018). arXiv:1811.09794 <http://arxiv.org/abs/1811.09794>
- [10] Yilun Du, Joshua Meier, Jerry Ma, Rob Fergus, and Alexander Rives. 2020. Energy-based models for atomic-resolution protein conformations. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26–30, 2020*. OpenReview.net. https://openreview.net/forum?id=S1e_9xrFvS
- [11] Raphael R. Eguchi, Namrata Anand, Christian A. Choe, and Po-Ssu Huang. 2020. IG-VAE: Generative Modeling of Immunoglobulin Proteins by Direct 3D Coordinate Generation. *bioRxiv* (2020). <https://doi.org/10.1101/2020.08.07.242347> arXiv:<https://www.biorxiv.org/content/early/2020/08/10/2020.08.07.242347.full.pdf>
- [12] Hiroyuki Fukuda and Kentaro Tomii. 2018. Deep Neural Network for Protein Contact Prediction by Weighting Sequences in a Multiple Sequence Alignment. *bioRxiv* (2018), 331926. <https://doi.org/10.1101/331926>
- [13] Yujuan Gao, Sheng Wang, Minghua Deng, and Jinbo Xu. 2018. RaptorX-Angle: real-value prediction of protein backbone dihedral angles through a hybrid method of clustering and deep learning. *BMC Bioinformatics* 19, 4 (2018), 100. <https://doi.org/10.1186/s12859-018-2065-x>
- [14] Vladimir Gligorijevic, P. Douglas Renfrew, Tomasz Kosciolke, Julia Koehler Leman, Kyunghyun Cho, Tommi Vatanen, Daniel Berenberg, Bryn Taylor, Ian M. Fisk, Ramiik J. Xavier, Rob Knight, and Richard Bonneau. 2019. Structure-Based Function Prediction using Graph Convolutional Networks. *bioRxiv* (2019), 786236. <https://doi.org/10.1101/786236>
- [15] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7–9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1412.6980>
- [16] Johannes Klicpera, Janek Groß, and Stephan Günnemann. 2020. Directional Message Passing for Molecular Graphs. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26–30, 2020*. OpenReview.net. <https://openreview.net/forum?id=B1eWbxStPH>
- [17] Andriy Kryshchak, Torsten Schwede, Maya Topf, Krzysztof Fidelis, and John Moulton. 2019. Critical assessment of methods of protein structure prediction (CASP)—Round XIII. *Proteins: Structure, Function, and Bioinformatics* 87, 12 (2019), 1011–1020. <https://doi.org/10.1002/prot.25823>
- [18] Michael Levitt and Arieh Warshel. 1975. Computer simulation of protein folding. *Nature* 253, 5494 (1975), 694–698. <https://doi.org/10.1038/253694a0>
- [19] Guohao Li, Matthias Müller, Ali K. Thabet, and Bernard Ghanem. 2019. DeepGCNs: Can GCNs Go As Deep As CNNs?. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 – November 2, 2019*. IEEE, 9266–9275. <https://doi.org/10.1109/ICCV.2019.00936>
- [20] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard S. Zemel. 2016. Gated Graph Sequence Neural Networks. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2–4, 2016, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1511.05493>
- [21] Zhong Li, Yuele Lin, Arne Elofsson, and Yuhua Yao. 2020. Protein Contact Map Prediction Based on ResNet and DenseNet. *BioMed Research International* 2020 (apr 2020), 1–12. <https://doi.org/10.1155/2020/7584968>
- [22] Elman Mansimov, Omar Mahmood, Seokho Kang, and Kyunghyun Cho. 2019. Molecular geometry prediction using a deep generative graph neural network. *CoRR abs/1904.00314* (2019). arXiv:1904.00314 <http://arxiv.org/abs/1904.00314>
- [23] Michal J. Pietal, Janusz M. Bujnicki, and Lukasz P. Kozlowski. 2015. GDFuzz3D: a method for protein 3D structure reconstruction from contact maps, based on a non-Euclidean distance function. *Bioinformatics* 31, 21 (06 2015), 3499–3505. <https://doi.org/10.1093/bioinformatics/btv390> arXiv:<https://academic.oup.com/bioinformatics/article-pdf/31/21/3499/17122233/btv390.pdf>
- [24] Michal J. Pietal, Janusz M. Bujnicki, and Lukasz P. Kozlowski. 2015. GDFuzz3D: a method for protein 3D structure reconstruction from contact maps, based on a non-Euclidean distance function. *Bioinform.* 31, 21 (2015), 3499–3505. <https://doi.org/10.1093/bioinformatics/btv390>
- [25] G.N. Ramachandran, C. Ramakrishnan, and V. Sasisekharan. 1963. Stereochemistry of polypeptide chain configurations. *Journal of Molecular Biology* 7, 1 (1963), 95–99. [https://doi.org/10.1016/S0022-2836\(63\)80023-6](https://doi.org/10.1016/S0022-2836(63)80023-6)
- [26] Schrödinger, LLC. 2015. The PyMOL Molecular Graphics System, Version 1.8. (November 2015).
- [27] Stefan Seemayer, Markus Gruber, and Johannes Söding. 2014. CCMpred—fast and precise prediction of protein residue–residue contacts from correlated mutations. *Bioinformatics* 30, 21 (07 2014), 3128–3130. <https://doi.org/10.1093/bioinformatics/btu500> arXiv:<https://academic.oup.com/bioinformatics/article-pdf/30/21/3128/17147204/btu500.pdf>
- [28] Andrew W. Senior, Richard Evans, John Jumper, James Kirkpatrick, Laurent Sifre, Tim Green, Chongli Qin, Augustin Zidek, Alexander W.R. Nelson, Alex Bridgland, Hugo Penadones, Stig Petersen, Karen Simonyan, Steve Crossan, Pushmeet Kohli, David T. Jones, David Silver, Koray Kavukcuoglu, and Demis Hassabis. 2020. Improved protein structure prediction using potentials from deep learning. *Nature* 577, 7792 (2020), 706–710. <https://doi.org/10.1038/s41586-019-1923-7>
- [29] Sheng Wang, Wei Li, Shiwang Liu, and Jinbo Xu. 2016. RaptorX-Property: a web server for protein structure property prediction. *Nucleic Acids Research* 44, W1 (04 2016), W430–W435. <https://doi.org/10.1093/nar/gkw306> arXiv:<https://academic.oup.com/nar/article-pdf/44/W1/W430/18787392/gkw306.pdf>
- [30] Sheng Wang, Siqi Sun, Zhen Li, Renyu Zhang, and Jinbo Xu. 2017. Accurate De Novo Prediction of Protein Contact Map by Ultra-Deep Learning Model. *PLoS Comput. Biol.* 13, 1 (2017). <https://doi.org/10.1371/journal.pcbi.1005324>
- [31] Qi Wu, Zhenling Peng, Ivan Anishchenko, Qian Cong, David Baker, and Jianyi Yang. 2020. Protein contact prediction using metagenome sequence data and residual neural networks. *Bioinformatics (Oxford, England)* 36, 1 (2020), 41–48. <https://doi.org/10.1093/bioinformatics/btz477>
- [32] Jinbo Xu. 2019. Distance-based protein folding powered by deep learning. *Proceedings of the National Academy of Sciences* 116, 34 (2019), 16856–16865. <https://doi.org/10.1073/pnas.1821309116>
- [33] Jianyi Yang, Ivan Anishchenko, Hahnbeom Park, Zhenling Peng, Sergey Ovchinnikov, and David Baker. 2020. Improved protein structure prediction using predicted interresidue orientations. *Proceedings of the National Academy of Sciences* 117, 3 (2020), 1496–1503. <https://doi.org/10.1073/pnas.1914677117> arXiv:<https://www.pnas.org/content/117/3/1496.full.pdf>
- [34] Jianyi Yang, Ivan Anishchenko, Hahnbeom Park, Zhenling Peng, Sergey Ovchinnikov, and David Baker. 2020. Improved protein structure prediction using predicted interresidue orientations. *Proceedings of the National Academy of Sciences* 117, 3 (2020), 1496–1503.
- [35] Jie Zhou, Ganqu Cui, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, and Maosong Sun. 2018. Graph Neural Networks: A Review of Methods and Applications. *CoRR abs/1812.08434* (2018). arXiv:1812.08434 <http://arxiv.org/abs/1812.08434>

A SUPPLEMENTAL MATERIAL: PRELIMINARIES

As illustrated in Fig. 10, protein backbone structure consists of distance based geometric representation and dihedral angle representation. We proceed to detail the geometric representations of protein graphs.

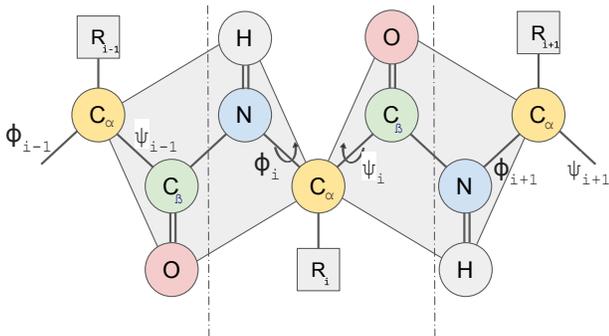


Figure 10: Protein backbone representation.

A.1 Distance-based geometric graph representation

We denoted edge attributes E' as distance geometric representation for protein graph. [18] introduced a chain of pseudoatoms placed at C_α positions to replace the protein main chain model. We adapted this simplification of the protein geometry and use distance matrix between C_α to represent the protein geometry. Thus, the distance matrix E' of C_α determines the overall shape of the backbone structure.

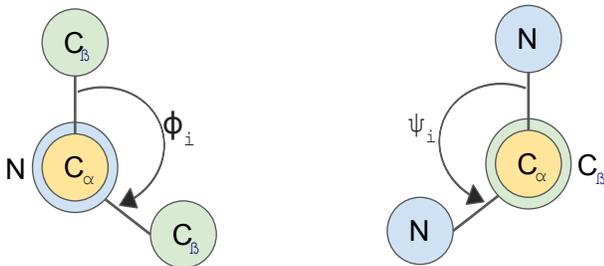


Figure 11: Dihedral angles Phi and Psi (a) The Phi ϕ torsion angle measures the rotation of the N- C_α bond (b) The Psi ψ torsion angle measures the rotation of the C_α -CO bond.

A.2 Dihedral angle geometric graph representation

The distance matrix E' itself cannot fully characterize the overall backbone structure due to the rotation of chemical bonds around

C_α : the N- C_α bond and the C_α -CO bond [25]. These two bonds attached to the C_α can rotate freely to form the unique folding patterns within protein. As shown in Fig. 10, the backbone dihedral angles Phi (ϕ) and Psi (ψ) are in sequence order on either side of C_α to represent the rotation of these bonds. Thus Phi (ϕ) and Psi (ψ) can be denoted as node attributes F' to evaluate the rotation patterns' rise around C_α . For i^{th} amino acid, the dihedral angle ϕ_i is the dihedral angle between the plane $C_\beta^{i-1}, N^i, C_\alpha^i$ and the plane $N^i, C_\alpha^i, C_\beta^i$ as shown in Fig. 11(a). The dihedral angle ψ_i is the dihedral angle between the plane $N^i, C_\alpha^i, C_\beta^i$ and the plane $C_\alpha^i, C_\beta^i, N^{i+1}$ as shown in Fig. 11(b). As each C_α associates with two torsion angles (ϕ, ψ), we denote the dihedral angle geometric representation as node attributes matrix F'_1, F'_2 .

B DETAILS RELATED TO EXPERIMENT SECTION

B.1 Multi-attributes graph representation features

Based on our formulation, each protein is represented by a graph $G(\mathcal{V}, \mathcal{E}, E, F)$ composed of nodes V presenting residues and edges E represented by the protein sequence. Thus the two types of features used in this project are named node features and edge features accordingly. Node feature matrix F are the properties of the single residue, including position-specific scoring matrix [5], predicted secondary structure, and solvent accessibility predictions. Edge feature matrix E are the features that contain pairwise information, such as co-evolution information [27], and distance potential [6],[30]. The extracted node features and edge features are used for the model described above. For edge translation path, the node features F_i and F_j are transformed by 1D convolution then concatenated into $E_{i,j}$ and $E_{j,i}$ to be used as feature map for edge attributes generation. For node translation path, the node feature matrix F and the pairwise distance matrix E' from the edge translation path output are used as the input for node attributes generation.

B.2 Training parameters

The parameters used in PG-GNN are presented in this section. All experiments are conducted on a 64-bit machine with Nvidia GPU (RTX 2080 Ti). For edge translation path, the number of edge translation blocks $S_E = 10$. For each edge translation block, 4 convolution layers per block N was applied sequentially. In parallel, one edge-to-edge layer was applied to input of the edge translation block. For both convolution and edge-to-edge translations ELU activation function was used after each convolution layer. For node translation path, the number of blocks in message passing is $S_V = 6$ with ReLU activation function in each block. The network was trained using Adam [15] with an initial learning rate of 0.00013.

C ADDITIONAL TABLES RELATED TO RESULTS SECTION

Table 2 shows the average MAE for ψ prediction.

Table 2: Performance comparison with different node translation path (ψ)

Datasets	GCN(k=5)	GCN(k=10)	GCN(k=20)	MPNN
CASP11	81.451	80.500	80.592	42.003
CASP13	81.864	83.178	84.556	46.944
CAMEO	81.863	84.185	80.602	44.110
membrane	82.229	81.423	84.559	35.713
PDB25	82.224	83.177	85.316	37.688