
Learning Complete Protein Representation by Deep Coupling of Sequence and Structure

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Learning effective representations is crucial for understanding proteins and their
2 biological functions. Recent advancements in language models and graph neural
3 networks have enabled protein models to leverage primary or tertiary structure
4 information to learn representations. However, the lack of practical methods to
5 deeply co-model the relationships between protein sequences and structures has
6 led to suboptimal embeddings. In this work, we propose CoupleNet, a network that
7 couples protein sequence and structure to obtain informative protein representations.
8 CoupleNet incorporates multiple levels of features in proteins, including the residue
9 identities and positions for sequences, as well as geometric representations for
10 tertiary structures. We construct two types of graphs to model the extracted
11 sequential features and structural geometries, achieving completeness on these
12 graphs, respectively, and perform convolution on nodes and edges simultaneously
13 to obtain superior embeddings. Experimental results on a range of tasks, such as
14 protein fold classification and function prediction, demonstrate that our proposed
15 model outperforms the state-of-the-art methods by large margins.

16 1 Introduction

17 Proteins are the fundamental building blocks of life and play essential roles in a diversity of ap-
18 plications, from therapeutics to materials. They are composed of 20 different basic amino acids,
19 which are lined by peptide bonds and form a sequence. The one-dimensional (1D) sequence of a
20 protein determines its structure, which in turn determines its biochemical function [40]. Due to recent
21 progress in protein sequencing [34], massive numbers of protein sequences are now available. For
22 example, the UniProt [3] database contains over 200 million protein sequences with annotations,
23 *e.g.*, gene ontology (GO) terms, similar proteins, family and domains. Notably, the development of
24 large-scale language models (LMs) in natural language processing has substantially benefited protein
25 research owing to similarities between human language and protein sequences [16, 27]. For instance,
26 models like ProtTrans [14] and ESM-series [39, 33] in learning protein representations have proven
27 successful utility of pre-training protein LMs with self-supervision to process protein sequences.

28 Thanks to the recent significant progress made by AlphaFold2 [30] in three-dimensional (3D) structure
29 prediction, a large number of protein structures from their sequence data are now made available. The
30 latest release of AlphaFold protein structure database [43] provides broad coverage of UniProt [3].
31 Recently proposed structure-based protein encoders become to utilize geometric features [25, 24,
32 53], *e.g.*, ProNet [47] learns representations of proteins with 3D structures at different levels, like the
33 amino acid, backbone or all-atom levels. There also exists a group of methods that build graph neural
34 networks and LMs (LSTMs or attention models) to process sequence and structure [53, 50, 19], for
35 example, GearNet [53] encodes sequential and spatial features by alternating node and edge message
36 passing on protein residue graphs.

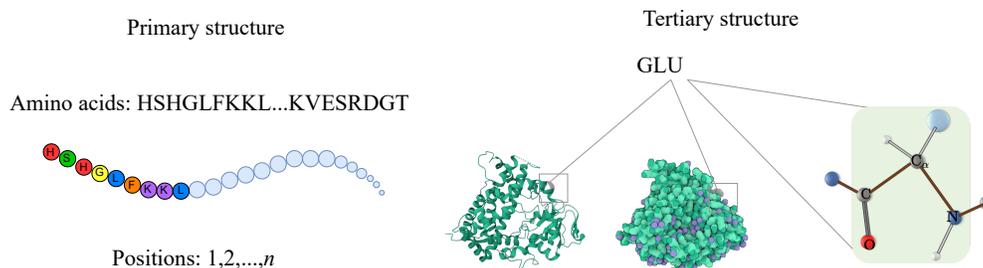


Figure 1: Illustration of the protein sequence and structure. 1) The primary structure comprises n amino acids. 2) The tertiary structure with atom arrangement in Euclidean space is presented, where each atom has a specific 3D coordinate. Amino acids have fixed backbone atoms (C_{α} , C, N, O) and side-chain atoms that vary depending on the residue types. GLU: Glutamic acid. Complete geometries can be obtained based on these coordinates. The sequence and structure provide different information types and data categories.

37 The 1D sequence and 3D structure of a protein provide different types of information, in detail, as
 38 shown in Figure 1, compared with the 1D sequential order and amino acids in peptide chains, the
 39 tertiary structure provides 3D coordinates of each atom in protein residues, which allow them to
 40 perform precise functions. Although a protein’s sequence determines its structure, various works
 41 have demonstrated the effectiveness of learning from either sequence or structure [33, 25]. However,
 42 rich constraints between the sequence and structure of a protein, which may be critical for protein
 43 tasks [4], have yet to be fully explored. Most protein sequence-structure modeling methods cannot
 44 deeply integrate the information behind sequence and structure for the reason that they tend to fuse
 45 representations together, extracted from sequence and structure encoders, respectively, by message
 46 passing mechanism [8] or by simple concatenation operations.

47 In this work, we aim to learn protein representations by deeply coupling the protein sequences and
 48 structures. Considering the relative positions of residues in the sequence and the spatial arrangement
 49 of atoms in the Euclidean space, the proposed CoupleNet constructs two categories of graphs for
 50 them, respectively. The complete representations are obtained at the amino acid and backbone
 51 levels on the two graphs, which are used as node and edge features to learn the final graph-level
 52 representations. Rather than concatenating sequence and structure representations, we take advantage
 53 of graph convolutions, performing node and edge convolutions simultaneously. The contributions of
 54 this paper are threefold:

- 55 • We propose a novel two-graph-based approach for representing the sequence and the 3D
 56 geometric structure of a protein, which is an effective way to guarantee completeness.
- 57 • We propose CoupleNet, a model that performs convolutions on nodes and edges of graphs
 58 to effectively integrate protein sequence and structure. This can better model the node-edge
 59 relationships and utilize the intrinsic associations between sequences and structures.
- 60 • Practically, the proposed model is verified by obtaining new state-of-the-art experimental
 61 results compared with current mainstream protein representation learning methods on a
 62 range of tasks, including protein fold classification, enzyme reaction classification, GO term
 63 prediction, domain prediction, and enzyme commission number prediction.

64 2 Related Work

65 **Protein Representation Learning** Protein representation learning has become an active and promis-
 66 ing direction in biology, which is essential to various downstream tasks in protein science. Because
 67 of the different levels of protein structures, existing methods mainly fall into three categories: protein
 68 LMs for sequences, structure models for geometry, and hybrid methods for both of them. As proteins
 69 are sequences of amino acids, considering their similarities with human languages, UniRep [1],
 70 UDSMProt [42] and SeqVec [23] use LSTM or its variants to learn sequence representations and
 71 long-range dependencies. TAPE [37] benchmarks a group of protein models, *e.g.*, 1D CNN, LSTM,
 72 and Transformer by various tasks. Elnaggar *et al.* [14] have trained six successful transformer variants

73 on billions of amino acid sequences, like ProtBert, and ProtT5. Similarly, ESM-series [39, 38, 33]
 74 employs a transformer architecture and a masked language modeling strategy to train robust represen-
 75 tations based on large-scale databases. Besides the protein sequence, as we have stated before, the
 76 3D geometric structure is vital to enhance protein representations. Most methods commonly seek to
 77 encode the spatial information of protein structures by convolutional neural networks (CNNs) [11],
 78 or graph neural networks [19, 2, 29]. For instance, SPROF [7] employs distance maps to predict
 79 protein sequence profiles, and IEConv [25] introduces a convolution operator to capture all relevant
 80 structural levels of a protein. GVP-GNN [29] designs the geometric vector perceptrons (GVP) for
 81 learning both scalar and vector features in an equivariant and invariant manner, Guo *et al.* [21]
 82 adopt SE(3)-invariant features as the model inputs and reconstruct gradients over 3D coordinates to
 83 avoid the usage of complicated SE(3)-equivariant models. ProNet [47] learns hierarchical protein
 84 representations at multiple tertiary structure levels of granularity. Moreover, CDConv [15] proposes
 85 continuous-discrete convolution using irregular and regular approaches to model the geometry and
 86 sequence structures. Some protein learning methods model the multi-level of structures at the same
 87 time [53, 6, 15], except for the primary structure and the tertiary structure, the second refers to the
 88 3D form of local segments of proteins (e.g., α -helix, β -strand), the quaternary is a protein multimer
 89 comprising multiple polypeptides, for example, PromptProtein [48] adopts a prompt-guided multi-task
 90 learning strategy for different protein structures with specific pre-training tasks. While previous
 91 works have attempted to combine protein sequence and structure, we focus on profoundly integrat-
 92 ing them by specifically designing two types of graphs respectively and conducting convolutions
 93 simultaneously to learn protein representations.

94 **Complete Message Passing Mechanism** ComENet [46] proposes rotation angles and spherical
 95 coordinates to fulfil the global completeness of 3D information on molecular graphs. By incorporating
 96 these designed geometric representations into the message passing scheme [18], the complete
 97 representation for a whole 3D graph is eventually yielded [47]. Unlike these methods, we couple
 98 sequence and structure via corresponding graphs and different geometric representations to obtain
 99 completeness representations.

100 3 Method

101 3.1 Preliminaries

102 **Notations** We represent a 3D graph as $G = (\mathcal{V}, \mathcal{E}, \mathcal{P})$, where $\mathcal{V} = \{v_i\}_{i=1, \dots, n}$ and $\mathcal{E} =$
 103 $\{\varepsilon_{ij}\}_{i, j=1, \dots, n}$ denote the vertex and edge sets with n nodes in total, respectively, and $\mathcal{P} =$
 104 $\{P_i\}_{i=1, \dots, n}$ is the set of position matrices, where $P_i \in \mathbb{R}^{k_i \times 3}$ represents the position matrix
 105 for node v_i . We treat each amino acid as a graph node for a protein, then k_i depends on the number
 106 of atoms in the i -th amino acid. The node feature matrix is $X = [\mathbf{x}_i]_{i=1, \dots, n}$, where $\mathbf{x}_i \in \mathbb{R}^{d_v}$ is
 107 the feature vector of node v_i . The edge feature matrix is $E = [\mathbf{e}_{ij}]_{i, j=1, \dots, n}$, where $\mathbf{e}_{ij} \in \mathbb{R}^{d_\varepsilon}$ is the
 108 feature vector of edge ε_{ij} . d_v and d_ε denote the dimensions of feature vectors \mathbf{x}_i and \mathbf{e}_{ij} .

109 **Invariance and Equivariance** We consider affine transformations that preserve the distance
 110 between any two points, *i.e.*, the isometric group SE(3) in the Euclidean space. This is called
 111 the symmetry group, and it turns out that SE(3) is the special Euclidean group that includes 3D
 112 translations and the 3D rotation group SO(3) [17, 12]. The matrix form of SE(3) is provided in
 113 Appendix A.1.

114 Given the function $f : \mathbb{R}^m \rightarrow \mathbb{R}^{m'}$, assuming the given symmetry group G acts on \mathbb{R}^m and $\mathbb{R}^{m'}$,
 115 then f is G-equivariant if,

$$f(T_g \mathbf{x}) = S_g f(\mathbf{x}), \forall \mathbf{x} \in \mathbb{R}^m, g \in G \quad (1)$$

116 where T_g and S_g are the transformations. For the SE(3) group, when $m' = 1$, the output of f is a
 117 scalar, we have

$$f(T_g \mathbf{x}) = f(\mathbf{x}), \forall \mathbf{x} \in \mathbb{R}^m, g \in G \quad (2)$$

118 thus f is SE(3)-invariant.

119 **Complete Geometric Representations** A geometric transformation $\mathcal{F}(\cdot)$ is complete if two 3D
 120 graphs $G^1 = (\mathcal{V}, \mathcal{E}, \mathcal{P}^1)$ and $G^2 = (\mathcal{V}, \mathcal{E}, \mathcal{P}^2)$, there exists $T_g \in \text{SE}(3)$ such that the representations

$$\mathcal{F}(G^1) = \mathcal{F}(G^2) \iff P_i^1 = T_g(P_i^2), \text{ for } i = 1, \dots, n \quad (3)$$

121 The operation T_g would not change the 3D conformation of a 3D graph [46]. Positions can generate
 122 geometric representations, which can also be recovered from them.

123 **Message Passing Paradigm** Message passing mechanism is mainly applied in graph convolutional
 124 networks (GCNs) [32], which follows an iterative scheme of updating node representations based on
 125 the feature aggregation from nearby nodes.

$$\begin{aligned} \mathbf{h}_i^{(0)} &= \text{BN}(\text{FC}(\mathbf{x}_i)), \\ \mathbf{u}_i^{(l)} &= f_{\text{Agg}}^{(l)}(\mathbf{h}_j^{(l-1)} | v_j \in \mathcal{N}(v_i)), \\ \mathbf{h}_i^{(l)} &= f_{\text{Update}}^{(l)}(\mathbf{h}_j^{(l-1)}, \mathbf{u}_i^{(l)}) \end{aligned} \quad (4)$$

126 where $\text{FC}(\cdot)$ and $\text{BN}(\cdot)$ mean the linear transformation and batch normalization respectively. $\mathcal{N}(v_i)$
 127 denotes the neighbours of node v_i . $f_{\text{Agg}}^{(l)}$ and $f_{\text{Update}}^{(l)}$ are aggregation and transformation functions at
 128 the l -th layer, which are permutation invariant and equivariant of node representations.

129 3.2 Sequence-Structure Graph Construction

130 Specifically, we represent each amino acid as a node,
 131 considering the residue types and their positions $i =$
 132 $1, 2, \dots, n$ (See Figure 1) in the sequence, we define
 133 the sequential graph primarily on the sequence,
 134 if $\|i - j\| < l$, the edge ε_{ij} exists, where l is a hyper-
 135 parameter. Besides the sequential graph, we predefine a
 136 radius r , and build the radius graph, and there exists an
 137 edge between node v_i and v_j if $\|P_{i, C_\alpha} - P_{j, C_\alpha}\| < r$,
 138 where P_{i, C_α} denotes the 3D position of C_α in the i -th
 139 residue.

140 Firstly, we design a base approach called CoupleNet_{aa}
 141 that only uses the C_α positions of the structures. In-
 142 spired by Ingraham *et al.* [28], we construct a local
 143 coordinate system (LCS) for each residue, as shown in
 144 Figure 2.

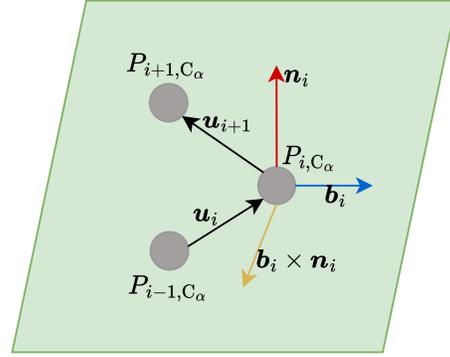


Figure 2: The local coordinate system.

$$\mathbf{Q}_i = [\mathbf{b}_i \quad \mathbf{n}_i \quad \mathbf{b}_i \times \mathbf{n}_i] \quad (5)$$

145 where $\mathbf{u}_i = \frac{P_{i, C_\alpha} - P_{i-1, C_\alpha}}{\|P_{i, C_\alpha} - P_{i-1, C_\alpha}\|}$, $\mathbf{b}_i = \frac{\mathbf{u}_i - \mathbf{u}_{i+1}}{\|\mathbf{u}_i - \mathbf{u}_{i+1}\|}$, $\mathbf{n}_i = \frac{\mathbf{u}_i \times \mathbf{u}_{i+1}}{\|\mathbf{u}_i \times \mathbf{u}_{i+1}\|}$. Then we can get the geometric
 146 representations at the amino acid level of a protein 3D graph,

$$\mathcal{F}(G)_{ij, aa} = (\|P_{i, C_\alpha} - P_{j, C_\alpha}\|, \mathbf{Q}_i^T \cdot \frac{P_{i, C_\alpha} - P_{j, C_\alpha}}{\|P_{i, C_\alpha} - P_{j, C_\alpha}\|}, \mathbf{Q}_i^T \cdot \mathbf{Q}_j) \quad (6)$$

147 where \cdot is the matrix multiplication, this implementation is $\text{SE}(3)$ -equivariant and obtains complete
 148 representations at the amino acid level; as if we have \mathbf{Q}_i , the LCS \mathbf{Q}_j can be easily obtained by
 149 $\mathcal{F}(G)_{ij, aa}$.

150 For a node v_i , the node features $\mathbf{x}_{i, aa}$ in the baseline approach is the concatenation of the one-hot
 151 embeddings of the amino acid types and the physicochemical properties of each residue, namely, a
 152 steric parameter, hydrophobicity, volume, polarizability, isoelectric point, helix probability and sheet
 153 probability [51, 22], which provide quantitative insights into the biochemical nature of each amino
 154 acid. And $\mathcal{F}(G)_{ij, aa}$ is set as edge features for CoupleNet_{aa}.

155 Secondly, we consider all backbone atoms C_α, C, N, O in CoupleNet. In detail, the peptide bond
 156 exhibits partial double-bond character due to resonance [20], indicating that the three non-hydrogen
 157 atoms comprising the bond (the carbonyl oxygen, carbonyl carbon, and amide nitrogen) are coplanar,

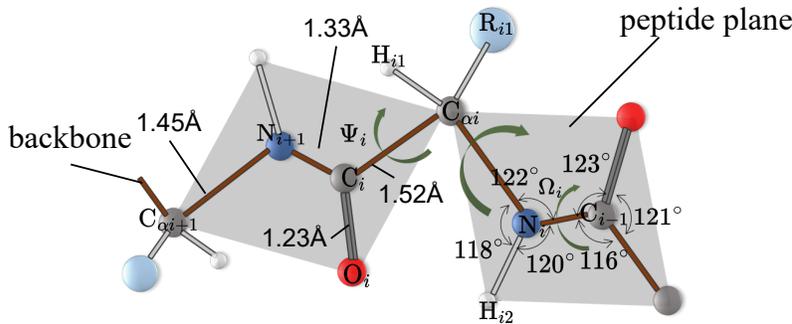


Figure 3: The polypeptide chain depicting the characteristic backbone bond lengths, angles, and torsion angles (Ψ_i, Φ_i, Ω_i). The planar peptide groups are denoted as shaded gray regions, indicating that the peptide plane differs from the geometric plane calculated based on the 3D positions.

158 as shown in Figure 3. There is some rotation about the connection. The $N_i - C_{\alpha i}$ and $C_{\alpha i} - C_i$
 159 bonds, are the two bonds in the basic repeating unit of the polypeptide backbone. These single bonds
 160 allow unrestricted rotation until sterically restricted by side chains [35, 45]. Since the coordinates of
 161 C_{α} can be obtained as we have the complete representations at the amino acid level, the coordinates
 162 of other backbone atoms based on these rigid bond lengths and angles are able to be determined with
 163 the remaining degree of the backbone torsion angles Φ_i, Ψ_i, Ω_i . The omega torsion angle around
 164 the C – N peptide bond is typically restricted to nearly 180° (trans) but can approach 0° (cis) in
 165 rare instances. Other than the bond lengths and angles presented in Figure 3, all the H bond lengths
 166 measure approximately 1 Å.

167 For the sequential graph, we compute the sine and cosine values of Φ_i, Ψ_i, Ω_i for each amino acid i ,
 168 and use them as another part of nodes features for node v_i .

$$\mathbf{x}_i = \mathbf{x}_{i,aa} \| ((\sin \wedge \cos)(\Phi_i, \Psi_i, \Omega_i)) \quad (7)$$

169 where $\|$ denotes concatenation. There is no isolated node for the designed graph, which means
 170 the backbone atoms can be determined one by one along the polypeptide chain based on the posi-
 171 tions of C_{α} and these three backbone dihedral angles. Therefore, the existing presentations
 172 $[\mathcal{F}(G)_{ij,aa}]_{i,j=1,\dots,n}$ and $[\mathbf{x}_i]_{i=1,\dots,n}$ are complete at the backbone level for the sequential graph.
 173

174 For the radius graph, we want to get the positions of back-
 175 bone atoms in any two amino acids i and j . Inspired by
 176 trRosetta [52], the relative rotation and distance are com-
 177 puted including the distance ($d_{ij,C_{\beta}}$), three dihedral angles
 178 ($\omega_{ij}, \theta_{ij}, \theta_{ji}$) and two planar angles ($\varphi_{ij}, \varphi_{ji}$), as shown in
 179 Figure 4, where $d_{ij,C_{\beta}} = d_{ji,C_{\beta}}, \omega_{ij} = \omega_{ji}$, but θ and φ
 180 values depend on the order of residues. These interresidue
 181 geometries define the relative locations of the backbone
 182 atoms of two residues in all their details [52], because the
 183 torsion angles of $N_i - C_{\alpha i}$ and $C_{\alpha i} - C_i$ do not influ-
 184 ence their positions. Therefore, these six geometries are
 185 complete for amino acids at the backbone level for the
 186 radius graph. The graph edges contain the relative spati-
 187 al information between any two neighboring amino acids
 188 $e_{ij} = \mathcal{F}(G)_{ij,aa} \| \mathcal{F}(G)_{ij,bb}$, where

$$\mathcal{F}(G)_{ij,bb} = (d_{ij,C_{\beta}}, (\sin \wedge \cos)(\omega_{ij}, \theta_{ij}, \varphi_{ij})) \quad (8)$$

189 The designed node and edge features, \mathbf{x}_i and e_{ij} , for the sequential and radius graphs, provide a
 190 new perspective to represent protein sequences and structures. Such integration can bring better
 191 performance for the following graph-based learning tasks.

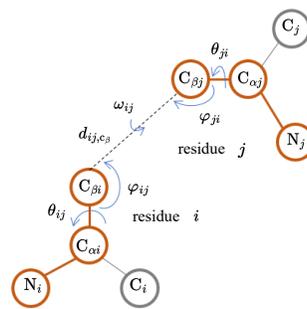


Figure 4: Interresidue geometries including angles and distances.

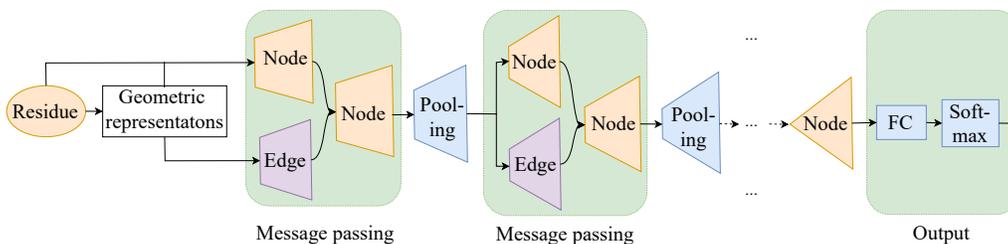


Figure 5: An illustration of CoupleNet.

192 3.3 Sequence-Structure Graph Convolution

193 Inspired by the message passing paradigm and continuous-discrete convolution [15], sequences
 194 and structures are encoded successfully together by convolutions. To deeply couple sequences and
 195 structures of proteins and encode them jointly, we employ convolution to embed them simultaneously,
 196 exploring their relationships to generate comprehensive and effective embeddings. Different from
 197 previous works, we innovatively construct two categories of graphs for sequence and structure and
 198 design various sequential and structural representations to achieve completeness on them at the amino
 199 acid and backbone levels. We then convolve node and edge features with the help of the message
 200 passing mechanism.

201 In order to implement convolution on nodes and edges simultaneously between sequence and
 202 we set ε_{ij} to exist if the following conditions are satisfied

$$\|i - j\| < l \quad \text{and} \quad \|P_{i,C\alpha} - P_{j,C\alpha}\| < r \quad (9)$$

203 The existing node and edge feature matrices (X, E) are complete representations of a protein 3D
 204 graph to reconstruct its backbone atom positions. Compared with the equation Eq. 4, the proposed
 205 CoupleNet first apply a $FC(\cdot)$ layer and a $BN(\cdot)$ layer to the node features to obtain the initial
 206 encoded representation. Then the $f_{\text{Agg}}^{(l)}$ is applied to gather neighboring features of nodes and
 207 edges by convolution, where $\sigma(\cdot)$ is the activation function. We use the dropout and add a residual
 208 connection from the previous layer as $f_{\text{Update}}^{(l)}$. For the consideration that the spatial arrangement and
 209 tight positioning of specific amino acids, which may be spaced widely apart on the linear polypeptide,
 210 are necessary for proteins to operate as intended [10], l is set to be a relatively large number, see
 211 Appendix A.2 for details.

$$\begin{aligned} \mathbf{h}_i^{(0)} &= \text{BN}(\text{FC}(\mathbf{x}_i)), \\ \mathbf{u}_i^{(l)} &= \sigma(\text{BN}(\sum_{v_j \in \mathcal{N}(v_i)} W e_{ij} \mathbf{h}_j^{(l-1)})), \\ \mathbf{h}_i^{(l)} &= \mathbf{h}_i^{(l-1)} + \text{Dropout}(\mathbf{u}_i^{(l)}) \end{aligned} \quad (10)$$

212 3.4 Model Architecture

213 Building upon the sequence-structure graph convolution, we build the CoupleNet, as shown in
 214 Figure 5. The inputs to the graph are the calculated sequential and structural representations (X, E).
 215 Following the existing protein graph models [15, 25, 47], our CoupleNet employs graph pooling
 216 layers to obtain deeply encoded, graph-level representations. After pooling, due to the decrease
 217 in nodes, we increase the predefined radius r to include more neighbors. The message passing
 218 mechanism only executes on nodes for the consideration of reducing model complexity. Another
 219 reason is that representations of sequences and structures have already been coupled by equation
 220 Eq. 4. A detailed description of the model architecture is provided in Appendix A.2.

Table 1: Accuracy (%) on fold classification and enzyme reaction classification. [*] means the results are taken from [15]. The best and suboptimal results are shown in bold and underline.

Input	Method	Fold Classification			Enzyme
		Fold	SuperFamily	Family	Reaction
Sequence	CNN [41]*	11.3	13.4	53.4	51.7
	ResNet [37]*	10.1	7.21	23.5	24.1
	LSTM [37]*	6.41	4.33	18.1	11.0
	Transformer [37]*	9.22	8.81	40.4	26.6
Structure	GCN [32]*	16.8	21.3	82.8	67.3
	GAT [44]*	12.4	16.5	72.7	55.6
	3DCNN_MQA [11]*	31.6	45.4	92.5	72.2
	IEConv (atom level) [25]*	45.0	69.7	98.9	87.2
Sequence-Structure	GraphQA [2]*	23.7	32.5	84.4	60.8
	GVP [29]*	16.0	22.5	83.8	65.5
	ProNet-Amino Acid [47]	51.5	69.9	99.0	86.0
	ProNet-Backbone [47]	52.7	70.3	99.3	86.4
	ProNet-All-Atom [47]	52.1	69.0	99.0	85.6
	IEConv (residue level) [25]*	47.6	70.2	99.2	87.2
	GearNet [53]	28.4	42.6	95.3	79.4
	GearNet-IEConv [53]	42.3	64.1	99.1	83.7
	GearNet-Edge [53]	44.0	66.7	99.1	86.6
	GearNet-Edge-IEConv [53]	48.3	70.3	99.5	85.3
	CDCConv [15]	<u>56.7</u>	<u>77.7</u>	<u>99.6</u>	<u>88.5</u>
CoupleNet (Proposed)	60.6	82.1	99.7	89.0	

221 4 Experiments

222 4.1 Datasets and Settings

223 The models are trained with the Adam optimizer [31] using the PyTorch and PyTorch Geometric
 224 libraries. Detailed descriptions of the datasets and experimental settings are provided in Appendix A.3.
 225 Following the tasks in IEConv [25], GearNet [53] and CDCConv [15], here, we evaluate the CoupleNet
 226 on four protein tasks: protein fold classification, enzyme reaction classification, GO term prediction
 227 and enzyme commission (EC) number prediction.

228 **Fold Classification** Protein fold is to predict the fold class label given a protein, which is crucial for
 229 understanding how protein structure and protein evolution interact [26]. In total, this dataset contains
 230 16, 712 proteins with 1, 195 fold classes. There are three test sets available, Fold: Training excludes
 231 proteins from the same superfamily. Superfamily: Training does not include proteins from the same
 232 family. Family: Proteins from the same family are included in the training.

233 **Enzyme Reaction Classification** Reaction categorization aims to predict a protein’s class of
 234 enzyme-catalyzed reactions, according to all four levels of the EC number [49, 36]. Following the
 235 setting in [25], this dataset has 37, 248 proteins from 384 four-level EC numbers [5].

236 **GO Term Prediction** The goal of GO term prediction is to foretell whether a protein is related
 237 to a certain GO term. Following [19], these proteins are organized into three ontologies: molecular
 238 function (MF), biological process (BP), and cellular component (CC), which are hierarchically
 239 connected, functional classes. MF describes activities that occur at the molecular level, BP represents
 240 the larger processes, and CC describes the parts of a cell or its extracellular environment [3].

241 **EC Number Prediction** This task seeks to predict the 538 EC numbers from the third level and
 242 fourth levels of different proteins [19], which describe their catalysis of biochemical reactions.

Table 2: F_{\max} on GO term and EC number prediction. [*] means the results are taken from [15]. The best and suboptimal results are shown in bold and underline.

Category	Method	GO-BP	GO-MF	GO-CC	EC
Sequence	CNN [41]*	0.244	0.354	0.287	0.545
	ResNet [37]*	0.280	0.405	0.304	0.605
	LSTM [37]*	0.225	0.321	0.283	0.425
	Transformer [37]*	0.264	0.211	0.405	0.238
Structure	GCN [32]*	0.252	0.195	0.329	0.320
	GAT [44]*	0.284	0.317	0.385	0.368
	3DCNN_MQA [11]*	0.240	0.147	0.305	0.077
Sequence-Structure	GraphQA [2]*	0.308	0.329	0.413	0.509
	GVP [29]*	0.326	0.426	0.420	0.489
	IEConv (residue level) [25]*	0.421	0.624	0.431	-
	GearNet [53]	0.356	0.503	0.414	0.730
	GearNet-IEConv [53]	0.381	0.563	0.422	0.800
	GearNet-Edge [53]	0.403	0.580	0.450	0.810
	GearNet-Edge-IEConv [53]	0.400	0.581	0.430	0.810
	CDCConv [15]	<u>0.453</u>	<u>0.654</u>	<u>0.479</u>	<u>0.820</u>
CoupleNet (Proposed)	0.467	0.669	0.494	0.866	

243 4.2 Baselines

244 We compare our proposed method with existing protein representation learning methods, which are
 245 classified into three categories based on their inputs, which could be a sequence (amino acid sequence),
 246 3D structure or both sequence and structure. 1) Sequence-based encoders, including CNN [41],
 247 ResNet [37], LSTM [37] and Transformer [37]. 2) Structure-based methods (GCN [32], GAT [44],
 248 3DCNN_MQA [11], IEConv (atom level) [25]). 3) Sequence-structure based models, *e.g.*, GVP [29],
 249 ProNet [47], GearNet [53], CDCConv [15], *etc.* GearNet-IEConv and GearNetEdge-IEConv [53] add
 250 the IEConv layer based on GearNet, which is found important in fold classification.

251 4.3 Results of Fold and Reaction Classification.

252 Table 1 provides the comparisons on the fold and enzyme reaction classification. The results are
 253 reported in terms of accuracy (%) for these two tasks. From this table, we can see that the proposed
 254 model CoupleNet achieves the best performance across all four test sets on the fold and enzyme
 255 reaction classification compared with recent state-of-the-art methods. Especially on the Fold and
 256 SuperFamily test sets, CoupleNet improves the results by about 4%, showing that CoupleNet is
 257 proficient at learning the mapping between protein sequences, structures and functions. Moreover,
 258 CDCConv [15] ranks second among these methods, both CDCConv and our method are implemented
 259 by sequence-structure convolution. This phenomenon illustrates that deeply coupling sequences
 260 and structures of proteins is conducive to learning better protein embeddings. And our proposed
 261 CoupleNet model utilizes complete geometric representations and the specially designed message
 262 passing mechanism, achieving new state-of-the-art results.

263 4.4 Results of GO Term and EC Prediction

264 We follow the split method in [19, 53] to guarantee that the test set only comprises PDB chains with
 265 sequence identity no higher than 95% to the training set for GO term and EC number prediction.
 266 Table 2 compares different protein modeling methods on GO term prediction and EC number
 267 prediction. The results are reported in terms of F_{\max} , which considers both precision and recall for
 268 evaluation, the equation of F_{\max} is provided in Appendix A.4. The proposed model, CoupleNet
 269 yields the highest F_{\max} across these four test sets of two tasks, outperforming other state-of-the-art
 270 models. This indicates CoupleNet can effectively predict the functions, locations, and enzymatic
 271 activities of proteins. These results once again illustrate the importance of deeply coupled sequences

Table 3: Ablation of our proposed method

Method	Fold Classification			Enzyme Reaction	GO			EC
	Fold	Superfamily	Family		BP	MF	CC	
CoupleNet	60.6	82.1	99.7	89.0	0.467	0.669	0.494	0.866
CoupleNet _{aa}	57.8	78.7	99.6	88.6	0.458	0.660	0.484	0.851
w/o Φ, Ψ, Ω	60.3	81.3	99.6	88.7	0.463	0.666	0.490	0.862
w/o $d, \omega, \theta, \varphi$	60.4	81.5	99.7	88.9	0.461	0.666	0.488	0.864

272 and structures. The improvements of CoupleNet over the suboptimal CDCConv [15] model indicate
273 the advanced modeling power of CoupleNet.

274 We employ different cutoff splits following [19,
275 15], which means that the proteins in the test set
276 are divided into groups that have, respectively,
277 30%, 40%, 50%, 70%, and 95% similarity to
278 the training set for GO term and EC number pre-
279 diction, as shown in Figure 6 and Appendix A.5.
280 From Figure 6, we can see that our proposed
281 model CoupleNet achieves the highest F_{\max}
282 scores across all cutoffs, especially when the
283 cutoffs are at 30% to 50%. There is a larger mar-
284 gin compared with GearNet, GearNet-Edge [53]
285 and CDCConv [15]. This demonstrates that Cou-
286 pleNet, which utilizes complete geometric repre-
287 sentations, is more robust, especially when there
288 is a low similarity between the training and test
289 sets.

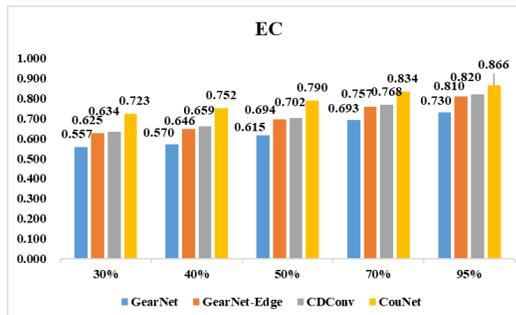


Figure 6: F_{\max} on EC number prediction under different cutoffs.

290 4.5 Ablation Study

291 Table 3 presents an ablation study of the proposed CoupleNet model on the four protein tasks. We
292 examined the impact of removing the backbone torsion angles (w/o Φ, Ψ, Ω) and removing the
293 interresidue geometric structure representations (w/o $d, \omega, \theta, \varphi$). The former is designed for the
294 sequential graph, and the latter is for the radius graph to achieve completeness at the protein backbone
295 level. However, we combine the two types of graphs together to enhance the relationships between
296 sequence and structure. From Table 3, we can also find that these complete geometries provide
297 complementary information to amino acid position features, with one of their removals leading to
298 minor performance drops for the reason that they both provide complete geometries from different
299 perspectives. Removing Φ, Ψ, Ω causes larger performance degradation compared with removing
300 $d, \omega, \theta, \varphi$. Such comparisons indicate that the backbone dihedral angles may have more effects
301 on learning protein representations in these experimental settings. Compared with CoupleNet_{aa},
302 CoupleNet achieves significant improvements on the four tasks, demonstrating the importance of
303 complete structural representations at the backbone level in learning protein embeddings.

304 5 Conclusions and Limitations

305 In this work, we propose CoupleNet, a novel protein representation learning method that deeply fuses
306 protein sequences and multi-level structures by conducting convolution on graph nodes and edges
307 simultaneously. We design the sequential graph and the radius graph, achieving completeness on
308 them at different protein structure levels. Our approach achieves new state-of-the-art results on the
309 protein tasks, which demonstrates the superiority of our proposed method. A limitation is that the
310 detailed inter-relationships between sequence and structures remain to be explored and uncovered.
311 We leave such research for future work.

312 While our model can enable advanced protein analyses and provide effective representations, there
313 may exist broader impacts and harmful activities. The representations could potentially be misused,
314 *e.g.*, for designing harmful molecules or proteins.

References

- 316 [1] Ethan C. Alley et al. “Unified rational protein engineering with sequence-based deep represen-
317 tation learning”. In: *Nature Methods* (2019).
- 318 [2] Federico Baldassarre et al. “GraphQA: protein model quality assessment using graph convolu-
319 tional networks.” In: *Bioinformatics* (2020).
- 320 [3] Alex Bateman. “UniProt: A worldwide hub of protein knowledge”. In: *Nucleic Acids Research*
321 (2019).
- 322 [4] Tristan Bepler and Bonnie Berger. “Learning the protein language: Evolution, structure, and
323 function”. In: *Cell systems* 12.6 (2021), pp. 654–669.
- 324 [5] Helen M Berman et al. “The protein data bank”. In: *Nucleic acids research* 28.1 (2000),
325 pp. 235–242.
- 326 [6] Can Chen et al. “Structure-aware protein self-supervised learning”. In: *Bioinformatics* 39.4
327 (2023), btad189.
- 328 [7] Sheng Chen et al. “To Improve Protein Sequence Profile Prediction through Image Captioning
329 on Pairwise Residue Distance Map”. In: *Journal of Chemical Information and Modeling*
330 (2020).
- 331 [8] Yihong Chen et al. “Refactor gnns: Revisiting factorisation-based models from a message-
332 passing perspective”. In: *Advances in Neural Information Processing Systems* 35 (2022),
333 pp. 16138–16150.
- 334 [9] G Marius Clore and Angela M Gronenborn. “NMR structure determination of proteins and
335 protein complexes larger than 20 kDa”. In: *Current opinion in chemical biology* 2.5 (1998),
336 pp. 564–570.
- 337 [10] Srinivasan Damodaran. “Amino acids, peptides and proteins”. In: *Fennema’s food chemistry* 4
338 (2008), pp. 425–439.
- 339 [11] Georgy Derevyanko et al. “Deep convolutional networks for quality assessment of protein
340 folds”. In: *Bioinformatics* 34.23 (2018), pp. 4046–4053.
- 341 [12] Weitao Du et al. “SE (3) Equivariant Graph Neural Networks with Complete Local Frames”.
342 In: *International Conference on Machine Learning*. PMLR. 2022, pp. 5583–5608.
- 343 [13] Arun Kumar Dubey and Vanita Jain. “Comparative study of convolution neural network’s relu
344 and leaky-relu activation functions”. In: *Applications of Computing, Automation and Wireless
345 Systems in Electrical Engineering: Proceedings of MARC 2018*. Springer. 2019, pp. 873–880.
- 346 [14] Ahmed Elnaggar et al. “ProfTrans: Towards Cracking the Language of Lifes Code Through
347 Self-Supervised Deep Learning and High Performance Computing”. In: *IEEE Transactions on
348 Pattern Analysis and Machine Intelligence* (2021).
- 349 [15] Hehe Fan et al. “Continuous-Discrete Convolution for Geometry-Sequence Modeling in
350 Proteins”. In: *The Eleventh International Conference on Learning Representations*. 2023.
- 351 [16] Noelia Ferruz and Birte Höcker. “Controllable protein design with language models”. In:
352 *Nature Machine Intelligence* (2022), pp. 1–12.
- 353 [17] Fabian Fuchs et al. “Se (3)-transformers: 3d roto-translation equivariant attention networks”.
354 In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 1970–1981.
- 355 [18] Justin Gilmer et al. “Neural message passing for quantum chemistry”. In: *International
356 conference on machine learning*. PMLR. 2017, pp. 1263–1272.
- 357 [19] Vladimir Gligorijević et al. “Structure-based protein function prediction using graph convolu-
358 tional networks”. In: *Nature communications* 12.1 (2021), p. 3168.
- 359 [20] ER HARD GROSS and JOHANNES MEIENHOFER. “The Peptide Bond”. In: *Major Methods
360 of Peptide Bond Formation: The Peptides Analysis, Synthesis, Biology, Vol. 1* 1 (2014), p. 1.
- 361 [21] Yuzhi Guo et al. “Self-supervised pre-training for protein embeddings using tertiary structures”.
362 In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 36. 6. 2022, pp. 6801–
363 6809.
- 364 [22] Jack Hanson et al. “Improving prediction of protein secondary structure, backbone angles,
365 solvent accessibility and contact numbers by using predicted contact maps and an ensemble
366 of recurrent and residual convolutional neural networks”. In: *Bioinformatics* 35.14 (2019),
367 pp. 2403–2410.
- 368 [23] Michael Heinzinger et al. “Modeling the language of life – Deep Learning Protein Sequences”.
369 In: *bioRxiv* (2019).

- 370 [24] Pedro Hermosilla and Timo Ropinski. “Contrastive representation learning for 3d protein
371 structures”. In: *arXiv preprint arXiv:2205.15675* (2022).
- 372 [25] Pedro Hermosilla et al. “Intrinsic-Extrinsic Convolution and Pooling for Learning on 3D
373 Protein Structures”. In: *International Conference on Learning Representations* (2021).
- 374 [26] Jie Hou, Badri Adhikari, and Jianlin Cheng. “DeepSF: deep convolutional neural network for
375 mapping protein sequences to folds”. In: *Bioinformatics* 34.8 (2018), pp. 1295–1303.
- 376 [27] Bozhen Hu et al. “Protein Language Models and Structure Prediction: Connection and Pro-
377 gression”. In: *arXiv preprint arXiv:2211.16742* (2022).
- 378 [28] John Ingraham et al. “Generative models for graph-based protein design”. In: *Advances in
379 neural information processing systems* 32 (2019).
- 380 [29] Bowen Jing et al. “Learning from Protein Structure with Geometric Vector Perceptrons”. In:
381 *Learning* (2020).
- 382 [30] John Jumper et al. “Highly accurate protein structure prediction with AlphaFold”. In: *Nature*
383 596.7873 (2021), pp. 583–589.
- 384 [31] Diederik P Kingma and Jimmy Ba. “Adam: A method for stochastic optimization”. In: *arXiv
385 preprint arXiv:1412.6980* (2014).
- 386 [32] Thomas N Kipf and Max Welling. “Semi-supervised classification with graph convolutional
387 networks”. In: *arXiv preprint arXiv:1609.02907* (2016).
- 388 [33] Zeming Lin et al. “Language models of protein sequences at the scale of evolution enable
389 accurate structure prediction”. In: *BioRxiv* (2022).
- 390 [34] Bin Ma. “Novor: real-time peptide de novo sequencing software.” In: *Journal of the American
391 Society for Mass Spectrometry* (2015).
- 392 [35] David L Nelson, Albert L Lehninger, and Michael M Cox. *Lehninger principles of biochemistry*.
393 Macmillan, 2008.
- 394 [36] Marina V Omelchenko et al. “Non-homologous isofunctional enzymes: a systematic analysis
395 of alternative solutions in enzyme evolution”. In: *Biology direct* 5 (2010), pp. 1–20.
- 396 [37] Roshan Rao et al. “Evaluating protein transfer learning with TAPE”. In: *Advances in neural
397 information processing systems* 32 (2019).
- 398 [38] Roshan M Rao et al. “MSA transformer”. In: *International Conference on Machine Learning*.
399 PMLR. 2021, pp. 8844–8856.
- 400 [39] Alexander Rives et al. “Biological Structure and Function Emerge from Scaling Unsupervised
401 Learning to 250 Million Protein Sequences”. In: *Proceedings of the National Academy of
402 Sciences of the United States of America* (2019).
- 403 [40] Andrew W Senior et al. “Improved protein structure prediction using potentials from deep
404 learning”. In: *Nature* 577.7792 (2020), pp. 706–710.
- 405 [41] Amir Shانهsazzadeh, David Belanger, and David Dohan. “Is transfer learning necessary for
406 protein landscape prediction?” In: *arXiv preprint arXiv:2011.03443* (2020).
- 407 [42] Nils Strodthoff et al. “UDSMProt: universal deep sequence models for protein classification”.
408 In: *Bioinformatics* 36.8 (Jan. 2020), pp. 2401–2409. ISSN: 1367-4803. DOI: 10 . 1093 /
409 bioinformatics/btaa003.
- 410 [43] Mihaly Varadi et al. “AlphaFold Protein Structure Database: massively expanding the structural
411 coverage of protein-sequence space with high-accuracy models”. In: *Nucleic acids research*
412 50.D1 (2022), pp. D439–D444.
- 413 [44] Petar Velickovic et al. “Graph attention networks”. In: *stat* 1050.20 (2017), pp. 10–48550.
- 414 [45] K Peter C Vollhardt and Neil E Schore. *Organic chemistry: structure and function*. Macmillan,
415 2003.
- 416 [46] Limei Wang et al. “ComENet: Towards Complete and Efficient Message Passing for 3D
417 Molecular Graphs”. In: *arXiv preprint arXiv:2206.08515* (2022).
- 418 [47] Limei Wang et al. “Learning Hierarchical Protein Representations via Complete 3D Graph
419 Networks”. In: *The Eleventh International Conference on Learning Representations*. 2023.
- 420 [48] Zeyuan Wang et al. “Multi-level Protein Structure Pre-training via Prompt Learning”. In: *The
421 Eleventh International Conference on Learning Representations*.
- 422 [49] Edwin C Webb et al. *Enzyme nomenclature 1992. Recommendations of the Nomenclature Com-
423 mittee of the International Union of Biochemistry and Molecular Biology on the Nomenclature
424 and Classification of Enzymes*. Ed. 6. Academic Press, 1992.

- 425 [50] Fang Wu, Dragomir Radev, and Jinbo Xu. “When Geometric Deep Learning Meets Pretrained
426 Protein Language Models”. In: *bioRxiv* (2023), pp. 2023–01.
- 427 [51] Gang Xu, Qinghua Wang, and Jianpeng Ma. “OPUS-Rota4: a gradient-based protein side-chain
428 modeling framework assisted by deep learning-based predictors”. In: *Briefings in Bioinformat-
429 ics* 23.1 (2022), bbab529.
- 430 [52] Jianyi Yang et al. “Improved protein structure prediction using predicted interresidue ori-
431 entations”. In: *Proceedings of the National Academy of Sciences* 117.3 (2020), pp. 1496–
432 1503.
- 433 [53] Zuobai Zhang et al. “Protein representation learning by geometric structure pretraining”. In:
434 *International Conference on Learning Representations*. 2023.