# RETRIEVAL-AUGMENTED LANGUAGE MODEL FOR KNOWLEDGE-AWARE PROTEIN ENCODING

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Protein language models often struggle to capture the biological functions encoded within protein sequences due to their lack of factual knowledge (e.g., gene descriptions of proteins). Existing solutions leverage protein knowledge graphs (PKGs), using knowledge as auxiliary encoding objectives. However, none of them explored the direct injection of correlated knowledge into protein language models, and task-oriented knowledge integration during fine-tuning, making them suffer from insufficient knowledge exploitation and catastrophic forgetting of pretrained knowledge. The root cause is that they fail to align PKGs with downstream tasks, forcing their knowledge modeling to adapt to the knowledge-isolated nature of these tasks. To tackle these limitations, we propose a novel knowledge retriever that can accurately predict gene descriptions for new proteins in downstream tasks and thus align them with PKGs. On this basis, we propose Knowledge-aware retrieval-augmented protein language model (Kara), achieving the first unified and direct integration of PKGs and protein language models. Using the knowledge retriever, both the pre-training and fine-tuning stages can incorporate knowledge through a unified modeling process, where contextualized virtual tokens enable token-level integration of high-order knowledge. Moreover, structure-based regularization is introduced to inject function similarity into protein representations, and unify the pre-training and fine-tuning optimization objectives. Experimental results show that Kara consistently outperforms existing knowledge-enhanced models in 6 representative tasks, achieving on average 5.1% improvements.

031 032

033

004

010 011

012

013

014

015

016

017

018

019

020

021

024

025

026

027

028

029

#### 1 INTRODUCTION

Proteins are essential for understanding biological processes and recent advances in artificial intelligence led to growing interest in learning generalized vector representations of proteins (Hu et al., 2024). By viewing amino acids as language tokens, protein language models (PLMs) such as ESM (Lin et al., 2023), ProteinBert (Brandes et al., 2022), and ProtBert (Ahmed et al., 2020) have proven highly valuable in various application tasks such as drug discovery (Hoang et al., 2024) and function prediction (Xu et al., 2024; Shaw et al., 2024). However, as pointed out by Kalifa et al. (2024); Zhou et al. (2023); Zhang et al. (2022), lacking factual knowledge (*e.g.*, gene descriptions) makes them struggle to capture intricate biological function encoded within protein sequences.

Existing solutions leverage protein knowledge graphs (PKGs) that describe the relationships between proteins and gene ontology (GO) entities with biological relations (Chen et al., 2023b). These
models use protein sequences and associated GO annotations as complementary encoding objectives
to infuse knowledge information. For example, OntoProtein (Zhang et al., 2022) uses knowledge
embedding objective (*i.e.*, TransE (Bordes et al., 2013)) to optimize the alignment between the protein representations and associated GO entity representations. KeAP (Zhou et al., 2023) integrates
GO entity representations into the masked token prediction of protein sequences through a crossattention mechanism. Despite their effectiveness, unfortunately, they still have several limitations.

Limitations. 1) Implicitly embed knowledge information. Existing methods use knowledge only
 as encoding objectives to supervise the pre-training of the model, assuming that knowledge infor mation can be well embedded within model parameters. However, as highlighted by Kandpal et al.
 (2023), LMs often struggle to precisely embed knowledge, particularly long-tail knowledge. Storing knowledge within model parameters also makes them unable to adapt to knowledge graph updates

054 (e.g., adding new knowledge), which further diminishes their usability. 2) Overlook the structure 055 information. Existing methods treat each knowledge triplet (*i.e.*, (protein, relation, GO)) inde-056 pendently. However, the neighboring GO entities of a protein are often correlated, and the high-order 057 connections between proteins (e.g., proteins linked to a GO entity through similar relations) can 058 provide additional insights into their functional similarities. Ignoring the structural relevance makes existing methods fail to fully exploit knowledge information within PKGs. 3) Inconsistent knowledge modeling. Existing methods incorporate knowledge modeling during pre-training but ignore 060 it during fine-tuning, leading to inconsistent optimization objectives between these two stages. This 061 inconsistency can cause the knowledge learned during pre-training to be catastrophically forgotten 062 during fine-tuning (Lee et al., 2020), making it difficult to transfer to downstream tasks. Overall, the 063 root cause of these limitations is that proteins in downstream tasks often fall outside the PKG, re-064 straining the use of knowledge during fine-tuning. Existing methods fail to align knowledge graphs 065 with downstream tasks, forcing their knowledge modeling to adapt to the knowledge-isolated nature 066 of these tasks (e.g., knowledge cannot be directly used as part of the input for protein encoding). 067

**Proposed Work.** To tackle these limitations, we propose 068 Kara, a Knowledge-aware retrieval-augmented protein 069 language model, achieving the first unified and direct integration of PKGs and protein language models. As 071 the core of Kara, we propose a knowledge retriever that can accurately predict potential gene descriptions for new 073 proteins and thus align them with PKGs. This align-074 ment allows the pre-training and fine-tuning stages of 075 Kara to be enhanced through a unified knowledge modeling process, and seamlessly adapt to knowledge updates. 076 By employing contextualized virtual tokens, we achieve 077 token-level information fusion between protein sequence and knowledge. Specifically, we categorized the virtual 079 tokens into knowledge tokens and structure tokens, enabling the direct injection of both knowledge informa-081 tion and high-order structure information into protein rep-082 resentations. To unify the optimization objectives, we 083 incorporate structure-based regularization into both two



Figure 1: Performance in downstream tasks. S-, M-, and L-Contact are the short-range, medium-range, and long-range contact prediction. PPI is the protein-protein interaction prediction.

stages, injecting function similarities into protein representations and helping the pre-trained knowl edge to be effectively transferred to downstream tasks.

As shown in Figure 1, experiments in 6 representative downstream tasks demonstrate the effectiveness of Kara. It consistently outperforms powerful baselines (*i.e.*, KeAP and ESM-2) across all the tasks. For instance, Kara exceeds the state-of-the-art knowledge-enhanced model KeAP by 11.6% in the long-range contact prediction and by 10.3% in the protein homology detection, highlighting Kara as a better paradigm for integrating protein knowledge graphs into protein language models.

092

#### 2 PRELIMINARY

094 **Protein Knowledge Graphs.** A protein knowledge graph (PKG) is  $G = \{V_p, V_{qo}, R, F\}$ , where  $V_p$ is the protein set and  $V_{ao}$  is the gene ontology (GO) entity set. R is the set of relations among proteins 096 and GO entities. The knowledge set F consists of two kinds of triplets: (p, r, g) which describes the properties of proteins, and  $(g_1, r, g_2)$  which describes the relationships between GO entities. Each 098 protein  $p \in V_p$  has an amino acid sequence s. Each GO entity  $q \in V_{qq}$  includes a text description  $t_q$ explaining the gene's function. Similarly, each relation  $r \in R$  comes with a text description  $t_r$ . We 099 first generate pre-trained embeddings of items in PKG and store them in vector databases for further 100 usage. Specifically, relation r and GO entity q are encoded based on their text descriptions using a 101 frozen PubMedBERT model (Gu et al., 2021), resulting in relation embedding r and GO embedding 102 g. Protein p is encoded based on its amino acid sequence via a frozen ProtBert model (Ahmed et al., 103 2020), resulting in protein embedding p. These stored embeddings will be further used to construct 104 virtual tokens in Kara. Following previous works, we use the ProteinKG25 knowledge graph (Zhang 105 et al., 2022). Detailed introduction of ProteinKG25 can be found in Appendix A. 106

- **Problem Formulation.** Given a PKG G, we aim to pre-train a knowledge-aware protein language model f so that for each protein with amino acid sequence s, we can generate its knowledge-
  - 2



Figure 2: **Overall architecture**. During pre-training, Kara directly integrates knowledge information via contextualized virtual tokens and structure-based regularization. During fine-tuning, the knowledge information can be similarly integrated into protein representations through a knowledge retriever, which can align new proteins in downstream tasks with the protein knowledge graph.

integrated vector representation as  $\tilde{\mathbf{p}} = f(G, s)$ . In Kara, f consists of a protein encoder, a knowledge projector, a protein projector, and a knowledge retriever. We use the ProtBert model (Ahmed et al., 2020) as the backbone of the protein encoder, which is the same as previous works (Zhou et al., 2023) for a fair comparison. By fine-tuning f on task-specific data, we further verify its capabilities to generalize pre-trained knowledge to downstream tasks (*e.g.*, protein homology detection).

#### 3 METHODOLOGIES

As shown in Figure 2, with a knowledge retriever to align new proteins with the protein knowl-133 edge graph, Kara can uniformly integrate knowledge information during both the pre-training and 134 fine-tuning stages <sup>1</sup>. Specifically, the contextualized virtual tokens allow Kara to directly inject the 135 associated knowledge information and high-order structure information of a protein into its repre-136 sentations. During pre-training, masked language modeling (MLM) helps the protein encoder learn 137 to fuse the information of protein sequences and structured knowledge at the token level. During 138 fine-tuning, downstream task modeling helps the protein encoder learn to extract task-specific use-139 ful knowledge from PKGs via virtual tokens. Additionally, based on the high-order connectivity 140 between proteins, structure-based regularization is incorporated during the two stages to unify their 141 optimization objectives and inject function similarities into protein representations. We detail each 142 part of Kara in the following and summarize important notations used in this paper in Table 1.

143 144

145

120

121

122

123 124

125

126

127

128

129 130

131 132

#### 3.1 PRE-TRAINING STAGE

#### 146 3.1.1 CONTEXTUALIZED VIRTUAL TOKENS

147 Existing protein language models struggle to encode knowledge information since 1) knowledge 148 in PKG is interconnected, providing the context of proteins based on the graph structure, but lan-149 guage models are only designed to encode sequential data, limiting their ability to capture graph 150 information; and 2) PKGs contain multi-modal information (e.g., amino acid sequences and GO 151 text descriptions), and protein language models can only encode amino acid sequences, failing to achieve effective multi-modal information fusion. As shown in Figure 2 A.1, we tackle the above 152 challenges by introducing contextualized virtual tokens. By summarizing the associated knowledge 153 of a protein as knowledge virtual tokens and summarizing its high-order structure as structure virtual 154 tokens, Kara can directly inject the knowledge and graph information into protein representations. 155 These virtual tokens are then concatenated with the amino acid token sequences as the knowledge 156 context, so that each amino acid can query them to integrate helpful knowledge information, en-157 abling effective token-level multi-modal information fusion. Specifically, for each protein  $p_i \in V_p$ , 158 we extract its one-hop GO entities with relations  $\mathcal{N}_1(p_i) = \{(r_i, g_i) | (p_i, r_i, g_i) \in F\}$  as its knowl-159 edge, and use its two-hop proteins  $\mathcal{N}_2(p_i) = \{p_j | (p_j, r_i, g_i) \in F; (r_i, g_i) \in \mathcal{N}_1(p_i)\}$  as its structure

160 161

<sup>&</sup>lt;sup>1</sup>Note that the knowledge retriever is only used during fine-tuning, ensuring that no data from downstream tasks can be leaked into the pre-training stage.

Table 1: Important notations and descriptions.						
Notation	Description					
G	A protein knowledge graph.					
$V_p, V_{go}, R$	Protein set, GO entity set, and relation set in G.					
$\overline{F}$	Set of triplets ( <i>i.e.</i> , knowledge) in G.					
$p_i, r_j, g_k$	A protein, a GO entity, and a relation.					
$s_i,s_i^m$	The amino acid sequence of protein $p_i$ , and each amino acid in $s_i$ .					
$\mathbf{p}_i, \mathbf{r}_j, \mathbf{g}_k$	Stored pre-trained embeddings of protein $p_i$ , relation $r_j$ , and GO entity $g_k$ (see Section 2).					
$\mathbf{v}_i^k, \mathbf{v}_i^p$	Knowledge virtual token and structure virtual token of protein $p_i$ .					
$\mathbf{S}_i, \mathbf{S}_i^L$	Input embedding sequence of protein $p_i$ , embedding sequence at the L-th layer.					
$ ilde{\mathbf{p}}_i$	Encoded embedding of protein $p_i$ by Kara.					
$\mathbf{g}_m^{go}, \mathbf{g}_m^{prot}$	Neighboring GO entity embedding and neighboring protein embedding of GO entity $g_m$ .					
$\mathbf{q}_n$	Query embedding corresponds to new protein $p_n$ .					
$\mathbf{ ilde{r}}_m$	Query embedding corresponds to relation $r_m$ .					
$\mathbf{c}_m$	Candidate embedding corresponds to GO entity $g_m$ .					
$\mathbb{S}(\cdot)$	Score function.					
$\mathtt{MLP}(\cdot)$	Trainable multi-layer perceptron.					
$N_1(p_i)$	One-hop GO entities with relations of protein $p_i$					
$N_2(p_i)$	Two-hop connected proteins of protein $p_i$ .					
$N_{go}(g_m)$	One-hop neighboring GO entities of GO entity $g_m$ .					
$N_{prot}(g_m)$	One-hop neighboring proteins of GO entity $g_m$ .					
$\mathcal{E}(r_m)$	Candidate GO entity set corresponding to relation $r_m$ .					

context. The knowledge virtual token of protein  $p_i$  is then constructed as

$$\mathbf{v}_{i}^{k} = \frac{1}{|\mathcal{N}_{1}(p_{i})|} \sum_{(r_{i},g_{i})\in\mathcal{N}_{1}(p_{i})} \mathsf{MLP}_{knowledge}([\mathbf{r}_{i}:\mathbf{g}_{i}]), \tag{1}$$

where  $\mathbf{r}_i$  and  $\mathbf{g}_i$  are respectively the pre-trained embeddings of relation  $r_i$  and GO entity  $g_i$  (see Section 2). [:] is the concatenation operation. MLP<sub>knowledge</sub> is a trainable multi-layer perceptron used to project text-modal information into a uniform semantic space. Similarly, to incorporate the structure information of  $p_i$ , we construct its structure virtual token as

$$\mathbf{v}_{i}^{p} = \frac{1}{|\mathcal{N}_{2}(p_{i})|} \sum_{p_{j} \in \mathcal{N}_{2}(p_{i})} \text{MLP}_{structure}(\mathbf{p}_{j}),$$
(2)

where  $\mathbf{p}_j$  is the pre-trained embedding of protein  $p_j$ . MLP<sub>structure</sub> is another trainable multi-layer perceptron used to project the amino acid sequence-modal information. We then construct the input embedding sequence for the protein encoder by concatenating virtual tokens with amino acid tokens. Given the amino acid sequence  $s_i = [s_i^1, s_i^2, ..., s_i^{|s_i|}]$  of protein  $p_i$ , where  $s_i^m$  represents an amino acid, we lookup the embedding vocabulary of protein encoder to initialize the input embedding sequence as  $\mathbf{S}_i = [\mathbf{s}_i^1, \mathbf{s}_i^2, ..., \mathbf{s}_i^{|s_i|}] \in \mathbb{R}^{|s_i| \times d}$ , then concatenate it as

$$\mathbf{S}_i \leftarrow [\mathbf{v}_i^k, \mathbf{v}_i^p, \mathbf{S}_i] \in \mathbb{R}^{(2+|s_i|) \times d},\tag{3}$$

where  $|s_i|$  is the length of amino acid sequence  $s_i$ , and d is the dimension of embeddings. During inference, any related knowledge updates can be perceived by constructing these virtual tokens.

#### 3.1.2 KNOWLEDGE-GUIDED PRE-TRAINING

The pre-training of Kara has two purposes: 1) achieving effective information fusion of the contextualized virtual tokens (*i.e.*, knowledge and structure information) and the amino acid tokens (*i.e.*, protein information); and 2) integrating the knowledge-based relevance (*i.e.*, function similarities) among proteins into their representations. For the first purpose, we introduce knowledge-guided masked language modeling, allowing each amino acid to query the virtual tokens to extract helpful knowledge information for restoring masked tokens, which achieves token-level information fusion at each layer of the protein encoder. Specifically, given the input embedding sequence  $S_i$ , we use a 15% probability to mask each amino acid token (*i.e.*, replace the amino acid embedding as the embedding of special token '[MASK]'). Then the masked embedding sequence is encoded by the
 Transformer component (Vaswani et al., 2017) within the protein encoder as follows:

$$\tilde{\mathbf{S}}_{i}^{l} = \mathrm{LN}(\mathbf{S}_{i}^{l} + \mathrm{MHA}(\mathbf{S}_{i}^{l})), \tag{4}$$

219 220 221

241

242

243

263

$$\mathbf{S}_{i}^{(l+1)} = \mathrm{LN}(\tilde{\mathbf{S}}_{i}^{l} + \mathrm{MLP}(\tilde{\mathbf{S}}_{i}^{l})), \tag{5}$$

where  $S_i^0$  is initiated by  $S_i$ . LN denotes the layer-norm unit and MHA denotes the multi-head attention unit. After modeling the correlations among virtual tokens and amino acid tokens layer by layer, we leverage cross-entropy loss  $\mathcal{L}_{MLM}$  on the last layer token embeddings (*i.e.*,  $S_i^L$ , where L is the number of Transformer layers in protein encoder) to estimate the masked tokens.

While the aforementioned masked language modeling achieves token-level multi-modal knowledge 226 infusion, we further introduce a sequence-level regularization based on graph connectivity between 227 proteins, integrating biological function similarities into their representations. As we mentioned 228 before, each protein  $p_i \in \mathcal{N}_2(p_i)$  is two-hop connected with  $p_i$  in graph structure. This high-order 229 connectivity indicates that  $p_i$  and  $p_j$  share the same knowledge  $(r_i, g_i)$  and thus should be similar 230 in their biological functions. Therefore, each pair  $(p_i, p_j \in \mathcal{N}_2(p_i))$  can be regarded as positive pair 231 that we hope their embeddings are closer in semantic space (e.g., A9JR22 and A9JR44 in Figure 2), 232 and  $(p_i, p_k \notin \mathcal{N}_2(p_i))$  can be regarded as negative pair (e.g., A9JR22 and O14910). Specifically, in 233 Kara, we generate the sequence-level embedding of protein  $p_i$  as  $\tilde{\mathbf{p}}_i = \text{MEAN}(\mathbf{S}_i^L[2:])$ , where MEAN 234 is the mean-pooling operation, and  $\mathbf{S}_i^L[2:]$  is the last layer token embeddings except the virtual 235 tokens. Then, we apply the margin loss on sequence-level protein embeddings to ensure high-order 236 connected protein  $p_i$  is closer to  $p_i$  than other proteins in semantic space.

$$\mathcal{L}_{\text{reg}} = -\frac{1}{|\mathcal{N}_2(p_i)|} \sum_{p_j \in \mathcal{N}_2(p_i)} \text{MAX}(0, \text{sim}(\tilde{\mathbf{p}}_i, \tilde{\mathbf{p}}_j)) - \text{sim}(\tilde{\mathbf{p}}_i, \tilde{\mathbf{p}}_k) + \gamma), \tag{6}$$

where sim indicates the similarity function (*e.g.*, cosine similarity). We finally pre-train the parameters within the protein encoder, knowledge projector, and structure projector by jointly optimizing  $\mathcal{L}_{\text{MLM}}$  and  $\mathcal{L}_{\text{reg}}$ . These three components are then used to handle downstream tasks.

#### 244 3.2 FINE-TUNING STAGE 245

#### 246 3.2.1 KNOWLEDGE RETRIEVER

247 Proteins in downstream tasks often fail outside the PKGs (Zhou et al., 2023), restraining the use of 248 knowledge during fine-tuning. Existing methods thus incorporate knowledge modeling solely during 249 pre-training, leaving the fine-tuning process only guided by task-specific objectives. However, this 250 strategy has several limitations. 1) The optimization objectives of the pre-training and fine-tuning 251 stages are inconsistent (*i.e.*, one is knowledge-guided while the other is knowledge-isolated), caus-252 ing the pre-training knowledge to be catastrophically forgotten during fine-tuning (Lee et al., 2020). 253 2) Without PKGs during fine-tuning, these models fail to explicitly extract helpful knowledge for 254 downstream tasks, leading to unsatisfactory performance. 3) Knowledge graphs are consistently up-255 dated (e.g., correcting obsolete knowledge). Existing models cannot adapt to these updates without undergoing retraining. To tackle these challenges, we propose a knowledge retriever that can accu-256 rately predict potential knowledge for new proteins, and thus align them with PKGs. This allows the 257 pre-training and fine-tuning stages to directly integrate with knowledge through a unified modeling 258 process, thus unifying the optimization objectives and seamlessly adapting to knowledge updates. 259

**Generating Candidate Embeddings.** We regard the GO entities in protein knowledge graphs as retrieval candidates. To achieve more accurate and stable retrieval, we integrate the neighboring structure information of each GO entity  $g_m$  and generate its candidate embedding as

$$\mathbf{c}_m = \mathsf{MLP}_{aggregation}([\mathsf{MLP}_G(\mathbf{g}_m) : \mathsf{MLP}_G(\mathbf{g}_m^{go}) : \mathsf{MLP}_P(\mathbf{g}_m^{prot})]), \tag{7}$$

where  $\mathbf{g}_m$  is the stored embedding of  $g_m$ . We use  $\mathbf{g}_m^{go}$  to incorporate the information of neighboring GO entities of  $g_m$ , defined as  $\mathbf{g}_m^{go} = \frac{1}{|\mathcal{N}_{go}(g_m)|} \sum_{g_k \in \mathcal{N}_{go}(g_m)} \mathbf{g}_k$ . Similarly,  $\mathbf{g}_m^{prot}$  is used to incorporate the information of  $g_m$ 's neighboring proteins, defined as  $\mathbf{g}_m^{prot} = \frac{1}{|\mathcal{N}_{prot}(g_m)|} \sum_{p_k \in \mathcal{N}_{prot}(g_m)} \mathbf{p}_k$ .  $\mathcal{N}_{go}(g_m)$  and  $\mathcal{N}_{prot}(g_m)$  are respectively the 1-hop neighboring GO entities and 1-hop neighboring proteins of  $g_m$ . All of  $\text{MLP}_{aggregation}$ ,  $\text{MLP}_G$ , and  $\text{MLP}_P$  are trainable multi-layer perceptrons. 270 **Retrieval Process.** For each new protein  $p_n$ , we use a frozen ProtBert model to generate its query 271 embedding as  $\mathbf{q}_n = \text{MLP}_P(\text{MEAN}(ProtBert(s_n)))$  where  $s_n$  is the amino acid sequence of  $p_n$ . 272 Intuitively, we can traverse the relation set R and the GO entity set  $V_{go}$  to find potential knowledge 273 for  $p_n$ . However, the time consumption of this strategy is unacceptable because of the large size 274 of  $V_{ao}$  (*i.e.*, 47K in ProteinKG25). Fortunately, we observed that each relation only connects with several specific GO entities in PKGs, inspiring us to reduce the retrieval complexity by finding 275 relation-GO combinations. Specifically, for relation  $r_m \in R$ , we construct its candidate GO entity 276 set as  $\mathcal{E}(r_m) = \{g_m | (p_x, r_m, g_m) \in F\}$ . During retrieval, we traverse each  $r_m \in R$  and use 277 each of its corresponding candidate GO entity  $g_m \in \mathcal{E}(r_m)$  to construct the candidate knowledge 278  $(p_n, r_m, g_m)$ . Then we use the TransE objective Bordes et al. (2013) to score  $(p_n, r_m, g_m)$  as 279

280 281

287

288

295

298

301 302  $\mathbb{S}(p_n, r_m, g_m) = ||\mathbf{q}_n + \tilde{\mathbf{r}}_m - \mathbf{c}_m||_1,$ (8)

where  $\tilde{\mathbf{r}}_m = MLP_{rel}(\mathbf{r}_m)$ . Finally, we rank all the candidate knowledge based on their scores, and 282 then add the top-K candidate knowledge into G to align new protein  $p_n$  with the knowledge graph. 283

284 **Training Strategy.** We use triplets  $(p_i, r_i, g_i) \in F$  as valid knowledge and by minimizing a marginbased ranking criterion, we hope that valid knowledge can receive lower scores than invalid knowl-286 edge. The training objective is defined as

$$\mathcal{L}_{margin} = \text{MAX}(0, \mathbb{S}(p_i, r_i, g_i) - \mathbb{S}(p_i, r_i, g_j) + \gamma), \tag{9}$$

289 where MAX is the maximum operation and  $\gamma$  is a hyper-parameter used to control the distance be-290 tween valid and invalid knowledge.  $(p_i, r_i, g_j) \notin F$  is invalid knowledge constructed by perturbing 291  $g_i$  in  $(p_i, r_i, g_i)$  with a random GO entity  $g_i$ . Since the retrieval process needs to match information 292 from different modalities (*i.e.*, text descriptions and amino acid sequences), we further propose a 293 cross-modal matching loss to unify the semantic space of embeddings from different modalities as

$$\mathcal{L}_{match} = \mathrm{MAX}(0, ||\mathrm{MLP}_G(\mathbf{g}_i) - \mathrm{MLP}_P(\mathbf{g}_i^{prot})||_1 - ||\mathrm{MLP}_G(\mathbf{g}_i) - \mathrm{MLP}_P(\mathbf{g}_j^{prot})||_1 + \gamma), \quad (10)$$

296 where  $\mathbf{g}_{i}^{prot}$  is the neighboring protein embedding of a randomly sampled GO entity  $g_{j}$ . This loss 297 forces the text modality information  $MLP_G(\mathbf{g}_i)$  of  $g_i$  is closer to its corresponding neighboring protein information  $MLP_P(\mathbf{g}_i^{prot})$  (*i.e.*, amino acid sequence modality) than other protein information 299  $MLP_P(\mathbf{g}_i^{prot})$ . After jointly optimizing  $\mathcal{L}_{margin}$  and  $\mathcal{L}_{match}$ , the knowledge retriever can accurately 300 predict the potential knowledge for new proteins, enabling its effective alignment with PKGs.

#### 3.2.2 TASK-ORIENTED FINE-TUNING 303

304 After being aligned with PKGs, new proteins can be uniformly encoded with the enhancement of knowledge following Equations (1)-(5), and any related knowledge updates will be perceived when 306 constructing virtual tokens, as they can access the latest version of the PKG to extract relevant 307 knowledge and structures. Then, the downstream task objectives will be used to fine-tune Kara, 308 enabling the protein encoder to extract task-specific useful knowledge from PKGs via virtual tokens. Note that for each new protein  $p_n$ , we exclude other new proteins from  $\mathcal{N}_1(p_n)$  when constructing structure virtual token  $\mathbf{v}_n^p$ , to avoid noises. 310

311 Moreover, the structure-based regularization can also be seamlessly adapted to the fine-tuning stage. 312 This brings two advantages. 1) Downstream tasks usually lack sufficient training data (Rao et al., 313 2019). The regularization term can introduce biological function similarities among new proteins 314 as an auxiliary optimization objective, thus effectively avoiding over-fitting. 2) By using this regu-315 larization as a unified optimization objective of pre-training and fine-tuning, pre-trained knowledge can avoid being catastrophically forgotten and thus effectively transfer to downstream tasks. 316

317 **Complexity.** Compared with vanilla protein language models, the additional time complexity of 318 Kara only stems from the virtual tokens and the retrieval process. The two virtual tokens let the 319 encoding complexity become  $O((|S|+2)^2d)$  from  $O(|S|^2d)$ , where |S| is the length of amino 320 acid sequences. Thanks to the proposed strategy of finding relation-GO combinations, the time 321 complexity of retrieving potential knowledge for a new protein is only  $O(|R|k_{max})$ , where |R| is the size of the relation set in the PKG and  $k_{max}$  is the maximum size of the candidate GO entity 322 sets for relations.  $k_{max}$  is much smaller than the size of the GO entity set in the protein knowledge 323 graph (e.g., In proteinKG25,  $k_{max}$  is about 2K and the size of the GO entity set is 47K).

Table 2: Performance comparisons in the amino acid contact prediction task, where seq indicates 325 the number of amino acids between two selected amino acids. P@L, P@L/2, and P@L/5 denote 326 the precision calculated upon top L (i.e., L most likely contacts), top L/2, and top L/5 predictions, 327 respectively. The best results are **bolded** and the second best results are underlined. 328

	6	$6 \le seq \le 12$			$12 \le seq \le 24$			$24 \le seq$		
Models	P@L	P@L/2	P@L/5	P@L	P@L/2	P@L/5	P@L	P@L/2	P@L/5	
LSTM	0.26	0.36	0.49	0.20	0.26	0.34	0.20	0.23	0.27	
ResNet	0.25	0.34	0.46	0.28	0.25	0.35	0.10	0.13	0.17	
Transformer	0.28	0.35	0.46	0.19	0.25	0.33	0.17	0.20	0.24	
ProtBert	0.30	0.40	0.52	0.27	0.35	0.47	0.20	0.26	0.34	
ESM-1b	0.38	0.48	0.62	0.33	0.43	0.56	0.26	0.34	0.45	
ESM-2	0.40	0.50	0.62	0.35	0.44	0.56	0.27	0.35	0.45	
OntoProtein	0.37	0.46	0.57	0.32	0.40	0.50	0.24	0.31	0.39	
KeAP	<u>0.41</u>	<u>0.51</u>	<u>0.63</u>	<u>0.36</u>	<u>0.45</u>	0.54	0.28	0.35	0.43	
Kara	0.45	0.55	0.65	0.39	0.48	0.59	0.31	0.39	0.48	

#### 336 337 338

339

341

324

#### **EXPERIMENTS AND ANALYSES** 4

340 We evaluate the generalization ability of Kara in 6 representative downstream tasks, including amino acid contact prediction, homology detection, stability prediction, protein-protein interaction identi-342 fication, binding affinity prediction, and semantic similarity inference. Ablation studies, hyperparameter studies, and analysis of the knowledge retriever are also provided. The detailed task 343 descriptions are provided in Appendix B. Experimental settings and implementation details are pro-344 vided in Appendix C. We run each experiment independently three times and report the average re-345 sults. Codes and datasets are at https://anonymous.4open.science/r/Kara-1DB8/. 346

#### 347 348

#### 4.1 AMINO ACID CONTACT PREDICTION

349 **Overview.** This task aims to predict whether two amino acids within a protein are in contact, which 350 is a token-level classification task (Rao et al., 2019). Following Zhou et al. (2023), we use variants 351 of LSTM, ResNet, and Transformer proposed by the TAPE (tasks assessing protein embeddings) 352 benchmark (Rao et al., 2019), pre-trained language models ProtBert (Ahmed et al., 2020), ESM-1b 353 (Rives et al., 2021), and typical knowledge-enhanced model OntoProtein (Zhang et al., 2022) as 354 baselines. The state-of-the-art knowledge-enhanced model KeAP (Zhou et al., 2023) and the recent 355 powerful protein language model ESM-2-30t (Lin et al., 2023) are also used for comparison.

356 **Results.** As shown in Table 2, Kara outperforms existing models by large margins in all short-357  $(6 \le seq \le 12)$ , medium-  $(12 \le seq \le 24)$ , and long-range  $(24 \le seq)$  contact predictions, 358 achieving on average 9.5% and 11.0% improvements in the P@L and P@L/2 metrics. Compared 359 with the state-of-the-art knowledge-enhanced model KeAP, Kara consistently surpasses it, especially 360 in challenging long-range predictions. This is due to Kara's use of contextualized virtual tokens, 361 which allows each amino acid token to explicitly extract task-oriented knowledge information from 362 protein knowledge graphs, thus producing knowledge-contextualized token embeddings with more information. However, KeAP fails to incorporate knowledge during the fine-tuning stage.

364 365

366

#### 4.2 PROTEIN-PROTEIN INTERACTION IDENTIFICATION

367 Overview. The protein-protein interaction (PPI) identification task aims to predict the interaction 368 types of protein pairs, which is a sequence-level multi-label classification problem. Our experiments 369 are performed on three widely-used datasets SHS27K (Chen et al., 2019), SHS148K (Chen et al., 2019), and STRING (Lv et al., 2021), where 7 types of interactions are included. Following Zhang 370 et al. (2022), we use DPPI (Hashemifar et al., 2018), DNNPPI (Li et al., 2018), PIPR (Chen et al., 371 2019), and GNN-PPI (Lv et al., 2021) as four baselines. The LM baselines include ProtBert, ESM-372 1b, and ESM-2. The knowledge-enhanced model baselines include KeAP and OntoProtein. 373

374 **Results.** Experimental results are shown in Table 3. We can see that Kara outperforms baselines 375 on nearly all datasets, highlighting its effectiveness in accurately understanding the relationships between protein sequences. An interesting observation is that the performance gains of KeAP com-376 pared with OntoProtein are very small on the STRING dataset. It is suggested in Zhou et al. (2023) 377 that this is because the large number of fine-tuning data in the STRING dataset reduces the impact of Table 3: Performance comparisons in the protein-protein interaction identification task. BFS (breadth-first search) and DFS (depth-first search) indicate the strategies used to generate test sets on three datasets. We use the F1 score as the evaluation metric.

	. We use the fif beore us the evaluation metho.										
			SHS27K			SHS148K			STRING		
	Models	BFS	DFS	Avg	BFS	DFS	Avg	BFS	DFS	Avg	
	DNN-PPI	48.09	54.34	51.22	57.40	58.42	57.91	53.05	64.94	59.00	
	DPPI	41.43	46.12	43.77	52.12	52.03	52.08	56.68	66.82	61.75	
	PIPR	44.48	57.80	51.14	61.83	63.98	62.91	55.65	67.45	61.55	
	GNN-PPI	63.81	74.72	69.27	71.37	82.67	77.02	78.37	91.07	84.72	
	ProtBert	70.94	73.36	72.15	70.32	78.86	74.59	67.61	87.44	77.53	
	ESM-1b	74.92	78.83	76.88	77.49	82.13	79.31	78.54	88.59	83.57	
	ESM-2	75.05	79.55	77.30	77.19	83.34	80.26	81.32	89.19	85.30	
	OntoProtein	72.26	78.89	75.58	75.23	77.52	76.38	76.71	91.45	84.08	
	KeAP	<u>78.58</u>	77.54	<u>78.06</u>	77.22	<u>84.74</u>	<u>80.98</u>	<u>81.44</u>	89.77	<u>85.61</u>	
-	Kara	81.18	<u>78.85</u>	80.01	79.62	86.02	82.82	82.73	92.46	87.59	

knowledge modeling in pre-training. In contrast, Kara incorporates knowledge modeling in both the

387 388 389

378

382

384 385 386

390 391

392

pre-training and fine-tuning stages, thus avoiding catastrophically forgetting pre-trained knowledge.

## 393 394

395

#### 4.3 HOMOLOGY DETECTION AND STABILITY PREDICTION

397 Overview. Homology detection aims to predict the re-398 mote homology of protein, which is a sequence-level classification task. We follow the datasets and experimental 399 settings of Hou et al. (2018), and ask the model to pre-400 dict the right fold type of protein from 1,195 different 401 types. We report average accuracy on the fold-level held-402 out set. Stability prediction aims to predict the intrinsic 403 stability of a protein, which is a sequence-level regres-404 sion task. Following Rocklin et al. (2017), we use Spear-405 man's rank correlation scores to evaluate the model per-406 formance. The same baselines are used as in Table 2. 407

Table 4: Performance comparisons inthe protein homology detection and sta-bility prediction tasks.

Models	Homology	Stability
LSTM	0.26	0.69
ResNet	0.17	0.73
Transformer	0.21	0.73
ProtBert	0.29	0.78
ESM-1b	0.11	0.77
ESM-2	0.13	0.80
OntoProtein	0.24	0.75
KeAP	<u>0.29</u>	<u>0.80</u>
Kara	0.32	0.83

**Results.** As illustrated in Table 4, existing knowledgeenhanced models (*i.e.*, OntoProtein and KeAP) cannot outperform traditional language models in this task. Pre-

vious works (Zhang et al., 2022) attributed this failure to the lack of sequence-level objectives during
 pre-training. Instead, using structure-based regularization, Kara incorporates the knowledge-based
 relevance (*i.e.*, function similarity) among proteins as a unified sequence-level optimization objective in both pre-training and fine-tuning stages, thus achieving better performance.

414 415

416

#### 4.4 PROTEIN-PROTEIN BINDING AFFINITY PREDICTION

417 **Overview.** This task aims to map each pair of proteins to a 418 real value to indicate their binding affinity changes, which is a 419 sequence-level regression task. Following Unsal et al. (2022), we 420 use Bayesian ridge regression to the element-wise multiplication 421 of protein embeddings for predicting the binding affinity. The 422 SKEMPI dataset (Moal & Fernández-Recio, 2012) is used and the 423 performance is reported based on the mean square error of 10-fold cross-validation. We use the same baselines as recent works (Zhou 424 et al., 2023), additionally with KeAP and ESM-2. 425

Results. As shown in Table 5, all of the existing knowledgeenhanced models fail to outperform ESM-1b. This is because protein structure features play a vital role in this task (Unsal et al.,

2022), and the existing models overlook the modeling of protein structures but ESM-1b can achieve
it via its network architecture. Kara can achieve competitive performance with ESM-1b because the
protein knowledge graph contains the description of the structure properties of proteins, and Kara can directly inject such knowledge information into protein embeddings via the virtual tokens.

 Table 5: Performance comparisons in the protein-protein binding affinity prediction.

	*
Models	Affinity $(\downarrow)$
PIPR	0.63
ProtBert	0.58
ESM-1b	<u>0.50</u>
ESM-2	0.50
OntoProtein	0.59
KeAP	0.52
Kara	0.50

ruble 7. Ablation study and performance of variants.					
Tasks	Concate ( $6 \le seq \le 12$ )	PPI (STRING)	Homology	Stability	Affinity (↓
w/o contextualized virtual tokens	0.42	85.16	0.28	0.81	0.55
w/o structure-based regularizations	0.43	86.49	0.30	0.80	0.52
Retrieval based on the protein sequence similarities	0.43	85.33	0.29	0.79	0.57
Kara	0.45	87.59	0.32	0.83	0.50

Table 7: Ablation study and performance of variants

#### SEMANTIC SIMILARITY INFERENCE 4.5

440 **Overview.** This task evaluates models' ability to extract the Table 6: Performance in the se-441 biomolecular functional similarity among proteins. Following 442 Unsal et al. (2022), we use biological process (BP) and cellu-443 lar component (CC) to divide protein attributes into two groups 444 and calculate the Lin similarity in each group as the groundtruth similarity. We then calculate the Manhattan similarity be-445 tween protein embeddings as the prediction. The Spearman's 446 rank correlation between these similarities is calculated as the 447 metric. We include another powerful protein language model 448 MSA Transformer (Rao et al., 2021) as baseline. 449

mantic similarity inference task.					
Models	BP	CC			
MSA Transformer	0.31	0.30			
ProtBert	0.35	0.36			
ESM-1b	0.42	0.37			
ESM-2	0.41	0.39			
OntoProtein	0.36	0.36			
KeAP	0.41	0.40			
Kara	<u>0.41</u>	0.41			

**Results.** Table 6 shows that Kara outperforms existing knowledge-enhanced models on both BP and 450 CC. This can be attributed to the explicit incorporation of the information of GO entities in Kara, 451 which describes the functionality of proteins. Kara is unable to outperform ESM-1b on BP may be 452 because of the larger number of parameters of ESM-1b. However, it can still outperform the larger 453 model ESM-1b on CC, indicating its effectiveness in explicitly incorporating GO entity information. 454

455 456

432

439

#### 4.6 ANALYSIS OF KARA

457 Ablation and Variants. In Table 7 we study the performance contribution of each component in 458 Kara. We can see that all of the virtual tokens, structure-based regularization, and knowledge re-459 triever are essential for achieving good performance. Specifically, removing contextualized virtual 460 tokens makes Kara unable to incorporate knowledge explicitly, and thus significantly degrades its 461 performance in the protein-protein binding affinity prediction task which requires the property understanding of proteins. After removing structure-based regularization, Kara fails to integrate func-462 tion similarities into sequence-level protein embeddings and thus results in performance degradation 463 in sequence-level tasks, such as homology detection and stability prediction. 464

465 To assess the effectiveness of our proposed knowledge retriever, we compare it to a variant that uses a protein similarity-based retriever. In this variant, we use the frozen ProtBert model to calculate 466 embedding similarities between new proteins and those in the PKG, selecting the top-K similar 467 proteins and using their embeddings as virtual tokens. However, this approach does not outperform 468 Kara. The reason is that similarity-based retrievers struggle to accurately predict associated knowl-469 edge (*i.e.*, gene descriptions) for proteins, but proteins with similar sequences can have different 470 functions, so this approach may introduce irrelevant protein information as noise during encoding. 471

Hyper-parameter Analysis. During pre-training, we use 472 the ground-truth knowledge graph structure to construct the 473 virtual tokens. However, in the fine-tuning stage, because 474 the new proteins are not included in the protein knowledge 475 graph, we need to use the knowledge retriever to predict its 476 top-K potential knowledge to construct the virtual tokens 477 for fine-tuning and inference, where K is a hyper-parameter 478 used to control the amount of predicted potential knowl-479 edge incorporated. Because the predicted potential knowl-480 edge can bring additional information but also inevitable 481 noise, in this part we study how K affects the performance



Figure 3: Performance of Kara with different numbers of knowledge K.

482 of Kara. As shown in Figure 3, the performance improves across different tasks when K increases from 0 to 1, showcasing the value of incorporating knowledge into protein representations. As K483 continues to increase, performance fluctuates due to the introduction of noise from additional knowl-484 edge. Nevertheless, it still outperforms the variant without knowledge (*i.e.*, K=0), demonstrating 485 Kara's ability to effectively extract useful knowledge for downstream tasks.

# 486 4.7 ANALYSIS OF KNOWLEDGE RETRIEVER

Ablation Study. The accurate retrieval of the knowledge retriever is extremely important for Kara's performance in downstream tasks. Therefore, here we analyze how different components and hyper-parameters affect the retrieval performance of the knowledge retriever. As we mentioned before, the knowledge retriever is trained on the provide the second secon

Table 8:	Ablation	study	results	of	the
knowledg	ge retrieve	r.			

Metrics	P@1	P@5
Without structure information	0.681	0.669
Without cross-modal matching	0.733	0.721
Without relation-GO combinations	0.649	0.538
Original	0.821	0.795

the ProteinKG25 knowledge graph and we use the randomly sampled 2,000 proteins as the test set to select the best model. During the evaluation, for each test protein  $p_t$  we first traverse each relation  $r \in R$  to construct query pairs  $(p_t, r, ?)$ , and then use the knowledge retriever model to score the corresponding candidate knowledge  $(p_t, r, g_i^r)$ , where  $g_i^r$  is the candidate GO entity from  $\mathcal{E}(r)$ . After traversing all the relations, we rank candidate knowledge based on their scores and calculate the Precision@n (P@n) metric to evaluate the retrieval performance, which indicates how much knowledge on the top-n ranked candidates is valid (*i.e.*, exists in the protein knowledge graph).

500 Hyper-parameter Analysis. In Table 8, without structure infor-501 mation means that we remove the neighbor information in can-502 didate GO embeddings (Equation equation 7), and without crossmodal matching means that the knowledge retriever is only opti-504 mized based on  $\mathcal{L}_{margin}$ . We can see that both of these two components are beneficial to the retrieval performance. Without relation-505 GO combinations means that for each relation, we use the whole 506 GO entity set as the candidates during retrieval. The worse perfor-507 mance of this variant shows that relation-GO combination strategy 508 can not only reduce the retrieval time consumption, but also help 509 to filter out irrelevant GO candidates and thus improve the retrieval 510 accuracy. As shown in Figure 4, we can see that the higher neighbor 511 sampling number helps to achieve better retrieval performance. 512



Figure 4: Performance of knowledge retriever with different neighbor sampling numbers.

#### 513 514 5 RELATED WORK

515 Protein representation learning has attracted much attention due to the rapid development of pretrained language models. Existing works treat amino acid sequences as token sequences, and train 516 the language model with either supervision signal (Bepler & Berger, 2019) or self-supervised pre-517 training objective (Alley et al., 2019; Rao et al., 2019; Xiao et al., 2021; Ahmed et al., 2020; Un-518 sal et al., 2022; Lin et al., 2023; Brandes et al., 2022). However, these approaches ignore factual 519 knowledge (e.g., gene descriptions of proteins), resulting in inferior representations. Recently, On-520 toProtein (Zhang et al., 2022) is the first attempt to incorporate knowledge graph by proposing a 521 hybrid encoder. KeAP (Zhou et al., 2023) further extends it by performing token-level knowledge 522 exploration via cross-attention module. However, both of them are limited by ignoring knowledge 523 graph structure and task-oriented knowledge modeling. Very recently, GOProteinGNN (Kalifa et al., 524 2024) explores the benefit of graph structure. However, it still suffers from inconsistent optimiza-525 tion objectives and fails to consider the high-order relationships among proteins. Instead, Kara can explicitly inject high-order knowledge during both the pre-training and fine-tuning stages. 526

Some models explore incorporating information from other modalities to improve their ability to
learn protein representations (Chen et al., 2023a). For example, Otter-Knowledge (Lam et al., 2023)
designs knowledge graphs for not only proteins but broadly biomedical concepts. ProtST (Xu et al., 2023) infers protein representations from biomedical texts, but with no graph structure. Our model
captures text descriptions together with knowledge graphs for high-order knowledge incorporation.

532

534

#### 6 CONCLUSION AND FUTURE WORK

We develop a retrieval-augmented language model (Kara) for knowledge-aware protein representation learning, achieving the first unified and direct integration of protein knowledge graphs and protein language models, while considering the high-order relationships within knowledge graphs. Experimental results demonstrate the effectiveness of Kara and its superiority in 6 downstream tasks. A promising future direction is integrating other modalities, such as 3D structures, with knowledge graphs to develop multi-modal, knowledge-aware protein language models.

## 540 REPRODUCIBILITY STATEMENT

541 542

Here we detail the efforts that we have made to ensure reproducibility of this work. As shown in
Section 4, we provide the anonymous link where the source code of Kara and source data (including
both the ProteinKG25 knowledge graph and downstream task datasets) are downloadable. In Appendix B, we provide detailed descriptions of the experimental settings and data processing steps
for each downstream task. In Appendix C, we provide detailed descriptions of the experimental
environment, backbone selection, hyper-parameter settings, and implementation details (including
all of the pre-training and fine-tuning stages, as well as the knowledge retriever). We also provide
the official links to pre-trained models and datasets that we used in Kara.

550 551

552

565

566

567

570

571

572

573

577

578

579

580

585

586

587

References

- Elnaggar Ahmed, M Heinzinger, C Dallago, G Rihawi, Y Wang, L Jones, T Gibbs, T Feher, C Angerer, S Martin, et al. Prottrans: towards cracking the language of life's code through self-supervised deep learning and high performance computing. *bioRxiv*, 2020.
- Ethan C Alley, Grigory Khimulya, Surojit Biswas, Mohammed AlQuraishi, and George M Church. Unified rational protein engineering with sequence-based deep representation learning. *Nature methods*, 16(12):1315–1322, 2019.
- Mohammed AlQuraishi. Proteinnet: a standardized data set for machine learning of protein structure. *BMC bioinformatics*, 20:1–10, 2019.
- Tristan Bepler and Bonnie Berger. Learning protein sequence embeddings using information from structure. In *International Conference on Learning Representations*, 2019.
  - Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. 2013.
- Nadav Brandes, Dan Ofer, Yam Peleg, Nadav Rappoport, and Michal Linial. Proteinbert: a universal deep-learning model of protein sequence and function. *Bioinformatics*, 38(8):2102–2110, 2022.
  - Can Chen, Yingxue Zhang, Xue Liu, and Mark Coates. Bidirectional learning for offline modelbased biological sequence design. In *International Conference on Machine Learning*, pp. 5351– 5366, 2023a.
- Jiaoyan Chen, Hang Dong, Janna Hastings, Ernesto Jiménez-Ruiz, Vanessa López, Pierre Monnin,
  Catia Pesquita, Petr Škoda, and Valentina Tamma. Knowledge graphs for the life sciences: Recent
  developments, challenges and opportunities. *arXiv preprint arXiv:2309.17255*, 2023b.
  - Muhao Chen, Chelsea J-T Ju, Guangyu Zhou, Xuelu Chen, Tianran Zhang, Kai-Wei Chang, Carlo Zaniolo, and Wei Wang. Multifaceted protein–protein interaction prediction based on siamese residual rcnn. *Bioinformatics*, 35(14):i305–i314, 2019.
- Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1): 1–23, 2021.
  - Somaye Hashemifar, Behnam Neyshabur, Aly A Khan, and Jinbo Xu. Predicting protein–protein interactions through sequence-based deep learning. *Bioinformatics*, 34(17):i802–i810, 2018.
- Thanh Lam Hoang, Marco Luca Sbodio, Marcos Martinez Galindo, Mykhaylo Zayats, Raul Fernandez-Diaz, Victor Valls, Gabriele Picco, Cesar Berrospi, and Vanessa Lopez. Knowledge enhanced representation learning for drug discovery. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 10544–10552, 2024.
- Jie Hou, Badri Adhikari, and Jianlin Cheng. Deepsf: deep convolutional neural network for mapping protein sequences to folds. *Bioinformatics*, 34(8):1295–1303, 2018.

- Bozhen Hu, Cheng Tan, Lirong Wu, Jiangbin Zheng, Jun Xia, Zhangyang Gao, Zicheng Liu, Fandi
  Wu, Guijun Zhang, and Stan Z Li. Advances of deep learning in protein science: A comprehensive
  survey. *arXiv preprint arXiv:2403.05314*, 2024.
- 598 Dan Kalifa, Uriel Singer, and Kira Radinsky. Goproteingnn: Leveraging protein knowledge graphs 599 for protein representation learning. *arXiv preprint arXiv:2408.00057*, 2024.
- Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. Large language models struggle to learn long-tail knowledge. In *International Conference on Machine Learning*, pp. 15696–15707, 2023.
- David E Kim, Frank DiMaio, Ray Yu-Ruei Wang, Yifan Song, and David Baker. One contact
   for every twelve residues allows robust and accurate topology-level protein structure modeling.
   *Proteins: Structure, Function, and Bioinformatics*, 82:208–218, 2014.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- Hoang Thanh Lam, Marco Luca Sbodio, Marcos Martínez Galindo, Mykhaylo Zayats, Raul
   Fernandez-Diaz, Victor Valls, Gabriele Picco, Cesar Berrospi Ramis, and Vanessa Lopez. Otter knowledge: benchmarks of multimodal knowledge graph representation learning from different
   sources for drug discovery. *arXiv preprint arXiv:2306.12802*, 2023.
- Cheolhyoung Lee, Kyunghyun Cho, and Wanmo Kang. Mixout: Effective regularization to finetune large-scale pretrained language models. In *International Conference on Learning Representations*, 2020.
- Hang Li, Xiu-Jun Gong, Hua Yu, and Chang Zhou. Deep neural network based predictions of protein
   interactions using primary sequences. *Molecules*, 23(8):1923, 2018.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- Jun-Jie Liu, Natalia Orlova, Benjamin L Oakes, Enbo Ma, Hannah B Spinner, Katherine LM Baney,
   Jonathan Chuck, Dan Tan, Gavin J Knott, Lucas B Harrington, et al. Casx enzymes comprise a
   distinct family of rna-guided genome editors. *Nature*, 566(7743):218–223, 2019.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.
- Guofeng Lv, Zhiqiang Hu, Yanguang Bi, and Shaoting Zhang. Learning unknown from correlations: Graph neural network for inter-novel-protein interaction prediction. *arXiv preprint* arXiv:2105.06709, 2021.
- Iain H. Moal and Juan Fernández-Recio. Skempi: a structural kinetic and energetic database of mutant protein interactions and its use in empirical models. *Bioinformatics*, 28(20):2600–2607, 2012.
- John Moult, Krzysztof Fidelis, Andriy Kryshtafovych, Torsten Schwede, and Anna Tramontano.
   Critical assessment of methods of protein structure prediction (casp)—round xii. *Proteins: Structure, Function, and Bioinformatics*, 86:7–15, 2018.
- Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter
   Abbeel, and Yun Song. Evaluating protein transfer learning with tape. 2019.
- Roshan M Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu, and Alexander Rives. Msa transformer. In *International Conference on Machine Learning*, pp. 8844–8856, 2021.
- 646
   647 Danny Reidenbach. Evosbdd: Latent evolution for accurate and efficient structure-based drug design. In *ICLR 2024 Workshop on Machine Learning for Genomics Explorations*, 2024.

633

640

- Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15):e2016239118, 2021.
- Gabriel J Rocklin, Tamuka M Chidyausiku, Inna Goreshnik, Alex Ford, Scott Houliston, Alexander
  Lemak, Lauren Carter, Rashmi Ravichandran, Vikram K Mulligan, Aaron Chevalier, et al. Global
  analysis of protein folding using massively parallel design, synthesis, and testing. *Science*, 357 (6347):168–175, 2017.
- Peter Shaw, Bhaskar Gurram, David Belanger, Andreea Gane, Maxwell L Bileschi, Lucy J Colwell,
  Kristina Toutanova, and Ankur P Parikh. Protex: A retrieval-augmented approach for protein
  function prediction. *bioRxiv*, pp. 2024–05, 2024.
  - Serbulent Unsal, Heval Atas, Muammer Albayrak, Kemal Turhan, Aybar C Acar, and Tunca Doğan. Learning functional properties of proteins with language models. *Nature Machine Intelligence*, 4 (3):227–245, 2022.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
   Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Infor *mation Processing Systems*, 2017.
  - Yijia Xiao, Jiezhong Qiu, Ziang Li, Chang-Yu Hsieh, and Jie Tang. Modeling protein using largescale pretrain language model. *arXiv preprint arXiv:2108.07435*, 2021.
- Minghao Xu, Xinyu Yuan, Santiago Miret, and Jian Tang. Protst: Multi-modality learning of protein sequences and biomedical texts. In *International Conference on Machine Learning*, pp. 38749–38767, 2023.
- Yaoyao Xu, Xuxi Chen, Tong Wang, Huan He, Tianlong Chen, and Manolis Kellis. Demystify the
  secret function in protein sequence via conditional diffusion models. In *ICLR 2024 Workshop on Generative and Experimental Perspectives for Biomolecular Design*, 2024.
- <sup>676</sup>
   <sup>677</sup>
   <sup>678</sup>
   <sup>679</sup>
   <sup>671</sup>
   <sup>672</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>678</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>671</sup>
   <sup>672</sup>
   <sup>672</sup>
   <sup>673</sup>
   <sup>674</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>678</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>670</sup>
   <sup>671</sup>
   <sup>671</sup>
   <sup>672</sup>
   <sup>672</sup>
   <sup>673</sup>
   <sup>673</sup>
   <sup>674</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>677</sup>
   <sup>678</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>670</sup>
   <sup>671</sup>
   <sup>672</sup>
   <sup>672</sup>
   <sup>673</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
- Hong-Yu Zhou, Yunxiang Fu, Zhicheng Zhang, Bian Cheng, and Yizhou Yu. Protein representation
   learning via knowledge enhanced primary structure reasoning. In *International Conference on Learning Representations*, 2023.
- 683 684

685 686

660

661

662

663

667

668

687 688

- 689 690
- 691
- 692 693
- 694
- 695
- 696 697
- 698
- 699
- 700
- '01

# 702 A DATASET DESCRIPTION

704 We train the Kara using the ProteinKG25 knowledge graph (Zhang et al., 2022), consistent with 705 previous knowledge-enhanced models to achieve a fair comparison. ProteinKG25 includes about 706 4.5 million triplets describing relationships between protein and gene ontology (GO) entities, and 100K triplets describing relationships between GO entities. There are 31 kinds of relations, 600K 708 proteins, and 50K GO entities in ProteinKG25. Each GO entity in ProteinKG25 can be a molecule, a 709 cellular component, or a biological process, and each protein in ProteinKG25 has an average of 8.64 710 relations. Following the strategy provided by Zhou et al. (2023), we removed proteins appearing in the datasets of downstream tasks to avoid data leakage. The raw data of ProteinKG25 can be found 711 in https://www.zjukg.org/project/ProteinKG25/. 712

713 714

715

## **B** DOWNSTREAM TASK DEFINITIONS

Amino Acid Contact Prediction. This is a pairwise token-level matching task, where each pair of input amino acids  $(s^i, s^j)$  from a protein sequence s is mapped to a label  $y_{i,j} \in \{0, 1\}$ , indicating whether they are in contact or not (< 8Å apart). Accurate contact maps can facilitate robust modeling of full 3D protein structure (Kim et al., 2014). Following previous works (Zhou et al., 2023), we use data that comes from ProteinNet (AlQuraishi, 2019) and report precision on the ProteinNet CASP12 test set, which is a standard metric reported in CASP (Moult et al., 2018).

722 Protein-protein Interaction Identification. This is a pairwise sequence-level classification task. 723 Given a pair of proteins  $(p_i, p_j)$ , the model aims to predict the interaction types  $y_{i,j}$  between them. 724 Similar to previous works (Zhou et al., 2023), 7 types of interactions are included in our experiments, 725 which are reaction, binding, post-translational modifications, activation, inhibition, catalysis, and 726 expression. Each protein pair may belong to several types simultaneously so this is a multi-label 727 classification problem. We use three widely-used datasets SHS27K (Chen et al., 2019), SHS148K 728 (Chen et al., 2019), and STRING (Lv et al., 2021) in our experiments, where SHS27K and SHS148K 729 can be regarded as two subsets of STRING, which remove proteins with no more than 50 amino acids or  $\geq 40\%$  sequence identity. The F1 score is used as the evaluation metric for this task. 730

**Homology Detection.** This is a sequence-level classification task where each input protein p is mapped to a label  $y \in \{1, 2, ..., 1195\}$  based on its representation generated by protein language models, which indicates its possible protein fold. This task requires the evolutionary understanding of proteins and thus is valuable in microbiology and medicine (*e.g.*, discover new CAS enzymes (Liu et al., 2019)). We follow the previous works and use data from Hou et al. (2018). By holding out entire evolutionary groups from the training set, the model is required to generalize across evolutionary gaps. Same as Hou et al. (2018), we report accuracy on the fold-level heldout set.

**Stability Prediction.** This is a sequence-level regression task. Each input protein p is mapped as a number  $y \in \mathbb{R}$ , which represents the most extreme conditions under which the protein maintains its structure above a concentration threshold, serving as a proxy for its intrinsic stability. Measuring the stability of proteins is important for finding top candidates from expensive protein engineering experiments (Rao et al., 2019). We use the data provided by Rocklin et al. (2017), where the training set includes proteins from four rounds of experimental design, while the test set contains proteins that are Hamming distance-1 neighbors of the top candidates. We report the Spearman's rank correlation scores on the test set to evaluate the model performance.

- 746 Protein-protein Binding Affinity Prediction. This is a pairwise sequence-level regression task 747 that maps each pair of proteins  $(p_i, p_j)$  as a real value  $y \in \mathbb{R}$ , indicating the binding affinity changes between them. This task evaluates how well a protein representation can predict changes in bind-748 ing affinity resulting from protein mutations, thus being valuable for many downstream applications 749 such as drug design (Reidenbach, 2024). Following Unsal et al. (2022), we use Bayesian ridge re-750 gression to the element-wise multiplication of protein embeddings for predicting the binding affin-751 ity. The SKEMPI dataset (Moal & Fernández-Recio, 2012) is used and the performance is reported 752 based on the mean square error of 10-fold cross-validation. 753
- 754 Semantic Similarity Inference. This is a pairwise sequence-level regression task, which evaluates
   755 how well protein language models can capture information about biomolecular functional similarity
   between proteins. In this task, we emphasize the biological process (BP) and cellular component

 Table 9: Hyper-parameter settings for different downstream tasks.

14010	// 11/JP01 J				or annorone	downou oun cuono.
Tasks	Train Steps	Batch Size	K	$\mathcal{L}_{reg}$	Learning Rate	Gradient Accumulation Step
Contact	30,000	1	5	False	3e-5	8
Homology	2,200	2	1	True	4e-5	16
Stability	4,800	5	5	True	1e-5	16

(CC) categories similar to previous works (Unsal et al., 2022). We first use BP and CC to divide protein attributes into two groups and calculate the Lin similarity in each group as the ground-truth similarity. We then calculate the Manhattan similarity between protein embeddings as the prediction. The Spearman's rank correlation between these similarities is calculated as the metric.

## C EXPERIMENTAL DETAILS

756

763

764

765

766 767

768 769

808

**Experimental Settings.** Same as previous knowledge-enhanced protein language models such as 770 KeAP and OntoProtein, we use the ProtBert model  $^2$  as the backbone of the protein encoder within 771 Kara for a fair comparison. The text descriptions of GO entities and relations are encoded by the 772 PubMedBert model<sup>3</sup>, which is also consistent with previous works. While generating the pre-773 trained embeddings of items in the protein knowledge graph (see Section 2), we represent each item 774 as averaging the embeddings of its amino acid or word tokens. Our model is implemented with 775 Python and we refer to the official code released by Zhou et al. (2023) to implement the downstream 776 task experiments. All tasks use standard datasets and metrics, consistent with previous works, to 777 ensure a fair comparison. Note that since the train/valid/test set splittings of SHS27K, SHS148K, 778 and STRING datasets are not provided, we use the official code released by Lv et al. (2021) to 779 split each dataset with three different random seeds, and the average performance of each dataset is reported. All the experiments are conducted on NVIDIA A40 with 48 GB memory.

781 **Pre-training Implementation Details.** In the pre-training stage, we set the protein encoder within 782 Kara (*i.e.*, a PortBert model) as full-parameter trainable similar to previous works (Zhang et al., 783 2022). We only use proteins and knowledge preserved in the ProteinKG25 knowledge graph to pre-784 train Kara, where the maximum token length is set as 1024 for proteins and 512 for text descriptions. 785 For each protein, we randomly select 10 knowledge and 10 high-order connected proteins respectively from  $\mathcal{N}_1$  and  $\mathcal{N}_2$  to construct its virtual tokens. The margin  $\gamma$  is set as 5 and the number of 786 negative samples is set as 2. We set the batch size to 4 with the maximum number of update steps 787 to 10,000, and the gradient accumulation step to 16. The learning rate is set as 1e-6 and we use 788 AdamW (Loshchilov & Hutter, 2017) for optimization. The weight decay is set as 1e-2. 789

790 Knowledge Retriever Implementation Details. In the knowledge retriever, we set the sampling 791 number of neighbors during the candidate embedding generation as 100. Similar to the pre-training stage, the maximum token length is 1024 for proteins and 512 for text descriptions. To train the 792 knowledge retriever, we randomly sample 2,000 proteins as well as their associated knowledge from 793 the ProteinKG25 knowledge graph as the test set, and the remaining proteins are used as training 794 data. The best knowledge retriever model is selected based on the Precision@5 metric on the test set. We train the knowledge retriever with the Adam optimizer (Kingma & Ba, 2015). The number of 796 training epochs is set as 500 with the batch size as 100, and we use the early stopping strategy with 797 a patience of 5. The learning rate is set as 1e-3 and the negative sampling number is set as 20. The 798 margin  $\gamma$  is also set as 5. Note that we only train the parameters within MLPs and the embeddings 799 of items in the protein knowledge graph are frozen, thus making our knowledge retriever seamlessly 800 generalize to knowledge updates. During inference, we rank all the candidate knowledge for a new 801 protein based on their scores S (lower is better), and then select top-K knowledge to add to the protein knowledge graph, where  $K \in \{1, 5, 50, 100\}$ . 802

Fine-tuning Implementation Details. In the fine-tuning stage, we freeze the knowledge projector  $MLP_{knowledge}$  and the structure projector  $MLP_{structure}$ , and only optimize the parameters within the protein encoder for downstream tasks. Note that the protein-protein interaction identification, the protein-protein binding affinity prediction, and the semantic similarity inference tasks do not need fine-tuning and we directly use the pre-trained Kara to encode proteins for these tasks. For the

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/Rostlab/prot\_bert

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/microsoft/BiomedNLP-BiomedBERT-base-uncased-abstract-fulltext

structure-based regularization term, we still set the margin  $\gamma$  as 5 and the number of negative samples as 2. Different downstream tasks require various fine-tuning hyper-parameters and we summarize them in Table 9. Additionally, we follow the implementations in GNN-PPI (Lv et al., 2021) for PPI prediction, where the number of epochs is 600 and batch size is 2048. The learning rate is set as 1e-3 for the SHS27K dataset and 1e-4 for the SHS148K and STRING datasets. We follow the implementations in PROBE (Unsal et al., 2022) for the binding affinity prediction and semantic similarity inference tasks.