

Retrieved Sequence Augmentation for Protein Representation Learning

Chang Ma¹ Haiteng Zhao² Lin Zheng¹ Jiayi Xin¹ Qintong Li¹ Lijun Wu³
Zhihong Deng² Yang Lu⁴ Qi Liu¹ Sheng Wang⁵ Lingpeng Kong¹

¹The University of Hong Kong ²Peking University ³Shanghai Artificial Intelligence Laboratory
⁴University of Waterloo ⁵University of Washington
{cma, lpk}@cs.hku.hk

Abstract

Protein Language Models traditionally depend on Multiple Sequence Alignments (MSA) to incorporate evolutionary knowledge. However, MSA-based approaches suffer from substantial computational overhead and generally underperform in generalizing to de novo proteins. This study reevaluates the role of MSA, proposing it as a retrieval augmentation method and questioning the necessity of sequence alignment. We show that a simple alternative, Retrieved Sequence Augmentation (RSA), can enhance protein representation learning without the need for alignment and cumbersome preprocessing. RSA surpasses MSA Transformer by an average of 5% in both structural and property prediction tasks while being 373 times faster. Additionally, RSA demonstrates enhanced transferability for predicting de novo proteins. This methodology addresses a critical need for efficiency in protein prediction and can be rapidly employed to identify homologous sequences, improve representation learning, and enhance the capacity of Large Language Models to interpret protein structures.¹

1 Introduction

Proteins are fundamental yet complex components of life. They exhibit a diverse range of functions within organisms. The enigmatic characteristic of these macromolecules originates from the intricate interplay between their sequences, structures, and functions, which is influenced jointly by physics and evolution (Sadowski and Jones, 2009). Protein language models (Elnaggar et al., 2020; Jumper et al., 2021; Lin et al., 2022) capture the co-occurrence probability of amino acids observed in nature, thus encapsulating structural and evolutionary information within the resulting representations. While this approach has demonstrated its effectiveness (Elnaggar et al., 2021; Jumper et al.,

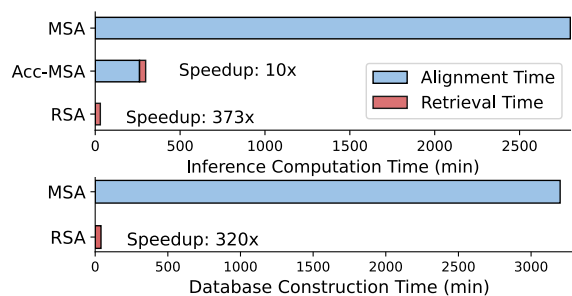


Figure 1: (Upper) Illustration of inference speed up by RSA compared to MSA on secondary structure prediction dataset with 8678 sequences. Accelerated MSA (Acc-MSA) refers to MSA built from sequences retrieved by our RSA retriever. (Lower) Illustration of speed up by RSA retrieval compared to MSA on database construction of 10000 protein sequences.

2021; Lin et al., 2022; Nijkamp et al., 2022; Rives et al., 2019), the evolutionary knowledge that can be extracted from a single sequence remains insufficient (Hu et al., 2022; Lin et al., 2022).

In order to compensate for this limitation, multiple sequence alignments (MSA) (Yanofsky et al. 1964; Altschuh et al. 1988; De Juan et al. 2013; Jumper et al. 2021) have been extensively used as a foundational protein feature engineering technique to extract protein evolutionary information in protein models (Rao et al., 2021; Jumper et al., 2021; Abramson et al., 2024). MSA draws on the evolutionary principle that functional constraints of species govern the mutation rate, which in turn drives the convergence of sequences. Therefore, key residues at functional sites tend to be conserved across protein families. MSA primarily aligns these conserved regions across homologous proteins to identify critical functional residues, such as substrate binding sites (Kunji and Robinson, 2006). Traditional approaches (Yanofsky et al., 1964; Altschuh et al., 1988) such as Potts Model (Balakrishnan et al., 2011), directly extracts statistical features from MSA for structural prediction. In recent studies (Jumper et al., 2021; Rao et al., 2021; Yang et al., 2020; Ju et al., 2021), mod-

¹Code and data are available at this [repo](#).

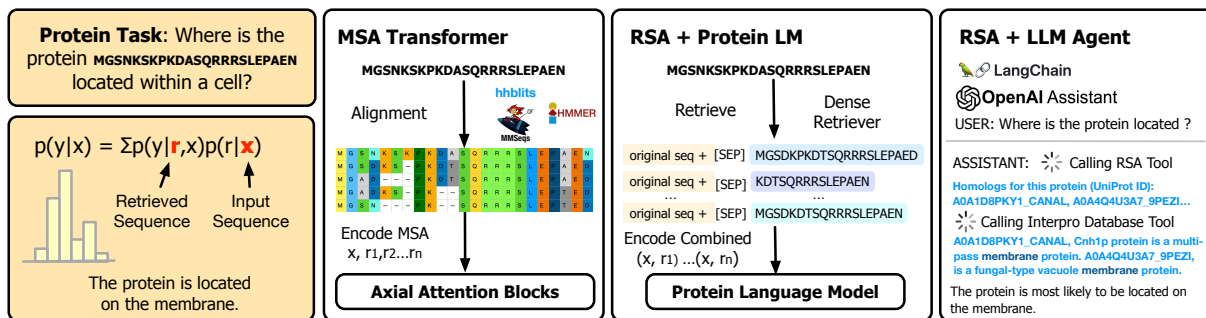


Figure 2: **Comparison between MSA Transformer and RSA.** MSA Transformer aligns query to the protein database and use axial attention to encode MSA feature. RSA could enhance protein language models by encoding both retrieved and original sequences. RSA could enhance LLM as a tool. Both MSA Transformer and RSA fall within the retrieval framework; however, RSA doesn't require the alignment process.

els like MSA Transformer (Rao et al., 2021) and AlphaFold (Jumper et al., 2021; Abramson et al., 2024) employ language models for predictions using MSAs as input feature. Despite being a vital component for state-of-the-art models, MSA carries a NP-Complete computational cost that scales with $O(L^N)$ (Wang and Jiang, 1994), where L represents the length of each sequence and N denotes the number of sequences examined. Even with acceleration techniques, MSA remains computationally intensive. For example, HHblits (Remmert et al., 2012) requires 10 seconds for a single iteration search on Pfam using 64 CPUs.

This motivates us to investigate alternatives to traditional alignment by addressing two research questions: (1) *Is alignment truly necessary for protein language models, and (2) is there a cost-efficient substitute for MSA?*

To answer these questions, we revisit MSA from a data-centric point of view and prove through theoretical analysis that it can be viewed as a retrieval-augmentation method (Goyal et al., 2022; Guu et al., 2020a; Khandelwal et al., 2019; Wang et al., 2022). We argue that MSA is retrieval through alignment. Retrieval-augmentation employs a large-scale memory of sequences as the knowledge base and utilizes multiple related input sequences instead of the single input to connect to the required knowledge. This approach offers the potential for more interpretable and modular knowledge capture (Guu et al., 2020b). It also enables rapid generalization to new domains (Khandelwal et al., 2019; Basu et al., 2022). Furthermore, we show that alignment is not essential as long as we have a strong sequence encoder, i.e. a transformer-based protein language model. This finding echoes previous research (Bhattacharya et al., 2020) that co-evolution patterns found through alignment could be cap-

tured with a single layer of attention without shared parameters across sequences. Since MSA is not indispensable and works mainly by enhancing protein language models as a retrieval-augmentation method, more efficient alignment-free retrievers can naturally serve as a substitute for MSA.

To this end, we explore Retrieved Sequence Augmentation (RSA) methods as a general framework to enhance protein representations. Specifically, RSA employs a pre-trained dense sequence retriever in search for protein sequences that are similar to the query sequence both in terms of homology as well as structure. By encoding retrieved sequences alongside the original protein, the model incorporates external knowledge and transfers it to new domains. Our assessment of this method consists of comprehensive experiments conducted across seven distinct tasks, including protein structure, function, evolution, and engineering, which require diverse knowledge. Using a vast database of approximately 40 million protein sequences, we show that a retrieval-based approach leveraging this data consistently outperforms state-of-the-art methods. Moreover, RSA employs retrieved sequences from dense retrievers without requiring an alignment process, thus resulting in a 373-fold speed-up and on-the-fly processing, as shown in Figure 1. Additionally, RSA without additional pretraining outperforms a pre-trained MSA Transformer in the downstream tasks, particularly for denovo proteins with few or no MSAs. It can be easily incorporated to augment any pre-trained protein language model, and be used as an efficient tool to boost the ability of large language model (LLM), e.g. GPT-4 (Achiam et al., 2023) to understand protein sequence. Consequently, we conclude that retrieval augmentation for proteins as a general framework can be a sound replacement for MSA in

terms of expressiveness, speed, and augmentation performance. Our contributions include:

- The investigation of retrieval-augmented protein language models and the proposal of the first alignment-free, efficient framework, RSA, for enhancing any protein representation model as well as LLM.
- The theoretical establishment of a unified framework reveal two insights: (1) MSA-augmented methods can be viewed as retrieval-augmented language models. Their performance can be explained by the injection of evolutionary knowledge. (2) The complex alignment process is not essential for deep protein language models.
- The demonstration that pre-trained dense retrievers offer greater efficiency and competitive efficacy in extracting homologous sequences and structurally similar sequences.

2 Augmenting Protein Representations with Retrieved Sequences – Is MSA necessary?

In this section, we rethink MSA-based models under a unified retrieval augmentation framework. We show that MSA sequences enhance representations in a similar way retrieved-augmentations do. Furthermore, we emphasize design elements that inspire our methodology for achieving increased efficiency and flexibility.

2.1 Background and Problem Statement

Given a protein $x = [o_1, o_2, \dots, o_L]$ comprising of L amino acids, the objective of a protein language model is to learn an embedding transferable to subsequent tasks, e.g. predicting properties of the sequence $p(y|x)$. The embedding is represented as $\text{Embed}(x) = [h_1, h_2, \dots, h_L]$, where $h_i \in \mathbb{R}^d$.

One approach to building an evolution-informed representation is to encode the input from a Multiple Sequence Alignment (MSA). An MSA includes several protein sequences, each a homolog of the query protein sequence—typically homologous proteins from species that are evolutionarily close. These proteins are aligned together such that each column in the alignment represents the evolutionary changes of an amino acid. In functionally important regions, amino acids tend to remain more stable, whereas in other regions, amino acids may undergo deletions, mutations, or insertions as evolution progresses. Figure 3 provides an illustration

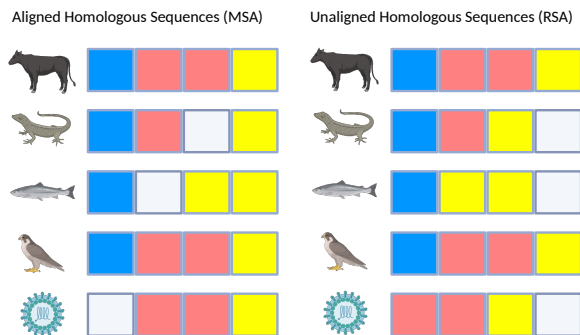


Figure 3: Illustrated difference of aligned and unaligned protein sequences. The white colour stands for the empty space in alignment "-".

for MSA. In our analyses, we consider MSA as N aligned protein homologs r_1, \dots, r_N . Prior studies (Yang et al., 2020; Ju et al., 2021) encode MSA as co-evolution statistics features $R_{1\dots N}$ and aggregate these features to derive the representation, while MSA Transformer (Rao et al., 2021; Jumper et al., 2021) perceives MSA as a matrix, employing axial attention to extract salient evolutionary traits. Here we also denote retrieved sequences as r_1, \dots, r_N and their features as $R_{1\dots N}$, though no alignment is performed on these sequences.

2.2 MSA is Retrieval through Alignment

Inspired by Guu et al. (2020b), we rethink state-of-the-art evolution augmentation methods under a new framework: *protein retrieval augmentation*. Specifically, we consider these methods as learning a downstream predictor $p(y|x)$ based on an aggregation of homologous protein representations $R_{1\dots N}$. From the view of retrieval, $p(y|x)$ is decomposed into two steps: *retrieve* and *predict*. For a given input x , the retrieve step first finds possibly helpful protein sequence r from a sequence corpus \mathcal{R} and then predict the output y conditioning on this retrieved sequence:

$$p(y|x) = \sum_{r \in \mathcal{R}} p(y|x, r)p(r|x) \approx \sum_{n=1}^N p(y|x, r_n)p(r_n|x). \quad (1)$$

The probability $p(r|x)$ denotes the possibility that r is sampled from the retriever given x . Intuitively it measures the similarity between the two sequences r and x . This framework also applies to the MSA-based augmentation methods. We explain in detail using a state-of-the-art MSA-augmentation model *MSA Transformer* (Rao et al., 2021) as an example. In MSA Transformer, the axial attention layers calculate self-attention both row-wise and column-wise. Column-wise attention is defined as follows,

Method	Retriever Form	Alignment Form	Weight λ_n	Aggregation Function
Existing Methods				
Potts Model	MSA	Aligned	—	—
Co-evolution Aggregator	MSA	Aligned	$\frac{1}{N}$	$\text{FFN}(\sum_{i=1}^N R_n(i)\lambda_n)$
MSA Transformer	MSA	Aligned	$\sigma(\frac{XW_Q(R_nW_K)^T}{N\sqrt{d}})$	$\text{FFN}(\sum_{i=1}^N R_n(i)\lambda_n)^\dagger$
Proposed Variants				
Unaligned MSA Augmentation	MSA	Not Aligned	$\sigma(-\ X - R_L\ _2)$	$\sum_{i=1}^N \text{FNN}(R_n(i))\lambda_n$
Accelerated MSA Transformer	Dense Retrieval	Aligned	$\sigma(\frac{XW_Q(R_nW_K)^T}{\lambda(N,d)})$	$\text{FFN}(\sum_{i=1}^N R_n(i)\lambda_n)$
Retrieval Sequence Augmentation	Dense Retrieval	Not Aligned	$\sigma(-\ X - R_n\ _2)$	$\sum_{i=1}^N \text{FNN}(\text{Embed}(x; r_n))\lambda_n$

Table 1: Protein Retrieval Augmentation methods decomposed along a different axis. We formulate the aggregation function in the sequence classification setting and use a feed-forward neural network $\text{FFN}(\cdot)$ to map representations to logits. The proposed variants vary in design axis from the existing methods. † Note that MSA Transformer performs the aggregation in each layer of axial attention, which differs from other variants.

given W_Q, W_K, W_V, W_O as the parameters in a typical attention function:

$$R_s(i) = \sum_{n=1}^N \sigma\left(\frac{R_s(i)W_Q(R_n(i)W_K)^T}{N\sqrt{d}}\right)R_n(i)W_VW_O, \quad (2)$$

where $R_n(i)$ denotes the i -th token representation of the n -th MSA sequence after performing the row-wise attention. Note that in MSA input, the first sequence r_1 is defined as the original sequence x . Then for a token prediction task, we define the i -th position output as y and the predicted distribution $p(y|x)$ can be expressed as:

$$\begin{aligned} p(y|x) &= \sum_{n=1}^N \sigma\left(\frac{R_1W_Q(R_nW_K)^T}{N\sqrt{d}}\right)(R_nW_VW_OW_y) \\ &= \sum_{n=1}^N p(y|x, r_n)\lambda_n = \sum_{n=1}^N p(y|x, r_n)p(r_n|x), \end{aligned} \quad (3)$$

where $\lambda_n = \sigma\left(\frac{R_1(i)W_Q(R_n(i)W_K)^T}{N\sqrt{d}}\right)$ is the weighting norm that represents the similarity of retrieved sequence r_n and original sequence x ; $p(y|x, r_n)$ is a predictor that maps the row-attention representation of r_n and x to label. Eq.3 gives a retrieval-augmentation view of MSA Transformer that essentially retrieves homologous sequences with multiple sequence alignment and aggregates representations of homologous sequences with regard to their sequence similarity. Taking one step further, we define a set of design dimensions to characterize the retrieving and aggregation processes. We introduce how popular models (Appendix E) and our proposed methods (§3) fall along them in Table 1. A detailed introduction of design details is available in Appendix D.

Our discussion and formulation so far reach the conclusion that retrieval augmentation serves as a comprehensive framework capable of extracting

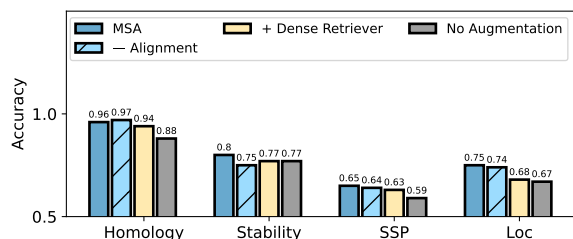


Figure 4: Comparison of variants on the retrieval form and the alignment form on property prediction tasks.

evolutionary knowledge, akin to multiple sequence alignment (MSA) augmentation methods. This underlines the prospects of retrieval sequence alignment (RSA) superseding MSA augmentations as an efficient and generalizable approach.

However, MSA-based methods claim a few advantages: the *alignment* process can help the model capture column-wise residue evolution; and the *MSA Retriever* uses a discrete, token-wise search criterion that ensures all retrieved sequences are homology. We introduce two variants to help challenge these claims: 1) **Unaligned MSA Augmentation (- Alignment)**, illustrated in Figure 8, uses the homologous sequences from MSA to augment representations without alignment and 2) **Accelerated MSA Transformer (+ Dense Retriever)** explores substituting the discrete retrieval process in MSA with a dense retriever. An empirical study on these models can be found in §2.3.

2.3 Do we still need alignment for proteins? An Empirical Analysis

It is commonly believed eliminating alignment could reduce expressiveness of proteins, as it highlights residue-wise mutations and compares across a protein family (Altschuh et al., 1988; De Juan et al., 2013). Bhattacharya et al. (2020) instead show that a single layer of attention suffices to

predict MSA-based statistics like pairwise residue co-evolution without shared parameters across the family. [Truong Jr and Bepler \(2023\)](#) also proposed using Transformers to represent co-evolution without alignment. Here we compare variants against MSA Transformer to further discuss the necessity of aligned feature when we have a strong protein language model as encoder. As shown in Figure 4, alignment does not consistently improve performance and unaligned variants achieve comparable performances on both homology and stability prediction. Additionally, a dense retriever competitively substitute aligner to find related sequences for retrieval augmentation. As alignment does not additionally improves performance when we have a strong protein language model, we could safely substitute MSA for dense retrieval augmentation methods.

3 RSA

Motivated by the potential of pre-trained retrievers to identify proteins that are homologous or geometric similar, we propose a pipeline, RSA (**R**etrieval **S**equences **A**ugmentation), to directly augment protein models on-the-fly. RSA follows the *retrieve-then-predict* framework in Eq. 1. It comprises of a neural sequence retriever $p(r|x)$, and a protein model that combines both original input and retrieved sequence to obtain prediction $p(y|x, r)$.

RSA Retriever is used for finding the sequences that are semantically close to the query. The similarity score $f(x, r)$ is defined as the negative L2 distance between the embedding G of the two sequences.

$$p(r|x) = \frac{\exp f(x, r)}{\sum_{r' \in \mathcal{R}} \exp f(x, r')}, \quad (4)$$

$$f(x, r) = -\|G(x) - G(r)\|_2$$

For protein retrieval, we aim to retrieve protein sequences that have similar structures or are homologous to the query sequence. Motivated by the high accuracy of k-nearest neighbor retrieval with ESM-1b ([Rives et al., 2019](#)) pre-trained embeddings (as shown in Table 2), we implement the embedding functions using a 34-layer ESM-1b encoder. We obtain sequence embeddings by performing average pooling over token embeddings. Note that finding the most similar proteins from a large-scale sequence database is computationally heavy. To accelerate retrieval, we use Faiss indexing ([Johnson](#)

Retrieval Task	Type	Recall	Precision
Pfam - Family	Homology	100	90.42
SCOPe - Fold	Structural	100	65.98
SCOPe - Superfamily	Structural	100	46.00
SCOPe - Family	Structural	100	24.71

Table 2: Accuracy for retrieving top 100 protein sequences with ESM1b embeddings. In dataset Pfam and SCOPe, we test whether retrieved proteins are of the same Family, Superfamily, or Fold as query protein.

[et al., 2019a](#)), which uses clustering and quantization to allow efficient similarity search.

Retrieval Augmented Protein Encoder Given a sequence x and a retrieved sequence r with length L and M respectively, the protein encoder combines x and r for prediction $p(y|x, r)$. To make our model applicable to any protein learning task, we need to augment both sequence-level representation and token-level representation (essential for structure prediction tasks). We concatenate the two sequences before input into the transformer encoder, which uses self-attention to aggregate global information from the retrieved sequence r into each token representation.

$$A = \sigma\left(\frac{(H_{[x;r]}W^Q)(H_{[x;r]}W^K)^T}{\sqrt{d}}\right), A = [A_x; A_r] \quad (5)$$

$$\text{Attn}(H_{[x;r]}) = (A_x H_x W^V + A_r H_r W^V) W^O$$

where $H_{[x;r]} = [h_1^x, h_2^x, \dots, h_L^x, h_1^r, \dots, h_M^r]$ denotes the input embedding of original and retrieved sequences. The output token representation h_i end-to-end learns to select and combine the representation of retrieved tokens. This can also be considered a soft version of MSA alignment. After computing for each pair of (x, r) , we aggregate them by weight $p(r|x)$ defined in Eq. 4.

Training For downstream finetuning, we maximize $p(y|x)$ by training on the retrieval augmented protein encoder. We freeze the retriever parameters during training. For a query sequence of length L with N retrieved proteins, suppose the length of retrieved proteins $L' \leq L$ the computation cost is N times the original model, $O(NL^2)$ for a transformer encoder layer, which is as efficient as MSA Transformer with $O(NL^2) + O(N^2L)$ complexity.

4 Experiments

4.1 General Setup

Downstream Task We evaluate RSA on seven downstream tasks: secondary structure predic-

tion (Klausen et al., 2019), contact prediction (AlQuraishi, 2019), remote homology prediction (Hou et al., 2018), subcellular localization prediction (Almagro Armenteros et al., 2017), stability prediction (Rocklin et al., 2017), protein-protein interaction (Pan et al., 2010) and structure prediction on CASP14 (Kryshtafovych et al., 2021). Please refer to Appendix Table 8 for more statistics of the datasets. The train-eval-test splits follow TAPE benchmark (Rao et al., 2019) for the first four tasks and PEER benchmark (Xu et al., 2022) for subcellular localization and protein-protein interaction.

Retriever and MSA Setup Limited by available computation resources, we build a database on Pfam (El-Gebali et al., 2018) sequences, which covers 77.2% of the UniProtKB (Apweiler et al., 2004) database and reaches the evolutionary scale. We generate ESM-1b pre-trained representations of 44 million sequences from Pfam-A and use Faiss (Johnson et al., 2019b) to build the retrieval index. For a fair comparison, the MSA datasets are also built on the Pfam database. We use HHblits (Remmert et al., 2012) to extract MSA, searching for 3 rounds with e-value threshold $1e-3$.

Baselines We apply our retrieval method to both pre-trained and from-scratch language models. Following Rao et al. (2019) and Rao et al. (2021), we compare our model with vanilla protein representation models, including LSTM (Liu, 2017), Transformers (Vaswani et al., 2017) and pre-trained models ESM-1b (Rives et al., 2019), ProtBERT (El-naggar et al., 2020). We also compare with state-of-the-art knowledge-augmentation models: Potts Model (Balakrishnan et al., 2011); MSA Transformer (Rao et al., 2021) injects evolutionary knowledge through MSA; OntoProtein (Zhang et al., 2022) uses gene ontology knowledge graph to augment protein representations and PMLM (He et al., 2021b) uses pair-wise pretraining to enhance co-evolution awareness.

Training and Evaluation To demonstrate RSA as a general method, we perform experiments both with a shallow transformer encoder, and a large pre-trained ProtBERT encoder. The Transformer model has 512 dimensions and 6 layers. Also, we combined our method with popular pre-trained protein folding architectures ESMFold and AlphaFold2. All self-reported models use the same truncation strategy and perform parameter searches on the learning rate, warm-up rate, and batch size.

4.2 Main Results

We show the result for downstream tasks in Table 3, including models with/without pretraining, and with/without knowledge augmentations. We form the following conclusion: Retrieval Sequence Augmentations perform on par with or even better than other knowledge-augmented methods without additional pre-training. Our method outperforms MSA Transformer on average by 5% and performs on par with PMLM on structure and evolution prediction tasks. Notably, both MSA Transformer and PMLM perform additional pre-training with augmentations, while our method uses no additional pre-training. From the results, we can see that RSA combined transformer model also improves by 10% than other shallow models. We also study retrieval sequence augmentations on pre-trained protein folding models in Table 4. Despite RSA was implemented without additional fine-tuning on folding models, we achieve a 2% improvement both on ESMFold and AlphaFold2.

4.3 Retrieval Augmentation for De Novo Proteins with Few Homologs

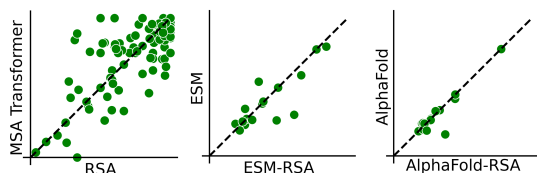
We test our model on a challenging problem in protein prediction, the prediction for proteins with few homologs, i.e. de novo (synthesized) proteins and orphan proteins (Fang et al., 2022; Wu et al., 2022). This task is especially difficult for MSA-based methods as alignment-based method often fails to generate MSA for these proteins, resulting in degraded performance. We test our model on 108 De Novo proteins from PDB (Berman et al., 2000) for the contact prediction task. It can be seen in Figure 5 that, RSA exceeds MSA transformer on 63.8% of data, demonstrating that RSA is more capable of locating augmentations for out-of-distribution proteins. We also test our model on the structure prediction task with 16 targets from CASP14-FM. CASP14-FM are considered more difficult because the absence of related templates requires the prediction methods to rely on de novo modeling techniques. We compare RSA augmented ESMFold and AlphaFold2 model with baselines in Figure 5, showing improved or competitive prediction on the majority of the targets. This results also show that our model surpasses MSA-based methods in transferring to unseen domains.

Method	Pretrain	Knowledge Pretrain	Knowledge Injection	SSP	Contact	Homology	Stability	Loc	PPI	Avg
Transformer	×	×	×	0.384	0.274	0.101	0.422	0.541	0.616	0.345
LSTM	×	×	×	<i>0.596</i>	<i>0.263</i>	<i>0.181</i>	<i>0.591</i>	<i>0.629</i>	<i>0.638</i>	0.404
RSA (Transformer backbone)	×	×	✓	0.541	0.332	0.346	0.602	0.591	0.700	0.518
ESM-1b	✓	×	×	<i>0.716</i>	<i>0.458</i>	<i>0.978</i>	<i>0.695</i>	<i>0.781</i>	<i>0.782</i>	0.668
ESM-2	✓	×	×	—	0.553	0.993	0.718	0.787	—	—
ProtBERT	✓	×	×	0.691	0.556	0.528	0.651	0.771	0.688	0.579
MSA Transformer (MSA N=1)	✓	✓	×	0.594	0.397	0.880	0.767	0.668	0.633	0.592
Gremlin (Balakrishnan et al., 2011)	×	×	✓	—	0.507	—	—	—	—	—
MSA Transformer	✓	✓	✓	0.654	0.618	0.958	0.796	0.694	0.751	0.672
OntoProtein (Zhang et al., 2022)	✓	×	✓	<i>0.68</i>	<i>0.40</i>	0.96	<i>0.75</i>	—	—	—
PMLM (He et al., 2021b)	✓	✓	×	0.728	0.717	<i>0.946</i>	—	—	—	—
RSA (ProtBERT backbone)	✓	×	✓	0.691	0.717	0.987	0.778	0.795	0.827	0.723

Table 3: Main Results for vanilla protein language models, knowledge-augmented baselines and our proposed RSA method. Note that *italized* result is reported by corresponding related work. The last column reports average result on all six tasks. For MSA Transformer and RSA, we all use 16 sequences (N=16) for augmentation. For Gremlin Potts model, we use the full MSA.

Methods	TM-Score	Percentage
ESMFold	0.678	
AlphaFold-single	0.335	
ESMFold-RSA	0.693	27.7%
AlphaFold-RSA	0.359	45.5%
AlphaFold-Full	0.747	
AlphaFold-Acc	0.551	19.7%

Table 4: Results for Structure Prediction on CASP14. Percentage represents the percentage of samples exceeding baselines.



PDB - De Novo Contact Prediction CASP14-FM - Structure Prediction (TM-score)

Figure 5: Prediction on proteins with few homologs, including contact prediction result on PDB de novo proteins and structure prediction result on CASP14-FM.

4.4 RSA as a Tool for Large Language Model

RSA can not only be used on small-scale representation learning model, it can also augment large language models, e.g. ChatGPT. Currently, even GPT4 model shows limited understanding of biological sequences. We follow ToolFormer (Schick et al., 2024) to equip RSA as a tool for GPT models, enabling LLM to query RSA and retrieve similar sequences as well as Pfam labels to improve understanding of the protein sequence. We benchmark RSA as a Tool on Gene Ontology tasks (Jensen et al., 2003). Results show that RSA as tool could improve protein understanding ability of LLMs for all tasks. We also show that RSA can be integrated

with other bioinformatics tools to build a LLM-agent for protein understanding in Appendix G.6.

Methods	CC	MF	BP	EC
GPT-3.5-Turbo	0.43	0.45	0.39	0.12
GPT-3.5-Turbo + RSA	0.60	0.45	0.58	0.37
GPT-4	0.54	0.50	0.37	0.54
GPT-4 + RSA	0.70	0.74	0.65	0.74

Table 5: Gene Ontology Results Using LLM (N=8)

4.5 Retrieval Speed

A severe speed bottleneck limits the use of previous MSA-based methods. We compare the computation time of RSA with MSA and an accelerated version of MSA as introduced in § 2.2. As shown in Figure 1, alignment time cost is much more intense than retrieval time. Even after reducing the number of alignment sequences to 500, accelerated MSA still need 270 min to build MSA. At the same time RSA only uses dense retrieval, and is accelerated 373 times. Also, MSA is limited by its cumbersome construction of retrieval HHM profile to perform HHM-HHM search. By contrast, RSA only needs to build the pre-trained features for the database, which can be accelerated with GPUs and batch forwarding. Results on a small database of 10000 proteins demonstrate a speedup of 320 times.

4.6 Ablation Study

Ablation on Retriever: Ablation on Retrieval Number Our study examines the effect of injected knowledge quantity for RSA and all retrieval baselines. The results are listed in Table 7. We select the Contact dataset because all baseline models are implemented on this dataset. RSA and all baselines perform consistently better as the retrieval

Tasks	MSA Transformer	Accelerated MSA Transformer	RSA
SSP	0.654	0.634	0.691
Contact	0.618	0.608	0.717
Homology	0.958	0.945	0.987
Stability	0.796	0.767	0.778
Loc	0.694	0.682	0.795
PPI	0.751	0.679	0.827

Table 6: Results for MSA Transformer and Accelerated MSA Transformer on downstream tasks. Accelerated MSA Transformer uses MSA built from dense retrieval.

number increases. Also, our model outperforms all baseline models for all augmentation numbers.

Methods	N=1	N=4	N=8	N=16	N=32	N= full
Potts Model	—	0.412	0.471	0.479	0.480	0.507
MSA Transformer	0.397	0.579	0.560	0.618	0.669	—
Accelerated MSA Transformer	0.397	0.524	0.538	0.608	0.654	—
RSA (ProtBERT backbone)	0.556	0.595	0.615	0.717	0.719	—

Table 7: The performance of RSA w.r.t. the number of retrieved sequences on contact prediction.

Ablation on Aggregation: We compare RSA with Accelerated MSA Transformer to evaluate whether our aggregation method is beneficial for learning protein representations. Note that only part of the retrieved sequences that are homologous are utilized after alignment. As shown in Table 6, the performance of the Accelerated MSA Transformer drops a lot compared to RSA. In contrast to MSA type aggregation, which is restricted by token alignment, our aggregation is more flexible and can accommodate proteins with variant knowledge.

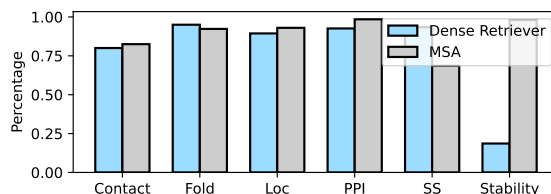
4.7 Retrieved Protein Interpretability

Dense Retrievers Find Homologous Sequences.

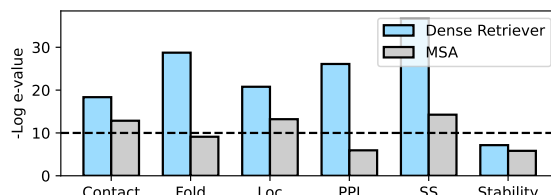
As illustrated in Figure 6(a), across all six datasets, our dense retriever retrieved a high percentage of homologous proteins that can be aligned to the original protein sequence, comparable to traditional MSA retrievers. We additionally plot each dataset’s negative log E-values distribution in Figure 6(b). Accordingly, dense retrieval show high potential for finding homologous sequences, which explains the ability of RSA to capture evolutionary knowledge.

RSA Retriever Find Structurally Similar Protein

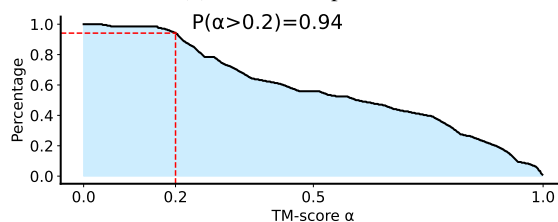
In Figure 6(c), we plot the TM scores between the RSA retrieved protein (structure obtained with ESMFold) and the origin protein on Protein-Net (AlQuraishi, 2019) test set. Most of the retrieved proteins TM-score exceed 0.2 (acceptable structural similarity) and about half are above 0.5 (high similarity), indicating dense retrieval is capable of finding proteins with structural knowledge.



(a) percentage comparison



(b) e-value comparison



(c) TM-score distribution

Figure 6: (a) Plot of the percentage of sequences that have found homologs on datasets for six tasks. (b) Plot of the $-\log(\text{E-values})$ of MSA and Dense Retriever obtained sequences. E-values of both methods are obtained with HHblits(Remmert et al., 2012). Sequences with $-\log \text{E-value} > 10$ are high-quality homologs. (c) Cumulative distribution of TM-scores for retrieved proteins.

4.8 Creating MSA with RSA

Despite the cumbersome computation, MSA is still widely used at present in SOTA models. In this section, we discuss the quality of MSA built by RSA, i.e. Accelerated MSA, a process 10 times faster. Table 6 illustrates that Accelerated MSA Transformer performs near to MSA Transformer (MSA N=16) for most datasets, except for Stability and PPI on which our retriever failed to find enough homologous sequences. Also, Accelerated MSA could be used as input for AlphaFold2 model, as shown in Table 4. However, the performance varies among samples, depending on retrieved sequence diversity, as further discussed in Appendix G.4.

5 Related Work

Retrieval-Augmented Language Models The integration of non-parametric retrieval and parametric models has been revolutionary for many problems (Kaplan et al., 2020; Guu et al., 2020b; He et al., 2021a; Borgeaud et al., 2021; Yogatama et al., 2021). Retrieval-augmentation introduces knowledge from memories and improve performance in

complex reasoning (Trivedi et al., 2022) and generalization (Khandelwal et al., 2019).

Protein Language Models To model and further understand the protein sequence data, language models are introduced to train on mass data (Heinzinger et al., 2019; Alley et al., 2019). Large scale pre-training enables language models to learn structural and evolutionary knowledge (El-naggar et al., 2021; Jumper et al., 2021; Lin et al., 2022). Despite these successes, many important applications still require MSAs and other external knowledge (Rao et al., 2021; Jumper et al., 2021; He et al., 2021b; Zhang et al., 2021; Ju et al., 2021; Rao et al., 2020). MSAs have been shown effective in improving representation learning, despite being extremely slow and costly in computation. Hu et al. (2022) and Hong et al. (2021) use dense retrieval to accelerate multiple sequence augmentation, while still dependent on alignment procedures. Recent work (Fang et al., 2022; Lin et al., 2022; Wu et al., 2022; Chowdhury et al., 2022) explores MSA-free language models though additional pre-training is involved. We take this step further to investigate retrieval-augmented protein language models. Another line of work improves MSA quality and generation speed by using generative models to augment and produce MSAs (Zhang et al., 2020; Zheng et al., 2024; Chen et al., 2024). We differ from this line of work as we discuss substitutes rather than augmentation of MSAs.

6 Conclusion

This work has highlighted the limitations inherent in traditional MSA-based approaches and proposed RSA as a substitute. Through extensive evaluation, we have demonstrated a significant improvement not only in the speed of processing—being more than 300 times faster than MSA—but also in enhancing predictive performance in downstream tasks with various models.

7 Limitation

One notable limitation of our method RSA is that it is highly dependent on high-quality pre-trained embeddings and the abundance of protein sequences. We found that our retriever tends to perform better in a database that has more protein sequences – that have not been screened by a clustering algorithm, like Uniclust30. This could be explained by our nearest neighbor retrieval technique which

often requires more similar sequences for augmentation. We also found different patterns in retrieval sequences from MSAs. Our retriever tends to show polarized retrieval quality, either finding many evolutionary close sequences or failing to find any homologous sequences. We believe this is due to the imbalanced training of pre-trained embeddings on different protein families and hope to mitigate this issue with further training on retrieval datasets.

We report other failed cases here for a more thorough view of our proposed method:

- Directly applying Accelerated MSAs to MSA-based pre-trained models often shows about 2-3% decrease on downstream performance than using original MSAs. This may be the natural gap between Acc-MSA and pre-training data. However, Accelerated MSAs are 10 times faster.
- The performance of RSA improves marginally with more sequences when $N > 16$. This is because we use the softmax distribution over L2 metrics to perform weighting, thereby assigning low weights to sequences further from the query.

We intend to further scale up our RSA method to larger protein databases and pre-train a retriever on abundant data in future work.

8 Acknowledgement

We would like to thank the anonymous reviewers for their valuable suggestions that greatly helped improve this work. This work is partially supported by the joint research scheme of the National Natural Science Foundation of China (NSFC) and the Research Grants Council (RGC) under grant number N_HKU714/21.

References

- Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. 2024. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, pages 1–3.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

- Ethan C Alley, Grigory Khimulya, Surojit Biswas, Mohammed AlQuraishi, and George M Church. 2019. Unified rational protein engineering with sequence-based deep representation learning. *Nature methods*, 16(12):1315–1322.
- José Juan Almagro Armenteros, Casper Kaae Sønderby, Søren Kaae Sønderby, Henrik Nielsen, and Ole Winther. 2017. Deeploc: prediction of protein subcellular localization using deep learning. *Bioinformatics*, 33(21):3387–3395.
- Mohammed AlQuraishi. 2019. Proteinnet: a standardized data set for machine learning of protein structure. *BMC bioinformatics*, 20(1):1–10.
- D Altschuh, T Vernet, P Berti, D Moras, and K Nagai. 1988. Coordinated amino acid changes in homologous protein families. *Protein Engineering, Design and Selection*, 2(3):193–199.
- Rolf Apweiler, Amos Bairoch, Cathy H Wu, Winona C Barker, Brigitte Boeckmann, Serenella Ferro, Elisabeth Gasteiger, Hongzhan Huang, Rodrigo Lopez, Michele Magrane, et al. 2004. Uniprot: the universal protein knowledgebase. *Nucleic acids research*, 32(suppl_1):D115–D119.
- Sivaraman Balakrishnan, Hetunandan Kamisetty, Jaime G Carbonell, Su-In Lee, and Christopher James Langmead. 2011. Learning generative models for protein fold families. *Proteins: Structure, Function, and Bioinformatics*, 79(4):1061–1078.
- Protein Data Bank. Rcsb pdb. 2022.
- Soumya Basu, Ankit Singh Rawat, and Manzil Zaheer. 2022. Generalization properties of retrieval-based models. *arXiv preprint arXiv:2210.02617*.
- Helen M Berman, John Westbrook, Zukang Feng, Gary Gilliland, Talapady N Bhat, Helge Weissig, Ilya N Shindyalov, and Philip E Bourne. 2000. The protein data bank. *Nucleic acids research*, 28(1):235–242.
- Nicholas Bhattacharya, Neil Thomas, Roshan Rao, Justas Dauparas, Peter K Koo, David Baker, Yun S Song, and Sergey Ovchinnikov. 2020. Single layers of attention suffice to predict protein contacts. *Biorxiv*, pages 2020–12.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2021. Improving language models by retrieving from trillions of tokens. *arXiv preprint arXiv:2112.04426*.
- Bo Chen, Zhilei Bei, Xingyi Cheng, Pan Li, Jie Tang, and Le Song. 2024. Msagpt: Neural prompting protein structure prediction via msa generative pre-training. *arXiv preprint arXiv:2406.05347*.
- Ratul Chowdhury, Nazim Bouatta, Surojit Biswas, Christina Floristean, Anant Kharkar, Koushik Roy, Charlotte Rochereau, Gustaf Ahdriz, Joanna Zhang, George M Church, et al. 2022. Single-sequence protein structure prediction using a language model and deep learning. *Nature Biotechnology*, 40(11):1617–1623.
- David De Juan, Florencio Pazos, and Alfonso Valencia. 2013. Emerging methods in protein co-evolution. *Nature Reviews Genetics*, 14(4):249–261.
- Sara El-Gebali, Jaina Mistry, Alex Bateman, Sean R Eddy, Aurélien Luciani, Simon C Potter, Matloob Qureshi, Lorna J Richardson, Gustavo A Salazar, Alfredo Smart, Erik L L Sonnhammer, Layla Hirsh, Lisanna Paladin, Damiano Piovesan, Silvio C E Tosatto, and Robert D Finn. 2018. [The Pfam protein families database in 2019](#). *Nucleic Acids Research*, 47(D1):D427–D432.
- Ahmed Elnaggar, Michael Heinzinger, Christian Dal-lago, Ghalia Rehawi, Yu Wang, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, Debsindhu Bhowmik, and Burkhard Rost. 2021. Prottrans: Towards cracking the language of life’s code through self-supervised learning. *bioRxiv*.
- Ahmed Elnaggar, Michael Heinzinger, Christian Dal-lago, Ghalia Rihawi, Yu Wang, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, et al. 2020. Prottrans: towards cracking the language of life’s code through self-supervised deep learning and high performance computing. *arXiv preprint arXiv:2007.06225*.
- Xiaomin Fang, Fan Wang, Lihang Liu, Jingzhou He, Dayong Lin, Yingfei Xiang, Xiaonan Zhang, Hua Wu, Hui Li, and Le Song. 2022. Helixfold-single: Msa-free protein structure prediction by using protein language model as an alternative. *arXiv preprint arXiv:2207.13921*.
- Anirudh Goyal, Abram L. Friesen, Andrea Banino, Theophane Weber, Nan Rosemary Ke, Adria Puigdomenech Badia, Arthur Guez, Mehdi Mirza, Ksenia Konyushkova, Michal Valko, Simon Osindero, Timothy Lillicrap, Nicolas Heess, and Charles Blundell. 2022. Retrieval-augmented reinforcement learning.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020a. Realm: Retrieval-augmented language model pre-training. *international conference on machine learning*.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020b. Retrieval augmented language model pre-training. In *International Conference on Machine Learning*, pages 3929–3938. PMLR.
- Junxian He, Graham Neubig, and Taylor Berg-Kirkpatrick. 2021a. Efficient nearest neighbor language models. *arXiv preprint arXiv:2109.04212*.
- Liang He, Shizhuo Zhang, Lijun Wu, Huanhuan Xia, Fusong Ju, He Zhang, Siyuan Liu, Yingce Xia, Jianwei Zhu, Pan Deng, et al. 2021b. Pre-training co-evolutionary protein representation via

- a pairwise masked language model. *arXiv preprint arXiv:2110.15527*.
- Michael Heinzinger, Ahmed Elnaggar, Yu Wang, Christian Dallago, Dmitrii Nechaev, Florian Matthes, and Burkhard Rost. 2019. Modeling aspects of the language of life through transfer-learning protein sequences. *BMC bioinformatics*, 20(1):1–17.
- Liang Hong, Siqi Sun, Liangzhen Zheng, Qingxiong Tan, and Yu Li. 2021. fastmsa: Accelerating multiple sequence alignment with dense retrieval on protein language. *bioRxiv*.
- Jie Hou, Badri Adhikari, and Jianlin Cheng. 2018. DeepSF: deep convolutional neural network for mapping protein sequences to folds. *Bioinformatics*, 34(8):1295–1303.
- Mingyang Hu, Fajie Yuan, Kevin K Yang, Fusong Ju, Jin Su, Hui Wang, Fei Yang, and Qiuyang Ding. 2022. Exploring evolution-aware & -free protein language models as protein function predictors. In *Advances in Neural Information Processing Systems*.
- Lars Juhl Jensen, Ramneek Gupta, H-H Staerfeldt, and Søren Brunak. 2003. Prediction of human protein function according to gene ontology categories. *Bioinformatics*, 19(5):635–642.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019a. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019b. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547.
- Fusong Ju, Jianwei Zhu, Bin Shao, Lupeng Kong, Tie-Yan Liu, Wei-Mou Zheng, and Dongbo Bu. 2021. Copulanet: Learning residue co-evolution directly from multiple sequence alignment for protein structure prediction. *Nature communications*, 12(1):1–9.
- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al. 2021. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Michael Lewis. 2019. Generalization through memorization: Nearest neighbor language models. *Learning*.
- Michael Schantz Klausen, Martin Closter Jespersen, Henrik Nielsen, Kamilla Kjaergaard Jensen, Vanessa Isabell Jurtz, Casper Kaae Soenderby, Morten Otto Alexander Sommer, Ole Winther, Morten Nielsen, Bent Petersen, et al. 2019. Netsurfp-2.0: Improved prediction of protein structural features by integrated deep learning. *Proteins: Structure, Function, and Bioinformatics*, 87(6):520–527.
- Andriy Kryshchak, Torsten Schwede, Maya Topf, Krzysztof Fidelis, and John Moult. 2021. Critical assessment of methods of protein structure prediction (casp)—round xiv. *Proteins: Structure, Function, and Bioinformatics*, 89(12):1607–1617.
- Edmund RS Kunji and Alan J Robinson. 2006. The conserved substrate binding site of mitochondrial carriers. *Biochimica et Biophysica Acta (BBA)-Bioenergetics*, 1757(9-10):1237–1248.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Sal Candido, et al. 2022. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv*.
- Xueliang Liu. 2017. Deep recurrent neural network for protein function prediction from sequence. *arXiv preprint arXiv:1701.08318*.
- Milot Mirdita, Konstantin Schütze, Yoshitaka Moriwaki, Lim Heo, Sergey Ovchinnikov, and Martin Steinegger. 2022. Colabfold: making protein folding accessible to all. *Nature methods*, 19(6):679–682.
- Erik Nijkamp, Jeffrey Ruffolo, Eli N Weinstein, Nikhil Naik, and Ali Madani. 2022. Progen2: exploring the boundaries of protein language models. *arXiv preprint arXiv:2206.13517*.
- Xiao-Yong Pan, Ya-Nan Zhang, and Hong-Bin Shen. 2010. Large-scale prediction of human protein-protein interactions from amino acid sequence based on latent topic features. *Journal of proteome research*, 9(10):4992–5001.
- Sudeep Pushpakom, Francesco Iorio, Patrick A Eyers, K Jane Escott, Shirley Hopper, Andrew Wells, Andrew Doig, Tim Williams, Joanna Latimer, Christine McNamee, et al. 2019. Drug repurposing: progress, challenges and recommendations. *Nature reviews Drug discovery*, 18(1):41–58.
- Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter Abbeel, and Yun Song. 2019. Evaluating protein transfer learning with tape. *Advances in neural information processing systems*, 32.
- Roshan Rao, Joshua Meier, Tom Sercu, Sergey Ovchinnikov, and Alexander Rives. 2020. Transformer protein language models are unsupervised structure learners. *Biorxiv*.
- Roshan M Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu, and Alexander Rives. 2021. Msa transformer. In *International Conference on Machine Learning*, pages 8844–8856. PMLR.

- Michael Remmert, Andreas Biegert, Andreas Hauser, and Johannes Söding. 2012. Hhblits: lightning-fast iterative protein sequence searching by hmm-hmm alignment. *Nature methods*, 9(2):173–175.
- Alexander Rives, Siddharth Goyal, Joshua Meier, Demi Guo, Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. 2019. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences of the United States of America*.
- Gabriel J Rocklin, Tamuka M Chidyausiku, Inna Greshnik, Alex Ford, Scott Houlston, Alexander Lemak, Lauren Carter, Rashmi Ravichandran, Vikram K Mulligan, Aaron Chevalier, et al. 2017. Global analysis of protein folding using massively parallel design, synthesis, and testing. *Science*, 357(6347):168–175.
- MI Sadowski and DT Jones. 2009. The sequence–structure relationship and protein function prediction. *Current opinion in structural biology*, 19(3):357–362.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2024. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36.
- Martin Steinegger, Markus Meier, Milot Mirdita, Harald Vöhringer, Stephan J Haunsberger, and Johannes Söding. 2019. Hh-suite3 for fast remote homology detection and deep protein annotation. *BMC bioinformatics*, 20(1):1–15.
- Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. *arXiv preprint arXiv:2212.10509*.
- Timothy Truong Jr and Tristan Bepler. 2023. Poet: A generative model of protein families as sequences-of-sequences. *Advances in Neural Information Processing Systems*, 36:77379–77415.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Dingmin Wang, Shengchao Liu, Hanchen Wang, Linfeng Song, Jian Tang, Song Le, Bernardo Cuenca Grau, and Qi Liu. 2022. Augmenting message passing by retrieving similar graphs.
- Lusheng Wang and Tao Jiang. 1994. On the complexity of multiple sequence alignment. *Journal of computational biology*, 1(4):337–348.
- Ruidong Wu, Fan Ding, Rui Wang, Rui Shen, Xiwen Zhang, Shitong Luo, Chenpeng Su, Zuofan Wu, Qi Xie, Bonnie Berger, et al. 2022. High-resolution de novo structure prediction from primary sequence. *BioRxiv*, pages 2022–07.
- Minghao Xu, Zuobai Zhang, Jiarui Lu, Zhaocheng Zhu, Yangtian Zhang, Chang Ma, Runcheng Liu, and Jian Tang. 2022. Peer: A comprehensive and multi-task benchmark for protein sequence understanding. *arXiv preprint arXiv:2206.02096*.
- Jianyi Yang, Ivan Anishchenko, Hahnbeom Park, Zhenling Peng, Sergey Ovchinnikov, and David Baker. 2020. Improved protein structure prediction using predicted interresidue orientations. *Proceedings of the National Academy of Sciences*, 117(3):1496–1503.
- Charles Yanofsky, Virginia Horn, and Deanna Thorpe. 1964. Protein structure relationships revealed by mutational analysis. *Science*, 146(3651):1593–1594.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.
- Jian Ye, Scott McGinnis, and Thomas L Madden. 2006. Blast: improvements for better sequence analysis. *Nucleic acids research*, 34(suppl_2):W6–W9.
- Dani Yogatama, Cyprien de Masson d’Autume, and Lingpeng Kong. 2021. Adaptive semiparametric language models. *Transactions of the Association for Computational Linguistics*, 9:362–373.
- Chengxin Zhang, Wei Zheng, SM Mortuza, Yang Li, and Yang Zhang. 2020. Deepmsa: constructing deep multiple sequence alignment to improve contact prediction and fold-recognition for distant-homology proteins. *Bioinformatics*, 36(7):2105–2112.
- He Zhang, Fusong Ju, Jianwei Zhu, Liang He, Bin Shao, Nanning Zheng, and Tie-Yan Liu. 2021. Co-evolution transformer for protein contact prediction. *Advances in Neural Information Processing Systems*, 34:14252–14263.
- Ningyu Zhang, Zhen Bi, Xiaozhuan Liang, Siyuan Cheng, Haosen Hong, Shumin Deng, Jiazhong Lian, Qiang Zhang, and Huajun Chen. 2022. Ontoprotein: Protein pretraining with gene ontology embedding. *arXiv preprint arXiv:2201.11147*.
- Wei Zheng, Qiqige Wuyun, Yang Li, Chengxin Zhang, P Lydia Freddolino, and Yang Zhang. 2024. Improving deep learning protein monomer and complex structure prediction using deepmsa2 with huge metagenomics data. *Nature Methods*, 21(2):279–289.

Hong-Yu Zhou, Yunxiang Fu, Zhicheng Zhang, Bian Cheng, and Yizhou Yu. 2023. Protein representation learning via knowledge enhanced primary structure reasoning. In *The Eleventh International Conference on Learning Representations*.

Contents

1	Introduction	1
2	Augmenting Protein Representations with Retrieved Sequences – Is MSA necessary?	3
2.1	Background and Problem Statement	3
2.2	MSA is Retrieval through Alignment	3
2.3	Do we still need alignment for proteins? An Empirical Analysis	4
3	RSA	5
4	Experiments	5
4.1	General Setup	5
4.2	Main Results	6
4.3	Retrieval Augmentation for De Novo Proteins with Few Homologs	6
4.4	RSA as a Tool for Large Language Model	7
4.5	Retrieval Speed	7
4.6	Ablation Study	7
4.7	Retrieved Protein Interpretability	8
4.8	Creating MSA with RSA	8
5	Related Work	8
6	Conclusion	9
7	Limitation	9
8	Acknowledgement	9
A	Limitations and Failed Case Analysis	15
B	Broader Impact and Potential Risks	15
C	A Brief Recap on Proteins	16
D	Details Introduction of Retrieval-augmentation Designs	17
E	Overview of Previous Protein Representation Augmentation Methods	17
F	Experiment Setups	18
F.1	In-depth Introduction to Protein Tasks	18
F.2	HHblits Settings	19
F.3	Model Hyperparameters	20
F.4	RSA and Variants Implementation Details	20
F.4.1	Retriever Implementation Details	20
F.4.2	ProtBERT-RSA Architecture and Implementation	20
F.4.3	RSA for Protein Folding	21
F.4.4	Accelerated MSA	21

G	Supplementary Experiment Analysis	22
G.1	Comparison of the Running time between RSA vs MSA	22
G.2	Case Study	22
G.3	Domain Adaptation Analysis	24
G.4	Comparison of Accelerated MSA vs MSA quality	24
G.5	Interpretability of RSA	26
G.6	ProteinChat: RSA Empowers ChatGPT on Protein Understanding	26

A Limitations and Failed Case Analysis

One notable limitation of our method RSA is that it is highly dependent on high-quality pre-trained embeddings and the abundance of protein sequences. We found that our retriever tends to perform better in a database that has more protein sequences – that have not been screened by a clustering algorithm, like Uniclust30. This could be explained by our nearest neighbor retrieval technique which often requires more similar sequences for augmentation. We also found different patterns in retrieval sequences from MSAs. Our retriever tends to show polarized retrieval quality, either finding many evolutionary close sequences or failing to find any homologous sequences. We believe this is due to the imbalanced training of pre-trained embeddings on different protein families and hope to mitigate this issue with further training on retrieval datasets.

We report other failed cases here for a more thorough view of our proposed method:

- Directly applying Accelerated MSAs to MSA-based pre-trained models often shows about 2-3% decrease on downstream performance than using original MSAs. However, Accelerated MSAs are 10 times faster.
- The performance of RSA improves marginally with more sequences when $N > 16$. This is because we use the softmax distribution over L2 metrics to perform weighting, thereby assigning low weights to sequences further from the query.
- We found that in protein folding tasks, performing Average Pooling on ESMFold/AlphaFold shows worse zero-shot performance than Max Pooling with a scoring model. This is due to the misalignment of protein structures and simple weighting could result in averaging the structures of proteins with different angles of view.

B Broader Impact and Potential Risks

In this section, we discuss the broader impact of RSA in terms of protein representation learning, de novo protein understanding, as well as the potential application to large language models.

RSA for Protein Representation Learning Developing efficient protein representation learning methods will significantly improve the ability to analyze complex protein structures, functions, and interactions. This would lead to a more comprehensive understanding of biological processes at the molecular level, consequently boosting advancements in the fields of bioinformatics and computational biology. In this paper, we propose RSA as an efficient and effective protein representation learning methods, which will spur the development of protein representation learning methods. Notably, our method requires no alignment methods. The traditional alignment process in MSA often requires mass CPU engines mostly available to academics. Our method on the other hand only requires a small memory GPU like 3090Ti and we will publicize our retrieval index, promoting democratic research in this field.

RSA for De Novo Protein Understanding We have shown in our work that RSA could perform De Novo Protein Understanding. This is particularly important for drug repurposing and virtual screening tasks (Pushpakom et al., 2019) for drug discovery. This can contribute to the development of personalized medicine by facilitating the identification of disease-specific protein biomarkers and selecting molecular cures for various diseases. However, de novo protein understanding often relies on newly-designed protein databases, which may include sensitive information about individuals, such as their genetic makeup, or violates intellectual property rights. Ensuring the privacy and security of this data is critical to prevent misuse and protect individual rights

RSA as Tool for Large Language Models In addition to the potential impacts in the field of biology, our method could also improve the ability of Large Language Models in biological sequence understanding. Currently, large language models like ChatGPT show difficulty in understanding protein sequences. We showcase how RSA could improve this ability with the combination of retrieval and chain of thought. This application is valuable in education and training, as users could rapidly learn about proteins through chat models, which help educate the next generation of researchers in bioinformatics, computational biology, and related fields. This will lead to a more skilled workforce in the life sciences.

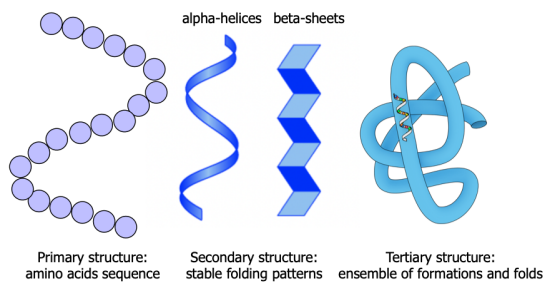


Figure 7: Illustrated explanation of protein levels of structures, primary structure, secondary structure and tertiary structure.

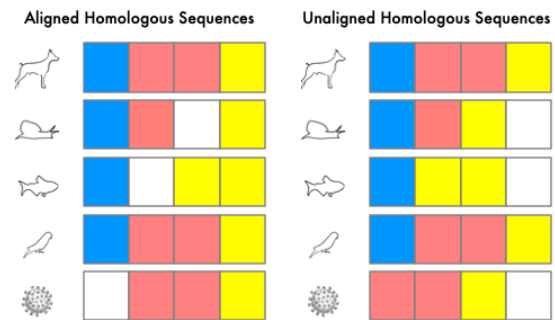


Figure 8: Illustrated difference of aligned and unaligned protein sequences. The white colour stands for the empty space in alignment "-".

C A Brief Recap on Proteins

Proteins are the end products of the decoding process that starts with the information in cellular DNA. As workhorses of the cell, proteins compose structural and motor elements in the cell, and they serve as the catalysts for virtually every biochemical reaction that occurs in living things. This incredible array of functions derives from a startlingly simple code that specifies a hugely diverse set of structures. In fact, each gene in cellular DNA contains the code for a unique protein structure. Not only are these proteins assembled with different amino acid sequences, but they also are held together by different bonds and folded into a variety of three-dimensional structures. The folded shape, or conformation, depends directly on the linear amino acid sequence of the protein. In fact, this phenomenon is denoted as the *sequence-structure-function paradigm*. Here we will emphasize four key concepts in protein understanding.

1. What are proteins made of ?

Amino acids. Within a protein, multiple amino acids are linked together by peptide bonds, thereby forming a long chain. There are 22 alpha-amino acids, from which proteins are composed. We model these amino acids in a similar way in NLP, as tokens. A tokenizer breaks the protein sequences into amino acid tokens that could be modeled by protein language models.

2. Protein structures

There are four levels of structures in protein, as illustrated in Figure 7:

- Primary structure: amino acids sequence
- Secondary structure: stable folding patterns, including Alpha Helix, Beta Sheet.
- Tertiary structure: ensemble of formations and folds in a single linear chain of amino acids
- macromolecules with multiple polypeptide chains or subunits

Predicting protein structure is an important and difficult task. In this work, we also perform experiments on three tasks – secondary structure prediction, protein contact prediction (tertiary structure), and protein folding (tertiary structure), with increasing task difficulty.

3. Protein Homology Protein homology is defined as shared ancestry in the evolutionary history of life. There exists different kinds of homology, including orthologous homology that may be similar

function proteins across species (human and mice α -goblin), and paralogous homology that is the result of mutations (human α -goblin and β -goblin). Homologies result in conservative parts in protein sequences, or leads to similar structures and functions.

4. Multiple Sequence Alignments A method used to determine conservative regions and find homologous sequences. An illustration (Figure 8) is given here to show how sequences are aligned. Aligned tokens may include the original amino acid, substitution, and deletions. The traditional way to generate MSA is using dynamic programming, with $O(L^N)$ complexity. Temporary methods use HMM-HMM alignment, as well as other acceleration methods. HH-Suite3 (Steinegger et al., 2019) reports a time complexity of $O(NL^2)$, which is still costly when performing alignment on a large database.

D Details Introduction of Retrieval-augmentation Designs

We introduce 4 design dimensions to distinguish RSA against MSA-based methods as well as discuss how we could design more efficient substitutes.

- **Retriever Form** indicates the retriever type used. Multiple Sequence Alignment is a discrete retrieval method that uses alignment (Ye et al., 2006) to find homologous sequences. Dense retrieval (Johnson et al., 2019b) has been introduced to accelerate discrete sequence retrieval.
- **Alignment Form** indicates whether retrieved sequences are aligned.
- **Weight Form** is the aggregation weight of homologous sequences, as the $p(r_n|x)$ in Eq. 3. Here we denote this weight as λ_n . Traditionally, aggregation methods consider different homologous sequences to be similarly important and use average weighting. MSA Transformer uses a weighted pooling method though the weights of λ_n use global attention and are dependent on all homologous sequences.
- **Aggregation Function** is how the representations of homologous sequences are aggregated to the original sequence to form downstream prediction, as in $p(y|x, r)$. For example, considering the sequence classification problem, a fully connected layer maps representations to logits. The retrieval augmentation probabilistic form first maps each representation to logits $p(y|x, r_n)$ and then linearly weight the logits with λ_n in Eq. 3.

Here retriever and alignment are the main bottlenecks of retrieval augmentation methods. The aggregation function and weight form are mainly dependent on model architecture and we focus on the first two dimensions in this paper.

E Overview of Previous Protein Representation Augmentation Methods

Below we introduce several state-of-the-art evolution augmentation methods for protein representation learning. These methods rely on MSA as input to extract representations. We use x to denote a target protein and its MSA containing N homologous proteins. We consider MSAs as N aligned protein homologs r_1, \dots, r_N . These studies (Yang et al., 2020; Ju et al., 2021) encode MSA as co-evolution statistics features $R_{1\dots N}$ and aggregate these features to derive the representation, while MSA Transformer (Rao et al., 2021; Jumper et al., 2021) perceives MSA as a matrix, employing axial attention to extract salient evolutionary traits. A unified view of these variants is available in Table 1 and §3.2 in the main paper.

Potts Model (Balakrishnan et al., 2011). This line of research fits a Markov Random Field to the underlying MSA with likelihood maximization. This approach is different from other protein representation learning methods as it only learns a pairwise score for residues contact prediction. We will focus on other methods that augment protein representations that can be used for diverse downstream predictions.

Co-evolution Aggregator (Yang et al., 2020; Ju et al., 2021). One way to build an evolution informed representation is to use a MSA encoder to obtain the co-evolution related statistics. By applying MSA encoder on the n -th homologous protein in the MSA, we can get a total of $L \times d$ embeddings R_n , each position is a d channel one-hot embedding indicating the amino acid type. We use w_n to denote the weight

from R_n when computing the token representation h_i :

$$h_i = \frac{1}{M_{eff}} \sum_{n=1}^N w_n R_n(i), \quad (6)$$

where $M_{eff} = \sum_{n=1}^N w_n$ and $w_n = \frac{1}{N}$. For contact prediction, pair co-evolution representation are computed in a similar way from the hadamard product:

$$h_{ij} = \frac{1}{M_{eff}} \sum_{n=1}^N w_n R_n(i) \otimes R_n(j). \quad (7)$$

Ensembling Over MSA (Rao et al., 2020). This approach aligns and ensembles representations of homologous sequences. Consider the encoder extract the same token representations for unaligned and aligned sequences. The ensembled token representation is:

$$h_i = \frac{1}{N} \sum_{n=1}^N R_n(i), h_{ij} = \frac{1}{N} \sum_{n=1}^N \sigma\left(\frac{R_n(i)W_Q(R_n(j)W_K)^T}{N\sqrt{d}}\right). \quad (8)$$

MSA Transformer (Rao et al., 2021) In each transformer layer, a tied row attention encoder extracts the dense representation R_n , then a column attention encoder

$$R_s(i) = \sum_{n=1}^N \sigma\left(\frac{R_s(i)W_Q(R_n(i)W_K)^T}{N\sqrt{d}}\right)R_n(i)W_V. \quad (9)$$

Knowledge Graph Augmentation (Zhang et al., 2022; Zhou et al., 2023). This line of research aims at incorporating factual knowledge in protein representations. Different from MSA-based methods that draw evolution knowledge from raw protein sequences, these methods are dependent on protein knowledge graphs that have been annotated by experts, therefore we only provide comparisons with these models in experimental studies and don't incorporate them into our unified framework.

F Experiment Setups

F.1 In-depth Introduction to Protein Tasks

Table 8: Overview for datasets in downstream tasks

Task Name	Dataset source	#train sequences	#test sequences
Secondary Structure Prediction	NetSurfP-2.0 (Klausen et al., 2019)	8,678	513
Contact Prediction	ProteinNet (AlQuraishi, 2019)	25,299	40
Remote Homology Prediction	DeepSF (Hou et al., 2018)	12,312	718
Stability Prediction	Rocklin's Dataset (Rocklin et al., 2017)	53,571	12,851
Subcellular Localization	DeepLoc (Almagro Armenteros et al., 2017)	8,945	2,768
Protein Protein Interaction	Pan's Dataset (Pan et al., 2010)	6,844	227
Protein Folding	CASP14 (Kryshtafovych et al., 2021)	-	65

Secondary structure prediction (SSP)

Task Formulation: 8-class classification $o_i \mapsto \{0, 1, \dots, 7\}$

Task Description: Secondary structure prediction aims to predict the secondary structure of proteins, which indicates the local structures. This task predicts an 8-class label for each token, indicating which local structure this amino acid belongs to.

Task Impact: This task helps to determine whether a model captures protein local structure.

Contact prediction (Contact):

Task Formulation: 2-class classification $(o_i, o_j) \mapsto \{0, 1\}$

Task Description: Contact prediction predicts the medium-range and long-range (distance >6) residue-residue contact, which measures the ability of models to capture global tertiary structures.

Task Impact: This task helps to determine whether a model captures protein tertiary structure. The assessment of this task focuses specifically on medium- and long-range interactions due to their crucial importance in the protein folding process.

Homology prediction (Homology):

Task Formulation: 1195-class classification $x \mapsto \{0, 1 \dots 1194\}$

Task Description: Homology prediction aims to predict the fold label of any given protein, which indicates the evolutionary relationship of proteins.

Task Impact: Protein fold classification is important for both functional analysis and evaluating evolutionary knowledge.

Stability prediction (Stability):

Task Formulation: regression $x \mapsto \mathbb{R}$

Task Description: Stability prediction is a protein engineering task, which measures the change in stability w.r.t. residue mutations.

Task Impact: Evaluate the ability of models to predict protein function as well as evaluate the ability of models to understand mutations, which is crucial for drug discovery and protein engineering.

Subcellular Localization (Loc):

Task Formulation: regression $x \mapsto \{0, 1, \dots, 7\}$

Task Description: Subcellular localization refers to the process of determining the specific location or compartment within a cell where a particular molecule or protein resides. This information is essential for understanding the function and behavior of molecules or proteins, as their subcellular locations often dictate their roles in cellular processes, interactions with other molecules, and influence on cellular functions. For example, proteins on the cell membrane generally have signaling and regulatory functions.

Task Impact: This task is closely related to protein functions and roles in biological processes.

Protein-Protein Interaction (PPI):

Task Formulation: two-class classification $(x_1, x_2) \mapsto \{0, 1\}$

Task Description: Protein-protein interaction predicts whether two proteins interact with each other.

Task Impact: This task is crucial for protein function understanding and drug discovery.

Protein Folding (Fold):

Task Formulation: $x \mapsto S$, where S is the 3d-structure of protein, including all coordinates of atoms.

Task Description: Protein Folding predicts the structure of protein sequences.

Task Impact: This task is known to be challenging, and requires elaborated knowledge of protein local and global structure to make atomic predictions.

Dataset Details: We report test results on CASP14 public available targets. We also remove all sequences over 800 tokens due to the computation memory limit. The reported targets are: T1024, T1025, T1026, T1027, T1028, T1029, T1030, T1031, T1032, T1033, T1034, T1035, T1036s1, T1037, T1038, T1039, T1040, T1041, T1042, T1043, T1045s1, T1045s2, T1046s1, T1046s2, T1047s1, T1047s2, T1048, T1049, T1050, T1051, T1053, T1054, T1055, T1056, T1057, T1058, T1059, T1060s2, T1060s3, T1062, T1063, T1064, T1065s1, T1065s2, T1066s1, T1066s2, T1067, T1068, T1069s1, T1069s2, T1070, T1071, T1072s1, T1072s2, T1073, T1074, T1075, T1076, T1077, T1078, T1079, T1082, T1083, T1084, T1085, T1086, T1087, T1088, T1089, T1090, T1092, T1093, T1094, T1095, T1096, T1098, T1099, T1100, T1101. The blue targets are from CASP14-FM set.

Table 8 gives the details of the datasets for these tasks.

De Novo Contact Prediction: We follow Chowdhury et al. (2022) to curate a de novo dataset of 108 proteins from Protein Data Bank (Bank). These proteins are originally designed de novo using computationally parametrized energy functions and are well-suited for out-of-domain tests. Note that different from orphan dataset, MSA can be built for this dataset, though showing a decline in quality.

F.2 HHblits Settings

For MSA datasets, We use HHblits (Remmert et al., 2012) to perform alignment. The commands for MSA dataset construction is:

```
hhblits -cpu $CPU_NUM -i $INPUT_FILE -d $DATABASE_DIR -oa3m $OUTPUT_FILE -n 1 -e 0.001
```

We also use HHblits to calculate E-value and determine whether we found homologous sequences in Figure 5 and §5.7 in the main paper. The commands for protein E-value calculation is:

```
hhalgn -i query.fasta -d retrieved.fasta -o output.aln -e 0.001
```

F.3 Model Hyperparameters

All self-reported models use the same truncation strategy and perform parameter searches on the learning rate among $[3e-8, 3e-6, 3e-5, 3e-4, 1e-3]$, warm-up rate among $[0, 0.08]$, seed among $[111, 222, 333, 444, 555, 666]$, and batch size among $[1, 2, 4, 8, 16]$. For evaluation, we choose the best-performing model on the validation set and perform prediction on the test set. The best performing hyperparameters could be found in the file:

```
./RSA-code/scripts/$MODEL_NAME/run_$TASK_NAME.sh
```

Also, code with download instructions for dataset and retrieval index is available in the supplementary.

F.4 RSA and Variants Implementation Details

F.4.1 Retriever Implementation Details

First, we calculate the ESM-1b embeddings of the 44 million sequences in Pfam-A 32.0. We use 16 V100 GPUs to calculate the embeddings in a day. A GPU as small as 3090 Ti would be enough, though it would take longer. Then, we adopt Faiss (Johnson et al., 2019b) indexing to accelerate the retrieval process by clustering the pre-trained dense vectors. In our implementation, we use the Inverted file with Product Quantizer encoding Indexing and set the size of quantized vectors to 64, the number of centroids to 4096, and the number of probes to 8. The construction of the Faiss index takes roughly 30 minutes using 0.5% randomly selected protein embeddings for index training. All embeddings as well as their id are subsequently added to the index.

During retrieval, for each query sequence, we first use ESM-1b to calculate its embedding, and then using this embedding, we query faiss to find the top N nearest neighbor of this embedding, getting the distance and sequence id of retrieved sequences. L2 distances are used to measure sequence similarity.

F.4.2 ProtBERT-RSA Architecture and Implementation

Here we provide the details for ProtBERT-RSA Architecture. An illustration of this process is also available in Figure 9. Note that in Step 2 retrieval of Faiss index could be further accelerated with GPU. In Step 4, the predictions of pairwise augmentation could be accelerated with batching on GPU, concurrently predicting k augmented sequences at the same time.

However, for large pre-trained models and when k is very large, the batch computation may exceed memory limit. In this case, we provide implementation for gradient accumulation, which calculates loss and gradients for individual prediction (predictions _{i}) and sum up the gradients with gradient accumulation. This implementation is a convex upperbound for the original loss function and we have validated its stability. This could also be implemented in batch size n , where each backward iteration calculates k/n retrieval augmentations, achieving trade-off between inference speed and memory limit.

```
Given query sequence $query, retrieval database $Faiss_Index, sequence database $Pfam, the number of retrieval $k, ProtBERT model $Model, and label $y.  
Step 1. embedding = ESM_1b(query)  
Step 2. distances, ids = Faiss_Index.retrieve(embedding, k)  
       retrieved_seqs = Pfam[ids]  
Step 3. predictions_i = Model([query, retrieved_seq]), i=1,2,..k  
Step 4. prediction = sum(predictions_i * softmax(distance_i))  
Step 5. loss = loss_function(prediction, y), perform training
```

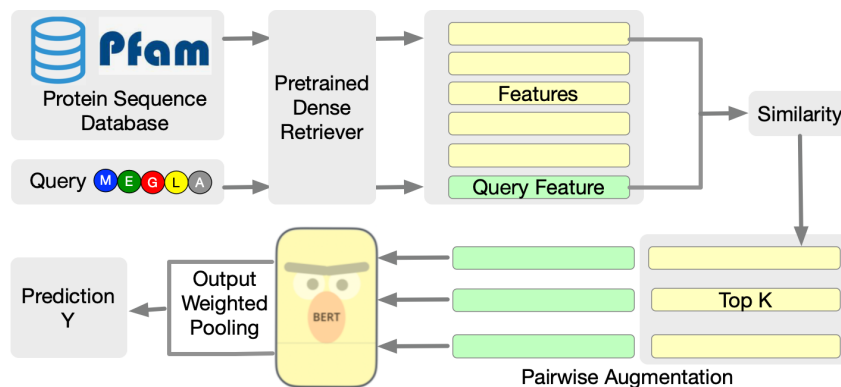


Figure 9: Detailed illustration of ProtBERT-RSA architecture.

F.4.3 RSA for Protein Folding

The major difference of RSA prediction for protein folding from other tasks is that we use a ranker to choose the final prediction rather than using weighted pooling. This is due to the misalignment of protein structures and simple weighting could result in averaging the structures of proteins with different angles of view. We train the ranker together with pTM-score loss (Lin et al., 2022) and contrastive loss on a subset of 1000 randomly chosen proteins from Protein Data Bank. These proteins are distinct from CASP14 test set. The ranker takes in the original structure prediction of the protein sequence and the k augmented predictions, and generate the highest ranking prediction as the final result. As current protein folding models are very large, we only provide zero-shot testing results on these pre-trained models, without further finetuning on our pipeline.

```
Given query sequence $query, retrieval database $Faiss_Index, sequence database
  $Pfam, the number of retrieval $k, Folding model $Model, Ranking model $Ranker
  and label $y.
```

```
Step 1. embedding = ESM_1b(query)
```

```
Step 2. distances, ids = Faiss_Index.retrieve(embedding, k)
  retrieved_seqs = Pfam[ids]
```

```
Step 3. predictions_i = Model([query, retrieved_seq]), i=1,2,..k
```

```
Step 4. prediction = Ranker(predictions_i), i=1,2,..k
```

Due to the different model architectures of ESMFold and AlphaFold, we explain in details the inference pipeline of Model([query, retrieved]).

ESMFold-RSA ESMFold is a **single** sequence protein folding model that consists of a protein representation model and a folding trunk based on the extracted representation. As illustrated in Figure 10(a), we concatenate query sequence with retrieved sequence and input them into the representation encoder. The encoder combines information from both query and retrieved sequence into query embedding via self-attention. Then we could use the pre-trained folding trunk to predict the structure of the query sequence. This pipeline could also be accelerated with batch prediction.

AlphaFold-RSA Different from ESMFold, AlphaFold encoder takes both single sequence representation and pairwise representation as input. Therefore, as shown in Figure 10(b), we generate the retrieved structure encoding with AlphaFold based on retrieved sequences, then we generate the structure of the query sequence based on the combination of single and pair representation. Note that we removed the template and MSA input in AlphaFold to ablation the effect of RSA.

F.4.4 Accelerated MSA

Accelerated MSA variant explores 165 substituting the discrete retrieval process in MSA with a dense retriever. We implement this method by first retrieving 500 sequences and then aligning these sequences with JackHMMer tool. Note that for most tasks we retrieve 500 sequences before alignment, as MSA Transformer can't take in many sequences. The command for aligning is:

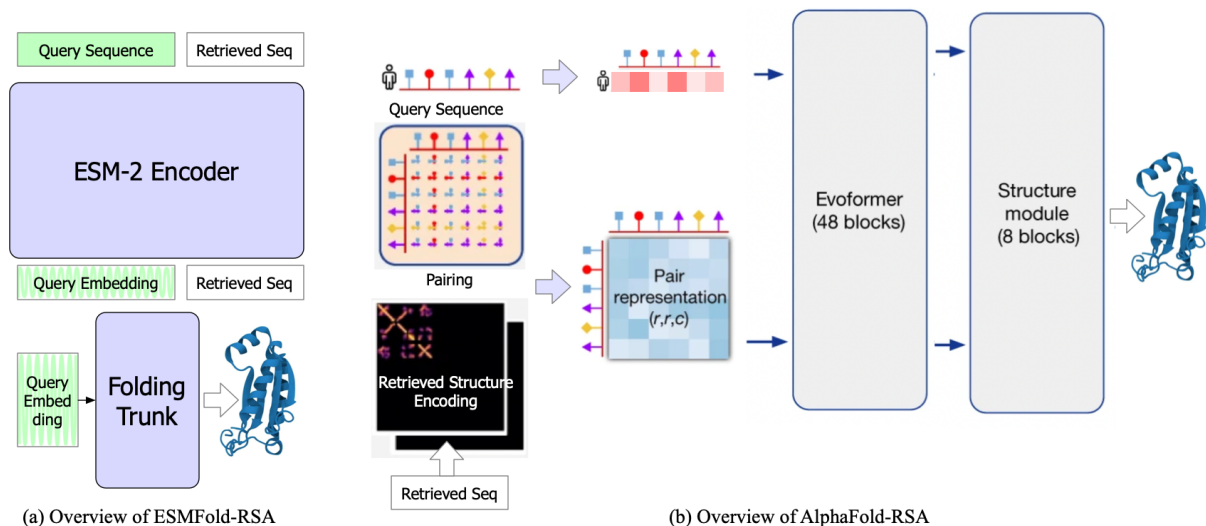


Figure 10: Illustration of the inference pipeline of RSA for Protein Folding

```
./jackhmmer -E 10.0 -A $aligned_file query.fasta retrieved.fasta
```

G Supplementary Experiment Analysis

G.1 Comparison of the Running time between RSA vs MSA

A severe speed bottleneck limits the use of previous MSA-based methods. In this part, we add analysis on database construction time as well as give details for inference time calculation. We calculate the total time used in each retrieval inference by summing: *alignment time* and *retrieval time*, as shown in Figure 11. Alignment time is the time used when finding MSA sequences through alignment and aligning found sequences with HHblits. Retrieval time is the time used during dense retrieval, including calculating the embedding of the query sequence with GPU. It is notable from the figure that alignment itself is a computationally costly procedure.

Also, MSA is limited by its cumbersome construction of retrieval HHM profile to perform HHM-HHM search. We follow the MSA custom database construction process in HHblits and compare with the construction time for RSA on a single V100 GPU (batch size=1) on a database of 10000 protein sequences. As shown in Figure 12, our method use only 10 minutes to finish the construction, though building a profile requires more than 3200 minutes.

G.2 Case Study

We cherry-picked one example of ProtBERT and ProtBERT-RSA on homology prediction (1195 class classification task) to showcase the interpretability as well as give intuition on our method. As shown in Figure 13, our method takes the original sequence as well as 16 retrieved sequences for prediction. After weighted summing of all predicted results, the prediction of probability on ground truth label increase and gives the correct prediction. We checked the most highly weighted (top 5) retrieved sequences, all five proteins are Colicins, which is a family under Toxins' membrane translocation domains. We can see from the case that weighting by distance helps the model focus on more similar retrieved instances.

We also provide two case studies on how RSA improves ESMFold. For target T1055, a DNA polymerase processivity factor, RSA retrieves *A0A1A8WBQ9_9APIC*, *A0A1Y4NGW6_9FIRM*, *A0A4V4NFM9_9ASCO*, *A0A1D3TXL7_9FIRM*, *A0A0V0QX86_PSEPJ*, *A9KN76_LACP7*, *A0A162CB07_9CRUS*, *A0A369KX60_9PROT,SKI2_SCHPO*, and the highest ranking augmentation prediction is from (T1055, A0A1A8WBQ9_9APIC). A0A1A8WBQ9_9APIC is a Merozoite surface protein. Merozoite surface protein 7 (MSP7) is a protein of the malaria parasite that has been found to be associated with processed fragments from the MSP1 protein in a complex involved in red blood cell

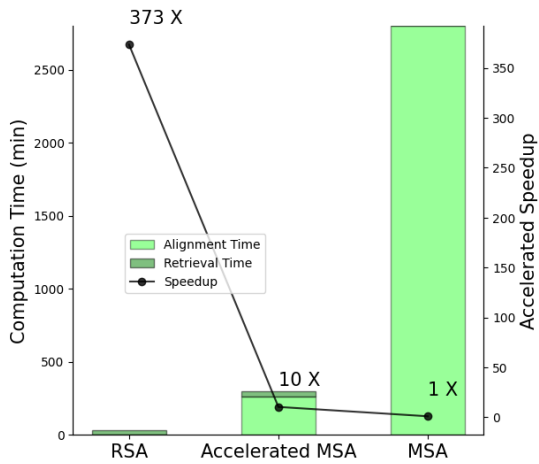


Figure 11: Illustration of speed up by RSA retrieval compared to MSA on secondary structure prediction dataset with 8678 sequences. Accelerated MSA refers to the MSA Transformer with MSA sequences retrieved by our RSA retriever.

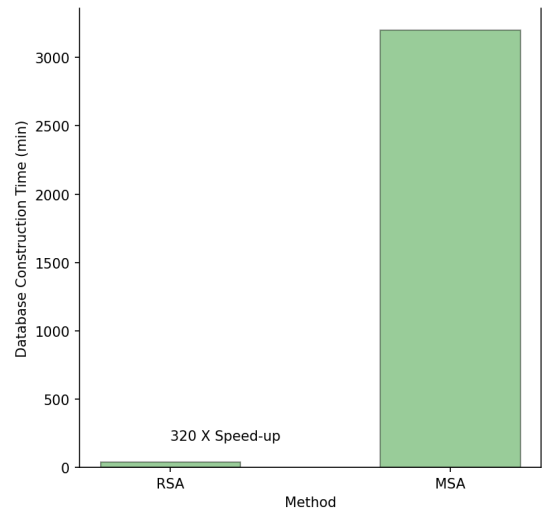


Figure 12: Illustration of speed up by RSA retrieval compared to MSA on database construction of 10000 protein sequences.

Method	Input	Prediction (P_{label} = Toxins' membrane translocation domains)
ProtBERT	sequence for d1cola_	0.0008
ProtBERT-RSA	d1cola_ and D2TV62_CITRI(weight=0.12) A0A3N1IY76_9ENTR(weight=0.10), C3K4R4_PSEFS(weight=0.07), A0A380QRI8_YERRU(weight=0.07), A0A6M8U9Q6_9GAMM(weight=0.07), B4F067_PROMH(weight=0.07), Q9I4Y4_PSEAE(weight=0.06), A0A0Q4MWP4_9GAMM(weight=0.06), C0AY95_9GAMM(weight=0.06), B2VE54_ERWT9(weight=0.05), A0A4P7L2K9_9GAMM(weight=0.05), A0A3N1J581_9ENTR(weight=0.05), A0A427K289_9GAMM(weight=0.05), A0A0Q4MTM7_9GAMM(weight=0.05), M1SHS7_MORM0(weight=0.04), B1VJ71_PROMH(weight=0.02),	P(d1cola_, D2TV62_CITRI) = 0.1315 P(d1cola_, A0A3N1IY76_9ENTR) = 0.1033 P(d1cola_, C3K4R4_PSEFS) = 0.2034 P(d1cola_, A0A380QRI8_YERRU) = 0.1034 P(d1cola_, A0A6M8U9Q6_9GAMM) = 0.1038 P(d1cola_, B4F067_PROMH) = 0.2132 P(d1cola_, Q9I4Y4_PSEAE) = 0.0003 P(d1cola_, A0A0Q4MWP4_9GAMM) = 0.1132 P(d1cola_, C0AY95_9GAMM) = 0.1938 P(d1cola_, A0A4P7L2K9_9GAMM) = 0.1211 P(d1cola_, A0A3N1J581_9ENTR) = 0.1034 P(d1cola_, A0A427K289_9GAMM) = 0.2309 P(d1cola_, A0A0Q4MTM7_9GAMM) = 0.1257 P(d1cola_, M1SHS7_MORM0) = 0.0000 P(d1cola_, B1VJ71_PROMH) = 0.0017 Prediction = 0.1063

Figure 13: Case study on homology prediction.

invasion. A0A1A8WBQ9_9APIC is a Merozoite surface protein C-terminal domain-containing protein that is related to DNA polymerase processivity factor through its requirement of a host factor, *E. coli* thioredoxin, in order to carry out its function. They also show similar structures with a TM-score of 0.42.

For target T1039, a virion RNA polymerase of crAss-like phage, RSA retrieves *A0A078ATM6_STYLE*, *A0A1D8P931_9FLAO*, *A0A363CW97_9PROT*, *D7JGI7_9BACT*, *A0A0B3VPN2_9FIRM*, *A0A1E4TQ27_PACTA*, *A0A1M6KY55_9FLAO*, *A0A1X7R9D3_9SACH*, *A0A0R1SCS6_9LACO*, *A0A367GMII_9SPHI*, *A0A2N1F639_9FLAO*, *A0A0D6TLE8_9FLAO*, *A0A3N4NFZ1_9FLAO*, *A0A1D2VEI9_9ASCO*, *A0A1L7I7H7_9FLAO*, *A0A1R0FA92_9RHIZ*. The highest ranking augmentation prediction is from (**T1039**, **A0A078ATM6_STYLE**). *A0A078ATM6_STYLE* is a COMM domain-containing protein 1. It has no distinct functional relationship with T1039, though the second chain of this protein has a similar structure to T1039, with a TM-score of 0.34.

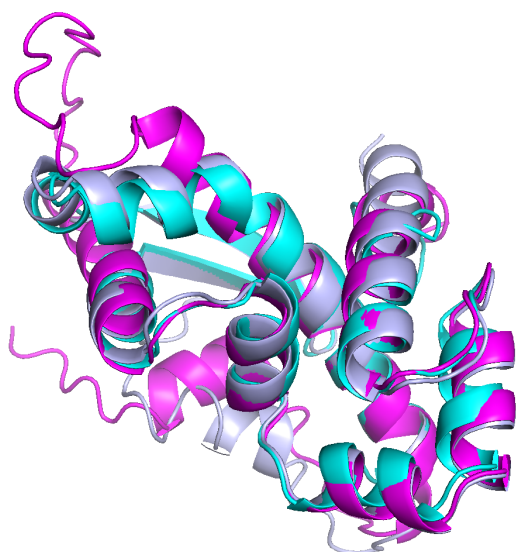


Figure 14: Structure Prediction for T1055, Cyan is the color for Ground truth. Pink is the color for ESMFold. Light purple is the color for ESMFold-RSA. The TM-score for ESMFold is 0.70, and the TM-score for ESMFold-RSA is 0.91.

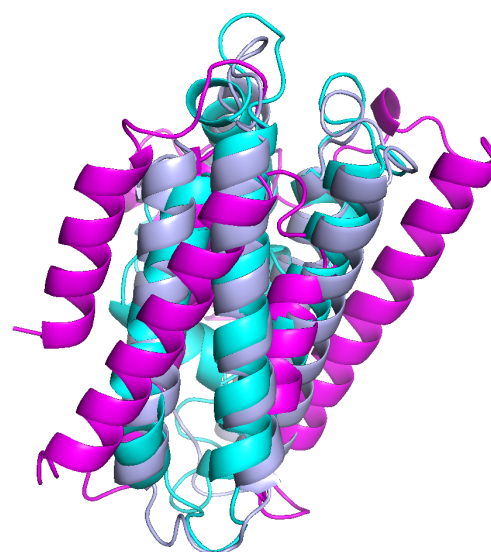


Figure 15: Structure Prediction for T1039, The TM-score for ESMFold is 0.61, and the TM-score for ESMFold-RSA is 0.29

G.3 Domain Adaptation Analysis

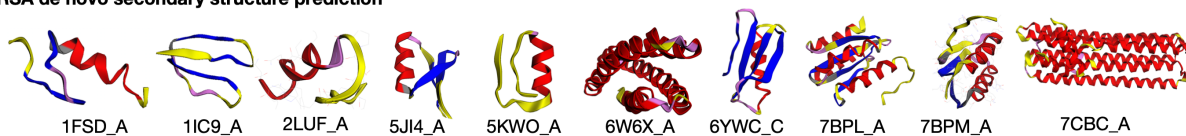
In this section, we perform additional analysis on the domain adaptation ability on secondary structure prediction tasks. We perform training on NetSurfP-2.0(Klausen et al., 2019) training set and test on two datasets with domain gaps. On CASP12, RSA marginally outperforms other baselines, as shown in Table 8. We also test on 10 de novo proteins (6YWC, 2LUF, 7BPM, 7BPL, 7CBC, 1FSD, 1IC9, 5JI4, 5KWO, 6W6X). Since we didn't find secondary structure labels for these proteins, we provide visualization in Figure 16, which shows that our model has an obvious overhead over MSA Transformer on predicting geometric components.

G.4 Comparison of Accelerated MSA vs MSA quality

Accelerated MSA performs worse than original MSA when directly applied to MSA Transformer, as well as AlphaFold. In this section, we showcase successful and failed cases in AlphaFold and compare the coverage of two kinds of MSA.

As shown in Figure 19, AlphaFold prediction is closely correlated to the coverage of MSA sequence. On cases where dense retriever fails to find a wide coverage of homologous sequences, AlphaFold performances drop starkly. Note that the MSA is implemented as ColabFold (Mirdita et al., 2022), using

RSA de novo secondary structure prediction



MSA Transformer de novo secondary structure prediction

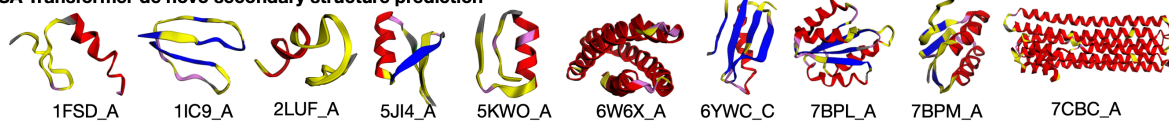


Figure 16: Prediction of Secondary Structure on De Novo Dataset. Each color corresponds to a different secondary structure.

Table 9: The domain adaptation performance of models on CASP12 secondary structure prediction.

Method	CASP12
ProtBERT	0.628
MSA Transformer	0.621
Accelerated MSA Transformer	0.620
RSA (ProtBERT backbone)	0.631

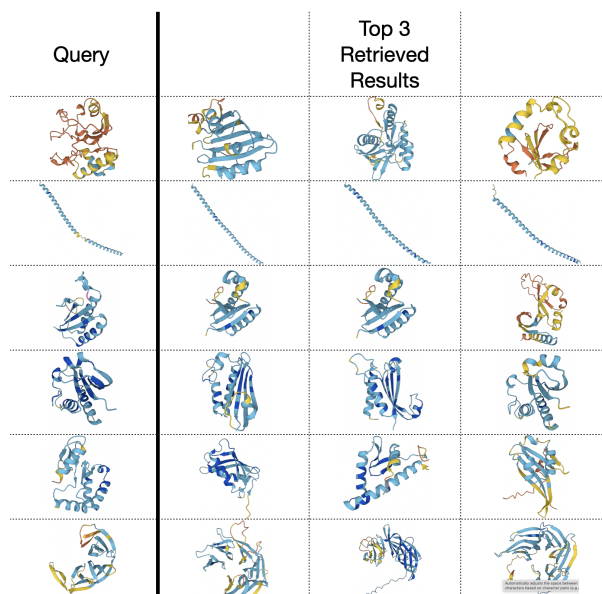


Figure 17: Query and Retrieved Sequence Structures

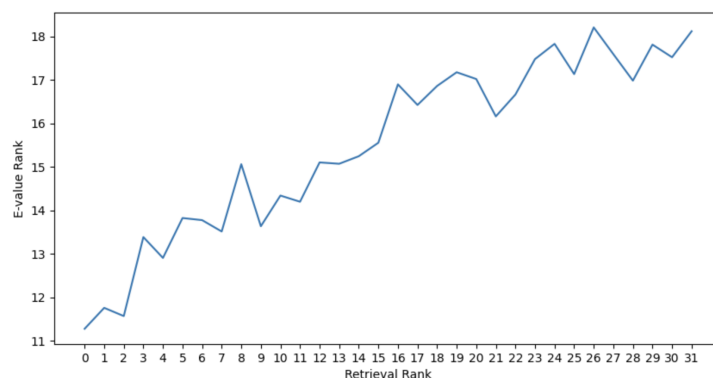


Figure 18: E-value rank against dense retrieval rank on in CB513 dataset.

Uniclust30 for MSA building, while our retriever database has a smaller coverage, using only Pfam database. Also we build accelerated MSA based on only top-500 sequences from retrieval.

G.5 Interpretability of RSA

In addition to analysis on interpretability in §5.7 in the main paper, we provide further analysis of the interpretability of RSA in terms of homology and structures.

Retrieval rank does not necessarily corresponds to the sequence closest to the query sequence token-wise. As shown in Figure 18, we calculate and rank the E-value of Top-32 retrieved protein sequences in CB513 dataset. We then calculate the average rank for the 1st, 2nd,... 32nd proteins in the dataset. It shows that the top-1 protein only has an average rank of 11, indicating that the retrieval rank does not necessarily corresponds to the sequence closest to the query sequence token-wise. Therefore, using dense retrieval, our retrieved results are diverse in sequences, though close to the query sequences in properties.

Visualization on Retrieval Structural Similarity As shown in Figure 17, we random picked a few more examples to illustrate the structural similarity between query protein and retrieval proteins. From the visualization, we can see that retrieved proteins exhibit similar structure or structure components, which could be used to boost structural knowledge.

G.6 ProteinChat: RSA Empowers ChatGPT on Protein Understanding

In the section, we point out a direction for our future work. RSA can not only be used on small-scale representation learning model, it can also augment large language models, e.g. ChatGPT. Currently, even GPT4 model shows limited understanding of biological sequences, especially long sequences like proteins. Galactica (Taylor et al., 2022) pretrains LLM on scientific text and biological sequences, however, calibration on the nature of these sequences are ubiquitous. Here, we propose a new method to enable ChatGPT to understand protein sequences, **ProteinChat**, by directly applying the protein dense retriever to ChatGPT. We implement ProteinChat by binding ChatGPT with four tools,

- RSA retriever, that finds homologous sequences to the query protein.
- Annotator, that annotates the homologous sequences in a protein online database, i.e. InterPro².
- Wikipedia Search, that facilitates more elaborated enquiry into the features of these proteins.
- PubMed Search, that facilitates elaborated enquiry into related scientific articles.

ProteinChat is implemented similar to ReAct agent (Yao et al., 2022), which views ChatGPT as an agent that uses chain-of-thoughts to break down a difficult task into parts that could be accomplished by available tools. Here we show several demos of ProteinChat.

Demo 1 (using Interpro + Wiki Search + RSA):

²<https://www.ebi.ac.uk/interpro/>

Question: Which family does this protein MSLNAEQKAKVVLEHGSSAHDGTGSTEVQVAL LTLRINDLQKHFLEHKKDHHSRRGLLRMVSQRRKLLDYLKRNISKYTDLIQSLGLRK belongs to? (Ground truth: Belongs to the universal ribosomal protein uS15 family.)

Thought: I need to find out the functional and evolutionary annotations of this protein.

Action: Protein Retrieval

Action Input: MSLNAEQKAKVVLEHGSSAHDGTGSTEVQVALLTLRINDLQKHFLEHKKDHHSRRGLLRMVSQRRKLLDYLKRNISKYTDLIQSLGLRK

Observation: Here are 4 homologs for this protein separated by ",":B4D3S0_9BACT, RS15_HAMD5, K6ZQA3_9ALTE, A0A1H6FC85_9GAMM

Thought: I need to use the entry labels to search the InterPro database.

Action: InterPro Annotation, Action Input: B4D3S0_9BACT

Observation: This is a 30S ribosomal protein S15. It is found in organism Chthoniobacter flavus Ellin428. It has molecular function that structural constituent of ribosome. It has a biological process translation. It has cellular component ribosome.

Thought: I need to search Wikipedia to find out which family this protein belongs to.

Action: Wikipedia Search, Action Input: 30S ribosomal protein S15

Observation: Page: **Ribosomal protein** Summary: A ribosomal protein (r-protein or rProtein) is any of the proteins that, in conjunction with rRNA, make up the ribosomal subunits involved in the cellular process of translation. E. coli, other bacteria and Archaea have a 30S small subunit and a 50S large subunit, whereas humans and yeasts have a 40S small subunit and a 60S large subunit.

Page: **40S ribosomal protein S30** Summary: 40S ribosomal protein S30 is a protein that in humans is encoded by the FAU gene.

Thought: I now know the answer.

Final Answer: The protein belongs to the 40S ribosomal protein S30 family.

Demo 2 (using Interpro + Wiki Search + RSA):

Question: Can this protein MPGVTVKDVNQQEFVRALAAFLKKSGLKVPPEWVDTVK-LAKHKELAPYDENWFYTRAASTARHLYLRGGAGVGSMTKIYGGQRNGVMPSHFSRG SKSVARRVLQALEGLKMVEKDQDGGKRLTPQGQRDLDRAGQVAAANKK binds with RNA? (Ground truth: Yes)

Thought: I need to find out what this protein does

Action: Protein Retrieval

Action Input: MPGVTVKDVNQQEFVRALAAFLKKSGLKVPPEWVDTVKLAKHKELAPYDENWFYTRAASTARHLYLRGGAGVGSMTKIYGGQRNGVMPSHFSRGSKSVARRVLQALEGLKMVEKDQDGGKRLTPQGQRDLDRAGQVAAANKK

Observation: Here are 4 homologs for this protein separated by ",":A0A315V8B8_GAMAF, A0A673GQT0_9TELE, A0A2K6SCP9_SAIBB, A0A674GDP2_TAEGU

Thought: I need to find out what this protein does

Action: InterPro Search, Action Input: A0A315V8B8_GAMAF

Observation: This is a 40S ribosomal protein S19. It is found in organism Gambusia affinis (Western mosquitofish).

Thought: I need to find out if this protein binds with RNA

Action: Wikipedia Search, Action Input: 40S ribosomal protein S19 binds RNA

Page: **40S ribosomal protein S19** Summary: 40S ribosomal protein S19 is a protein that in humans is encoded by the RPS19 gene. Page: **5S ribosomal RNA** Summary: The 5S ribosomal RNA (5S rRNA) is an approximately 120 nucleotide-long ribosomal RNA molecule with a mass of 40 kDa. It is a structural and functional component of the large subunit of the ribosome in all domains of life (bacteria, archaea, and eukaryotes), with the exception of mitochondrial ribosomes of fungi and animals. The designation 5S refers to the molecule's sedimentation velocity in an ultracentrifuge, which is measured in Svedberg units (S). Page: **Ribosomal protein** Summary: A ribosomal protein (r-protein or rProtein) is any of the proteins that, in conjunction with rRNA, make up the ribosomal

subunits involved in the cellular process of translation. E. coli, other bacteria and Archaea have a 30S small subunit and a 50S large subunit, whereas humans and yeasts have a 40S small subunit and a 60S large subunit.)

Thought: I now know the final answer.

Final Answer: Yes, 40S ribosomal protein S19 binds with 5S ribosomal RNA as part of the ribosomal subunits involved in the cellular process of translation

Demo 3 (using PubMed API + RSA):

Question: Identify the potential function of protein MSRTIFCTFLNKEAD-GLDFQLYPGELGKRIFNEISKEAWGQWMAKQTMLINEKKLNTMN-PDDRKLLEQEMVRFLFEGHDVHIDGYTPPEK and write a detailed protocol for measuring the function of this protein. (Ground truth: N/A)

Thought: I need to find what gene this protein is associated with

Action: Protein Retrieval

Action Input: MSRTIFCTFLNKEADGLDFQLYPGELGKRIFNEISKEAWGQWMAKQTM-LINEKKLNTMNPDDRKLLEQEMVRFLFEGHDVHIDGYTPPEK

Observation: Homologs of this protein has attributes: Iron_Traffic with importance 0.8876, E1_UFD with importance 0.1123

Thought: I need to find out more information about iron traffic

Action: PubMed Search, Action Input: iron traffic

Observation: Title: Moving Fe²⁺ from ferritin ion channels to catalytic OH centers depends on conserved protein cage carboxylates

Abstract: ... (too long, abbreviated)

Title: L-ferritin binding to scara5: a new iron traffic pathway potentially implicated in retinopathy

Abstract: ... (too long, abbreviated)

Title: Ceruloplasmin-ferroportin system of iron traffic in vertebrates

Abstract: ... (too long, abbreviated)

Thought: I now know enough to write a detailed protocol.

Final Answer: Based on the information gathered from the protein sequence and the relevant literature, it appears that the protein is potentially involved in iron trafficking. Iron trafficking is crucial for various cellular processes, and the ceruloplasmin-ferroportin system represents one of the major iron traffic pathways in vertebrates. Ferritin is an iron storage protein that has also been suggested as an iron carrier. The Scara5 receptor, which is specific for L-ferritin, has been found in various tissues, including the retina.

To measure the function of this protein, we can design a protocol that investigates the protein's role in iron transport and its interaction with other proteins involved in iron trafficking, such as ceruloplasmin, ferroportin, and Scara5. Here's a detailed protocol:

1. Protein expression and purification: Clone the gene encoding the protein of interest into an appropriate expression vector. Transform the vector into a suitable host, such as E. coli, and induce protein expression. Purify the expressed protein using affinity chromatography and other relevant purification techniques.

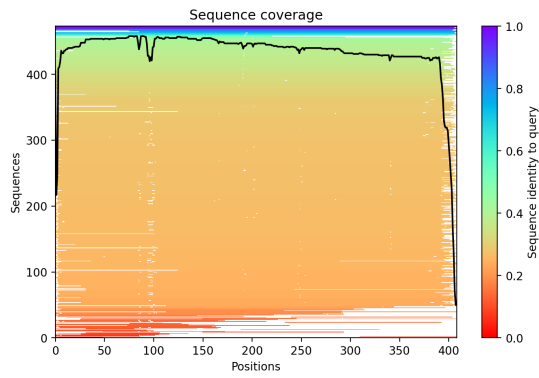
2. Iron-binding assay: Assess the ability of the purified protein to bind iron using a colorimetric iron-binding assay, such as the ferrozine assay. Measure the absorbance of the protein-iron complex at the appropriate wavelength to determine the iron-binding capacity of the protein.

3. Protein-protein interaction studies: Investigate the interaction between the protein of interest and other proteins involved in iron trafficking, such as ceruloplasmin, ferroportin, and Scara5. This can be done using techniques like co-immunoprecipitation, pull-down assays, or surface plasmon resonance (SPR). Cellular iron transport assay: Transfect cells with a plasmid encoding the protein of interest and perform an iron transport assay using radiolabeled iron (e.g., ⁵⁹Fe). Measure the rate of iron uptake and release by the cells expressing the protein of interest compared to control cells.

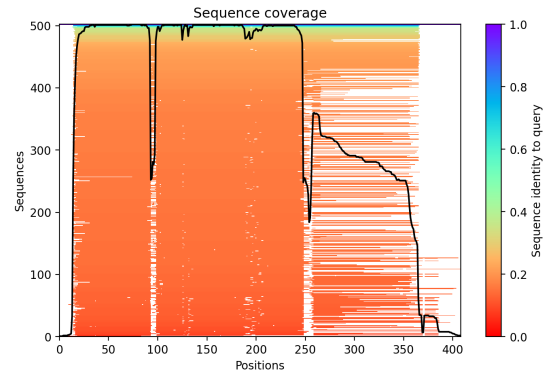
4. Functional studies in animal models: Assess the in vivo function of the protein by generating a knockout or transgenic animal model. Evaluate the iron status of the animals (e.g., serum iron, transferrin saturation, and ferritin levels) and examine the expression of other proteins involved in iron trafficking (e.g., ceruloplasmin, ferroportin, Scara5) in the knockout or transgenic animals compared to control animals.

5. Histological analysis: Examine the localization and distribution of the protein in various tissues, especially the retina, using immunohistochemistry. Determine the impact of the protein on retinal iron levels and assess the potential involvement of the protein in retinopathy.

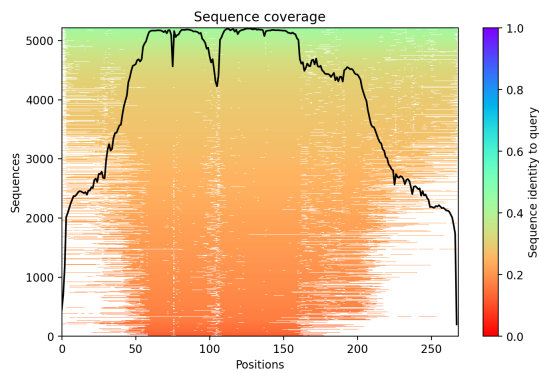
6. Data analysis: Analyze the data obtained from the various assays and experiments to draw conclusions about the function of the protein in iron trafficking and its potential role in retinopathy.



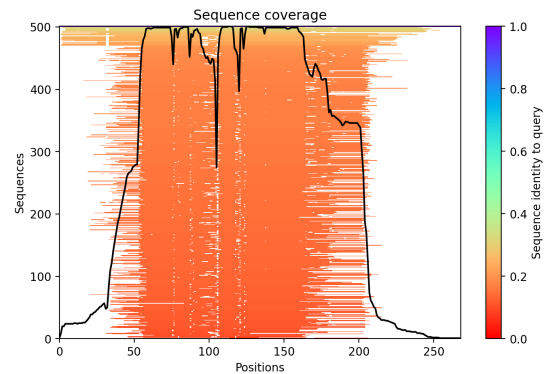
(a) T1024 MSA, AlphaFold prediction TM-score: 0.9842



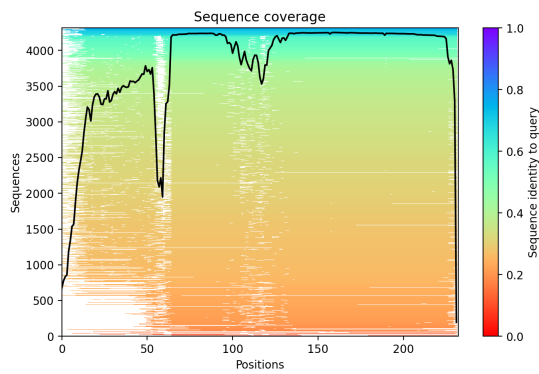
(b) T1024 Acc-MSA, AlphaFold prediction TM-score: 0.9897



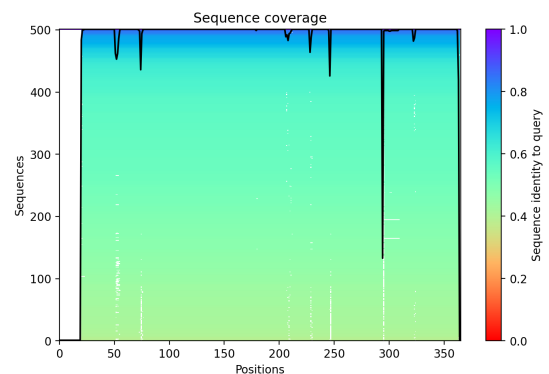
(c) T1025 MSA, AlphaFold prediction TM-score: 0.9229



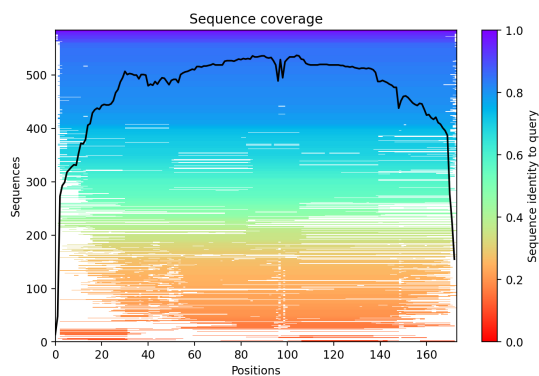
(d) T1025 Acc-MSA, AlphaFold prediction TM-score: 0.9204



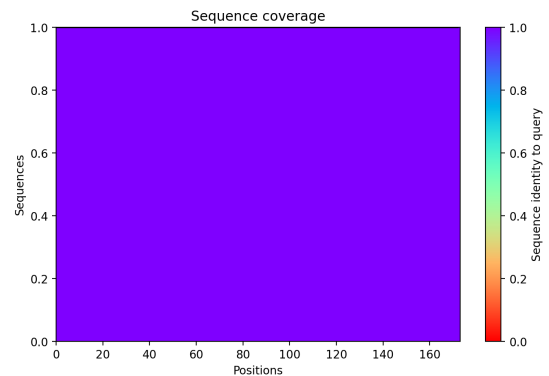
(e) T1047s1 MSA, AlphaFold prediction TM-score: 0.5020



(f) T1047s1 Acc-MSA, AlphaFold prediction TM-score: 0.4214



(g) T1045s2 MSA, AlphaFold prediction TM-score: 0.9356



(h) T1045s2 Acc-MSA, AlphaFold prediction TM-score: 0.2759

Figure 19: Visualization of the coverage rate of Accelerated MSA VS MSA.