

Improved the Protein Complex Prediction with Protein Language Models

Bo Chen^{1†}, Ziwei Xie^{2†}, Jiezhong Qiu³, Zhaofeng Ye³, Jinbo Xu^{2*} and Jie Tang^{1*}

¹Department of Computer Science and Technology, Tsinghua University, Beijing, China.

²Toyota Technological Institute at Chicago, Chicago, IL 60637, USA.

³Tencent, Shenzhen, China.

*Corresponding author(s). E-mail(s): jinboxu@gmail.com;
jietang@tsinghua.edu.cn;

Contributing authors: cb21@mails.tsinghua.edu.cn;
zwxie@ttic.edu; jiezhongqiu@tencent.com;
michaelzfy@tencent.com;

†These authors contributed equally to this work.

Abstract

AlphaFold-Multimer has greatly improved protein complex structure prediction, but its accuracy also depends on the quality of the multiple sequence alignment (MSA) formed by the interacting homologs (i.e., interologs) of the complex under prediction. Here we propose a novel method, denoted as ESMPair, that can identify interologs of a complex by making use of protein language models (PLMs). We show that ESMPair can generate better interologs than the default MSA generation method in AlphaFold-Multimer. Our method results in better complex structure prediction than AlphaFold-Multimer by a large margin (+10.7% in terms of the Top-5 best DockQ), especially when the predicted complex structures have low confidence. We further show that by combining several MSA generation methods, we may yield even better complex structure prediction accuracy than AlphaFold-Multimer (+22% in terms of the Top-5 best DockQ). We systematically analyze the impact factors of our algorithm and find out the diversity of MSA of

32 interologs significantly affects the prediction accuracy. Moreover, we show
33 that ESMPair performs particularly well on complexes in eucaryotes.

34 **Keywords:** Protein Complex Structure Prediction, Protein Language Model

35 1 Introduction

36 Most proteins function in a form of protein complexes[1–5]. Consequently,
37 obtaining accurate protein complex structures is vital to understanding how
38 a protein functions at the atom level. Experimental methods, such as X-ray
39 crystallography and cryo-electron microscopy, are costly and low-throughput,
40 and require intensive efforts to prepare samples for structure determination.
41 The computational methods, termed as protein complex prediction (PCP)
42 or protein-protein docking, is an attractive alternative for solving complex
43 structures. PCP takes sequences and/or the unbound structures of individ-
44 ual protein chains as inputs and then predicts the bound complex structures.
45 Traditional computational methods often rely on the global search paradigm,
46 such as fast-Fourier transform based methods like ClusPro [6], PIPER [7], and
47 ZDOCK [8] and Monte Carlo sampling-based methods like RosettaDock [9],
48 have been widely used in practice. These methods exhaustively search the con-
49 formation space of a complex, and optimize score functions to obtain the final
50 structures. Since the conformation space is large, these methods have to make
51 restrictive constraints on the search space in order to obtain results within
52 a reasonable amount of time. Typical constraints include reducing the search
53 resolutions, making the input monomers rigid bodies, and using score func-
54 tions that can be quickly evaluated [7, 8]. As a result, global search methods
55 have relatively low prediction accuracy and are used with more computation-
56 ally intensive local refinement methods to obtain higher resolution predictions
57 [10]. To date, PCP is still a fundamental and longstanding challenge in com-
58 putational structural biology [11, 12]. Various methods have been proposed
59 for PCP, but with limited accuracy. When only sequences are given as inputs,
60 PCP is even harder because the unbound structures of individual chains and
61 auxiliary information on the complex interfaces are unavailable.

62 In the last decades, deep learning has enabled substantial progress in quite
63 a few computational structural biology tasks, such as protein contact [13–15],
64 tertiary structure prediction [16–18], and cryo-electron microscopy structure
65 determination [19, 20]. Among these, co-evolution analysis based contact
66 prediction [18, 21, 22] and structure prediction [23, 24] have made sub-
67 stantial progress and demonstrated state-of-the-art accuracy for monomers.
68 These methods utilize the co-evolutionary information hidden in MSA to
69 infer inter-residue interactions or three-dimensional structures of the targets.
70 AlphaFold2 is the representative method, which has showed unparalleled accu-
71 racy in CASP14 [16]. Recently, AlphaFold-Multimer [25], a derived version

of AlphaFold2 for multimers, significantly outperforms prior protein complex prediction systems [6, 23, 24]. However, compared to the accuracy of AlphaFold2 [16] on folding monomers, the accuracy of AlphaFold-Multimer on predicting the protein complex structures is far from satisfactory. Its success rate is around 70% and the mean DockQ score is around 0.6 (medium quality judged by DockQ) [24]. The most important input feature to AlphaFold-Multimer is the multiple sequence alignment (MSA) [23, 24]. Compared with AlphaFold2 [16] that takes the MSA of a single protein as the input, AlphaFold-Multimer needs to build an MSA of interologs for protein complex structure prediction. However, how to construct such an MSA is still an open problem for heteromers. It requires the identification of interacting homologs in the MSAs of constituent single chains, which may be challenging since one species may have multiple sequences similar to the target sequence (paralogs). Several algorithms have been proposed to identify putative interologs from genome data, such as profiling co-evolved genes [26], and comparing phylogenetic trees [27]. Genome co-localization and species information are two commonly used heuristics to form interologs for co-evolution-based complex contact and structure prediction [25, 28]. Genome co-localization is based on the observation that, in bacteria, many interacting genes are coded in operons [29, 30] and are co-transcribed to perform their functions. However, this rule does not perform well for complexes in eukaryotes with a large number of paralogs, since it becomes more difficult to disambiguate correct interologs [28, 31]. The other phylogeny-based method for identifying interologs is first proposed in ComplexContact [28] and later similar ideas are adopted by AlphaFold-Multimer. This method first identifies groups of paralogs (sequences of the same species) from the MSA of each chain, then ranks the paralogs based on their sequence similarity to their corresponding primary chain, and last pairs sequences of the same species and with the same rank together. However, they are all hand-crafted approaches which merely take effects on the specific domains. In this paper, we instead investigate general and automatic algorithms for constructing MSAs of interologs for heterodimers effectively.

Representation learning via pre-training techniques has been prevailing in different applications [34–37]. Inspired by this, protein language models [38–40] (PLMs) have surged as the main regime for protein representation learning built on a large amount of protein sequences, which benefits downstream tasks, such as contact prediction [15, 39], remote homology detection [41, 42] and mutation effect prediction [43]. PLMs can comprehensively capture the biological constraints and co-evolutionary information encoded in the sequence, which is a plausible interpretation for their impressive performance on various downstream tasks than canonical methods relying on dedicated hand-crafted traits. To this, a natural question arises: ***Can we leverage the co-evolutionary information featured by PLMs to build effective interologs?***

In this paper, we mainly focus on *ab-initio* protein complex structure prediction, i.e., predicting the complex structure without prior information on the binding interfaces of the target complex. To our best knowledge, we are

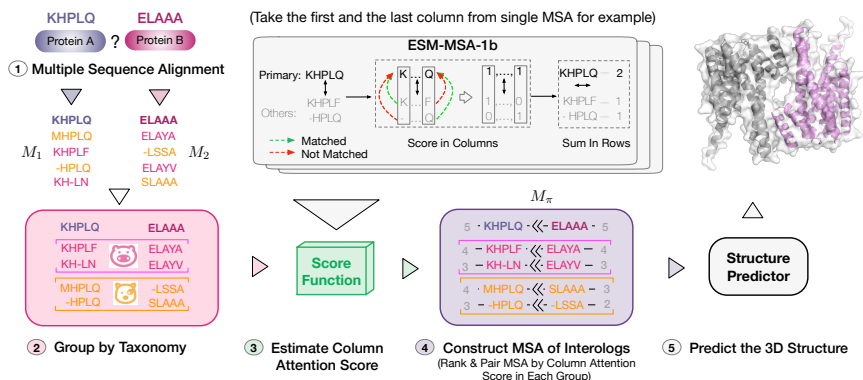


Fig. 1 Schematic illustration of ESMPair that builds interologs as the input to AlphaFold-Multimer. Given a pair of query sequences as input: 1) we first search the UniProt database [32] with JackHMMER [33] to generate the MSA for each query sequence, 2) sequences of the same taxonomy rank are grouped into the same cluster, 3) ESM-MSA-1b is applied to estimate the column attention score (ColAttn_score) between each sequence homolog of MSA with the query sequence. 4) One interolog is obtained by directly concatenating two matched sequence homologs. We match two sequence homologs of the same taxonomy group with similar attention scores from the two query sequences, 5) AlphaFold-Multimer takes the interolog MSA as input to predict the complex structure.

117 the first to propose a simple yet effective MSA pairing algorithm that uses
 118 the immediate output from protein language models to form joint MSAs, i.e.,
 119 MSA of interologs. In particular, we leverage column-wise attention scores from
 120 ESM-MSA-1b [39] to identify and pair homologs from MSAs of constituent
 121 single chains, coined as ESMPair. We conduct extensive experiments on three
 122 test sets, i.e., pConf70, pConf80, and DockQ49. Compared with previous
 123 methods, ESMPair achieves state-of-the-art structure prediction accuracy on
 124 heterodimers (+10.7%, +7.3%, and +3.7% in terms of the Top-5 best DockQ
 125 score over AlphaFold-Multimer on three test sets, respectively). Moreover, we
 126 find out that the ensemble strategies, which combine ESMPair with other MSA
 127 pairing methods, significantly improve the structure prediction accuracy over
 128 the standard single strategy. We further analyze the performance of complexes
 129 from eukaryotes, bacteria, and archaea, and find out ESMPair performs the
 130 best on eukaryotes for which identifying interologs is quite difficult [28, 31].
 131 Most strikingly, on a few targets where one of the constituent chains is from
 132 eukaryotes while the other is from bacteria, ESMPair considerably outper-
 133 forms other baselines (+25% in overall performance over AlphaFold-Multimer),
 134 which strongly demonstrates that the PLM-enhanced MSA pairing method
 135 is robust for targets from different superkingdoms. Then we posit that the
 136 diversity of interologs has a significant positive correlation with the predic-
 137 tion accuracy. Lastly, we explore other approaches that utilize the output of
 138 ESM-MSA-1b. For example, we take the cosine-similarity score between the
 139 sequence embeddings as the metric to build interologs, which performs on par
 140 with the default protocol used in AlphaFold-Multimer. Generally, ESMPair is

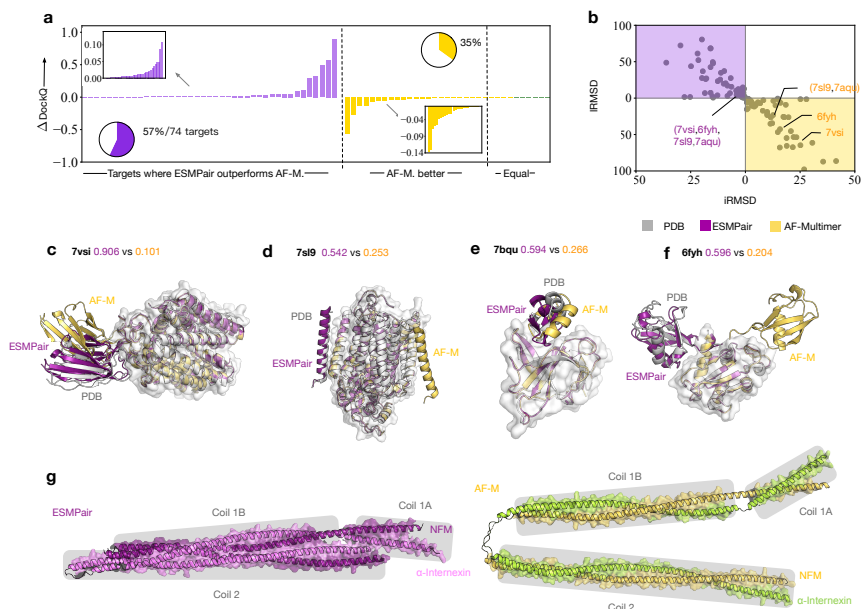


Fig. 2 Comparing ESMPair and AF-Multimer predictions on newly-release targets (a-f) and unresolved cases (g). a-f, The comparisons are made on the 74 targets whose release data is later on 2018-4-30. a, The bar charts show the relative performance gap between ESMPair and AF-Multimer on three categories: ESMPair outperforms AF-M.; AF-M. outperforms ESMPair; Equal performance. b, The interface and ligand RMSD distributions of predicted structures via ESMPair (Purple) and AF-M (Yellow). c-f, Four cases are further visualized. Among this, the ligand orientation are wrongly-predicted via AF-M. on 7VSI and 7AQU, while the binding site are wrongly-predicted by AF-M. on 7SL9 and 6FYH. g, The intermediate filament protein NFM-INA heterodimer structure predicted via ESMPair shows a four-helix bundle. The gray boxes show the interacting motifs of coil 1A, coil 1B and coil 2 of the two proteins.

141 the first simple yet effective algorithm that incorporates the strength of PLMs
 142 into tackling the issues of identifying MSA of interologs. We believe ESMPair
 143 will facilitate the fields of protein structure prediction which highly resorts to
 144 the co-evolution information hidden in MSA.

145 2 Results

146 In this section, we first briefly outline the framework of ESMPair for protein
 147 complex prediction (Section 2.1). Then, we discuss our proposed methods
 148 obtain better complex prediction accuracy than previous MSA pairing methods
 149 (Section 2.2). We find out the ensemble strategy showcase the excellent performance
 150 that the default single strategy (Section 2.3). We further quantitatively
 151 analyze several key factors and hyperparameters that may impact the performance
 152 of our method, and also explore the capability of different measurements
 153 to distinguish acceptable predictions from unacceptable ones (Section 2.4).

Table 1 DockQ scores and success rate of PLM-enhanced pairing methods and baselines. We report the average of Top-5 Best DockQ score, Top-1 Best DockQ score, and Success Rate (DockQ \geq 0.23) (%) on pConf70, Quality49, and pConf80 test sets. For one test target, we predicted 5 different structures using the five AlphaFold-Multimer models. Subscript in red represents the performance gain of our method over the default MSA pairing strategy in Alphafold-Multimer (%).

Methods	pConf70			Quality49			pConf80		
	Top5	Top1	SR	Top5	Top1	SR	Top5	Top1	SR
Non-Pairing Methods									
Block	0.199	0.179	30.4	0.212	0.194	49.0	0.351	0.319	51.2
Baseline Pairing Methods									
Genome	0.215	0.182	33.7	0.219	0.195	49.0	0.377	0.346	54.7
AF-Multimer	0.234	0.203	42.4	0.247	0.219	58.0	0.408	0.369	62.5
PLM-enhanced Pairing Methods									
InterLocalCos	0.218	0.180	33.7	0.236	0.210	52.3	0.389	0.353	56.5
InterGlobalCos	0.224	0.182	35.9	0.229	0.206	52.9	0.391	0.350	57.1
IntraCos	0.235	0.199	37.0	0.251	0.219	54.8	0.400	0.362	58.3
ESMPair	0.259	0.214	42.4	0.265	0.235	58.7	0.423	0.378	63.1
	(+10.7)	(+5.4)	(+0.0)	(+7.3)	(+7.3)	(+1.2)	(+3.7)	(+2.4)	(+1.0)

2.1 ESMPair overview

The overall framework of ESMPair is illustrated in Fig. 1 with the details in Methods. In complex structure prediction, predictors such as AlphaFold-Multimer make use of inter-chain co-evolutionary signals by pairing sequences between MSA of constituent single chains of the query complex. Formally, given a query heterodimer, we obtain individual MSAs of its two constituent chains, denoted as $M_1 \in \mathcal{A}^{N_1 \times C_1}$ and $M_2 \in \mathcal{A}^{N_2 \times C_2}$, where \mathcal{A} is the alphabet used by PLM, N_1 and N_2 are the number of the sequences in MSAs M_1 and M_2 , and C_1 and C_2 are the sequence length. The MSA pairing pipeline aims at designing a matching or an injection $\pi : [N_1] \rightarrow [N_2]$ between MSAs from each chain to build the MSA of interologs, dubbed as $M_\pi \in \mathcal{A}^{N \times (C_1 + C_2)}$, where N is the number of the sequence in the joint MSA. In practice, the MSA of interologs M_π is a collection of the concatenated sequence $\{\text{concat}(M_1[i], M_2[\pi(i)]) : i \in \mathcal{P}\}$, where \mathcal{P} is the indices of the sequences from M_1 that can be paired with any sequences from M_2 according to the matching pattern π . Then MSA of interologs is taken by predictors as input to predict the structure of the query heterodimer. Our aim is to leverage the superiority of PLMs to explore an effective matching strategy π that facilitates the protein complex structure prediction.

173 2.2 ESMPair outperforms other MSA pairing methods 174 on heterodimer predictions

175 **Overall evaluation.** For each test target we predict five 3D structures using
176 Alphafold-Multimer’s 5 models and then report the average of Top- k ($k=1, 5$)
177 Best DockQ score of the predicted structures and the corresponding success
178 rate (SR) in Table 1. Our method outperforms the other methods. To be spe-
179 cific, our method outperforms the AF-Multimer’s default MSA pairing strategy
180 on all three test sets (0.259 vs. 0.234 on pConf70, 0.423 vs. 0.406 on pConf80,
181 and 0.265 vs. 0.242 on Quality49, in term of Top-5 DockQ score). Our exper-
182 imental results confirm that our proposed column-wise attention based MSA
183 pairing method, denoted as ESMPair, is better than 1) the sequence similarity-
184 based method used in AF-Multimer, and 2) the cosine similarity-based method
185 based on the mixed noisy residue embedding, i.e., ESMPair(InterLocalCos),
186 ESMPair(InterGlobalCos), and ESMPair(IntraCos). Hereinafter, we abbrevi-
187 ate them as IntraLocalCos, InterGlobalCos, and InterCos.

188 Among all the MSA pairing methods, block diagonalization performs the
189 worst (-30% compared with ESMPair in terms of the average of Top-5 best
190 DockQ). The result indicates that the inter-chain co-evolutionary information
191 helps with complex structure prediction. Among MSA pairing baselines, AF-
192 Multimer surpasses genetic co-localization by a large margin (+12.8% Top-5
193 DockQ). All the proposed PLM-enhanced pairing methods substantially out-
194 perform the block diagonalization and the genetic-based methods. Even though
195 AF-Multimer may have overly optimistic performance using the default pair-
196 ing method since the training MSAs are built using it, IntraCos MSA pairing
197 method performs on a par with AF-Multimer, and ESMPair further exceeds
198 it by a large margin (+4.2~10.7% Top-5 DockQ score over three test sets).

199 **Intra-ranking methods are superior to inter-ranking ones both in
200 effectiveness and scalability.** From Table. 1, we can also see inter-
201 ranking methods like InterLocalCos and InterGlobalCos underperform the
202 intra-ranking ones, i.e. IntraCos and ESMPair. We speculate that as ESM-
203 MSA-1b pre-trains in the monomer data, it fails to directly capture the
204 underlying correlations across the constituent chains in the complex. Besides
205 heterodimers, when it extends to predict the structure of multimer with more
206 than two chains, intra-ranking strategies are the self-contained methods that
207 only need to rank the MSAs in each single chain, and then match MSA of the
208 same rank with other chains to build effective interologs with time complexity
209 of $O(N)$, where N is the depth of MSA. While the inter-pairing strategies suf-
210 fer from the exponential growth of combinations with increasing interacting
211 chains with the time complexity $O(N^r)$, where r is the number of chains in the
212 multimer. Thus, intra-ranking methods are more time-efficient and scalable
213 than inter-ranking ones.

214 **ESMPair performs better on low pConf targets.** As shown in Table. 1,
215 the performance gap between ESMPair and AF-Multimer becomes narrower
216 on pConf80 than on pConf70, with improvement ratio from 3.7% to 10.7%. To

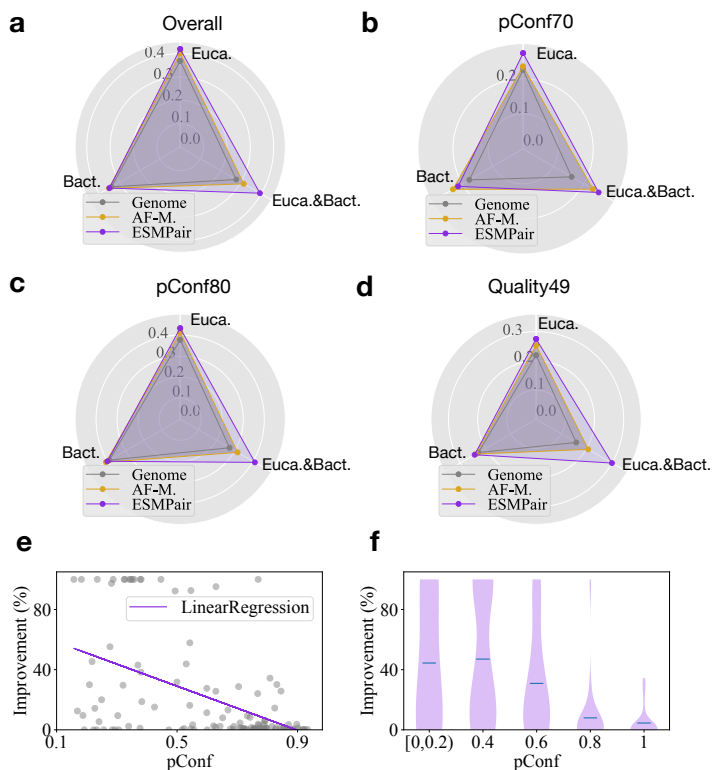


Fig. 3 Comparison of the prediction performance on different domains. a-d, We compare the DockQ score among ESMPair, AF-Multimer, and Genome on Eucaryote, Bacteria, and Eucaryote&Bacteria domains. The Euca.&Bact. is a special domain means the two constituent chains in the heterodimer belong to the two domains respectively. Specifically, the heterodimers of our dataset are from Eucaryotes, Bacteria, Viruses, Archaea, Eucaryotes;Bacteria respectively. We group the data from Bateria, Viruses, and Archaea as the Bateria domain. In all test sets, ESMPair significantly outperforms other two baselines on the Eucaryote targets. e-f, The negative correlations between the relative improvement between ESMPair and AF-Multimer.

217 take an in-depth analysis, we quantitatively analyze the correlations between
218 the predicted confidence score (pConf) estimated by AF-Multimer and the
219 performance gap of the average of Top-5 Best DockQ score between ESM-
220 Pair and AF-Multimer on Quality49, as illustrated in Fig. 3(e-f). The relative
221 improvement is negatively correlated (Pearson Correlation Coefficient is -0.49)
222 with the predicted confidence score. When pConf is less than 0.2, the relative
223 improvements even achieve 100%, while when pConf is more than 0.8, ESMPair
224 performs nearly on par with AF-Multimer. This is because AF-Multimer can
225 do well on a relatively easier target, it is very challenging to further improve it.

226 **ESMPair has the higher prediction accuracy on eucaryote targets.**
227 We further compare the DockQ distribution of ESMPair, AF-Multimer, and
228 Genome on three kingdoms, i.e. Eucaryote, Bacteria, and Eucaryote&Bacteria,

229 as shown in Fig. 3(a-d), we can see that ESMPair rivals the other two MSA
230 pairing methods on the Eucaryotes data by a large margin (0.420 for ESMPair
231 , 0.402 for AF-Multimer, and 0.369 for Genome on the overall data). As we all
232 know that it is notoriously difficult to identify homologous protein sequences
233 for the Eucaryotes data, ESMPair has a desirable property to build effective
234 interologs on the Eucaryotes. While in the Bacteria data, three strategies have
235 similar performance (around 0.35 on the whole data). Most strikingly, we find
236 ESMPair has an extraordinary performance on the Euca.&Bact. data over the
237 other two methods (0.394 for ESMPair, 0.314 for AF-Multimer, and 0.277 for
238 Genome on the overall data). We further check the performance gap for each
239 target from the Euca.&Bact. data. ESMPair performs significantly better on
240 the three out of six targets, 0.443 (ESMPair) versus 0.013 (AF-Multimer) on
241 5D6J, 0.289 versus 0.201 on 6B03, and 0.864 versus 0.854 on 7AYE. Besides,
242 ESMPair performs on par with AF-Multimer on the other three targets. These
243 results shed light on the robustness of protein language models. As PLMs are
244 pre-trained on billions of protein data [38–40], it can break the bottleneck
245 that other hand-crafted MSA pairing methods, such as genetic-based methods,
246 phylogeny-based methods, etc, which merely take effect in the specific domain.
247 While our proposed PLMs-enhanced methods can identify the co-evolutionary
248 signals effectively to build MSA of interologs across different superkingdoms.

249 **ESMPair outperforms AF-Multimer on the most of newly-released**
250 **targets.** We further select 74 targets that AF-Multimer does not train on [25],
251 i.e., the targets whose release date is later than 2018-4-30 from the test dataset.
252 Then we compare the performance of predicted structures on these targets
253 between ESMPair and AF-Multimer in Fig. 2. From Fig. 2(a), ESMPair out-
254 performs AF-Multimer on the most of targets (57%) with a relative larger
255 performance gap, while AF-Multimer outperforms ESMPair on fewer targets
256 (35%) with a relative lower gap. We further plot the distributions between
257 interface RMSD and ligand RMSD of predicted structures via ESMPair and
258 AF-Multimer in Fig. 2(b). The holistic distributions predicted by ESMPair
259 are closer to the origin of coordinates than that predicted by AF-Multimer,
260 which strongly proved ESMPair is superior to AF-Multimer on the predictions
261 of newly-released targets.

262 Furthermore, we show why ESMPair performs better than AF-Multimer
263 by analysing four PDB targets, 7VSI, 7AQU, 6FYH, and 7VSI. in Fig. 2(c-f).
264 Among these, 7VSI and 6FYH have a larger predicted iRMSD and IRMSD
265 variance by AF-Multimer, because AF-Multimer predicts the wrong binding
266 sites. While AF-Multimer predicts the right binding sites on 7SL9 and 7AQU
267 that have a smaller predicted iRMSD and IRMSD variance, it unfortunately
268 predicts the wrong ligand orientations. By contrast, our proposed ESMPair
269 correctly predicts the binding sites on the receptor and also places the ligand
270 in the approximately correct relative orientation.

271 To better illustrate the usage of ESMPair in predicting the protein com-
272 plexes without known resolved 3D structures, we inspected the intermediate
273 filament heterodimer formed between the neurofilament medium polypeptide

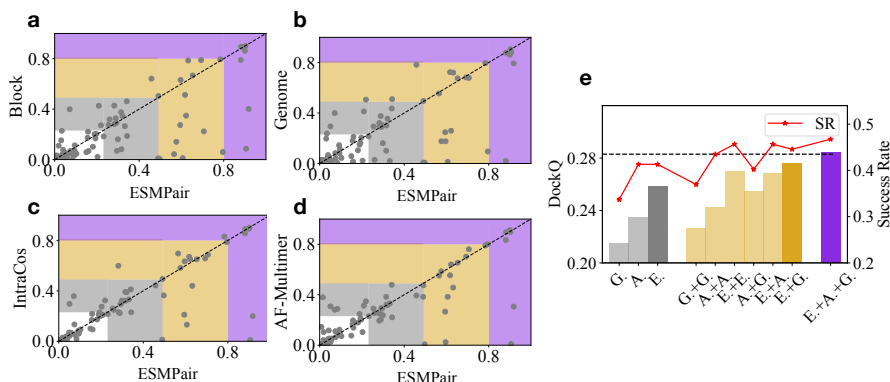


Fig. 4 The comparisons between ESMPair and four alternative MSA pairing approaches (a-d), and various ensemble strategy (e) on the targets from pConf70. (a-d), The coordinates of each point demonstrate the reported DockQ score of the target between ESMPair (x-axis) and other methods (y-axis). A point under the diagonal dash line implies ESMPair performs better than the compared method on this target. The highlight regions represent the incorrect (white), acceptable (grey), medium (yellow), and high-quality (purple) predicted models according to DockQ score. (e), The grey bars represent the performance of single strategies, where G. stands for Genome, A. is for AF-Multimer, and E. is for ESMPair. ESMPair is the best with 0.259 DockQ score and 42.4% Success Rate. The yellow bars show the ensemble performance of the two strategies. Among these, ESMPair + Genome performs the best with 0.277 DockQ score with 44.6% Success Rate. The purple bar implies the best performance about the ensemble of all the three strategies with 0.285 DockQ score with 46.8% Success Rate.

274 (NFM, UniProt ID P08553) and α -internexin (UniProt ID P46660), which is
 275 known to form an anti-parallel four-helix bundle[44, 45]. As shown in Fig. 2(g),
 276 both ESMPair and AF-Multimer correctly predict the three binding interfaces
 277 between the coil 1A, coil 1B and coil 2 motifs from NFM and α -internexin.
 278 However, ESMPair predicted the two coiled coils to pack as a four-helix bundle,
 279 which is consistent with the experimental evidences, while the AF-multimer
 280 predicted the two coiled coils to be separated. This case demonstrate the
 281 potential to apply ESMPair to model unresolved protein complexes.

282 2.3 Ensemble improves the prediction accuracy

283 From Fig. 4 (a-d), we found that different MSA pairing methods have their own
 284 advantages, even block diagonalization performs slightly better than ESMPair
 285 on about 30% targets, which implies that they can complement each other. To
 286 verify that, we combine ten models predicted by any two of the MSA paring
 287 methods, then we report the average of Top-5 Best DockQ score, as shown in
 288 Fig. 4 (e). The ensem strategies, i.e., the yellow and purple bars, significantly
 289 outperform the corresponding single strategy, i.e., the grey bars. Specifically,
 290 the performance of intra-ensemble strategies surpass the corresponding single
 291 strategy, for example, the DockQ score of ESMPair + ESMPair is 0.269 versus
 292 0.259 of ESMPair, which demonstrates that simply increasing the number of
 293 predictions of each model also benefits the structure prediction accuracy of

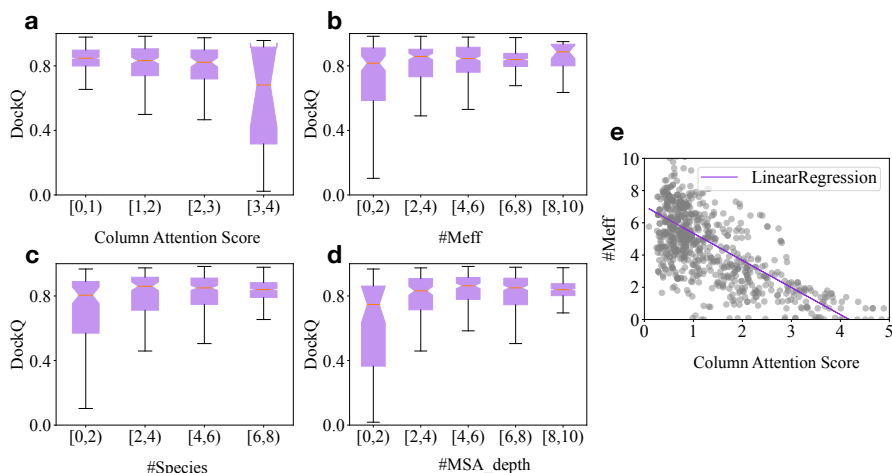


Fig. 5 Different factors affect the performance of structure prediction. The correlations between the average of Top-5 Best DockQ score (Y-axis) and (a) the column attention score predicted by ESM-MSA-1b, (b) the number of effective sequences measured by Meff, (c) the number of species, and (d) the depth of matched MSA. (e) The distribution of column attention score (X-axis) and the number of effective interologs in the paired MSA (Y-axis). The red curve is the visualization of the fitted linear regression model. The Pearson correlation coefficient is about -0.70, which strongly indicates that an increasing column attention score results in the decreasing number of effective interologs.

294 each target. Among the inter-ensemble strategies, ESMPair pluses any one of
295 the single strategy always have a better performance than the one without
296 ESMPair, for example, the SR of ESMPair + Genome is 44.6% versus 40.4%
297 of AF-Multimer + Genome. Finally, the ensemble of all three strategies, i.e.,
298 the purple bar, reaches the best performance with 0.285 DockQ score and
299 46.8% Success Rate, which motivates us that instead of merely using a single
300 strategy to build interologs, the ensemble MSA pairing strategy may be the
301 silver bullet to identify more effective interologs.

302 2.4 More analytic studies of ESMPair: key factors, 303 hyperparameters, and measurements to identify 304 high-quality predictions

305 In this part, we analytically and empirically investigate the inherent prop-
306 erties of ESMPair. Generally, we find out the diversity of the formed MSA
307 of interologs has a strong correlation with the performance of ESMPair.
308 Moreover, we study the effect of different layers of ESM-MSA-1b [39] on iden-
309 tifying homologs. Lastly, we demonstrate the predicted confidence score output
310 by AlphaFold-Multimer is a rational measurement to discriminate correct
311 predictions from incorrect ones.

312 **The diversity about MSA of interologs affects the predicted struc-
313 ture accuracy by ESMPair.** We investigate the connections between the

314 performance of ESMPair and some key factors of the formed MSA of interologs,
315 such as the column-wise attention score (i.e., ColAttn_score), the number of
316 effective sequences within MSA measured by Meff (i.e., #Meff), the number of
317 species (i.e., #Species), and the depth of MSA (i.e., MSA_Depth). To be spe-
318 cific, we predict 1,689 heterodimers sampled from PDB without filtering and
319 divide them into different regions according to the value of each factor. Notably,
320 for ColAttn_score, we average the score of each single chain in interolog, then
321 re-scaling it in the logarithm form, and then averaging ColAttn_score of all
322 interologs from the paired MSA as the final score of the target. For #Meff,
323 #Species, and MSA_Depth, we directly calculate the corresponding statistics
324 based on the interologs.

325 The correlations between DockQ score and each of above factors are illus-
326 trated in Fig. 5. #Meff, #Species, and MSA_Depth have a similar trend that
327 the predicted structure accuracy improves with the increasing of these factors.
328 It implies that MSA with more diversity represents the more co-evolutional
329 information that benefits structure predictions of AF-Multimer, which also
330 meets with previous insights[39]. Moreover, the increasing ColAttn_score
331 results in the decreasing structure prediction accuracy. Considering the self-
332 attention mechanism in the protein language model, given a sequence as the
333 query, the self-attention mechanism aims at identifying the sequence with high
334 homology affinity, i.e., the sequence with a high similarity score [15]. Therefore,
335 a large ColAttn_score indicates the MSA with a low #Meff, which potentially
336 results in an inaccurate structure prediction. To justify our speculation, we
337 explicitly characterize the dependency between ColAttn_score and #Meff, as
338 shown in Fig. 5 (e). ColAttn_score has shown a negative correlation to the
339 #Meff, with the Pearson correlation coefficient of -0.70, which elucidates that
340 a higher ColAttn_score reflects MSA with lower sequence diversity.

341 **ESMPair built on the last few transformer Layers has the bet-**
342 **ter performance.** As ESMPair leverages the column-wise attention output
343 by ESM-MSA-1b[39] to rank and match interologs, how do the column-wise
344 attention weight matrices by different transformer layers affect the efficacy of
345 ESMPair? To answer this, we use the DockQ score of predicted structures as
346 the metric to measure the quality of the input interologs built by ESMPair, as
347 shown in Supplement Fig. A2. ESMPair that based on the attention output of
348 layer 6 (0.258 DockQ score and 40.2% Success Rate), layer 7 (0.249 and 43.0%),
349 and AVG (0.262 and 42.2%) perform better than other layers. Overall, the AVG
350 aggregation of all the layers is relatively superior to others, thus we use AVG
351 as the default setting of ESMPair. What's more, ESMPair which built on the
352 last few layers (6-12th) identifies homologous sequences more precisely than
353 the former layers (1-5th). The phenomenon is consistent with the empirical
354 insights about how to effectively fine-tune the pre-trained language models in
355 the downstream tasks: the last few layers are the most task-specific, while the
356 former layers encode the general knowledge of the training data[46–48], thus
357 only aggregating latter layers may exploiting more homologous information
358 form MSAs. We leave this in future work.

359 **Predicted confidence score as an indicator to distinguish acceptable**
360 **models.** Practically, besides the substantial improved DockQ performance
361 through ESMPair, it is also vital to figure out how to identify the correct mod-
362 els (DockQ \geq 0.23) from incorrect ones [24]. To achieve this, we also predict all
363 the 1,689 heterodimers via ESMPair, then we apply: 1) the predicted Confid-
364 ence Score (pConf), 2) Interface pTM (ipTM), 3) predicted TM-score (pTM),
365 and 4) the number of contacts between residues from two chains (the distance
366 of C_β atoms in the residues from different chains within 8 Å) (Contacts) as the
367 metric to rank models, as shown in Supple. Fig. A1. From Fig. 1(a), we find
368 both pConf and ipTM are capable of distinguishing acceptable models from
369 unacceptable ones with AUC of 0.97. pTM has a worse performance with AUC
370 of 0.85, as pTM is used as the pessimistic predictor to measure the predicted
371 structure accuracy of each single chain, it ignores the interactions between
372 chains. Contacts merely count the number of interacting residues from dif-
373 ferent chains, which hardly indicates the accuracy of the predicted structure.
374 pConf and ipTM both consider the structure in both the single chain and inter-
375 faces, which are considerate indicators to validate the quality of the predicted
376 structure. We further quantify the interplays between pConf and DockQ score
377 of the predicted structure, as shown in Fig. 1(b), which further confirms the
378 strong correlations between pConf and the structure prediction accuracy.

379 3 Methods

380 In this part, we introduce the framework of our proposed PLMs-enhanced MSA
381 pairing method, i.e., ESMPair. Besides, we explore other promising alternative
382 pairing methods built on PLMs, such as InterGlobalCos, InterGlobalCos, and
383 IntraCos. The overall framework of ESMPair is illustrated in Fig. 1.

384 3.1 The PLM-enhanced MSA pairing pipeline

385 Previous efforts [38–40] have confirmed that protein language models (PLMs)
386 can characterize the co-evolutionary signals and biological structure con-
387 straints encoded in the protein sequence. Moreover, the MSA-based PLMs [15,
388 39] further explicitly capture the co-evolutionary information hidden in MSAs
389 via axial attention mechanisms [49, 50]. In light of this, we adopt the state-
390 of-the-art MSA-based PLM, i.e., ESM-MSA-1b [39], as the basis to explore
391 how to utilize them to build rational MSA of interologs to improve the protein
392 complex prediction based on AlphaFold-Mutimer [25].

Column Attention (ESMPair). The column attention weight matrix, which is calculated via each column of MSA via ESM-MSA-1b, can be treated as the metric to measure pairwise similarities between aligned residues in each column. Formally, for each chain, we have the MSA $M \in \mathcal{A}^{N \times C}$. The collections of column attention matrices are denoted as $\{A_{lhc} \in \mathbb{R}^{N \times N} : l \in [L], h \in [H], c \in [C]\}$, where L is the number of layers in PLM, H is the number of attention heads of each layer, and C is the sequence length, i.e., the number of residues of each sequence. We first symmetrize each column attention matrix,

and then aggregate the symmetrized matrices along the dimension of L , H and C to obtain the pairwise similarity matrix among the sequences of MSA, denoted as $S \in \mathbb{R}^{N \times N}$ (Eq.(1)). S is symmetric and its first row $S_1 \in \mathbb{R}^{1 \times N}$ can be viewed as measuring similarity scores between the query sequence and other sequences in the MSA,

$$S = \underset{l \in [L], h \in [H], c \in [C]}{\text{AGG}} \{A_{lhc} + (A_{lhc})^\top\}, \quad (1)$$

where \top represents the transpose operation and AGG is an entry-wise aggregation operator such as entry-wise mean operation $\text{MEAN}(\cdot)$, sum operator $\text{SUM}(\cdot)$, etc. Unless otherwise specified, AGG is specified as $\text{SUM}(\cdot)$ in this paper.

The MSA pairing strategy is specified as follows, for a query heterodimer, we first obtain S_1 of individual MSAs of constituent single chains. Then we group sequences from the MSA by their species, and rank sequences according to their similarity score of S_1 in each MSA, respectively. Finally, the sequences of each MSA with the same rank in the same species group are concatenated as interologs.

Cosine Similarity. The cosine similarity measurement has been thoroughly explored by pre-train language models [51, 52]. Intuitively, as PLMs generate residue-level embeddings for each sequence in the MSA, the sequence embedding can be directly obtained by aggregating all the residue embeddings in the sequence. Thus we can calculate the cosine similarity matrix between each sequence to measure their pairwise similarities.

To be more specific, we specify two MSA pairing strategies, i.e., Intra-ranking (IntraCos) and Inter-pairing, based on the cosine similarity measurement between sequence embeddings as follows:

Intra-ranking (IntraCos). Firstly, for all sequences from a given MSA $M \in \mathcal{A}^{N \times C}$, we obtain a collection of residue-level embedding $\{E_{ln} \in \mathbb{R}^{C \times d} : l \in [L], n \in [N]\}$, where d is the embedding dimension. For sequence $n \in [N]$, we can obtain its sequence-level embeddings $E_n = \text{AGG}_{l \in [L], c \in [C]}(E_{lnc})$ by aggregating over all layers L and all residues C , where $E_n \in \mathbb{R}^d$. Then we compute cosine similarities between the query sequence embedding, E_1 , and other sequence embeddings, $\{E_n, \text{ where } n \neq 1\}$, in the MSA to obtain the pairwise similarity score matrix (IntraCosScore) $S_1 \in \mathbb{R}^{1 \times N}$. After that, we build interologs like ESMPair does.

Inter-ranking. Instead of ranking sequences in each MSA and matching sequences of the same rank, here we directly compute the similarity score matrix between sequences from different MSAs. Formally, given two MSAs $M_1 \in \mathcal{A}^{N_1 \times C_1}$ and $M_2 \in \mathcal{A}^{N_2 \times C_2}$, we obtain two individual collections of sequence embeddings $\{E_n^{(1)} : n \in [N_1]\}$ and $\{E_n^{(2)} : n \in [N_2]\}$. The inter-chain cosine similarity matrix is denoted by $B \in \mathbb{R}^{N_1 \times N_2}$, where $B_{ij} = \cos(E_1[i], E_2[j])$. Without loss of generality, we assume $N_i \leq N_j$, we propose two algorithms to build interologs as follows:

- 429 1. **Global Maximization Optimization (InterGlobalCos).** We for-
430 malize the pairing problem as a maximum-weighted bipartite matching
431 problem. The weighted bipartite $G = (V, E)$ is constructed as follows:
432 sequences from individual MSAs of two chains form the set of vertices in
433 G , i.e., $V^{(1)} = \{M_i^{(1)} \in \mathcal{A}^{C_1} : i \in [N_1]\}$, $V^{(2)} = \{M_j^{(2)} \in \mathcal{A}^{C_2} : j \in [N_2]\}$,
434 and $V = V^{(1)} \cup V^{(2)}$. There are no edges among sequences from the same
435 chain MSA, thus $V^{(1)}$ and $V^{(2)}$ are two independent sets. There is an
436 edge e_{ij} between $M_i^{(1)}$ and $M_j^{(2)}$ if these two sequences are from the same
437 species; the weight associated with e_{ij} is B_{ij} . An optimal MSA match-
438 ing pattern can be obtained by Kuhn-Munkres (KM) algorithm[53] in the
439 polynomial time.
- 440 2. **Local Maximization Optimization (InterLocalCos).** KM algorithm
441 finds a global optimal solution. However, as suggested by [54], in each
442 species, the sequence that is most similar to the query sequence may be
443 more informative, while other sequences that are less similar may add
444 noises. Thus we propose a greedy algorithm that focuses on pairs that
445 have high similarity scores with the query sequence. We iteratively select
446 a pair of sequences (i, j) that have the largest score in B among sequences
447 that have not been selected before until reaching a pre-defined maximal
448 number of pairs.

449 **Complex structure prediction of heteromers with more than two**
450 **different chains.** The proposed methods, such as ESMPair and IntraCos,
451 can be easily extended to build MSA of interologs for heteromers with more
452 than two different chains. In practice, we can rank the MSAs in each query
453 sequence by the similarity matrix obtained by the corresponding metric, then
454 we match them of the same rank in each species to build effective interologs.

455 3.2 Settings

456 **Evaluation metric.** We evaluate the accuracy of predicted complex struc-
457 tures using DockQ [55], a widely-used metric in the computational structural
458 biology community. Specifically, for each protein complex target, we calculate
459 the highest DockQ score among its top- N predicted models selected by their
460 predicted confidences from AlphaFold-Multimer. We refer to this metric as the
461 best DockQ among top- N predictions.

462 **Datasets.** In order to investigate how improving pairing MSAs can improve
463 the performance of AlphaFold-Multimer, we construct a test set satisfying the
464 following criteria:

- 465 1. There are at least 100 sequences that can be paired given the species
466 constraints.
- 467 2. The two constituent chains of a heterodimeric target share $< 90\%$
468 sequence identity.

469 We select heterodimers consisting of chains with 20~1024 residues (due to
470 the constraint of ESM-MSA-1b and also ignore peptide-protein complex), and

471 the overall number of residues in a dimer is less than 1600 (due to GPU mem-
472 mory constraint) from Protein Data Bank (PDB), as accessed on March 3, 2022.
473 We use the default AlphaFold-Multimer MSA search setting to search the
474 UniProt database [32] with JackHMMER [33], which is used for MSA pairing.
475 We also search the Uniclust30 database [56] with HHblits [57], which is used for
476 monomers, i.e., block diagonal pairing. We further select those heterodimers
477 with at least 100 sequences that can be paired by AlphaFold-Multimer’s default
478 pairing strategy. We define two dimers as at most $x\%$ similar, if the maximum
479 sequence identity between their constituent monomers is no more than $x\%$.
480 Overall, we select 801 heterodimeric targets from PDB that are at most 40%
481 similar to any other targets in the dataset, and satisfy the aforementioned two
482 criteria. Then we use AlphaFold-Multimer (using the default MSA matching
483 algorithm) to predict their complex structures. Based on their predicted con-
484 fidence scores (pConf) or DockQ scores, 92 targets with their pConf less than
485 0.7 are denoted as the pConf70 test set. We select 0.7 as the low confidence
486 cutoff based on our fitted logistic regression models over 7,000 DockQ and
487 pConf pairs, because the conditional probability of the model having medium
488 or better quality given pConf equals 0.7 is slightly greater than 0.5 (around
489 0.6), while the probability is less than 0.5 if pConf equals 0.6. For more com-
490 parisons, we also select 0.8 as the cutoff, which results in the pConf80 test set
491 of 168 targets, and 155 targets with their predicted DockQ scores less than
492 0.49 are denoted as the DockQ49 test set.

493 **Baselines.** Several heuristic MSA pairing strategies have been developed for
494 protein complex contact and 3D structure prediction [17, 23].

495 *Phylogeny-based method.* The strategy is first proposed in ComplexCon-
496 tact [28] for complex contact prediction and is widely adopted by the
497 community. AlphaFold-Multimer employed a similar strategy. This strategy
498 first groups sequences in an MSA by their species and then ranks sequences of
499 the same species by their similarity to the query sequence. When there is more
500 than one sequence in a species group, it joins two sequences of the same rank
501 within the same species group to form an interolog. AlphaFold-Multimer uses
502 this strategy and shows state-of-the-art accuracy in complex structure predic-
503 tion [25]. Practically, we run the implementation code of Alphafold-Multimer
504 following the default setting of official repertory¹. Notably, we only evaluate the
505 unrelaxed model without the template information for the time efficiency[16].

506 *Genetic distances.* In bacteria, interacting genes sometimes are co-located
507 in operons and co-transcribed to form protein complexes [58]. Consequently,
508 we can detect interologs by the genetic distance of two genes. This strategy
509 pairs sequences of the same species based on the distances of their positions
510 in the contigs, which are retrieved from ENA. In our implementations, given
511 a sequence from the first chain, we pair it with the sequence from the second
512 chain that is closest to it in terms of genetic distance. If there are more than
513 one closest sequence, we select the one that has the lowest e-value to the query

¹<https://github.com/deepmind/alphafold>

514 sequence of the second chain; the e-value is calculated by the MSA search
515 algorithm used to construct the chain MSA.

516 *Block diagonalization.* This strategy pads each chain sequence with gaps to
517 the full length of the complex [23]. Therefore, each sequence in the constructed
518 joint MSA, except for the query sequence, will include non-gap tokens in
519 exactly one chain and gap tokens in other chains. By sorting sequences in the
520 joint MSA, we can make non-gap tokens to appear only in the diagonal blocks,
521 thus this strategy is termed as block diagonalization. In our implementations,
522 given a sequence from the first (second) chain, we append (prepend) non-gap
523 tokens to it until the number of non-gap tokens equals the length of the second
524 (first) chain.

525 **Running environment.** We conduct the experiments on an Enterprise Linux
526 Server with 56 Intel(R) Xeon(R) Gold 5120 CPU @ 2.20GHz, and a single
527 NVIDIA Tesla V100 SXM2 with 32GB memory size.

528 4 Conclusion & Discussion

529 This paper explores a series of simple yet effective MSA pairing algorithms
530 based on pre-trained protein language models (PLMs) for constructing effective
531 interologs. To our best knowledge, this is the first time that PLMs are used to
532 construct joint MSAs. Experimental results have confirmed the proposed ESM-
533 Pair significantly outperforms the state-of-the-art phylogeny-based protocol
534 adopted by AlphaFold-Multimer. What's more, ESMPair performs particularly
535 better on targets from eukaryotes which are hard to be predicted accurately
536 by AF-Multimer. We further confirm that, instead of using the conventional
537 single strategy to build interologs, the ensemble MSA pairing strategy can
538 largely improve the structure prediction accuracy. Generally, ESMPair has a
539 profound impact on biological applications depending on the high-quality MSA.
540 In the future, we will continue to explore more potential ways to leverage the
541 advantages of PLM in building and choosing MSA. We also looking forward to
542 applying our proposed methods to improve current MSA-based applications.

543 **Limitations.** In this paper, we merely consider how to build effective
544 interologs for heterodimers, which broadly benefits biological applications
545 depending on the high-quality MSA, such as the complex contact predic-
546 tion [59, 60], complex structure prediction discussed in this paper, etc.
547 However, there also have a large proportion of homodimers in biological assem-
548 blies. As it is trivial to build interologs for them, how to select high-quality
549 MSA for homodimers is a more challenging yet important question. Previ-
550 ous work [39, 54] has an empirical insight that instead of using the full MSA
551 searched from the protein sequence database, we can select a few high-quality
552 MSA following some promisings, such as using the MSA maximizing the
553 sequence diversity [39], or choosing the MSA owning the largest sequence sim-
554 ilarity with the primary sequence [54]. To date, few efforts have systematically
555 investigated the MSA-selection problem. We leave this for future work.

556 As we propose a series of MSA paring methods built on the output of
557 PLMs, the representation ability of the PLMs directly affects the performance
558 of our proposed methods. In this paper, we choose the state-of-the-art pro-
559 tein language model so far, i.e., ESM-MSA-1b [39], to support our algorithms.
560 However, it is always worth exploiting the potential correlations between differ-
561 ent PLM configurations and the performance of our proposed PLM-enhanced
562 methods to identify effective interologs.

563 Although ESMPair has advantages over the default strategy adopted by
564 AF-Multimer in identifying MSA of interologs, their success rate is similar.
565 After a deep analysis, we observe ESMPair outperforms AF-Multimer most
566 in acceptable cases ($\text{DockQ} \geq 0.23$), however it is notoriously difficult for
567 ESMPair to improve DockQ score of unacceptable cases to be acceptable
568 (Only 3% targets). As we follow the pipeline of the complex structure pre-
569 diction via AF-Multimer (Fig. 1), thus the limited ability of AF-Multimer
570 becomes the bottleneck of the performance of ESMPair. Nevertheless, the
571 above extensive experimental results have proved ESMPair consistently out-
572 performs AF-Multimer despite AF-Multimer having an inductive training
573 bias towards its default MSA pairing strategy. From the training process of
574 AF-Multimer, we know that the performance of structure prediction highly
575 depends on the quality of the input MSA. In light of this, we assume that if
576 AF-Multimer can fine-tune, or totally train from scratch based on ESMPair's
577 MSA pairing method, the accuracy of structure predictions may be further
578 improved. Moreover, compared with the conventional MSA pairing method
579 that only uses a single strategy to identify interologs, the ensemble strategy has
580 shown superior performance both in DockQ score and Success Rate without
581 fine-tuning AF-Multimer. We assure that the ensemble strategy proposes a new
582 perspective on how to comprehensively exploit the co-evolutionary patterns
583 among MSA, thus further having a wide impact on the biological algorithms
584 resorting to the input MSA.

585 **5 Author Contributions**

586 B.C proposed the main idea, conducted the main experiments, and wrote ini-
587 tial manuscript. Z.W.X, J.Z.Q, and Z.F.Y collected the experimental data,
588 designed experiments, and wrote the initial manuscript. J.B.X and J.T gave
589 the detailed instructions and refined the manuscript.

590 **6 Data Availability**

591 Data that involved in this work can be obtained from Github:
592 <https://github.com/allanchen95/ESMPair>.

593 **7 Code Availability**

594 The code of this study can be obtained from GitHub:
595 <https://github.com/allanchen95/ESMPair>.

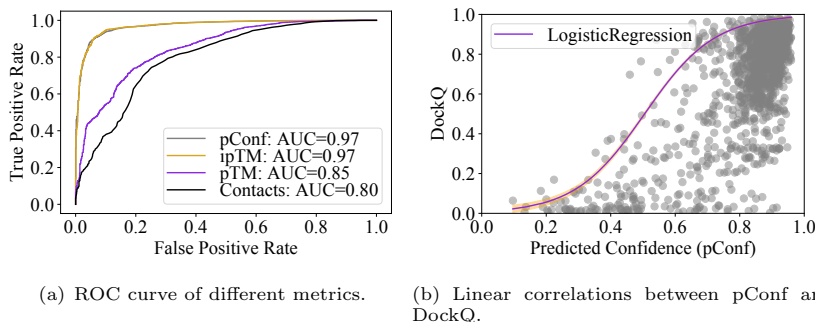


Fig. A1 Different metrics assessment. **a.** ROC curve of different metrics of distinguish acceptable cases ($\text{DockQ} \geq 0.23$) predicted by ESMPair. **b.** The distribution of predicted confidences (pConf, x-axis) and DockQ scores (left y-axis). And the conditional probability of the prediction having $\text{DockQ} \geq 0.23$ given pConf. The red curve is the visualization of the fitted logistic regression model.

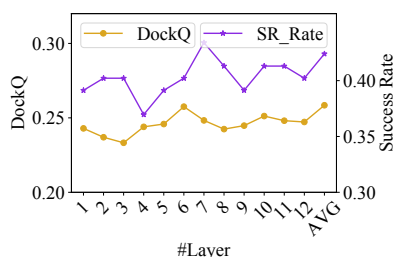


Fig. A2 The average of Top-5 Best DockQ scores of ESMPair based on the different layers of ESM-MSA-1b on the pConf70 dataset. AVG means that ESMPair is based on the column-wise attention matrix by averaging the one generated from all the twelve transformer layers.

596 Appendix A Supplement Material

597 **The Number of Effective Interlogs (Meff).** It counts the number of non-
598 redundant interlogs in an MSA, which measures the amount of homologous
599 information. Here we use the toolkit from RaptorX² to estimate the value of
600 Meff. Specifically, we set 70% sequence identity as the cutoff to judge if two
601 interlogs are redundant or not. If the number of interlogs (including itself)
602 similar to interlog i is n_i , then the weight of interlog i is $1/n_i$. Finally, Meff
603 is calculated by summing the weight of all interlogs.

604 **Supplement Experiments.** We conduct some additionally experiments
605 listed here.

²<https://github.com/j3xugit/RaptorX-3DModeling>

Table A1 The comparisons of Top-5 Best DockQ scores between *ESMPair* (EP) and *AF-Multimer* (AFM) on all test cases.

PDBID	EP	AFM	PDBID	EP	AFM	PDBID	EP	AFM
1AOK	0.343	0.221	1CXZ	0.012	0.012	1FAP	0.095	0.098
1GKA	0.565	0.552	1H2L	0.457	0.632	1H2M	0.458	0.394
1KMI	0.622	0.64	1MAF	0.322	0.316	1NW9	0.606	0.384
1S6C	0.232	0.23	1T6B	0.123	0.094	1UZX	0.576	0.586
1VGO	0.617	0.554	1WRD	0.79	0.784	1YHN	0.076	0.064
2AKA	0.04	0.015	2F8E	0.104	0.019	2H7Z	0.693	0.702
2VRW	0.918	0.306	3AU4	0.281	0.071	3BEG	0.209	0.313
3CPJ	0.283	0.25	3OUN	0.04	0.018	3PIN	0.074	0.111
3PZD	0.051	0.308	3TUF	0.654	0.653	3V2A	0.332	0.332
3V6B	0.796	0.801	3VF0	0.04	0.038	3VMF	0.26	0.254
3W5K	0.044	0.052	3ZN1	0.896	0.882	4CZZ	0.9	0.912
4GFT	0.15	0.076	4HUK	0.027	0.052	4JEG	0.282	0.194
4LC9	0.035	0.027	4N3B	0.29	0.295	4P3Y	0.008	0.007
4TU3	0.005	0.005	4WM0	0.291	0.255	4YBH	0.006	0.005
4YNO	0.092	0.084	4YOC	0.288	0.084	4YXC	0.156	0.166
5FOY	0.004	0.005	5H3J	0.238	0.299	5JJW	0.158	0.197
5JQY	0.341	0.479	5JZU	0.319	0.321	5KVM	0.189	0.248
5LN1	0.049	0.052	5N7E	0.048	0.04	5OJ6	0.49	0.384
5TAR	0.009	0.008	5TQB	0.104	0.203	5TVQ	0.007	0.007
5WHZ	0.043	0.117	5WWL	0.475	0.491	5XEQ	0.009	0.009
5XLN	0.315	0.383	5YY0	0.139	0.16	5Z51	0.024	0.024
5Z5K	0.058	0.053	6ABO	0.013	0.025	6B03	0.289	0.201
6DXO	0.197	0.203	6FYH	0.024	0.596	6HG4	0.044	0.044
6HTF	0.103	0.072	6JT1	0.295	0.292	6KIP	0.497	0.007
6LOJ	0.015	0.38	6M49	0.179	0.181	6N89	0.031	0.034
6POG	0.019	0.007	6QU1	0.037	0.039	6THL	0.01	0.01
6TTT	0.903	0.898	6UML	0.348	0.219	6W38	0.013	0.01
6WCW	0.193	0.214	6YT3	0.761	0.706	7A8T	0.016	0.004
7BQU	0.594	0.266	7CEG	0.007	0.006	7K2V	0.01	0.004
7LY5	0.008	0.145	7MRS	0.834	0.881	7RSI	0.034	0.034
7SL8	0.013	0.004	7VSI	0.01	0.906	1BQN	0.152	0.154
1KTZ	0.923	0.925	1M4U	0.019	0.019	1TXQ	0.025	0.025
2IWT	0.858	0.9	2Q2E	0.373	0.36	2QNA	0.504	0.506
2X1X	0.842	0.836	2XJY	0.707	0.697	2Y0I	0.734	0.745
2Y48	0.908	0.881	2ZUP	0.675	0.658	3AV0	0.484	0.376
3C5X	0.871	0.89	3D13	0.681	0.688	3EUJ	0.669	0.716
3LBX	0.672	0.677	3MCA	0.527	0.52	3N40	0.238	0.269
3NQU	0.853	0.854	3O1H	0.374	0.043	3OG6	0.835	0.824
3OJA	0.481	0.507	3ZYI	0.577	0.583	4DBG	0.814	0.822
4DSS	0.01	0.01	4F3L	0.163	0.163	4F7G	0.715	0.66
4LD3	0.922	0.917	4LJO	0.113	0.73	4OL0	0.586	0.582
4P2A	0.778	0.627	4PW9	0.678	0.732	4RGW	0.916	0.922
4RSI	0.741	0.747	4UN2	0.844	0.819	4WND	0.779	0.765
4XXB	0.872	0.88	4Y5O	0.639	0.612	5BQC	0.085	0.085
5C46	0.006	0.006	5C58	0.66	0.721	5CHL	0.733	0.76
5D6J	0.443	0.013	5KP6	0.856	0.851	5ME5	0.949	0.952
5NRO	0.826	0.834	5Y94	0.775	0.802	5W83	0.791	0.825
5YVI	0.917	0.912	5Z2W	0.787	0.805	5ZRZ	0.611	0.884
6AKM	0.244	0.813	6EC0	0.795	0.762	6EG0	0.69	0.666
6FKM	0.692	0.651	6G4J	0.773	0.714	6IRE	0.438	0.422
6IRT	0.502	0.5	6IWS	0.472	0.47	6JZE	0.932	0.934
6L5K	0.633	0.658	6LZ0	0.631	0.604	6OBP	0.466	0.454
6OD1	0.448	0.34	6Q00	0.905	0.903	6S0A	0.602	0.634
6SF1	0.88	0.879	6UUI	0.779	0.786	6WO1	0.846	0.809
6ZPH	0.265	0.264	7AYE	0.856	0.864	7DCR	0.785	0.775
7JW7	0.025	0.008	7KNT	0.64	0.647	7LVS	0.754	0.748
1A6U	0.334	0.335	1ARO	0.108	0.112	1CC1	0.329	0.33
1F45	0.445	0.457	1G4U	0.315	0.312	1H2A	0.32	0.32
1HTR	0.329	0.329	1I79	0.329	0.328	1JEQ	0.334	0.332
1KA9	0.331	0.327	1MHM	0.325	0.325	1NT2	0.257	0.258
1U0S	0.333	0.334	1V18	0.592	0.513	1WQ1	0.327	0.328
2I07	0.277	0.445	2Q5W	0.316	0.313	2QK7	0.04	0.04
2QSF	0.269	0.274	2RD7	0.269	0.275	3A2F	0.188	0.187
3C7N	0.175	0.177	3LQC	0.475	0.482	3NUH	0.338	0.366
3NVM	0.233	0.232	3P71	0.266	0.265	3SUS	0.007	0.835
3U73	0.289	0.291	3WCY	0.465	0.447	3WVN	0.019	0.019
4C9B	0.305	0.306	4CRW	0.504	0.485	4DEY	0.328	0.327
4EHP	0.491	0.429	4GMN	0.128	0.13	4HG6	0.353	0.37
4KHA	0.383	0.385	4MRT	0.331	0.333	4N6R	0.32	0.319
4RCA	0.303	0.241	4RS1	0.33	0.33	4U1C	0.321	0.321
4YC7	0.304	0.311	4YL8	0.537	0.531	4ZN3	0.442	0.329
5CM2	0.333	0.332	5HPK	0.44	0.433	5LOW	0.299	0.3
5OW0	0.324	0.323	5VPA	0.532	0.535	5YQZ	0.31	0.31
5YWR	0.794	0.803	6GK2	0.309	0.309	6INE	0.267	0.28
6MTL	0.203	0.156	6NDU	0.818	0.81	6OVM	0.281	0.284
6PFJ	0.324	0.325	6TX3	0.317	0.31	6UCC	0.825	0.821
6YXQ	0.481	0.477	7AX1	0.291	0.294	7BY2	0.401	0.423
7NKZ	0.335	0.335	7SL9	0.253	0.542			

References

- 606
- 607 [1] Jones, S., Thornton, J.M.: Principles of protein-protein interactions.
608 Proceedings of the National Academy of Sciences **93**(1), 13–20 (1996)
- 609 [2] Liddington, R.C.: Structural basis of protein-protein interactions. Protein-
610 Protein Interactions, 3–14 (2004)
- 611 [3] Sharan, R., Suthram, S., Kelley, R.M., Kuhn, T., McCuine, S., Uetz, P.,
612 Sittler, T., Karp, R.M., Ideker, T.: Conserved patterns of protein interac-
613 tion in multiple species. Proceedings of the National Academy of Sciences
614 **102**(6), 1974–1979 (2005)
- 615 [4] Tuller, T., Atar, S., Ruppim, E., Gurevich, M., Achiron, A.: Common and
616 specific signatures of gene expression and protein–protein interactions in
617 autoimmune diseases. Genes & Immunity **14**(2), 67–82 (2013)
- 618 [5] Pržulj, N., Malod-Dognin, N.: Network analytics in the age of big data.
619 Science **353**(6295), 123–124 (2016)
- 620 [6] Kozakov, D., Hall, D.R., Xia, B., Porter, K.A., Padhorny, D., Yueh, C.,
621 Beglov, D., Vajda, S.: The cluspro web server for protein–protein docking.
622 Nature protocols **12**(2), 255–278 (2017)
- 623 [7] Kozakov, D., Brenke, R., Comeau, S.R., Vajda, S.: Piper: an fft-based
624 protein docking program with pairwise potentials. Proteins: Structure,
625 Function, and Bioinformatics **65**(2), 392–406 (2006)
- 626 [8] Pierce, B.G., Wiehe, K., Hwang, H., Kim, B.-H., Vreven, T., Weng, Z.:
627 Zdock server: interactive docking prediction of protein–protein complexes
628 and symmetric multimers. Bioinformatics **30**(12), 1771–1773 (2014)
- 629 [9] Lyskov, S., Gray, J.J.: The rosettadock server for local protein–protein
630 docking. Nucleic acids research **36**(suppl.2), 233–238 (2008)
- 631 [10] Desta, I.T., Porter, K.A., Xia, B., Kozakov, D., Vajda, S.: Performance
632 and its limits in rigid body protein-protein docking. Structure **28**(9),
633 1071–1081 (2020)
- 634 [11] Keskin, O., Gursoy, A., Ma, B., Nussinov, R.: Principles of protein-
635 protein interactions: what are the preferred ways for proteins to interact?
636 Chemical reviews **108**(4), 1225–1244 (2008)
- 637 [12] Nooren, I.M., Thornton, J.M.: Diversity of protein–protein interactions.
638 The EMBO journal **22**(14), 3486–3492 (2003)
- 639 [13] Billings, W.M., Morris, C.J., Della Corte, D.: The whole is greater than its
640 parts: ensembling improves protein contact prediction. Scientific Reports

641 **11**(1), 1–7 (2021)

642 [14] Singh, J., Litfin, T., Singh, J., Paliwal, K., Zhou, Y.: Spot-contact-lm:
643 improving single-sequence-based prediction of protein contact map using
644 a transformer language model. *Bioinformatics* **38**(7), 1888–1894 (2022)

645 [15] Zhang, H., Ju, F., Zhu, J., He, L., Shao, B., Zheng, N., Liu, T.-Y.: Co-
646 evolution transformer for protein contact prediction. *Advances in Neural*
647 *Information Processing Systems* **34** (2021)

648 [16] Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger,
649 O., Tunyasuvunakool, K., Bates, R., Židek, A., Potapenko, A., *et al.*:
650 Highly accurate protein structure prediction with alphafold. *Nature*
651 **596**(7873), 583–589 (2021)

652 [17] Baek, M., DiMaio, F., Anishchenko, I., Dauparas, J., Ovchinnikov, S., Lee,
653 G.R., Wang, J., Cong, Q., Kinch, L.N., Schaeffer, R.D., *et al.*: Accurate
654 prediction of protein structures and interactions using a three-track neural
655 network. *Science* **373**(6557), 871–876 (2021)

656 [18] Roy, R.S., Quadir, F., Soltanikazemi, E., Cheng, J.: A deep dilated con-
657 volutional residual network for predicting interchain contacts of protein
658 homodimers. *Bioinformatics* **38**(7), 1904–1910 (2022)

659 [19] Si, D., Moritz, S.A., Pfab, J., Hou, J., Cao, R., Wang, L., Wu, T.,
660 Cheng, J.: Deep learning to predict protein backbone structure from
661 high-resolution cryo-em density maps. *Scientific reports* **10**(1), 1–22
662 (2020)

663 [20] Sanchez-Garcia, R., Gomez-Blanco, J., Cuervo, A., Carazo, J.M., Sorzano,
664 C.O.S., Vargas, J.: Deepenhancer: a deep learning solution for cryo-em
665 volume post-processing. *Communications biology* **4**(1), 1–8 (2021)

666 [21] Zhou, T.-m., Wang, S., Xu, J.: Deep learning reveals many more inter-
667 protein residue-residue contacts than direct coupling analysis. *bioRxiv*,
668 240754 (2018)

669 [22] Xie, Z., Xu, J.: Deep graph learning of inter-protein contacts. *Bioinfor-*
670 *matics* **38**(4), 947–953 (2022)

671 [23] Gao, M., Nakajima An, D., Parks, J.M., Skolnick, J.: Af2complex predicts
672 direct physical interactions in multimeric proteins with deep learning.
673 *Nature communications* **13**(1), 1–13 (2022)

674 [24] Bryant, P., Pozzati, G., Elofsson, A.: Improved prediction of protein-
675 protein interactions using alphafold2. *Nature communications* **13**(1), 1–11
676 (2022)

- 677 [25] Evans, R., O'Neill, M., Pritzel, A., Antropova, N., Senior, A.W., Green,
678 T., Žídek, A., Bates, R., Blackwell, S., Yim, J., et al.: Protein complex
679 prediction with alphafold-multimer. *BioRxiv* (2021)
- 680 [26] Pellegrini, M., Marcotte, E.M., Thompson, M.J., Eisenberg, D., Yeates,
681 T.O.: Assigning protein functions by comparative genome analysis: pro-
682 tein phylogenetic profiles. *Proceedings of the National Academy of*
683 *Sciences* **96**(8), 4285–4288 (1999)
- 684 [27] Juan, D., Pazos, F., Valencia, A.: High-confidence prediction of global
685 interactomes based on genome-wide coevolutionary networks. *Proceedings*
686 *of the National Academy of Sciences* **105**(3), 934–939 (2008)
- 687 [28] Zeng, H., Wang, S., Zhou, T., Zhao, F., Li, X., Wu, Q., Xu, J.: Com-
688 plexcontact: a web server for inter-protein contact prediction using deep
689 learning. *Nucleic acids research* **46**(W1), 432–437 (2018)
- 690 [29] Feinauer, C., Szurmant, H., Weigt, M., Pagnani, A.: Inter-protein
691 sequence co-evolution predicts known physical interactions in bacterial
692 ribosomes and the *trp* operon. *PloS one* **11**(2), 0149166 (2016)
- 693 [30] Ovchinnikov, S., Kamisetty, H., Baker, D.: Robust and accurate prediction
694 of residue–residue interactions across protein interfaces using evolutionary
695 information. *elife* **3**, 02030 (2014)
- 696 [31] Rodriguez-Rivas, J., Marsili, S., Juan, D., Valencia, A.: Conservation
697 of coevolving protein interfaces bridges prokaryote–eukaryote homologies
698 in the twilight zone. *Proceedings of the National Academy of Sciences*
699 **113**(52), 15018–15023 (2016)
- 700 [32] Apweiler, R., Bairoch, A., Wu, C.H., Barker, W.C., Boeckmann, B., Ferro,
701 S., Gasteiger, E., Huang, H., Lopez, R., Magrane, M., *et al.*: Uniprot:
702 the universal protein knowledgebase. *Nucleic acids research* **32**(suppl_1),
703 115–119 (2004)
- 704 [33] Johnson, L.S., Eddy, S.R., Portugaly, E.: Hidden markov model speed
705 heuristic and iterative hmm search procedure. *BMC bioinformatics* **11**(1),
706 1–8 (2010)
- 707 [34] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal,
708 P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., *et al.*: Language
709 models are few-shot learners. *Advances in neural information processing*
710 *systems* **33**, 1877–1901 (2020)
- 711 [35] Qiu, J., Chen, Q., Dong, Y., Zhang, J., Yang, H., Ding, M., Wang,
712 K., Tang, J.: Gcc: Graph contrastive coding for graph neural network
713 pre-training. In: *Proceedings of the 26th ACM SIGKDD International*

- 714 Conference on Knowledge Discovery & Data Mining, pp. 1150–1160
715 (2020)
- 716 [36] Devlin, J., Chang, M.-W., Lee, K., Toutanova, K.: Bert: Pre-training
717 of deep bidirectional transformers for language understanding. arXiv
718 preprint arXiv:1810.04805 (2018)
- 719 [37] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X.,
720 Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et
721 al.: An image is worth 16x16 words: Transformers for image recognition
722 at scale. arXiv preprint arXiv:2010.11929 (2020)
- 723 [38] Elnaggar, A., Heinzinger, M., Dallago, C., Rehawi, G., Wang, Y., Jones,
724 L., Gibbs, T., Feher, T., Angerer, C., Steinegger, M., et al.: Prot-
725 trans: towards cracking the language of life’s code through self-supervised
726 learning. bioRxiv, 2020–07 (2021)
- 727 [39] Rao, R.M., Liu, J., Verkuil, R., Meier, J., Canny, J., Abbeel, P., Sercu,
728 T., Rives, A.: Msa transformer. In: International Conference on Machine
729 Learning, pp. 8844–8856 (2021). PMLR
- 730 [40] Rives, A., Meier, J., Sercu, T., Goyal, S., Lin, Z., Liu, J., Guo, D., Ott,
731 M., Zitnick, C.L., Ma, J., et al.: Biological structure and function emerge
732 from scaling unsupervised learning to 250 million protein sequences.
733 Proceedings of the National Academy of Sciences **118**(15) (2021)
- 734 [41] Rao, R., Bhattacharya, N., Thomas, N., Duan, Y., Chen, P., Canny,
735 J., Abbeel, P., Song, Y.: Evaluating protein transfer learning with tape.
736 Advances in neural information processing systems **32** (2019)
- 737 [42] Vig, J., Madani, A., Varshney, L.R., Xiong, C., Rajani, N., *et al.*: Bertol-
738 ogy meets biology: Interpreting attention in protein language models. In:
739 International Conference on Learning Representations (2020)
- 740 [43] Meier, J., Rao, R., Verkuil, R., Liu, J., Sercu, T., Rives, A.: Language
741 models enable zero-shot prediction of the effects of mutations on protein
742 function. Advances in Neural Information Processing Systems **34** (2021)
- 743 [44] Khalil, M., Teunissen, C.E., Otto, M., Piehl, F., Sormani, M.P., Gat-
744 tringer, T., Barro, C., Kappos, L., Comabella, M., Fazekas, F., *et al.*:
745 Neurofilaments as biomarkers in neurological disorders. Nature Reviews
746 Neurology **14**(10), 577–589 (2018)
- 747 [45] Gonzalez-Lozano, M., Koopmans, F., Sullivan, P., Protze, J., Krause, G.,
748 Verhage, M., Li, K., Liu, F., Smit, A.: Stitching the synapse: Cross-linking
749 mass spectrometry into resolving synaptic protein interactions. Science
750 advances **6**(8), 5783 (2020)

- 751 [46] Durrani, N., Sajjad, H., Dalvi, F.: How transfer learning impacts linguistic
752 knowledge in deep nlp models? In: Findings of the Association for
753 Computational Linguistics: ACL-IJCNLP 2021, pp. 4947–4957 (2021)
- 754 [47] Merchant, A., Rahimtoroghi, E., Pavlick, E., Tenney, I.: What happens to
755 bert embeddings during fine-tuning? In: Proceedings of the Third Black-
756 boxNLP Workshop on Analyzing and Interpreting Neural Networks for
757 NLP, pp. 33–44 (2020)
- 758 [48] Fayyaz, M., Aghazadeh, E., Modarressi, A., Mohebbi, H., Pilehvar, M.T.:
759 Not all models localize linguistic knowledge in the same place: A layer-
760 wise probing on bertoids’ representations. In: Proceedings of the Fourth
761 BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks
762 for NLP, pp. 375–388 (2021)
- 763 [49] Ho, J., Kalchbrenner, N., Weissenborn, D., Salimans, T.: Axial atten-
764 tion in multidimensional transformers. arXiv preprint arXiv:1912.12180
765 (2019)
- 766 [50] Huang, Z., Wang, X., Huang, L., Huang, C., Wei, Y., Liu, W.: Ccnet:
767 Criss-cross attention for semantic segmentation. In: Proceedings of the
768 IEEE/CVF International Conference on Computer Vision, pp. 603–612
769 (2019)
- 770 [51] Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for
771 contrastive learning of visual representations. In: International Conference
772 on Machine Learning, pp. 1597–1607 (2020). PMLR
- 773 [52] Gao, T., Yao, X., Chen, D.: Simcse: Simple contrastive learning of sentence
774 embeddings. arXiv preprint arXiv:2104.08821 (2021)
- 775 [53] Munkres, J.: Algorithms for the assignment and transportation problems.
776 *Journal of the society for industrial and applied mathematics* **5**(1), 32–38
777 (1957)
- 778 [54] Si, Y., Yan, C.: Protein complex structure prediction powered by multi-
779 ple sequence alignment of interologs from multiple taxonomic ranks and
780 alphafold2. *bioRxiv* (2021)
- 781 [55] Basu, S., Wallner, B.: Dockq: a quality measure for protein-protein
782 docking models. *PloS one* **11**(8), 0161879 (2016)
- 783 [56] Mirdita, M., Von Den Driesch, L., Galiez, C., Martin, M.J., Söding, J.,
784 Steinegger, M.: Uniclust databases of clustered and deeply annotated pro-
785 tein sequences and alignments. *Nucleic acids research* **45**(D1), 170–176
786 (2017)

- 787 [57] Remmert, M., Biegert, A., Hauser, A., Söding, J.: Hhblits: lightning-
788 fast iterative protein sequence searching by hmm-hmm alignment. *Nature*
789 *methods* **9**(2), 173–175 (2012)
- 790 [58] Gueudré, T., Baldassi, C., Zamparo, M., Weigt, M., Pagnani, A.: Simul-
791 taneous identification of specifically interacting paralogs and interprotein
792 contacts by direct coupling analysis. *Proceedings of the National Academy*
793 *of Sciences* **113**(43), 12186–12191 (2016)
- 794 [59] Fukuda, H., Tomii, K.: Deepeca: an end-to-end learning framework for
795 protein contact prediction from a multiple sequence alignment. *BMC*
796 *bioinformatics* **21**(1), 1–15 (2020)
- 797 [60] Varnai, C., Burkoff, N.S., Wild, D.L.: Improving protein-protein inter-
798 action prediction using evolutionary information from low-quality msas.
799 *PLoS one* **12**(2), 0169356 (2017)