Linker-Tuning: Optimizing Continuous Prompts for Heterodimeric Protein Prediction

Anonymous Author(s) Affiliation Address email

Abstract

| 1 | Predicting the structure of interacting chains is crucial for understanding biological |
|----|--|
| 2 | systems and developing new drugs. Large-scale Pre-trained Protein Language |
| 3 | models (PLMs), such as ESM-2, have shown an impressive ability to extract |
| 4 | biologically meaningful representations for protein contact and structure prediction. |
| 5 | In this paper, we show that ESMFold, which has been successful in computing |
| 6 | accurate atomic structures for single-chain proteins, can be adapted to predict the |
| 7 | heterodimer structures in a lightweight manner. We propose Linker-tuning, which |
| 8 | learns a continuous prompt to connect the two chains in a dimer before running |
| 9 | it as a single sequence in ESMFold. Experiment results show that our method is |
| 10 | significantly better than the ESMFold-Linker baseline, with relative improvements |
| 11 | of +28.13% and +54.55% in DockQ score on the i.i.d heterodimer test set and |
| 12 | the out-of-distribution (OOD) test set HeteroTest2, respectively. Notably, on the |
| 13 | antibody heavy chain light chain (VH-VL) test set, our method successfully predicts |
| 14 | all the heavy chain light chain docking interfaces, with 46/68 medium-quality and |
| 15 | 22/68 high-quality predictions, while being $9 \times$ faster than AF-Multimer. |

16 **1** Introduction

Proteins are large biomolecules essential to life. They are sequences compromised of 20 types of 17 18 amino acids and fold into three-dimensional (3D) structures to carry out functions. Predicting the 3D structures of proteins from amino acid sequences is a long-standing challenge in computational 19 biology. It is important for the mechanical understanding of protein functions as well as for designing 20 new drugs. In 2021, AlphaFold2 (AF2) strikes a huge success in solving this challenge, achieving 21 near experimental accuracy on protein structure prediction [1]. However, this system heavily relies 22 on Multiple Sequence Alignments (MSAs) to extract the evolutionary information, but MSAs are not 23 always available or high quality, especially for orphan proteins and fast-evolving antibodies [2]. 24

Inspired by the success of transformer language models in the field of Natural Language Processing 25 (NLP), there is a line of work resorting to large-scale PLMs for protein structure prediction [2, 3, 4, 5]. 26 These PLM-based models, such as ESMFold [3], take only amino acid sequences as input, eliminating 27 the need for MSAs. Powered by PLMs, they show strong abilities in capturing protein structure 28 information [6, 7]. And they are able to predict protein 3D structures at the atomic level with high 29 accuracy while being an order of magnitude faster than AF2. However, these models are developed 30 for predicting the structures of single-chain proteins and it is not clear how to use them to predict 31 multi-chain protein structures. 32

To adapt these models for protein complex prediction, some researchers have proposed to use a poly-Glycine *linker* to join chains and input the linked sequences to the model to predict complex structures [8, 9]. The rationale is that the model should identify the linker segment as unstructured and fold the linked sequence in a similar way to multiple chains. Experimental result on AF2 shows that this approach is simple yet effective. However, for the PLM-based models, whether a linker is effective or not for protein complex prediction remains unexplored. In the work of ESMFold, they

effective or not for protein complex prediction remains unexplored. In the work of ESMFold, they briefly mention that they use a 25-residue poly-Glycine linker (denoted as G25 in the following) to

join different chains for a specific protein complex example [3]. But they do not test the performance

- of the linker systematically. Based on existing work, we would like to investigate the following
- 42 questions in this paper: 1) How well can a G25 linker perform on protein complex prediction? 2)
- ⁴³ Can we optimize the linker to achieve a better result? And how?

Viewing proteins as the language of life, linkers in fact are the same things as prompts in natural 44 language. Inspired by prompt engineering [10, 11] in NLP, we propose Linker-tuning, which is to 45 automatically learn a linker for the PLM-based model ESMFold on the task of heterodimeric protein 46 structure prediction. Our goal is to find a linker that can link the two chains of a heterodimer so the 47 structure prediction model can fold it in a similar way as a single-chain protein. How to best achieve 48 this goal, however, is non-trivial and remains under-explored for the complicated protein structure 49 prediction model. Through preliminary analysis, we find that it is better to place linker optimization 50 at the Folding Module instead of at the PLM, which is different from intuition. 51

Considering ESMFold is a model with large-scale pre-trained PLM ESM2 that scales up to 15B parameters, to accelerate the linker learning procedure, we train and select our model on a proxy task called *distogram prediction* [12], a task that aims to predict inter-residue distance bins in the 3D space for each pair of residues in a given protein. After training, we test our learned linker on the 3D structure prediction task on three datasets to investigate the generalization ability of our method.

- 57 In summary, our main contributions are as follows:
- We propose Linker-tuning, a lightweight adaptation method that automatically learns a
 linker in the continuous space to adapt the single-chain ESMFold for heterodimer structure
 prediction.
- We show that our method outperforms the ESMFold-Linker baseline by large margins on both contact and structure prediction tasks on the heterodimer test set.
- We find that our method generalizes well to predict heterodimers with low sequence similarity and antibody VH-VL complex.

65 2 Biological background

Linker In biology, linkers are short amino acid sequences created in nature to separate multiple domains in a single protein [13]. Biologists have found that linkers rich in Glycine act as independent units and do not affect the function of the individual proteins to which they attach [14, 15]. Therefore, we can use the Glycine-rich linker to join interacting chains to make it a single sequence, hoping it folds in the way they suppose to. Grounded in biological principles, we further extend the natural discrete linkers to virtual continuous linkers for better protein complex structure prediction.

Distogram and contact map The 3D structure of a protein is expressed as (x, y, z) coordinates of 72 the residues' atoms in the form of a pdb file [16]. The distance between two residues in a protein 73 3D structure is defined as the Euclidean distance between their C_{β} atoms (C_{α} for Glycine). Binning 74 all the inter-residue distances in a protein into k distance bins, we can obtain the distogram matrix 75 [12]. For a protein with L residues, the distogram d is an $L \times L$ matrix, with entry d_{ij} referring to 76 the distance category of residue i and j. In a coarser granularity, we can compute the contact map 77 $c \in \mathcal{R}^{L \times L}$, where $c_{ij} = 1$ means the distance between residue *i* and *j* is less than or equal to 8Å. 78 For protein complexes, we are especially interested in the inter-chain contact maps where the contacts 79 are formed by two residues from different chains. The inter-chain contact map reflects the interface 80 of interacting proteins, which is essential for predicting the 3D structure of the complex. 81

82 3 Related work

83 3.1 Protein structure prediction

Single-chain protein structure prediction In recent years, single-chain protein structure prediction
 has attracted increasing attention from researchers in the Artificial Intelligence (AI) community,

mainly due to the ground-breaking success of the deep learning model AF2. Deep learning based 86 protein structure prediction methods can be classified into two main categories: 1) MSA-based 87 methods, such as AF2, that take protein sequences and MSAs as input and predict 3D structures 88 [1, 17, 18]; 2) PLM-based methods, such as ESMFold, that take only protein sequences as input 89 and predict 3D structures [3, 2, 4, 5, 19, 20, 21, 22, 23]. PLM-based methods do not rely on MSAs, 90 which are time-consuming in searching homologs and not always available for some proteins like 91 92 orphan proteins. Instead, they adopt large-scale pre-trained PLMs to learn evolutionary and structural meaningful representations for 3D structure prediction. In this work, we build our method upon PLM-93 based methods. Specifically, we adopt ESMFold [3] as the backbone since its code and pre-trained 94 weights are all released and convenient to use. The overall architecture of ESMFold contains two 95 parts: 1) ESM2: a PLM pre-trained with masked language modeling objective and scales up to 15B 96 parameters; 2) Folding Module: contains Folding Trunk (similar to Evoformer in AF2) and Structure 97 Module (same as the one in AF2), which are responsible for structure folding. 98

Multi-chain protein structure prediction In biology, multi-chain proteins are protein complexes 99 formed by interacting single-chain proteins where the interactions are driven by the same physical 100 forces as protein folding [24]. Recently, there is a line of work repurposing single-chain AF2 for 101 protein complex structure prediction. The methods can be summarized into two main categories: 102 1) input-adapted methods that provide AF2 with pseudo-multimer inputs either by adding a large 103 number to the residue index between chains to indicate chain break [25, 26, 27, 28] or using a linker 104 to join chains [8, 9]; and 2) training-adapted methods that retrain AF2 on multimeric proteins, such 105 as AF-Multimer, the state-of-the-art (SOTA) method [29]. On the one hand, the two types of methods 106 either do not update any parameters, or update all parameters of the base model, while our method 107 108 falls in between, adding only a tiny number of extra parameters to the base model. On the other hand, existing work mainly focuses on the MSA-based method AF2, with little attention being paid to the 109 PLM-based methods. In this work, we focus on adapting the PLM-based methods for two-chain 110 protein structure prediction, which has not yet been explored. 111

112 **3.2 Prompt engineering**

In the NLP community, with the rise of large-scale pre-trained language models (LMs) such as GPT-3 113 [30], "pre-train, prompt, and predict" has become a prevalent paradigm to steer the LM to perform a 114 wide range of downstream tasks [10]. In this paradigm, the downstream tasks are reformulated in a 115 form that is similar to the LM pre-training task using a textual prompt [30, 31]. The key challenge in 116 prompt-based learning is to find the right prompt for a specific task, termed "prompt engineering". 117 There is a line of work that automatically search the right prompts for downstream tasks [32, 33]. In 118 particular, instead of natural language prompts, some researchers propose to use continuous prompts, 119 directly performing prompting in the embedding space of the LM [34, 11]. In their experiment, 120 continuous prompts achieve strong results in both language understanding and generation tasks. In 121 this work, we follow the idea of continuous prompting, searching for the linkers in the continuous 122 123 space.

124 4 Method: Linker-tuning

To adapt the single-chain model for multi-chain protein structure prediction, we propose a lightweight
 adaptation method called Linker-tuning and a novel weighted distogram loss. The basic idea of our
 method is to optimize linkers, i.e., prompts, in the embedding space of ESMFold.

128 4.1 Problem formulation

Continuous linker tuning of ESMFold for protein complex structure prediction is a continuous optimization problem. Our goal is to find a linker that maximizes the performance of ESMFold on protein complex prediction. To be specific, we first denote training data as $D_{train} = \{(x_1, y_1), ..., (x_n, y_n)\}$ where $x_i = (x_i^A, x_i^B)$ and x_i^A, x_i^B represent the amino acid sequences of two chains, y_i is the structure of protein x_i . For a specified linker length L, the linker optimization problem is defined as follows:

$$\boldsymbol{l}^* = \operatorname*{arg\,min}_{\boldsymbol{l}\in E_L} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(x_i, y_i, \boldsymbol{l}) \tag{1}$$



Figure 1: Overview of Linker-tuning method with ESMFold as backbone. (A) Training. Based on ESMFold (shown in blue colors), we add a linker embedding module E_L (shown in yellow colors) with linker length L. Given a protein with multiple chains, we add the linker specified in the linker embedding module between each chain before running it as a single chain through the ESMFold model. The model outputs a distogram with the linker part removed. We use a weighted distogram loss as the objective function to train the linker embedding module while freezing all the parameters in ESMFold. (B) Inference. After training, ESMFold with our linker embedding module can be treated as a whole black box model, denoted as ESMFold-Linker*. The input for this model is just protein sequences. And the model outputs a predicted distogram as well as all the atoms' 3D coordinates for the protein.

where l denotes a linker, $E_L \subset \mathcal{R}^{L \times d}$ denotes a specific embedding space with embedding dimension of d, $\mathcal{L}(x_i, y_i, l)$ denotes complex structure prediction loss w.r.t. protein (x_i, y_i) using linker l.

¹³⁶ Therefore, the linker optimization is placed at the task level instead of at the instance level.

137 4.2 Model architecture

Our method is implemented based on ESMFold, a PLM-based strong structure prediction model. 138 As shown in Figure 1, we place the continuous linker at Folding Module of ESMFold, which takes 139 both the sequence representation from ESM2 and the amino acid sequence as input. There are two 140 main reasons that motivate us to place the continuous linker at Folding Module instead of at ESM2. 141 142 First, we can utilize the pre-trained distogram head while avoiding backpropagating to the giant ESM2 model. If we put it on the ESM2 side, the combined depth of training will go up to 104 143 layers, making it easily suffer from gradient vanishing and exploding. Second, preliminary analysis 144 on inter-chain contact prediction (shown in Table 4) shows that using Folding Module on top of 145 ESM2-3B increases prediction precision dramatically over ESM2-3B while ESM2-3B just performs 146 slightly better ESM2-650M, implying that Folding Module is more sensitive to structure prediction 147 and easier to control. 148

We implement a plug-in linker embedding module, which contains $L \times d$ learnable parameters where d is the embedding dimension of Folding Module. During training, only the linker embedding module is trainable, while all the original parameters in ESMFold are frozen. Therefore, ESM2 is just a sequence feature extractor that generates features for Folding Module. As shown in Figure 1(A), we first use a poly-Glycine linker of the same length as the continuous linker to join different chains for the ESM2 input. Then we obtain the protein sequence representation and input it to Folding Module along with the chains connected by the continuous linker. Finally, the distogram head outputs a probability distribution $p_{ij}^D \in \mathcal{R}^{64}$ of each residue pair (i, j) on 64 distance bins, which is used for computing the loss function. After training, we view ESMFold and the linker embedding module as a whole and name it as ESMFold-Linker*. As shown in Figure 1(B), it can be used to predict the distograms as well as the 3D coordinates of all the residues for multi-chain protein sequences.

160 4.3 Weighted distogram loss

Intuitively, to predict the structure of a protein complex, we need to know two things: 1) the structures
of each chain, on which ESMFold has been trained; and 2) the interaction interface between chains,
which ESMFold has never seen before. Therefore, we propose to weight the intra-chain predictions
and inter-chain predictions differently, with a focus on learning better interface between chains.

Formally, let N_A, N_B be the number of residues in two chains in a protein complex, $N = N_A + N_B$ be the total number of residues in the protein complex. Let $y_{ij} \in \mathcal{R}^{64}$ denote the one-hot labels of the 3D space distance bins between residue pair (i, j) and $p_{ij} \in \mathcal{R}^{64}$ be the corresponding predicted probability. We define a weighted distogram loss for a protein complex as follows:

$$\mathcal{L}(x, y, \boldsymbol{l}) = \mathcal{L}_1(x^A, y^A) + \mathcal{L}_1(x^B, y^B) + \lambda \mathcal{L}_2(x, y, \boldsymbol{l})$$
(2)

where $\mathcal{L}_1(.,.)$ denotes the single-chain distogram loss given as follows:

$$\mathcal{L}_1(x_A, y_A) = -\frac{2}{N_A(N_A + 1)} \sum_{i=1}^{N_A} \sum_{j \ge i}^{N_A} \sum_{b=1}^{64} y_{ijb} log(p_{ijb}^D)$$
(3)

and $\mathcal{L}_2(x, y, l)$ denotes the inter-chain distogram loss defined as follows:

$$\mathcal{L}_{2}(x, y, \boldsymbol{l}) = -\frac{1}{N_{A}N_{B}} \sum_{i=1}^{N_{A}} \sum_{j=1}^{N_{B}} \sum_{b=1}^{64} y_{ijb} log(p_{ijb}^{D})$$
(4)

and $\lambda \ge 2$ is a hyperparameter controlling the attention we place on the interface of a protein complex.

¹⁷² In our method, we use the weighted distogram loss as the training objective and validation metric.

173 5 Experiments

174 5.1 Experiment setting

Datasets We mainly perform experiments on heteromers of two chains. For training, we use the 175 dataset from APOC [35], which contains heterodimers released in the Protein Data Bank (PDB) 176 before 2018-09-30. After filtering out similar sequences at a 40% sequence identity threshold, it 177 is split into train/valid/test¹ sets by CDPred [36]. We further filter out those proteins that contain 178 missing C_{β} coordinates (C_{α} for Glycine) in the pdb file. The resulting train/valid/test sample sizes are 179 2,946/193/172, respectively. The average number of residues in the test set is 367, with a maximum 180 of 998. Furthermore, we use the largest blind test set HeteroTest2² from CDPred, which contains 181 55 heterodimers released in PDB between 2021-09-01 to 2021-10-20 [36]. The average number of 182 residues is 505, with a maximum of 979. In addition, we use the antibody VH-VL test set from 183 XtrimoDock [37]. It contains 68 samples released in PDB after 2022-02-01. Each sample consists of 184 one heavy and one light chain, forming the fragment variable region (Fv), which is a critical part of 185 antigen binding. The average number of residues is 231, with a range of [223, 244]. 186

Models We use ESMFold-v1³ as our backbone model. ESMFold-v1 consists of a 3B ESM2 model and a 670M Folding Module, which is the largest yet publicly available ESMFold checkpoint. For the Linker-tuning method, the linker length L is set to 25, equal to the length of the manual poly-Glycine linker. So the plug-in linker embedding module contains 0.027M parameters. We initialize the linker embedding using the embedding of Glycine. During training, only the linker embedding module is trainable, while all the original parameters in ESMFold are frozen. The hyperparameter λ in the

¹https://github.com/BioinfoMachineLearning/CDPred/tree/main/example/training_datalists

²https://zenodo.org/record/6647564#.ZDWvMuxBxhE

³https://dl.fbaipublicfiles.com/fair-esm/models/esmfold_3B_v1.pt

weighted distogram loss is set to 4. We train the model on a single Nvidia A100 80GB GPU with
batch_size=1 and num_epoch=15. The protein sequences in the training set are cropped to 225
residues to fit in GPU memory using the multi-chain cropping algorithm from AF-Multimer [29].
The number of recycles is set to 1 during training to reduce computation. We use Adam optimizer
with a learning rate of 5e-4. We select the best model based on the validation weighted distogram
loss. During inference, the number of recycles is set to 3.

- 199 **Baselines** We compare our method with several baselines and one SOTA model as follows:
- ESMFold-Linker: ESMFold-v1 with chains joined by the G25 linker as input.
- ESMFold-Gap: ESMFold-v1 with residue_index_offset set to 512.
- AlphaFold-Linker [29]: AF2 with a 21 residue repeated Glycine-Glycine-Serine linker.
- HDOCK [38]: rigid docking with single chains predicted by AF2.
- AF-Multimer(v3 best) [29]: AF-Multimer contains five models that are trained on all protein structures released in PDB before 2021-09-30. We take the best prediction from the five AF-Multimer models.

Metrics For protein complex 3D structure prediction, we use DockQ [39] to evaluate the quality 207 of the predicted interfaces. As defined by Critical Assessment of PRediction Interactions (CAPRI), 208 interfaces with DockQ < 0.23 means incorrect prediction, interfaces with $0.23 \le DockQ < 0.49$ 209 means acceptable prediction, $0.49 \leq \text{DockQ} < 0.80$ means medium quality prediction, and $\text{DockQ} \geq$ 210 0.80 means high-quality prediction. To evaluate the whole predicted protein complex structure rather 211 than the interfaces, we adopt two commonly used global structure metrics, namely, Root Mean 212 Squared Deviation (RMSD), and Template-Modeling Score (TM-Score) [40]. Besides, we use the 213 top-k precision as an evaluation metric for inter-chain contact prediction. We set $k = N_s/5$, where 214 N_s is the minimum chain length for a given protein complex. 215

216 5.2 General heterodimer structure prediction

Table 1 shows the protein complex structure prediction results of our methods and the baselines on 217 the heterodimer test set and HeteroTest2. On the i.i.d. heterodimer test set, ESMFold-Linker achieves 218 a 0.32 DockQ score and a 0.76 TM-score on average. By optimizing the linker, our model, i.e., 219 ESMFold-Linker*, achieves a 0.36 DockQ score and a 0.79 TM-score on average on the same test 220 set, outperforming the ESMFold-Linker baseline by 13.61% and 3.28%, respectively. Interestingly, 221 the gain of the interface quality (13.61%) is much larger than the gain of the whole structure quality 222 (3.28%), indicating that our learned linker mainly improves the interfaces more than the overall 223 structures. We further improve the ESMFold-Linker* by incorporating a large chain break, which 224 adds a large number to the residue index in Folding Module. And the model ESMFold-Linker*-225 Gap achieves a 0.41 DockQ score and 0.80 TM-score, outperforming ESMFold-Linker by 28.13% 226 and 5.26%, respectively. On the OOD test set HeteroTest2, we observe similar results. ESMFold-227 Linker*-Gap surpasses ESMFold-Linker by 54.55% DockQ score and 4.82% TM-score, respectively, 228 suggesting that our learned linker can generalize well to OOD data. 229

Compared to AlphaFold-Linker, a model that takes linked sequences and MSAs as input, our best 230 model ESMFold-Linker*-Gap achieves similar DockQ scores on both test sets, with lower values in 231 RMSD. Meanwhile, it outperforms the classic docking method HDOCK with AF2 predicted chains 232 as input in terms of DockQ score and RMSD. Furthermore, we compare it with the SOTA model 233 AF-Multimer.⁴ From Table 1, we can see there is still a large gap between our method and the 234 AF-Multimer(v1 best) on HeteroTest2. There are three main reasons responsible for this gap: 1) The 235 base model for AF-Multimer is AF2, which is a model stronger than ESMFold in general, especially 236 for those proteins that have high-quality MSAs; 2) AF-Multimer is a fully fine-tuned version of AF2 237 on a larger protein complex structure dataset while our model is a prompt tuning method trained only 238 on the heterodimer dataset; 3) AF-Multimer ensembles five models, while we only use one model. 239 However, our method is able to predict some proteins that are hard for both ESMFold-Linker and 240 AF-Multimer. As shown in Figure 2, ESMFold-Linker* successfully predicts the interface of the 241

⁴We use AF-Multimer v1 here because of the overlapping training data of AF-Multimer v3 and HeteroTest2. Since AF-Multimer v1 contains the heterodimer test set in its training data, we do not report the performance.

| | F | Ieterodimer | test | | HeteroTest | 2 |
|----------------------------|------------|-------------|------------|------------|-------------|------------|
| | DockQ↑ | RMSD↓ | TM-score↑ | DockQ↑ | RMSD↓ | TM-score↑ |
| ESMEald Linkar | 0.32 | 10.76 | 0.76 | 0.11 | 20.10 | 0.62 |
| ESMIFOId-LIIIKEI | ± 0.34 | ± 8.68 | ± 0.19 | ± 0.20 | ± 10.31 | ± 0.19 |
| ESMEald Car | 0.34 | 10.37 | 0.77 | 0.11 | 20.17 | 0.63 |
| ESMFold-Gap | ± 0.35 | ± 8.89 | ± 0.19 | ± 0.21 | ± 11.70 | ±0.19 |
| AlmhaEald Linkar | 0.42 | 9.38 | 0.83 | 0.17 | 20.95 | 0.71 |
| AlphaFold-Linker | ± 0.40 | ± 9.46 | ± 0.17 | ± 0.32 | ± 11.93 | ± 0.18 |
| UDOCK | 0.36 | 9.74 | 0.81 | 0.15 | 19.49 | 0.68 |
| HDOCK | ± 0.38 | ± 8.66 | ± 0.17 | ± 0.29 | ± 11.72 | ± 0.18 |
| ESMEald Linkar*(aura) | 0.36 | 9.19 | 0.79 | 0.14 | 19.03 | 0.65 |
| ESMFold-Linker*(ours) | ± 0.35 | ± 8.04 | ±0.19 | ± 0.23 | ± 10.93 | ± 0.20 |
| ESMEald Linkart Con(avera) | 0.41 | 8.59 | 0.80 | 0.17 | 18.53 | 0.65 |
| ESMIFOID-LINKEr*-Gap(ours) | ± 0.35 | ± 8.39 | ± 0.19 | ± 0.25 | ± 11.27 | ± 0.20 |
| AE multimer (w1 heat) | | | | 0.30 | 15.07 | 0.73 |
| AF-mulumer(VI best) | | | | ± 0.35 | ± 11.78 | ± 0.20 |

Table 1: Structure prediction results on Heterodimer data.



Figure 2: Comparison of predicted structure quality and inference time of heterodimer **7D7F_AD** by ESMFold-Linker, ESMFold-Linker*(ours), and AF-Multimer(v3 best). 7D7F is a membrane protein comprising 917 residues in the A and D chains. Structures are drawn using Protein Imager [41]. Gray indicates the ground truth structure.

membrane protein 7D7F_AD with a DockQ score of 0.39 while ESMFold-Linker and AF-Multimer cannot predict the interface correctly.

244 5.3 Antibody heavy chain light chain docking

We further test our method on antibodies, an important type of protein in designing new drugs. 245 Particularly, we focus on the heavy chain and the light chain docking. Table 2 shows the structure 246 prediction results of MSA-free methods (first two methods) and MSA-based methods (last four 247 methods) on the VH-VL test set. As an MSA-free model, ESMFold-Linker predicts all the interfaces 248 successfully with an average DockQ score of 0.737, better than the classical docking method HDOCK. 249 But it still lacks behind AlphaFold-Linker. For the three linker-based models, the distributions of their 250 DockQ scores are shown in Figure 3. Equipped with the optimized linker, ESMFold-Linker* achieves 251 an average DockQ score of 0.753, with 7 more high-quality interface predictions than ESMFold-252 Linker and 5 more high-quality interface predictions than AlphaFold-Linker. This result indicates 253 that our learned linker trained on the general heterodimer dataset generalizes well to antibody data. 254 Although the interface prediction performance of our method still falls behind AF-Multimer(v3 best), 255 the gap in the DockO score is much smaller compared to the case in HeteroTest2. Besides, it is 256 quite close in TM-score to XtrimoDock [37], which is trained on an antibody-antigen dataset. Given 257 our method only requires sequences as input, it can be a potentially useful model in the scenario of 258 antibody design where the evolving antibody might not have MSAs. 259

| Table 2: Struct | ire prediction | results on | VH-VL |
|-----------------|----------------|------------|-------|
|-----------------|----------------|------------|-------|

| | DockQ↑ | RMSD↓ | TM-score↑ |
|---------------------------|--------|-------------|-----------|
| EQME-14 Links | 0.737 | 1.459 | 0.955 |
| ESMFold-Linker | ±0.084 | ±0.474 | ±0.019 |
| ESMEold Linker*(ours) | 0.753 | 1.388 | 0.959 |
| L'SIVITOIQ-LIIIKET (OUIS) | ±0.083 | ±0.498 | ±0.019 |
| UDOCK | 0.705 | 2.0318 | 0.926 |
| HDOCK | ±0.202 | ± 2.405 | ±0.101 |
| AlphaEold Linkor | 0.746 | 1.4068 | 0.957 |
| AlphaFold-Lilikei | ±0.089 | ±0.520 | ±0.021 |
| AE multimor (y2 host) | 0.779 | 1.287 | 0.963 |
| AF-Inutumet (V3 best) | ±0.091 | ±0.518 | ±0.020 |
| VtrimoDock | 0.775 | 1.264 | 0.965 |
| AumoDock | ±0.021 | ±0.572 | ±0.097 |



Figure 3: Boxplot of DockQ on VH-VL.

260 6 Analysis and Discussion

ESMFold-Linker* is 9× faster than AF-Multimer in inference

We report the structure inference time of the MSA-free methods 262 (ESMFold-Gap, ESMFold-Linker, and ESMFold-Linker*) and the 263 SOTA MSA-based model AF-Multimer on the VH-VL dataset using 264 A100 80G GPU. Table 3 shows the total model inference time on 265 the VH-VL test set, where AF-Multimer's time is only for one 266 model, excluding the time of MSA search. As shown in Table 3, on 267 the VH-VL test set with an average sequence length of 231, both 268 ESMFold-Linker and ESMFold-Linker* take 4 minutes to run the 269 inference, which is $9 \times$ faster than AF-Multimer. 270

Table 3: Inference time.

| | Time |
|-----------------|--------|
| ESMFold-Gap | 3 min |
| ESMFold-Linker | 4 min |
| ESMFold-Linker* | 4 min |
| AF-Multimer | 36 min |

Large chain break or linker, or both? We perform an ablation study on ESMFold with chain break 271 and linker to better understand the contribution of each operation. Table 4 shows the comparison 272 of inter-chain contact prediction precision of ESMFold-based methods on the heterodimer test set 273 and HeteroTest2.⁵ As shown in Table 4, it is hard to tell whether ESMFold-Linker or ESMFold-Gap 274 is better. However, combining the two (ESMFold-Linker-Gap) provides significant performance 275 gains over using either operation alone on both datasets. We observe similar effects in our method 276 when incorporating chain break with the optimized linker. Compared to using a chain break, the 277 major limitation of using a linker is that it increases the computation cost (shown in Table 3). But 278 we can enjoy the advantage of a large degree of freedom for improvement and better performance. 279 Empirically, combining the two gives a better performance than just using each of them. 280

| (%) | Het | erodimer te | st | Н | eteroTest2 | |
|---------------------------|----------|-------------|--------|----------|------------|--------|
| (10) | top Ns/5 | top Ns/2 | top Ns | top Ns/5 | top Ns/2 | top Ns |
| ESM2-650M-Linker | 12.02 | 9.89 | 8.33 | | | |
| ESM2-3B-Linker | 12.14 | 10.86 | 8.89 | | | |
| ESMFold-Linker | 49.88 | 47.04 | 40.64 | 23.00 | 18.92 | 13.72 |
| ESMFold-Gap | 51.15 | 48.13 | 40.82 | 22.09 | 18.21 | 13.08 |
| ESMFold-Linker*(ours) | 57.55 | 53.04 | 44.37 | 27.11 | 22.14 | 15.46 |
| ESMFold-Linker-Gap | 57.72 | 53.44 | 45.41 | 25.20 | 19.84 | 14.66 |
| ESMFold-Linker*-Gap(ours) | 60.40 | 56.27 | 48.00 | 28.00 | 23.69 | 17.26 |

|--|

The learned linker allows more chain twist while rarely interacting with the chains In Figure 4, we visualize the predicted contact maps of two proteins with the linker inside to understand how the

⁵The contact map probabilities are obtained from the predicted distogram probabilities by summing the probability mass in each distribution below 8.25Å.



Figure 4: Contact maps of viral proteins 7VYR_HL (A) and 7WPE_YZ (B).

linker interacts with the chains. The two proteins are 7VYR HL and 7WPE YZ, corresponding to a 283 good case (0.77 DockQ score) and a bad case (0.01 DockQ score) in our model ESMFold-Linker*. 284 As shown in Figure 4, both the G25 linker (middle) and our learned linker (right) seem to rarely 285 interact with the protein chains in both cases. This result indicates that ESMFold is able to recognize 286 the linker part as a disordered region and fold the connected sequences as multi-domain proteins. 287 Furthermore, there are more predicted contacts using the learned linker than using the G25 linker in 288 289 both cases. This result suggests that the learned linker allows the connecting chains to freely twist and rotate to recruit binding partners more than the manual linker. 290

Limitations Our method has some limitations. First, if the base model (ESMFold-v1) is not good at predicting a certain type of protein complexes, such as the heterodimers in HeteroTest2, adding an optimized linker can not make it a strong model for that type of data since the trainable parameter size is very small. Second, our method is tested on heterodimers, whether it generalizes to homodimers or multi-chain proteins is unknown. Third, the linker is only optimized at the Folding Module, while the linker at ESM2 remains constant. And the linker length is treated as a hyperparameter, which can be further optimized to improve performance and speed.

298 7 Conclusions and future work

The use of prompts in protein structure prediction models is not always clear due to the high 299 complexity of models and a general lack of biological knowledge for AI researchers. In this work, 300 we have proposed Linker-tuning, a prompt tuning method to adapt the single-chain pre-trained 301 ESMFold for heterodimer structure prediction. As proof-of-concept, we showcase that we can place 302 a soft prompt in ESMFold. The task is reformulated as a pre-trained task itself under the biological 303 prior. Experiments show that merely tuning a prompt on ESMFold can significantly improve the 304 predicted complex structure quality over the discrete prompt handcrafted with strong biological 305 insight. Hopefully, our work can inspire more work on AI for Protein Science. 306

There are two directions for future work. Firstly, we would like to extend our work to antibodyantigen structure prediction, a critical task with direct relevance to drug design. Secondly, we are going to explore structural-aware antibody design using our method since it is efficient and fast. By pursuing these directions, our objective is to make progressive contributions towards the development of effective drugs for disease treatment and pain relief.

312 **References**

- [1] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn
 Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure
 prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [2] Ruidong Wu, Fan Ding, Rui Wang, Rui Shen, Xiwen Zhang, Shitong Luo, Chenpeng Su, Zuofan Wu,
 Qi Xie, Bonnie Berger, et al. High-resolution de novo structure prediction from primary sequence. *BioRxiv*,
 pages 2022–07, 2022.
- [3] 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. Language models of protein sequences at the scale
 of evolution enable accurate structure prediction. *BioRxiv*, 2022.
- [4] Xiaomin Fang, Fan Wang, Lihang Liu, Jingzhou He, Dayong Lin, Yingfei Xiang, Xiaonan Zhang, Hua Wu,
 Hui Li, and Le Song. Helixfold-single: Msa-free protein structure prediction by using protein language
 model as an alternative. *arXiv preprint arXiv:2207.13921*, 2022.
- [5] Ratul Chowdhury, Nazim Bouatta, Surojit Biswas, Christina Floristean, Anant Kharkar, Koushik Roy,
 Charlotte Rochereau, Gustaf Ahdritz, Joanna Zhang, George M Church, et al. Single-sequence protein
 structure prediction using a language model and deep learning. *Nature Biotechnology*, 40(11):1617–1623,
 2022.
- [6] Roshan Rao, Joshua Meier, Tom Sercu, Sergey Ovchinnikov, and Alexander Rives. Transformer protein language models are unsupervised structure learners. In *International Conference on Learning Representations*, 2020.
- [7] Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle
 Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15):e2016239118, 2021.
- [8] Junsu Ko and Juyong Lee. Can alphafold2 predict protein-peptide complex structures accurately? *BioRxiv*, 2021.
- [9] Tomer Tsaban, Julia K Varga, Orly Avraham, Ziv Ben-Aharon, Alisa Khramushin, and Ora Schueler Furman. Harnessing protein folding neural networks for peptide–protein docking. *Nature communications*, 13(1):1–12, 2022.
- [10] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train,
 prompt, and predict: A systematic survey of prompting methods in natural language processing. ACM
 Computing Surveys, 55(9):1–35, 2023.
- [11] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. GPT
 understands, too. *arXiv preprint arXiv:2103.10385*, 2021.
- [12] Andrew W Senior, Richard Evans, John Jumper, James Kirkpatrick, Laurent Sifre, Tim Green, Chongli
 Qin, Augustin Žídek, Alexander WR Nelson, Alex Bridgland, et al. Improved protein structure prediction
 using potentials from deep learning. *Nature*, 577(7792):706–710, 2020.
- [13] Vishnu Priyanka Reddy Chichili, Veerendra Kumar, and Jayaraman Sivaraman. Linkers in the structural
 biology of protein–protein interactions. *Protein science*, 22(2):153–167, 2013.
- 14] Athena D Nagi and Lynne Regan. An inverse correlation between loop length and stability in a four-helixbundle protein. *Folding and Design*, 2(1):67–75, 1997.
- Janet E Deane, Daniel P Ryan, Margaret Sunde, Megan J Maher, J Mitchell Guss, Jane E Visvader, and
 Jacqueline M Matthews. Tandem lim domains provide synergistic binding in the lmo4: Ldb1 complex.
 The EMBO journal, 23(18):3589–3598, 2004.
- Stephen K Burley, Helen M Berman, Charmi Bhikadiya, Chunxiao Bi, Li Chen, Luigi Di Costanzo, Cole
 Christie, Jose M Duarte, Shuchismita Dutta, Zukang Feng, et al. Protein data bank: the single global
 archive for 3d macromolecular structure data. *Nucleic Acids Research*, 47(D1), 2018.
- [17] Minkyung Baek, Frank DiMaio, Ivan Anishchenko, Justas Dauparas, Sergey Ovchinnikov, Gyu Rie Lee,
 Jue Wang, Qian Cong, Lisa N Kinch, R Dustin Schaeffer, et al. Accurate prediction of protein structures
 and interactions using a three-track neural network. *Science*, 373(6557):871–876, 2021.

- [18] Jianyi Yang, Ivan Anishchenko, Hahnbeom Park, Zhenling Peng, Sergey Ovchinnikov, and David Baker.
 Improved protein structure prediction using predicted interresidue orientations. *Proceedings of the National Academy of Sciences*, 117(3):1496–1503, 2020.
- [19] Wenkai Wang, Zhenling Peng, and Jianyi Yang. Single-sequence protein structure prediction using
 supervised transformer protein language models. *Nature Computational Science*, 2(12):804–814, 2022.
- Yining Wang, Xumeng Gong, Shaochuan Li, Bing Yang, YiWu Sun, Chuan Shi, Hui Li, Yangang Wang,
 Cheng Yang, and Le Song. xtrimoabfold: De novo antibody structure prediction without msa. *arXiv preprint arXiv:2212.00735v3*, 2022.
- [21] Jinhua Zhu, Zhenyu He, Ziyao Li, Guolin Ke, and Linfeng Zhang. Uni-fold musse: De novo protein complex prediction with protein language models. *bioRxiv*, pages 2023–02, 2023.
- Brennan Abanades, Wing Ki Wong, Fergus Boyles, Guy Georges, Alexander Bujotzek, and Charlotte M
 Deane. Immunebuilder: Deep-learning models for predicting the structures of immune proteins. *bioRxiv*,
 pages 2022–11, 2022.
- Yining Wang, Xumeng Gong, Shaochuan Li, Bing Yang, Yiwu Sun, Yujie Luo, Hui Li, and Le Song. Fast de
 novo antibody structure prediction with atomic accuracy. *Cancer Research*, 83(7_Supplement):4296–4296,
 2023.
- Ozlem Keskin, Attila Gursoy, Buyong Ma, and Ruth Nussinov. Principles of protein- protein interactions:
 what are the preferred ways for proteins to interact? *Chemical reviews*, 108(4):1225–1244, 2008.
- Ian R Humphreys, Jimin Pei, Minkyung Baek, Aditya Krishnakumar, Ivan Anishchenko, Sergey Ovchin nikov, Jing Zhang, Travis J Ness, Sudeep Banjade, Saket R Bagde, et al. Computed structures of core
 eukaryotic protein complexes. *Science*, 374(6573):eabm4805, 2021.
- Patrick Bryant, Gabriele Pozzati, and Arne Elofsson. Improved prediction of protein-protein interactions
 using alphafold2. *Nature communications*, 13(1):1265, 2022.
- [27] Mu Gao, Davi Nakajima An, Jerry M Parks, and Jeffrey Skolnick. Af2complex predicts direct physical
 interactions in multimeric proteins with deep learning. *Nature communications*, 13(1):1744, 2022.
- [28] Milot Mirdita, Konstantin Schütze, Yoshitaka Moriwaki, Lim Heo, Sergey Ovchinnikov, and Martin
 Steinegger. Colabfold: making protein folding accessible to all. *Nature methods*, 19(6):679–682, 2022.
- Richard Evans, Michael O'Neill, Alexander Pritzel, Natasha Antropova, Andrew Senior, Tim Green,
 Augustin Žídek, Russ Bates, Sam Blackwell, Jason Yim, et al. Protein complex prediction with alphafold multimer. *BioRxiv*, pages 2021–10, 2022.
- [30] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.
 Advances in neural information processing systems, 33:1877–1901, 2020.
- [31] Timo Schick and Hinrich Schütze. Exploiting cloze-questions for few-shot text classification and natural
 language inference. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 255–269, 2021.
- [32] Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt: Eliciting
 knowledge from language models with automatically generated prompts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4222–4235, 2020.
- [33] Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot learners.
 In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the
 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages
 3816–3830, 2021.
- Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In
 Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4582–
 4597, 2021.
- [35] Mu Gao and Jeffrey Skolnick. Apoc: large-scale identification of similar protein pockets. *Bioinformatics*,
 29(5):597–604, 2013.

- [36] Zhiye Guo, Jian Liu, Jeffrey Skolnick, and Jianlin Cheng. Prediction of inter-chain distance maps of
 protein complexes with 2d attention-based deep neural networks. *Nature Communications*, 13(1):6963,
 2022.
- Yujie Luo, Shaochuan Li, Yiwu Sun, Ruijia Wang, Tingting Tang, Beiqi Hongdu, Xingyi Cheng, Chuan
 Shi, Hui Li, and Le Song. xtrimodock: Rigid protein docking via cross-modal representation learning and
 spectral algorithm. *bioRxiv*, pages 2023–02, 2023.
- [38] Yumeng Yan, Huanyu Tao, Jiahua He, and Sheng-You Huang. The hdock server for integrated protein–
 protein docking. *Nature protocols*, 15(5):1829–1852, 2020.
- [39] Sankar Basu and Björn Wallner. Dockq: a quality measure for protein-protein docking models. *PloS one*, 11(8):e0161879, 2016.
- [40] Yang Zhang and Jeffrey Skolnick. Scoring function for automated assessment of protein structure template
 quality. *Proteins: Structure, Function, and Bioinformatics*, 57(4):702–710, 2004.
- [41] Gianluca Tomasello, Ilaria Armenia, and Gianluca Molla. The protein imager: a full-featured online
 molecular viewer interface with server-side hq-rendering capabilities. *Bioinformatics*, 36(9):2909–2911,
 2020.